

**Common Lisp
Analytical Statistics Package:
User Manual**

**Scott D. Anderson, Adam Carlson,
David L. Westbrook, David M. Hart,
and Paul R. Cohen**

Computer Science Technical Report 93-55*

Experimental Knowledge Systems Laboratory
Department of Computer Science, Box 34610
Lederle Graduate Research Center
University of Massachusetts
Amherst, MA 01003-4610

* Replaces TR-90-85

Common Lisp Analytical Statistics Package: User Manual

**Experimental Knowledge Systems Laboratory
Computer Science Department, LGRC
University of Massachusetts
Box 34610
Amherst, Massachusetts 01003-4610**

**Scott D. Anderson
Adam Carlson
David L. Westbrook
David M. Hart
Paul R. Cohen**

with contributions by

David A. Fisher

CLASP Version 1.3.2
CLIP Version 1.4

February 2, 1994

Abstract

This document is the user manual for the Common Lisp Analytical Statistics Package (CLASP), a tool for visualizing and statistically analyzing data, and also for the Common Lisp Instrumentation Package (CLIP), a tool for automating data collection and experimentation.

Through statistics we attempt to separate those events which are simply due to chance from those which are related to other events in our world. CLASP provides the tools necessary for these tasks. CLASP and CLIP are implemented in Common Lisp, with the CLASP user interface implemented in CLIM, the Common Lisp Interface Manager. Along with describing the operation of CLASP and CLIP, this manual details the CLASP programming interface, an aid for building extensions to CLASP.

Contents

1	Introduction	1
1.1	History and Acknowledgments	1
1.2	Design	2
1.3	Getting CLIP/CLASP	2
1.4	Conventions	3
2	CLIP	5
2.1	Motivation	5
2.2	Implementation	6
2.2.1	Define-simulator Macro	6
2.2.2	Define-experiment Macro	7
2.2.3	Write-current-experiment-data Function	9
2.2.4	Run-experiment Function	9
2.2.5	Explicitly Stopping Trials	10
2.3	Clip Definition	10
2.3.1	Simple Clips	11
2.3.2	Clips with Components	11
2.3.3	Time-series Clips	11
2.3.4	Defclip Macro	11
2.4	Examples from a Simple Agent Simulator	13
2.4.1	Using Simple Clips to Collect Trial Summary Data	13
2.4.2	Using a Mapping Clip to Map Simple Clips Over Multiple Agents	14
2.4.3	Full Agent Simulator Experiment with Time-Series Clips	15
3	Graphical User Interface	21
3.1	CLASP Layout	21
3.2	CLASP Command Syntax	21
3.2.1	Basic syntax	21
3.2.2	Mapped arguments	22
3.3	Interacting with CLASP	22
3.3.1	Using the mouse	23
3.3.2	Using the keyboard	23
3.4	Examples of CLASP command usage	24
3.5	CLASP Command Menus	25
3.5.1	The File Menu	25
3.5.2	The Graph Menu	27
3.5.3	The Describe Menu	28

3.5.4	The Manipulate Menu	29
3.5.5	The Transform Menu	31
3.5.6	The Test Menu	32
3.5.7	The Sample Menu	33
4	Statistics	35
4.1	Descriptive Statistics	35
4.1.1	Location	35
4.1.2	Spread	36
4.1.3	Covariance	37
4.2	Test Statistics	37
4.2.1	Confidence Intervals	38
4.2.2	t Tests	39
4.2.3	Analysis of Variance	39
4.2.4	Linear Regression	42
4.2.5	Chi Square	43
5	Data Manipulation	45
5.1	Function Application	45
5.1.1	Formula Expressions	45
5.1.2	Smoothing	46
5.2	Partition Clauses	47
5.2.1	Partition Clause Logic Operators	48
5.2.2	Partition Clause Comparison Operators	48
5.2.3	Exhaustive Partitioning Operator	49
6	Functions	51
6.1	Descriptive Statistic Functions	51
6.1.1	Location	52
6.1.2	Spread	53
6.1.3	Correlation	54
6.1.4	Summary	56
6.2	Test Statistic Functions	56
6.2.1	t Tests	56
6.2.2	Confidence Intervals	57
6.2.3	Analysis of Variance	58
6.2.4	Linear Regression	61
6.2.5	Contingency Table Analysis	63
6.2.6	Smoothing	64
6.3	Data Manipulation Functions	64
6.3.1	Dataset Functions	65
6.3.2	Input/Output	67
6.3.3	Manipulation	68
6.4	Density and Combinatorial Functions	70
6.5	Sampling	75
6.6	Graphics	76

A	Computational Methods	77
A.1	Quantiles	77
A.1.1	Bisection	78
A.1.2	Cumulative	79
A.1.3	Discussion	80
A.2	Analysis of Variance	81
A.3	Linear Regression	85
A.4	Regression Analysis: Matrix Method	87
B	Clip Examples	91
B.1	More Experiment Definitions Using CLIP	91
B.1.1	Measuring an Airport Gate Maintenance Scheduler’s Performance	91
B.1.2	PHOENIX Real-Time-Knob Experiment	96
The PHOENIX System	96	
Identifying the Factors that Affect Performance	97	
Experiment Design	97	
RTK Experiment Clips	98	
B.1.3	Example from a Transportation Planning Simulation	106
	References	111
	Index	111

List of Figures

2.1	Simple Clip Example	15
2.2	Mapping Clip Example	16
2.3	Event-Driven Time-Series Clips	19
2.4	Periodic Time-Series Clips	19
4.1	One Way ANOVA Calculations	41
4.2	Two Way ANOVA Calculations	41
A.1	Bisection Method for Computing Quantiles	78
A.2	The Cumulative Method for Computing Quantiles	79
A.3	Interpolating in the Cumulative Method	80
A.4	Comparison of the Bisection and Cumulative Methods	81

Chapter 1

Introduction

When events occur in the world, people turn to statistics to analyze and describe them. This is especially true of the behavioral sciences, psychology and sociology. Through statistics we attempt to separate those events which are simply due to chance from those which are a result of, or are related to, other events in our world. Through careful analysis it is possible to understand a great deal about the state of the world. In AI, we have a new aspect of behavioral science, the study of the behavior of intelligent agents and computer programs in complex environments, and so we turn to statistical analysis to describe and analyze their behavior and the factors that predict their behavior. Two examples of such AI programs are PHOENIX, a simulator of agents that fight forest fires in Yellowstone National Park, and TRANSIM, a simulator of transportation plans in a global logistical problem. We expect that these simulated worlds will lend themselves to the same forms of analysis we use to study the natural world. Our goal is to smoothly connect the simulators and other computer programs with statistical analysis tools, written in Common Lisp and running in the same Lisp world, so that high quality statistical analysis will be conveniently available to help us understand and test the behavior of our programs.

Why use statistics to understand our programs? Why not just look at the code? As our programs become increasingly large and complex, studying the code becomes as effective in understanding overall behavior as studying biochemistry is in understanding animal behavior. The interaction of components and the response of the program to different circumstances of its execution make global behavior hard to predict. Furthermore, programs will act differently in different executions because of differences in their environments, and statistics can help us summarize, visualize, and analyze the overall behavior, abstracting away the details of particular executions. We implemented CLASP, the Common Lisp Analytical Statistics Package, to help AI researchers analyze the behavior of their programs.

1.1 History and Acknowledgments

This document describes the second implementation of CLASP. Many people helped develop CLASP since the idea was first conceived in the Experimental Knowledge Systems Laboratory (EKSL) in 1989. First, Paul Silvey developed our relational database technology, called the Relational Table Manager, or RTM. This gave us a good way of saving and processing our data in the same Lisp environment that generated it. It also led us to want to compute means and perform t-tests to analyze the differences of means. Soon, it became apparent that there was no end to the kinds of analyses that we might need to perform, given the

increasing sophistication of our models and experimental questions, and so the idea of a coherent statistical package was born.

David Fisher designed and implemented the statistical functions and graphical interface for the first implementation of CLASP. He was assisted by David Westbrook in many ways, but primarily in the TI Explorer windowing and graphics code. Paul Silvey assisted in integrating CLASP with RTM. Other people implemented specific statistical functions, including Sameen Fatima's implementation of Chi-square analysis of two-way contingency tables and Dorothy Mammen's implementation of log-linear analysis for higher dimensional contingency tables.

Even before the first implementation of CLASP, EKSL was collecting data from systems like PHOENIX and statistically analyzing them. Therefore, we quickly started implementing general methods for instrumenting PHOENIX, collecting data, and running experiments. This effort was begun by Scott Anderson and David Westbrook, and eventually evolved into what call CLIP, which is simply a set of macros to make it easier to run experiments.

The second implementation of CLASP is a joint effort by Scott Anderson, Adam Carlson and David Westbrook. Its goal is to make CLASP more convenient and generally available by simplifying the data representation, eliminating its dependence on the TI Explorer by using CLIM, the Common Lisp Interface Manager, and improving the user interface based on our experience with the first implementation.

1.2 Design

This manual roughly follows the design of CLIP and CLASP, and therefore breaks down into five sections. The first is a description of CLIP, our aids for automated data collection and experimentation. The other three sections are aspects of CLASP: its graphical user interface, its data manipulation functionality, and its statistical functionality. We designed CLASP so that these aspects are fairly modular: the graphical user interface makes it easy to call the data manipulation and statistical subsystems, but those functions are separate from it and from each other. The functions of those subsystems can be called directly from a user's program or from a Common Lisp evaluator. Consequently, the fifth section describes the actual Common Lisp functions. Indeed, individual functions, CLOS classes and methods are documented by automatically typesetting of the name, argument list, type, and documentation string, pulling this information directly out of the code.

1.3 Getting Clip/Clasp

CLIP/CLASP is available via anonymous ftp from *ftp.cs.umass.edu* in the directories `pub/eksl/clip` and `pub/eksl/clasp`.

CLIP requires Common Lisp and the Common Lisp Object System, CLOS. CLASP requires these also, and in addition requires the Common Lisp Interface Manager, or CLIM, and SCIGRAPH, a publicly available scientific graphing package written in Common Lisp and CLIM by Bolt, Beranek and Newman (BBN). SCIGRAPH is included in the distribution. CLIP/CLASP runs on a number of platforms under a variety of lisp implementations. Check the release notes for a detailed list of currently supported platforms.

If you have any problems, questions, or suggestions, contact us at:

`clasp-support@cs.umass.edu`

1.4 Conventions

Wherever lisp or CLASP objects appear in a body of text, the following conventions hold:

- Names of CLASP commands are printed in SMALL CAPS and prefixed with a semi-colon.
- Names of arguments to CLASP functions and commands appear in *italics*.
- Names of keywords are printed in `teletype` and are prefixed with a semi-colon.
- Names of all other lisp objects, including functions are printed in `teletype`.

Where sections of code are included, they will be in `teletype`. In examples of interactive CLASP sessions, san serif will be used and for anything the system prints and **bold san serif** will be used for user input.

Chapter 2

CLIP

2.1 Motivation

We collect information from software systems for many reasons. Sometimes, the very purpose of the system is to produce information. Other times, we collect data for debugging, feedback to the user—in general, for understanding the system’s behavior. Unfortunately, there seem to be few general tools available for this task. Instead, we often find ad hoc, domain-specific code, spread throughout a system, varying from component to component, despite strong similarity in the requirements for any data collection code.

CLIP, the Common Lisp Instrumentation Package, standardizes data collection. It is named by analogy with the “alligator clips” that connect diagnostic equipment to electrical components. By standardizing data collection, we can more easily add and delete instrumentation from our systems. We can reuse a piece of data collection code, called a “clip,” in other systems we are studying. We can use general tools to analyze the collected data. Instrumentation in a system becomes more easily used and understood, because the basic functionality of the system is separated from the functionality of the instrumentation.

We designed CLIP to be used for scientific experiments, in which a system is run under a number of different conditions to see whether and how the behavior changes as a function of the condition. Consequently, we view an experiment as a series of trials, varying the settings of experiment variables and collecting information. Generally speaking, an experiment comprises the following steps:

1. Creating clips and inserting them into the system to be studied. Generally, clip measures one particular aspect of a system, but you can also define clips that combine many measurements, including other clips. Often, the clips you need will already be defined.
2. Define the experiment, in terms of what system is to be run, any necessary initialization code, and the conditions under which it is to run (different input parameters or environmental conditions).
3. Run a series of trials, saving the output into a CLASPformat data file.* This format is described in the CLASP manual section Data Manipulation Functions, although it isn’t necessary to understand the file format. CLASP can read what CLIP writes.

*CLIP can also write data files in standard tab or space delimited format used by most other statistical packages, databases and spreadsheets.

4. Analyze the data using CLASP.

2.2 Implementation

CLIP provides five major forms to support experiments:

- `define-simulator`,
- `defclip`,
- `define-experiment`,
- `write-current-experiment-data`, and
- `run-experiment`

Each is described in detail below.

2.2.1 Define-simulator Macro

`define-simulator` (*name* *key* *system-name* *system-version* *reset-system* *start-system* *stop-system* *schedule-function* *deactivate-scheduled-function* *seconds-per-time-unit* *timestamp*) [Macro]

Define the interface to a simulation. The following options are recognized:

<code>:system-name</code>	string naming system	[2]
<code>:system-version</code>	function or form which handles the arguments of the experiment and returns a string which denotes the version of the system ^a	[1]
<code>:start-system</code>	function or form that handles the experiment variables and arguments of the experiment; this function is called during the experiment loop to begin execution of the system; when it returns the trial is considered to be completed	[2]
<code>:reset-system</code>	same as <code>:start-system</code> ; this function will be called during the experiment loop before the system is started;	[1]
<code>:stop-system</code>	same arg handling as <code>:start-system</code> ; this function can be used to execute code when a trial is shutdown; it is most useful when instrumenting multiprocessing systems where other processes need to be terminated	[1]
<code>:schedule-function</code>	function or form that handles the lambda list (<i>function time period name</i>) and optionally returns a data-structure which represents the event; used to provide access to the simulator's event-scheduling mechanism	[3]
<code>:deactivate-scheduled-function</code>	function or form that handles one arg, namely the data structure returned by the <code>schedule-function</code> function; used to provide access to the simulator's event-scheduling mechanism	[1]
<code>:timestamp</code>	a function or a list of (<function-name> <clip-name>) where function should return the current time in units specified by <code>:seconds-per-time-unit</code>	[3]
<code>:seconds-per-time-unit</code>	the number of seconds in 1 simulator time quantum (default is 1)	[3]

Keywords marked [1] are optional; Keywords marked [2] are required, and keywords marked [3] are required when using time-series clips

^aFunctions must accept the formal arguments specified in `define-experiment` and the actual arguments specified in a call to `run-experiment`. Forms can refer to the arguments lexically.

The `define-simulator` form is used to eliminate the need to specify redundant information when multiple experiments are to be defined for the same system. All of the options specified in this form can be overridden in a `define-experiment` form.

2.2.2 Define-experiment Macro

`define-experiment` (*name arguments &body body &aux documentation*) [Macro]

A `define-experiment` form sets up the environment in which the system is to be run. Options for `define-experiment` support the sequential nature of experimentation. That

is, for the purposes of alligator clips, we view an experiment as a series of trials, varying the settings of experiment variables and collecting information. The `define-experiment` macro allows code to be run at the beginning and end of the experiment, usually to set up data-recording or to do data analysis. It also allows code to be run at the beginning and end of each trial; often this is used to reset counters, clear data structures, or otherwise set up the data collection apparatus. Through the `:script` keyword a user may add specific code to run at particular times during the experiment. A user can also specify experiment variables and their associated sets of values; the experiment code will ensure that all combinations of the experiment variables are run. For example, specifying:

```
:variables ((rating-scheme '(:OPPORTUNISTIC :ORDER-BY-ORDER))
            (percentage from 100 downto 50 by 10))
```

describes an experiment with twelve conditions (six values of percentage times two kinds of rating schemes), and so the number of trials will be a multiple of twelve, with an equal number of trials being executed under each condition. The options accepted by `define-experiment` are:

<code>:simulator</code>	specify the name of the simulator definition to use
<code>:before-trial</code>	called with the current values of the experiment variables and arguments to specify operations performed before each trial (such as initializations)
<code>:after-trial</code>	similar to <code>:before-trial</code> , specifies operations to be performed <i>after</i> each trial. Summary data files for each trial can be written at this time using <code>write-current-experiment-data</code> .
<code>:before-experiment</code>	called with the arguments. The value should be a form that refers to the arguments or a function that handles the correct number of arguments. Used to initialize system and define experiment environment.
<code>:after-experiment</code>	Similar to <code>:before-experiment</code> , called with just the arguments. Used to close files, write any experiment summarization output and clean up after the experiment.
<code>:instrumentation</code>	is a list of names defined with <code>defclip</code> which will be enabled during the experiment. Use <code>write-current-experiment-data</code> in the after-trial code to write the data to the output file.
<code>:variables</code>	<code>{{(var exp) (var {loop-for-clause})}}</code> - define experiment variables
<code>:locals</code>	<code>{{(var exp) var}*}</code> - bind variables in the context of the experiment
<code>:script</code>	<code>{{(descriptive-name initial-time [next-time] form) — script-element-name}*}</code> <i>initial-time</i> can be a string, a number or form to evaluate. <i>next-time</i> is a time interval and should be a fixnum or form.

In addition, `define-experiment` accepts all the options accepted by `define-simulator`.

The following steps are generated from the `define-experiment` specification to run a typical experiment. After initializing the system using the `:before-experiment` form, it loops over the cross-product of the specified experiment variables, running one trial for each element in the cross-product. This is repeated for as many times as is specified in the `:repetitions` option. After the last trial, the `:after-experiment` form is run.

```
Run :BEFORE-EXPERIMENT code with args
LOOP
  Update experiment-variables for trial
  Run :BEFORE-TRIAL code with experiment-variables and args
  Run :RESET-SYSTEM code with experiment-variables and args
  Instantiate Script Events
  Reset and Enable all the instrumentation
  Run :START-SYSTEM code with experiment-variables and args
  <system runs (possible periodic and event based collection done)>
  Run :AFTER-TRIAL code (possible post-hoc collection done)
END LOOP when all trials completed
Run :AFTER-EXPERIMENT code with args
```

2.2.3 Write-current-experiment-data Function

`write-current-experiment-data` (*ℰkey separator format instrumentation stream*) [Function]

Causes each experiment instrumentation to write its data to *stream* or the output-files specified in the `run-experiment` call or in a `defclip`. *separator* should be a character which will be used to separate fields. It defaults to the value of `*data-separator-character*`. *format* should be one of `:CLASP` which means write a clasp native format data file, `:ASCII` which means write a standard separator delimited data file including column names or `:DATA-ONLY` which is the same as `:ASCII` except no column names are included. *format* defaults to the value of `*output-format*`. *instrumentation* can be used to specify a subset of the experiment's instrumentation to write to the data file.

2.2.4 Run-experiment Function

`run-experiment` (*experiment-name ℰkey args repetitions number-of-trials length-of-trial output-file error-file extra-header starting-trial-number*) [Function]

`Run-experiment` starts execution of the experiment named *experiment-name*. The *args* are passed on to the experiment. *output-file* is optional, but must be specified if `write-current-experiment-data` is called from within your experiment. *number-of-trials* can be used to specify an exact number of trials to run. If *number-of-trials* is not specified it will be calculated so as to vary all the experiment variables across all their values *repetitions* (default 1) times. *starting-trial-number* can be used to change this value to something other than one (1) which is useful for continuing partially completed experiments. *length-of-trial* can be specified for time-based simulators to put a limit on the maximum trial length. *error-file* can also be used to direct the error/debug output to a file. *extra-header* is written at the end of the header of the output file.

2.2.5 Explicitly Stopping Trials

The following functions can be called from user code to explicitly terminate a CLIP trial or experiment. Usually this is done from within the code of a script or as part of an error handler. They all take no arguments and assume that an experiment is in progress.

`shutdown-and-rerun-trial ()` [Function]

This will cause the current trial to be aborted (no data written) and restarted. This is a function that users should call when they have detected an error condition of some sort that renders the trial worthless, but rerunning the trial may work.

`shutdown-and-run-next-trial ()` [Function]

This will cause the current trial to be stopped (data will be written) and the next trial started. This is a function that users should call when want to normally shutdown a trial and collect and report the data from that trial.

`shutdown-experiment ()` [Function]

This will cause the current trial to be aborted (no data will be written) and will return the system to the state it was before the experiment began (by running the after-experiment code).

2.3 Clip Definition

There are basically only a small number of ways to instrument a software system. These are: adding special purpose code to collect data, interrogating global (usually objects with state) data structures or using advice.

The first way is to build in special purpose code to collect data about the state of the system while the system is running. We used this technique in the Phoenix_j testbed to keep information about the execution time of timeline entries and also message traffic patterns. The Transsim simulator also uses this technique when it keeps track of daily costs, ship locations, demon firing intervals. The key point here is that the only reason for adding the piece of code was to allow an experimenter to analyze the behavior of the system. The simulation itself does not use the information. This method increases the complexity and reduces the readability and maintainability of the software system. In code that is highly instrumented it is often difficult to determine what is intrinsic to the running of the system and what is used to instrument.

Another method involves interrogating global data structures to determine the state *post hoc*. Examples include most everything we collected about Phoenix (ie., fireline built, first plan tried, bulldozer digging times, etc.). This technique is fine if the information is hanging around after you have finished running the simulation, but this is not always the case. Also, there is some information that is time-sensitive, and therefore must be collected on a periodic or event-driven basis. Collecting this type of information often involves resorting to the first method – altering the code itself.

Alternatively, one can use the `advise` facility available in many lisp environments to non-intrusively collect information when functions are called. This is a nice modular approach, but requires a good deal of knowledge about the advise facility and the software system itself. Unfortunately, the `advise` facility is not standardized across Lisp implementations.

The `defclip` macro encapsulates the code necessary to instrument a software system and separates it from the system itself. It avoids the pitfalls of adding special purpose code

directly to the software system while at the same time allowing periodic and event-driven data collection. It also allows collection to be done by perusing global data structures.

2.3.1 Simple Clips

Simple clips have no components and collect data immediately prior to reporting it to the output file at `:after-trial` time. If they are defined with a `:schedule` or `:trigger-event` `defclip` option their default behavior is store all of the data collected during a trial and report a single value which is the mean of all the values collected.[†]

2.3.2 Clips with Components

Clips with components, as specified by the `:components` keyword, generate multiple columns in a data file each time they are reported. Depending on other options they may produce one column per component (*composite clips*) or multiple columns per component (*mapping clips*). Mapping clips are specified using the `:map-function` option to `defclip`. Clips with components are sometimes referred to as *super clips*. For a good example of clips with components and further discussion of their use, see Appendix B.1.1.

2.3.3 Time-series Clips

Clips that have `:components` and either the `:schedule` or `:trigger-event` option are *time-series clips*. They generate multiple data columns in the manner of component clips (which they are) and also multiple data rows. Each row corresponds to a single collection and is either triggered by a particular event or activated periodically. Since time-series clips generate multiple rows, they are generally written to a data file that is separate from the main experiment (summary) data file. The name of the data file associated with a time-series clip is specified using the `:output-file` option to `defclip`.

The `:schedule-function`, `:seconds-per-time-unit`, and `:timestamp` keywords to `define-simulator` must be specified for periodic time-series clips. Event-based time-series clips require only that the Common Lisp implementation provide some mechanism similar to the `advise` function.[‡]

2.3.4 Defclip Macro

`defclip` (*name args (options) &rest body*) [Macro]

A `defclip` form defines and names a particular data-collection instrument. Each clip collects a single variable which, together with other variables, make up a dataset. The dataset can then be analyzed by CLASP. The form may contain arbitrary code to gather data for the clip: access to global variables, function calls, and so forth. Naming the clips makes it easy to turn them on and off programatically or with menu selections.

options is a list of keywords and values which modify the main form implemented by the clip. *options* can be any of the following, all of which are optional:

[†]The default behavior if the collected values are non-numeric is to generate an error.

[‡]Actually, one can explicitly call the undocumented (until now), unexported 'clip::collect' function on a particular instrumentation and achieve the same effect as an event-based clip, but requires modifying the code.

<code>:schedule</code>	a list of keyword/value pairs; currently allowed are <code>:period</code> which provides the time interval between collections for time-series data, and <code>:start-time</code> which specifies the time of the first collection; if neither <code>:schedule</code> nor <code>:trigger-event</code> is specified, collection will be done immediately prior to the report being output to the data stream
<code>:trigger-event</code>	a function, list of functions or list of trigger specifications which trigger event-driven collection ; each trigger spec should of the form (<code><fname> [:BEFORE :AFTER] [:PREDICATE <fname>]</code>)
<code>:components</code>	a list of clips that are associated with this clip; ie., they are collected and reported as a group by collecting or reporting this clip
<code>:map-function</code>	provides a list of items to map the <code>:components</code> over ; this function should return the same arguments in the same order each time it is called
<code>:report-key</code>	allows overriding of the column header key which is written to the data stream; avoid using this for component clips unless you really know what you are doing as the default key will be generated to correctly differentiate between multiple invocations of the same clip with different parents; this string is used in a call to 'format'; for most component clips this string should handle the same arguments as the clip
<code>:initial-status</code>	whether the clip is enabled or disabled by default
<code>:report-function</code>	can be used to override the default report function; for expert users only
<code>:enable-function</code>	code ^a to set up data structures for the clip; runs once when the clip is turned on
<code>:disable-function</code>	code to remove data structures set up for the clip; runs once when the clip is turned off
<code>:reset-function</code>	code to reinitialize data structures at the beginning of each trial; for example, setting counters back to zero
<code>:display-function</code>	code for graphical display of the information
<code>:output-file</code>	used to specify the output file for time-series clips - merged with the pathname specified in the <code>run-experiment</code> call; should be a string or a function of one arg; the function is called with the name of the clip and should return a value suitable for use as the first argument to <code>merge-pathnames</code>

^a'code' at this level is specific to the user's system, and is assumed to be user-supplied Lisp code. For example, enabling a counter of events in a system may require creation and initialization of a counter variable.

Some examples of `defclip` usage:

```
A simple clip with code used to report a value:
(defclip number-of-dead-bulldozers ()
  (length (bds-that-died)))
```

An example showing a clip with components:

```
(defclip methods-each-bd (bulldozer)
  "Salient information for each instance of applying a recovery method:"
  (:components (trial-number
                 agent-name
                 method-type
                 failure-type
                 calculate-recovery-cost
                 method-succeeded-p
                 order-of-application)
   :map-function (gather-recovery-method-instances (name-of bulldozer)))
  ;; This code executes before the map-function is executed.
  (send (fire-system) :set-frame-system (name-of bulldozer)))

(defclip ii-projection-attempts-made ()
  "Clip for reporting the number of attempts that were made by
ii-projection during the last projection of the trial."
  (:reset-function (setf *ii-projection-attempts-made* nil))
  (car *ii-projection-attempts-made*))
```

2.4 Examples from a Simple Agent Simulator

The following three examples illustrate experiment control and data collection using CLIP. The agents in this simulator stochastically increment their state until they reach a final state. The first example defines the experimental interface to an agent simulator and two simple clips for collecting data. The second and third examples build on the first, defining more complicated clips for mapping a simple clip onto multiple agents and for collecting time-series data.

2.4.1 Using Simple Clips to Collect Trial Summary Data

This example collects trial summary data about overall agent-cost and task completion-time. Collection occurs at the end of each trial and is written out to a summary file in CLASP format. The `define-simulator` form designates methods for starting the simulation at the beginning of the experiment, resetting and reinitializing it if necessary, and stopping it after the final trial. The `define-experiment` form designates the simulator to use, in this case the one we have defined. It also specifies: a set of experiment variables and the experimental settings for each, a list of the clips we've defined under `:instrumentation`, initializations for global variables before each trial, and a function for writing out a row of data at the end of each trial.

For each trial, the trial number (assigned sequentially) and the value of each experiment variable are written out to the summary file along with the values of the specified clips. An example of the CLASP output file produced by this experiment is shown in Figure 2.1.

In this first example we have defined two simple clips. One collects the combined cost of all agents at the end of each trial and the other records the time at which the trial ends.

```
;;;-----
```

```

(define-simulator agent-sim-1
  :start-system (run-agent-simulation :reset nil)
  :reset-system (reset-agent-simulation)
  :stop-system stop-simulation)

;;;-----
;;; Clip Definitions

(defclip agents-cost ()
  ()
  (reduce #' + (find-agents) :key #'cost))

(defclip completion-time ()
  ()
  (current-time))

;;; *****
;;; The Experiment Definition

(define-experiment simple-agent-experiment-1 ()
  "A test experiment."
  :simulator agent-sim-1
  :variables ((transition-probability in '(.01 .1))
             (cost-factor from 1 to 3))
  :instrumentation (agents-cost completion-time)
  :before-trial (setf *transition-probability* transition-probability
                    *relative-cost* cost-factor)
  :after-trial (write-current-experiment-data))

#| Execute this to run the demo experiment.

(run-experiment 'simple-agent-experiment-1
  :output-file
  #+Explorer "ed-buffer:data.clasp")
|#

```

2.4.2 Using a Mapping Clip to Map Simple Clips Over Multiple Agents

This example defines a mapping clip that maps over all the agents at the end of each trial and returns the cost accrued by each. It produces a summary output file like that in Figure 2.2 recording after each trial an entry that includes trial number, each experiment variable, each of three agent's cost, and the completion time of the trial.

```

;;;-----
;;; Mapping clip to collect cost of each agent

(defclip all-agents-costs ()

```

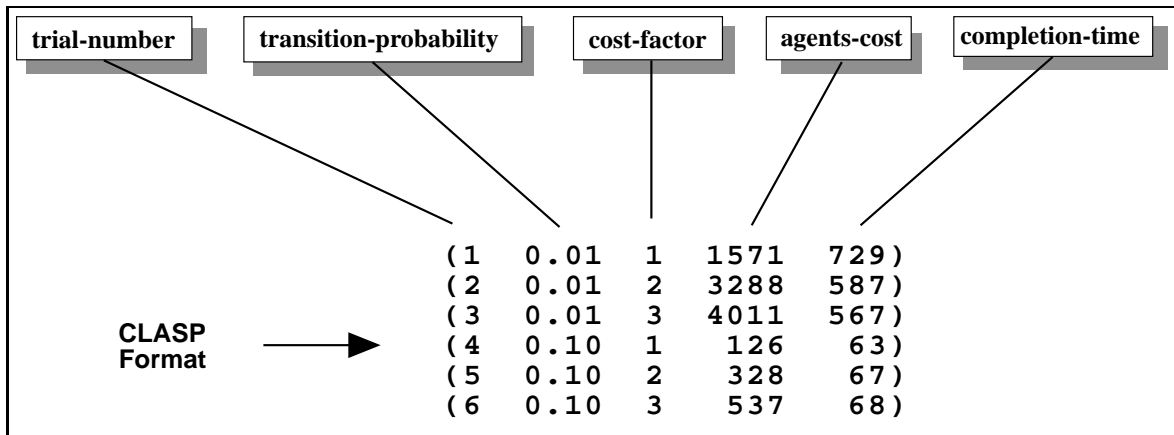



Figure 2.1: Simple Clip Example from the Agent Experiment.

```
(:map-function (find-agents)
 :components (each-agent-cost)))

(defclip each-agent-cost (agent)
 ()
 (cost agent))

;;; *****
;;; The Experiment Definition

(define-experiment simple-agent-experiment-2 ()
 "A test experiment."
 :simulator agent-sim-1
 :variables ((transition-probability in '(.01 .1))
 (cost-factor from 1 to 3))
 :instrumentation (agents-cost all-agents-costs completion-time)
 :before-trial (setf *transition-probability* transition-probability
 *relative-cost* cost-factor)
 :after-trial (write-current-experiment-data))

#| Execute this to run the experiment.

(run-experiment 'simple-agent-experiment-2
 :output-file
 #+Explorer "ed-buffer:data.clasp")

|#
```

2.4.3 Full Agent Simulator Experiment with Time-Series Clips

This example uses both periodic and event-driven time-series clips to collect data about each agent's state. To use time-series clips the experimenter must supply CLIP with enough information to schedule data collection *during* trials. This is done by specify-

trial-number	transition-probability	cost-factor	agents-cost	all-agents-costs :map (agent-1 agent-2 agent-3)	completion-time
(1	0.01	1	1370	390 622 358	622)
(2	0.01	2	2108	746 364 998	499)
(3	0.01	3	4938	1299 2772 867	924)
(4	0.1	1	148	36 59 53	59)
(5	0.1	2	280	130 80 70	65)
(6	0.1	3	627	192 258 177	86)

Figure 2.2: Mapping Clip Example from the Agent Experiment. Mapping clips are used to map one or more simple clips over multiple objects.

ing in `define-simulator` a `:schedule-function` that tells `defclip` how to schedule clip execution. The optional `:deactivate-scheduled-function` provides a function for un-scheduling clips (if necessary). `:seconds-per-time-unit` tells CLIP how to translate the time units of the simulator into seconds. The value of `:timestamp` is used to automatically record the time of collection in each row of a time-series data file.

This example illustrates three uses of clips: collecting summary data about the highest agent state, gathering periodic snapshots of all agent states, and recording particular aspects of the agents' changes of state. The periodic clip definition specifies an output file and a scheduling interval. Two event-driven clips are triggered in this example by the same function, `change-of-state-event-function`, which looks for any change of agent state. However, since their component clips are not identical, their output is routed to separate files (CLASP file format expects all rows to have the same number of columns and column names).

Four output files are produced by this experiment. As in the previous examples, a summary file records a row of data at the end of each trial. Each change in an agent's state triggers event-driven clips that record a row in a time-series output file, as in Figure 2.3. The event-driven clip definitions `change-of-state` and `change-of-state-pred` produce one output file each. The fourth output file also records time-series data, one row for each periodic snapshot of all agents' states (cf. Figure 2.4).

Collecting time-series data is more time-consuming than collecting summary data. In the first two examples, where only summary data was collected, trials were allowed to run to completion. In the third example periodic clips run every 12 minutes and event-driven clips run as frequently as agent's states change. We can limit the duration of trials by specifying a `:length-of-trial` value to `run-experiment`. This is useful for statistical tests on time-series data, such as cross-correlation, that expect the number of rows for each trial to be the same. Thus, in this example, specifying `:length-of-trial` to be 500 minutes ensures that the periodic data collection will produce the same number of rows for each trial.

```
;;;-----
(define-simulator agent-sim
  :start-system (run-agent-simulation :reset nil)
```

```

:reset-system (reset-agent-simulation)
:stop-system stop-simulation
;; a function that places functions to run on the queue of events.
:schedule-function (lambda (function time period name &rest options)
                    (declare (ignore name options))
                    (schedule-event function nil time period))
;; a function that removes functions from the queue of events.
:deactivate-scheduled-function unschedule-event
:seconds-per-time-unit 60
:timestamp current-time)

;;;-----
;;; Clip Definitions

;; This post-hoc clip produces two values.
(defclip highest-agent-state ()
  (:components (highest-state highest-agent))

  (loop
    with agents = (find-agents)
    with highest-agent = (first agents)
    with highest-state = (state highest-agent)
    for agent in (rest agents)
    for agent-state = (state agent) do
      (when (state< highest-state agent-state)
        (setf highest-state agent-state
              highest-agent agent))
      finally (return (values highest-state highest-agent))))

;;;-----
;;; Periodic collection

;; This clip invokes its component every 12 minutes.
(defclip periodic-agent-state-snapshot ()
  (:output-file "snapshot.clasp"
   :schedule (:period "12 minutes")
   :map-function (clip::find-instances 'agent)
   :components (each-agent-state-snapshot)))

;; Simple clip that returns the state of the agent.
(defclip each-agent-state-snapshot (agent)
  "Record the state at an agent."
  ()
  (state agent))

;;;-----
;;; Event-driven collection

```

```

;; This clip accepts the arguments passed to 'change-of-state-event-function'
;; and simply passes them through.
(defclip change-of-state (agent-name new-state)
  (:output-file "state-change.clasp"
   :trigger-event (change-of-state-event-function :BEFORE)
   :components (new-state agent-name))
  (values new-state agent-name))

;; This clip accepts no arguments and returns two values computed by other
;; functions.
(defclip change-of-state-pred ()
  (:output-file "state-change-pred.clasp"
   :trigger-event (change-of-state-event-function :AFTER)
   :components (fred barney))
  (values (compute-fred) (compute-barney)))

;;; *****
;;; The Experiment Definition

(define-experiment agent-experiment ()
  "A test experiment."
  :simulator agent-sim
  :variables ((transition-probability in '(.01 .1))
             (cost-factor from 1 to 5 by 2))
  :instrumentation (agents-cost
                   all-agents-costs
                   completion-time
                   highest-agent-state
                   change-of-state-pred
                   change-of-state
                   periodic-agent-state-snapshot)
  :before-trial (setf *transition-probability* transition-probability
                    *relative-cost* cost-factor)
  :after-trial (write-current-experiment-data))

;; Execute this to run the demo experiment.

(defun rexp ()
  (run-experiment 'agent-experiment
                 :output-file "data.clasp"
                 :length-of-trial "500 minutes"))

```

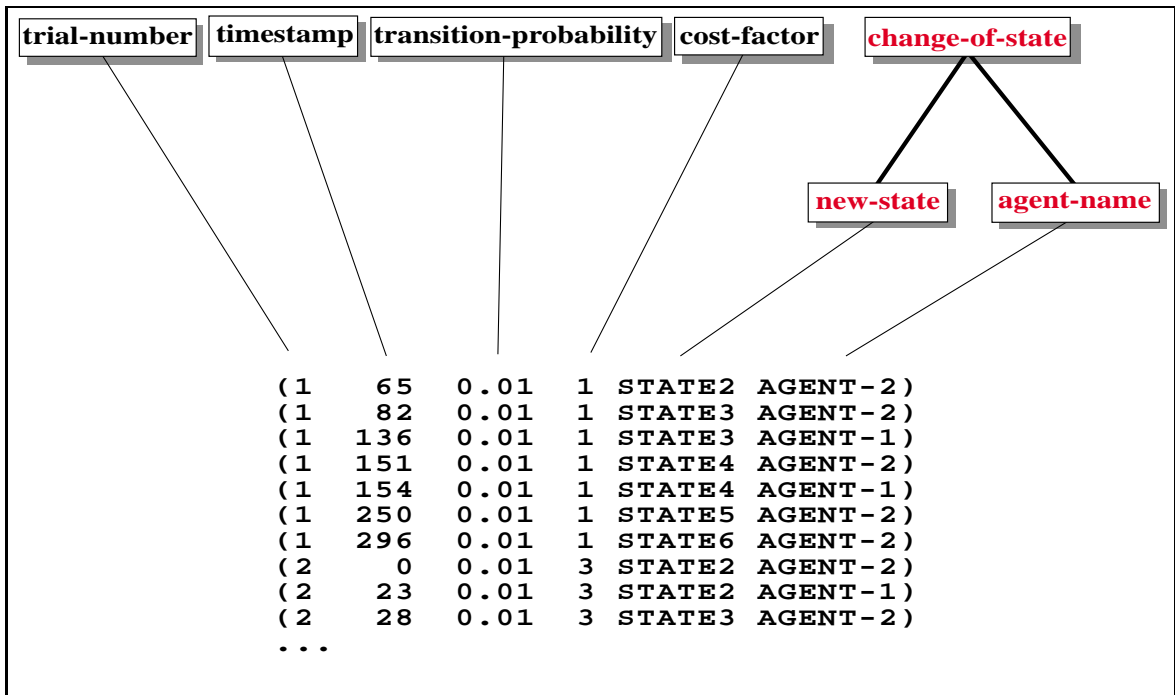


Figure 2.3: Event-Driven Time-Series Clips are used to collect data when an agent’s state changes in the Full Agent Experiment.

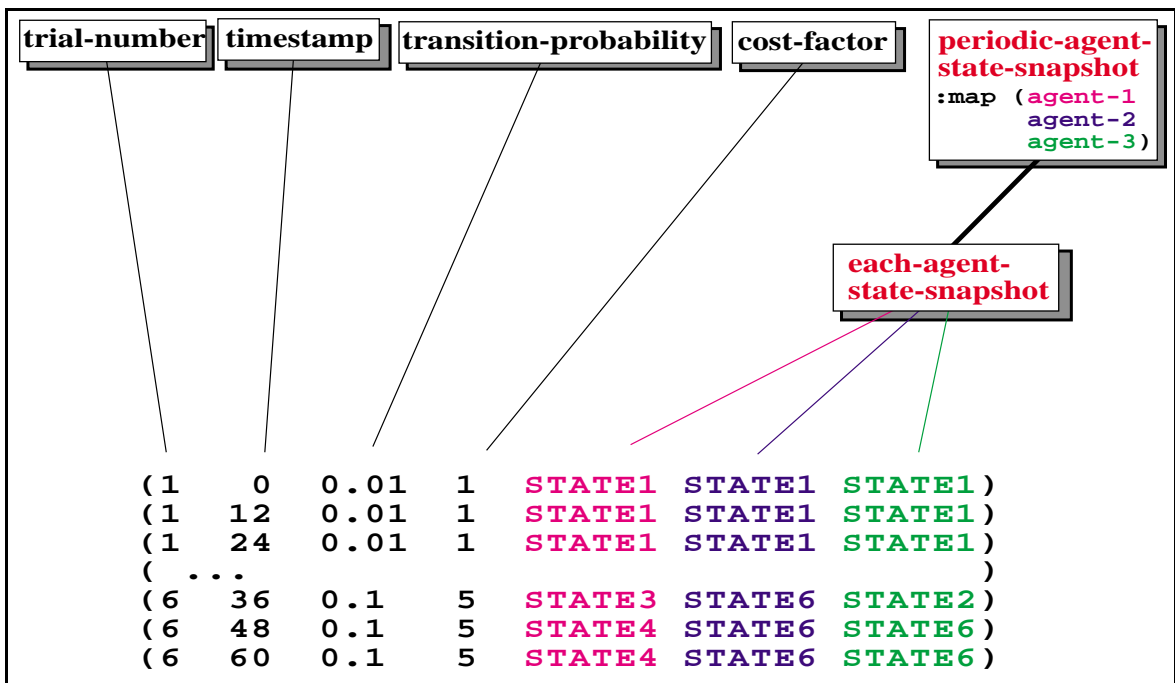


Figure 2.4: Periodic Time-Series Clips collect snapshots of each agent’s state at regular intervals during the Full Agent Experiment.

Chapter 3

Graphical User Interface

This chapter describes the layout and use of the CLASP graphical user interface. The interface is built using the Common Lisp Interface Manager (CLIM), so many CLIM features are available in CLASP.

The first section describes the layout of CLASP. The second section discusses CLASP command syntax. The next section explains how to interact with CLASP using the mouse and keyboard. Some examples of CLASP usage follow. Finally, each of the command menus are described.

3.1 Clasp Layout

The main window of CLASP is called the CLASP *frame*. This frame contains a number of sections called *panes*. Information in CLASP is displayed in one of four areas.

Menu Bar The CLASP menus will appear across the top of the CLASP frame. The menus are: File, Graph, Describe, Manipulate, Transform, Test and Sample. Clicking on a menu name will pop up the menu.

Datasets Display This contains the names of all Datasets and Variables currently loaded into CLASP. Dataset names are aligned to the left margin; variable names for each dataset are indented below the dataset name.

Results Display A mouse-sensitive icon for each CLASP result appears in this window.

Notebook CLASP commands and lisp expressions are entered here.

While the Menu Bar and Notebook are always in the CLASP frame, the Datasets Display and Results Display are under user control. They can either be panes at the left edge of the CLASP frame, or in separate windows. This is controlled through the `:PREFERENCES` command.

3.2 Clasp Command Syntax

3.2.1 Basic syntax

Commands have the form `:VERB arg1 arg2 ... argn`. `:VERB` is the name of the command and `arg1 arg2 ... argn` are the arguments it accepts. Note that when being entered, commands are prefixed by a colon character (`:`).

Some examples of command usage:

:Load Dataset (Pathname) **my-data.clasp**

:Rename Variable (Variable) **Var-1** (New name) **Numnodes**

:Merge Datasets (Datasets) **Dataset-1, Dataset-2**

3.2.2 Mapped arguments

Many CLASP commands give the user the option of entering multiple values for their arguments. When this happens, the command will be executed for each value of the argument it is given. This is called argument mapping. If a command has more than one argument that allows mapping, each argument must be given the same number of values and multiple values will be mapped over in parallel. That is, the command will be executed on all the first values of its arguments, then all the second values, etc.

Some examples of mapped arguments are:

:Mean (X(s)) **Var-1, Var-2, Var-3**

Will return the mean of Var-1, Var-2 and Var-3.

:Anova One Way Variables (Y(s)) **Y1, Y2** (X(s)) **X1, X2**

This will do a one-way anova of Y1 on X1 and a one-way anova of Y2 on X2.

:Correlation (Y(s)) **Y1, Y2** (X(s)) **X1**

This is illegal, X and Y are both mapping arguments, but they were given different numbers of inputs.

The following example shows an instance where not all the arguments of a command are mapped arguments. In :T TEST TWO SAMPLES, although the X and Y arguments map, the *Tails* argument doesn't. Whatever single value is entered for *Tails* will be used for every iteration of the other arguments. The prompt for an argument that may be mapped over will end in (s).

:T Test Two Samples (Y(s)) **Y1, Y2** (X(s)) **X1, X2** (Tails: [Both, Positive or Negative])
Positive

This will perform two 2 sample t-tests, one on X1 and Y1 and another on X2 and Y2. Both tests will be Positive tail tests.

3.3 Interacting with Clasp

At the command prompt (\Rightarrow), CLASP will accept either a CLASP command or a lisp expression. If a command is entered, its arguments will be prompted for and then the command will be run. If a lisp expression is typed, it will be evaluated and the value printed as in a normal lisp listener. When CLASP is looking for a particular kind of input, a variable, a string, etc. it is said to be *accepting* that kind of input. When accepting a particular type of input, CLASP makes it easy to enter that type of data by making objects of that type mouse selectable and also by doing completion on the names of those objects. For instance, when CLASP is accepting a variable as an argument to a statistical command, all CLASP variables will be selectable. In addition any sequence of numbers being displayed in the CLASP frame will be selectable, as will any lisp symbol which evaluates to a sequence of numbers.

3.3.1 Using the mouse

To enter a command, click on the menu in which the command resides and then click on the command. When a command has been selected and is waiting for arguments to be entered, any object which is a valid argument will become mouse selectable.

To select multiple objects, either for a mapped argument or a command which accepts multiple objects (such as `:MERGE DATASETS`), click on each object.

At the command prompt, clicking on a previously entered command or lisp expression will re-execute it. Clicking on a CLASP object (a Dataset, Variable or Result) will cause it to `:OPEN`, clicking on an open object will cause it to `:CLOSE`.

While entering a lisp expression, clicking on any other lisp expression will insert it at the cursor. Clicking on a CLASP variable will insert its value into the current expression as a quoted list.

3.3.2 Using the keyboard

As with the mouse, whenever CLASP is accepting a variable, the name of any CLASP variable may be typed. In addition, the name of any lisp symbol which has a sequence of numbers as its value may be typed. To enter multiple objects, separate them by commas.

CLASP also performs completion on typed input. When entering a command or the name of a variable, the Tab key will do partial completion, filling in as much information as possible until there are two possibilities which will be distinguished by the next character. For instance, since there are a number of sampling commands, all of which start with `:SAMPLE`, typing `:SAM` followed by a Tab will complete the word `:Sample`. The Space key completes the current word or object, and then prompts for the next piece of input (either the next word of a command, or the next argument.) The Return key indicates to CLASP that the input is complete and ready for processing. Returning to the `:SAMPLE` commands, if the user types `:SAMPLE U` (which uniquely determines the command `:SAMPLE UNIFORM`), then the three completion characters will behave differently. Space will complete the command and prompt the user for the first argument. Tab will complete the command, but it will not prompt for the first argument. Return will complete not only the command, but the first few arguments, which have default values, and then it will try and process the input line. In general, a good rule of thumb is: “When in doubt, hit the Space key, if that doesn’t do what you want, hit the Backspace key and go from there.”

Since CLASP uses the CLIM input editing facility, most of the normal CLIM editing keys will work. The following is a table of editing keys available at the CLASP command line.

When typing a lisp expression in CLASP, the bang key (!) can be used to reference CLASP variables. The expression `!variable-name` will substitute the value of the variable into the expression. If more than one variable has the same name,

`!variable-name!dataset-name`

will uniquely determine a variable, and its value will be used. If there are multiple variables with the same name, and the dataset is not specified, the matching variable belonging to the *most recently created dataset* will be used. When a final close parenthesis is entered in a lisp expression, the expression will immediately be evaluated – there is no need to hit Return.

A useful trick for editing a lisp expression from a previous line is:

Forward character	control-F
Backward character	control-B
Beginning of line	control-A
End of line	control-E
Delete next character	control-D
Delete previous character	Rubout
Kill to end of line	control-K
Transpose adjacent characters	control-T
Yank from kill ring	control-Y

Table 3.1: Keys used for editing input

1. Type an open parenthesis - This will tell CLASP that you are entering an expression.
2. Click on the line to edit - This will insert it into the current line.
3. Edit the line- Make whatever changes are desired using the keyboard editing characters described below.
4. Remove initial parenthesis - Go to the beginning of the line (you can use Control-A) and delete the extra open parenthesis (Control-D.)
5. Move to the end of the line (Control-E) - This will cause the new line to be entered.

3.4 Examples of Clasp command usage

Some examples of other CLASP commands follow. These commands make a new dataset, define a new variable for the dataset and rename the new variable.

:Make Dataset from Rows (Name) *Exper-1* (Data) **A, 13, 431; A, 10, 389; B, 11, 214; B, 10, 203 (Variable names) **Method, Time, Score****

Create a new dataset called *EXPER-1* with three variables, *METHOD*, *TIME*, and *SCORE*, and 4 rows of data. Note that in entering data, single items are separated by a comma, rows are separated by a semicolon.

:User Defined (Dataset) *exper-1* (Expression) (/ **score time)**

USER DEFINED is a command from the *Transform* menu which allows arbitrary transformations of variables. In this case, the command will add a new variable to *EXPER-1* which divides *SCORE* by *TIME*.

:Rename Variable (Variable) /-score-time** (New name) **points-per-second****

Commands from the *transform* menu compute a name for the new variable by using the expression that created it. Often the user will want to give the new variable a more semantically appropriate name. This command rename the newly created variable, */-SCORE-TIME* to *POINTS-PER-SECOND*.

3.5 Clasp Command Menus

The CLASP command menu pane has seven menus, each of which is described in this section. Clicking on a menu name will pop up the menu, left clicking on a menu item will select that item, right clicking on it will bring up help for that item.

3.5.1 The File Menu

File
Load Dataset
Import Dataset
Save Dataset
Export Dataset
Clear All
Clear Notebook
Save Notebook
Print
Preferences
Quit

File menu commands are used to manage datasets and the CLASP notebook, and to print results.

`:LOAD DATASET (Filename)` [*Command*]

Loads a dataset from a disk file in CLASP format.*

`:IMPORT DATASET (Filename) (Separator) (Include variable names)` [*Command*]

Loads a dataset from a disk file in columnar format. *Separator* is the character used to separate columns in the file. If *Include variable names* is Yes, then the first line of the file will be interpreted as variable names.†

`:SAVE DATASET (Filename)` [*Command*]

Saves a dataset to a disk file in CLASP format.

`:EXPORT DATASET (Dataset) (Filename) (Separator) (Include variable names)` [*Command*]

Saves a dataset to a disk file in columnar format. As with `:IMPORT DATASET`, *Separator* is the character used to separate columns. If *Include variable names* is Yes, then the variable names will be written to the first line of the file.

`:CLEAR ALL` [*Command*]

Clears the CLASP notebook and all datasets and results.

`:CLEAR NOTEBOOK` [*Command*]

Clears the CLASP notebook, but doesn't delete datasets and results.

`:SAVE NOTEBOOK` [*Command*]

Saves the current notebook to a postscript file.

`:PRINT (Object) (Filename)` [*Command*]

*CLASP file format is described in 6.3.2.

†For more information on importing and exporting files, see 6.3.2.

Object may be a CLASP dataset or the result of a CLASPcommand. The output is written in postscript format to *Filename*. Because of limitations in Clim 1.1's postscript support, printing graphs is difficult. The following procedure will produce reasonable postscript output for graphs.

1) Unless the graph will be printed with a color postscript printer, set the color of each dataset in the graph to black. This can be accomplished by clicking on the datasets in the graphs legends and using the window that pops up to set the dataset color. (Note, since SCIGRAPH displays graphs with a black background, the datasets will not be visible on screen.)

2) Print the graph to a file using the `:PRINT` command.

3) Run the postscript file through the sed script `ps-fix.sed`, to do this type

```
sed -f ps-fix.sed file.ps > new-file.ps
```

at the shell prompt. In the above example, `file.ps` is the name of the file that was created by the `:PRINT` command and `new-file.ps` is the name of the file created by sed. The file `ps-fix.sed` is available via anonymous ftp from the same location as CLASP.

`:PREFERENCES`

[*Command*]

Brings up a window in which the user can modify the behavior of CLASP. The window has three sections, the results display preferences, the error handling preferences and the datasets and results display preferences.

In the first section of the preferences window, the user can specify how a number of result types are displayed in CLASP. Each type will be followed by four options, `Display-In-Interactor`, `Display-In-Window`, `Iconify` and `NIL`. Clicking on an option will select it. If `Display-In-Window` is selected, then when results of that type are created, they will be displayed directly in the notebook. If `Display-In-Window` is selected, then results of that type will initially be displayed in a separate window and an icon will be placed in the notebook. If `Iconify` is selected, then an icon will be placed in the notebook and the result will not be displayed until the icon is opened. If `NIL` is selected, then nothing will be shown in the notebook. In all cases, the results window will always show all existing results. The default setting for all result types is `Display-In-Interactor`.

The available error handling preferences are `Use-Native-Debugger` and `Trap-Underflow-Errors`. When `Use-Native-Debugger` is not selected, CLASP captures lisp errors and prints a message in the notebook. When it is selected CLASP passes lisp errors down to the underlying lisp process. This can be useful when debugging code. When `Trap-Underflow-Errors` is not selected, CLASP will always convert floating point underflow errors to 0.0. Setting this parameter will cause CLASP to signal an error when a floating point underflow occurs. The default value for both error handling preferences is unselected.

The `Display datasets and results` preferences controls where the datasets and results lists are displayed. The `Pane` setting will display these lists directly in the CLASP frame, to the left of the notebook. The `Window` setting will cause them to be displayed in separate windows. This could be useful if when there are many datasets or variables with long names. Most window managers allow resizing of windows, so with this option, the user may set the datasets and results listings to a convenient

size. The final setting is None which will suppress displaying these lists altogether.

:QUIT [Command]

Exits from CLASP. It is up to the user to save whatever datasets were created or changed during the session. CLASP will not save these automatically, nor does it prompt the to save datasets before quitting.

3.5.2 The Graph Menu

Graph
Histogram
Scatter Plot
Line Plot
Row Line Plot
Regression Plot
Overlay Graphs

The Graph menu has commands which allow the creation of graphical views of data.

Graphs in CLASP are created by the SCIGRAPH system from BBN. The elements of these graphs are mouse sensitive. To change the style or color of a line or point, select the entry for that data from the legend. This will pop up a window which allows for editing of data attributes. For a more detailed description of SCIGRAPH and its functionality, see the SCIGRAPH documentation, available via anonymous ftp from cambridge.apple.com in the directory /pub/clim/clim-1-and-2/scigraph.

:HISTOGRAM $(X(s))$ (*Color by*) [Command]

Creates a histogram from the variable X . Unless *Color by* is None, the points in X are grouped by the *Color by* variable and each group is shown in a different color. To use the *Color by* option, both X and *Color by* must be CLASP variables and they must come from the same dataset. *Color by* is not a mapped variable, if multiple variables are selected for X , the same *Color by* option will be used for all of them.

:SCATTER PLOT $(Y(s))$ $(X(s))$ (*Color by*) [Command]

Creates a scatter plot of Y on X . Coloring is the same as for :HISTOGRAM, except that since there are two variable arguments, they must both be from the same dataset as *Color by* for coloring to be used.

:LINE PLOT $(Y(s))$ $(X(s))$ (*Color by*) [Command]

Creates a line plot of Y on X . Coloring is the same as for :SCATTER PLOT.

:ROW LINE PLOT $(Y(s))$ (*Color by*) [Command]

Creates a line of Y against 1, 2, 3 ... N, where N is the length of Y . Coloring is the same as for :HISTOGRAM.

:REGRESSION PLOT $(Y(s))$ $(X(s))$ [Command]

Creates a scatter plot of Y on X and overlays the regression line from the linear regression of Y on X .

:OVERLAY GRAPHS (*Graphs*) [Command]

Graphs is a sequence of graphs. :OVERLAY GRAPHS produces a single graph with each member of *Graphs* overlaid. :OVERLAY GRAPHS will attempt to choose high contrast colors so that points from different source graphs may be distinguished.

To produce colored graphs when the original data are not variables which are in the same dataset as the coloring variable, partition the data by hand, produce a single graph for each partition, and overlay the graphs.

3.5.3 The Describe Menu

Describe
Data Length
Mean
Variance
Standard Deviation
Minimum
Maximum
Range
Quantile
Median
Trimmed Mean
Mode
Interquartile Range
Statistical Summary
Covariance
Correlation
Cross Correlation
Autocorrelation

The Describe menu contains commands which produce descriptive statistics on data. Since most of these commands are described in section 4.1, only their arguments are given here.

:DATA LENGTH ($X(s)$)	[Command]
:MEAN ($X(s)$)	[Command]
:VARIANCE ($X(s)$)	[Command]
:STANDARD DEVIATION ($X(s)$)	[Command]
:MINIMUM ($X(s)$)	[Command]
:MAXIMUM ($X(s)$)	[Command]
:RANGE ($X(s)$)	[Command]
:QUANTILE ($X(s)$) (<i>Percentile</i>)	[Command]

Note that *Percentile* isn't a mapped argument. It must be a real number from 0 to 1 inclusive, and the same value will be mapped over all X 's.

:MEDIAN ($X(s)$)	[Command]
:TRIMMED MEAN ($X(s)$) (<i>Trimming factor</i>)	[Command]

Note that *Trimming factor* isn't a mapped argument. It must be a real number from 0 to .5, and the same value will be mapped over all X 's.

:MODE ($X(s)$)	[Command]
:INTERQUARTILE RANGE ($X(s)$)	[Command]

- `:STATISTICAL SUMMARY (X(s))` [Command]
 Displays a variety of commonly used descriptive statistics.
- `:COVARIANCE (Y(s) (X(s)))` [Command]
- `:CORRELATION (Y(s) (X(s)))` [Command]
- `:CROSS CORRELATION (Y(s) (X(s)) (Min lag) (Max lag))` [Command]
 Note that *Min lag* and *Max lag* are not mapped arguments.
- `:AUTOCORRELATION (Y(s) (X(s)) (Min lag) (Max lag))` [Command]
 Note that *Min lag* and *Max lag* are not mapped arguments.

3.5.4 The Manipulate Menu

Manipulate
Rename Dataset
Rename Variable
Delete Result
Make Dataset From Rows
Make Dataset From Columns
Add Variable To Dataset
Partition Dataset
Partition On
Merge Datasets
Open
Close
Describe

The Manipulate menu has commands for manipulating data within a dataset and deleting, renaming, opening, closing and describing data objects.

- `:RENAME DATASET (Dataset) (New name)` [Command]
 Changes the name of *Dataset* to *New name*.
- `:RENAME VARIABLE (Variable) (New name)` [Command]
 Changes the name of *Variable* to *New name*.
- `:DELETE RESULT (Object)` [Command]
 Deletes *Object*. *Object* can be a variable, a dataset or a result from a statistical command.
- `:MAKE DATASET FROM ROWS (Name) (Data) (Variable names)` [Command]
 Creates a new dataset called *Name*, using *Data* as the row-major data and *Variable names* as the names of the columns. To enter data, each row is a comma separated sequence, and semicolons separate the rows. For instance, to create a dataset called DS1, with variables V1, V2, V3 and V4, and containing the data:

```
2 4 2 3
5 7 4 3
```

Enter DS1 for *Name*, 2, 4, 2, 3; 5, 7, 4, 3 for *Data*, and V1, V2, V3, V4 for *Variable names*.

:MAKE DATASET FROM COLUMNS (*Name*) (*Data*) (*Variable names*) [Command]

Creates a new dataset called *Name*, using *Data* as the column-major data and *Variable names* as the names of the columns. To enter data, each row is a comma separated sequence, and semicolons separate the rows. For instance, to create a dataset called DS1, with variables V1, V2, V3 and V4, and containing the data:

```
2 4 2 3
5 7 4 3
```

Enter DS1 for *Name*, 2, 5; 4, 7; 2, 4; 3, 3 for *Data*, and V1, V2, V3, V4 for *Variable names*.

:ADD VARIABLE TO DATASET (*Dataset*) (*Data*) (*Name*) [Command]

Adds a new variable, called *Name* to *Dataset* using *Data*, which should be a comma separated sequence.

:PARTITION DATASET (*Dataset*) (*Partition clause*) (*Include variables*) [Command]

Partitions *Dataset* using *Partition clause*. *Include variables* should be a comma separated list of variables from the original dataset to include in the new dataset(s). The default is to include all variables. Partitioning is described in greater detail in section 5.

:PARTITION ON (*Dataset*) (*Partition variable*) (*Include variables*) [Command]

Partitions *Dataset* on *Partition variable*. *Include variables* is the same as for :PARTITION DATASET.

:MERGE DATASETS (*Datasets*) [Command]

Creates a new dataset by combining the rows of *Datasets*. Merging is described in greater detail in section 5.

:OPEN (*Object*) [Command]

Displays *Object* in a separate window. *Object* can be a dataset, a variable or a result.

:CLOSE (*Object*) [Command]

Closes the window that *Object* is being displayed in. *Object* can be a dataset, a variable or a result.

:DESCRIBE (*Object*) [Command]

Prints a history of *Object*'s creation. *Object* can be a dataset, a variable or a result.

3.5.5 The Transform Menu

Transform
User Defined
Sort
Index
Add a Constant
Natural Log
Log10
Recode Categorical
Reassign Bins
Smooth 4235h
Discrete Derivative

The Transform menu contains commands that allow for various transformations of data. Unlike most other commands, the Transform commands require that their arguments be actual CLASP variables. All the menu options create a new variable in the dataset to which the variable being transformed belongs. The name of the new variable will be derived from the transformation being performed, for instance, smoothing a variable called A will result in a variable called SMOOTH-OF-A.

:USER DEFINED (*Dataset*) (*Expression*) [*Command*]

Adds a new variable to *Dataset* by evaluating *Expression* for each row of *Dataset*. *Expression* may refer to variables in *Dataset* by name. For example, to calculate the processing time per node of a search, *Expression* might be (*/ nodesvisited totaltime*).

:SORT (*Variable*) [*Command*]

Sorts *Variable* in ascending order.

:INDEX (*Variable*) (*Index*) [*Command*]

Sorts *Variable* by *Index*. Both must be CLASP variables from the same dataset.

:ADD A CONSTANT (*Variable*) (*Number*) [*Command*]

Adds *Number* to *Variable*.

:NATURAL LOG (*Variable*) [*Command*]

Takes the natural log of *Variable*.

:LOG10 (*Variable*) [*Command*]

Takes the common log (base 10) of *Variable*.

:RECODE CATEGORICAL (*Variable*) (*Old values*) (*New values*) [*Command*]

Substitutes each member of *New values* for the corresponding member of *Old values* in *Variable*. This can be used to create an integer coding of a categorical variable.

:REASSIGN BINS (*Variable*) (*Bin limits*) (*Bin names*) [*Command*]

Allows the conversion of a continuous variable into a categorical variable. *Bin limits* are the upper limits on the bins, and *Bin names* are the names of the corresponding bins. The list of *Bin names* must be one element longer than the list of *Bin limits*. The last element of *Bin names* will name the bin which holds any values of *Variable* which are larger than the largest element of *Bin limits*.

:SMOOTH $4253h$ (*Variable*) [Command]

Apply a $4253h$ smoothing operator to *Variable*. This is described in greater detail in section 5.

:DISCRETE DERIVATIVE (*Variable*) [Command]

Take the pairwise difference of each value of *Variable* with the next value.

3.5.6 The Test Menu

Test
Confidence Interval Using z Statistic
Confidence Interval Using t Statistic
Confidence Interval of a Proportion
t-test One Sample
t-test Two Samples
t-test Matched Pairs
D test
Anova - One Way
Anova - Two Way
Chi-Square Counts
Chi-Square 2x2
Chi-Square RxC
Linear Regression - brief
Linear Regression - verbose
Multiple Linear Regression - verbose

The commands in the Test menu are described in section 4.2, and so only their arguments are given here.

:CONFIDENCE INTERVAL USING Z STATISTIC ($X(s)$) (*Confidence level*) [Command]

Confidence level is a real number from 0 to 1 inclusive. Note that it is not a mapping argument.

:CONFIDENCE INTERVAL USING T STATISTIC ($X(s)$) (*Confidence level*) [Command]

Confidence level is a real number from 0 to 1 inclusive. Note that it is not a mapping argument.

:CONFIDENCE INTERVAL OF A PROPORTION (*Successes*) (*Trials*) (*Confidence level*) [Command]

Successes is an integer ≥ 0 , *Trials* is an integer ≥ 1 , and *Confidence level* is a real number from 0 to 1 inclusive.

:T-TEST ONE SAMPLE ($X(s)$) (*Mean(s)*) (*Tails*) [Command]

Tails must be one of Both, Positive or Negative.

:T-TEST TWO SAMPLES ($Y(s)$) ($X(s)$) (*Tails*) [Command]

Tails must be one of Both, Positive or Negative.

:T-TEST MATCHED PAIRS ($Y(s)$) ($X(s)$) (*Tails*) [Command]

Tails must be one of Both, Positive or Negative.

- `:D TEST (Y(s)) (X(s)) (Tails)` [Command]
Tails must be one of `Both`, `Positive` or `Negative`.
- `:ANOVA - ONE WAY (Y(s)) (X(s))` [Command]
 Displays an anova table, a list of the group means, a plot of the group means, a scheffe table and an alternate calculation of the sum of squares total.
- `:ANOVA - TWO WAY (Y(s)) (X 1(s)) (X 2(s))` [Command]
 Displays an anova table, a table of the cell means, and row and column based means plots.
- `:CHI SQUARE COUNTS (A) (B) (C) (D)` [Command]
 Performs a chi square analysis on the contingency table
 A B
 C D.
- `:CHI SQUARE 2X2 (Y(s)) (X(s))` [Command]
 Displays the Chi Square statistic, the G statistic and contingency tables expressed in terms of frequency, percent of row totals, percent of column totals and expected values. *Y* and *X* must each be two-valued variables.
- `:CHI SQUARE RXC (Y(s)) (X(s))` [Command]
 Displays the Chi Square statistic, the G statistic and contingency tables expressed in terms of frequency, percent of row totals, percent of column totals and expected values.
- `:LINEAR REGRESSION - BRIEF (Y(s)) (X(s))` [Command]
 Displays the slope, intercept, r^2 , standard error of the slope and significance.
- `:LINEAR REGRESSION - VERBOSE (Y(s)) (X(s))` [Command]
 Displays all the information the `:LINEAR REGRESSION - BRIEF` does, and in addition, it displays the correlation, an anova table and a regression plot.
- `:MULTIPLE LINEAR REGRESSION - VERBOSE (Y) (X's)` [Command]
 Performs a multiple linear regression of *Y* on *X*'s, displaying the correlation, an anova table, the F statistic, the coefficients, betas and t statistic for each *X*, and a correlation matrix.

3.5.7 The Sample Menu

Sample
Sample from Uniform Distribution
Sample from Normal Distribution
Sample from Binomial Distribution
Sample from Poisson Distribution
Sample from Gamma Distribution

The Sample menu contains commands which take samples from various theoretical distributions. The variables in the dataset will be sampled from the specified distribution. In addition to the distribution parameters, all the sampling commands take a `Size` of

Sample argument, indicating how many values each variable should contain, and a Number of Samples argument indicating how many variables to create. The following descriptions only list the parameters of the distributions.

:SAMPLE FROM UNIFORM DISTRIBUTION (*Minimum*) (*Maximum*) [Command]

Will produce samples that are uniformly distributed between *Minimum* and *Maximum*, which must both be integers.

:SAMPLE FROM NORMAL DISTRIBUTION (*Mean*) (*Standard Deviation*) [Command]

Mean and *Standard Deviation* are both real numbers.

:SAMPLE FROM BINOMIAL DISTRIBUTION (*P*) (*N*) [Command]

P is a real number from 0 to 1, it is the probability of success. *N* is an integer ≥ 1 , it is the number of Bernoulli trials per sample.

:SAMPLE FROM POISSON DISTRIBUTION (*Mean*) [Command]

Mean is a real number.

:SAMPLE FROM GAMMA DISTRIBUTION (*Integer Order*) [Command]

Integer Order is an integer ≥ 1 .

Chapter 4

Statistics

Some computations on our data are intended merely to summarize, to distill the general trend or tendency, to help us see the signal by eliminating the noise. Other computations are meant to answer a question, to say whether this data is like that data or whether this model fits the data. The former computations are called descriptive statistics and the latter are called inferential or test statistics. In practice, the line between them is rather blurry. For example, fitting a model to the data describes its basic shape, but asking *whether* the model fits the data is a test, yet many statistical computations do both simultaneously. Nevertheless, CLASP's functions are categorized this way with respect to its menus and so is this manual, to help in organizing them.

4.1 Descriptive Statistics

Descriptive statistics are usually a function from the data to a single number. One example is the arithmetic mean, which we are all familiar with and which describes the data by a particular kind of central tendency or location. Some of the following statistics may be unfamiliar to you, but the basic goals are these: what is the data's location, how spread out is it, and how strong is its association with other data?

4.1.1 Location

The following statistics describe the "location" of the sample, usually the "center" of the sample, for some reasonable definition of the center.

Mean The mean is the average of a group of numbers, which describes the center of mass of the sample. This value is a good measure of the central tendency of a population; however, it is affected by a skewed distribution or when there are extreme outliers.

Median Like the mean, the median is another measure of central tendency; it indicates the value at the center of the distribution. Half of the scores will be above and half below. The median is less affected by outliers than the mean, since it is insensitive to how far the numbers are from the median. However, like the mean, the median is also affected by skewed distributions. The larger the difference between the median and the mean, the more likely that the distribution is not normal.

Quantile A generalization of the median is a quantile: given an fraction q , the q th quantile is a number x such that q of the data is greater than x and $1 - q$ is less than x . Clearly, the median is simply the $1/2$ quantile.

Trimmed Mean Discarding the highest and lowest values in a distribution and averaging the remaining values yields the trimmed mean. It is common to discard the highest and the lowest 25 percent of the values, but other trimming percentages may be used. Here is a sorted distribution:

(1 1 2 15 16 18 19 22 23 28 31 100)

The 25 percent trimmed mean is the mean of the middle six numbers. The trimmed mean is a robust alternative to the mean and is in some ways preferable to the median.

Mode The mode is the most common value in a distribution; for example, the mode of the distribution (1, 2, 2, 3, 4, 4, 4) is 4. If the data are continuous, real numbers, then the mode is apt to be uninformative because of the very low probability that two or more numbers will have exactly the same value. One solution is to map real-valued data into discrete numbers by rounding or sorting into bins for frequency histograms and so on, in which case the mode of a distribution depends on bin size. Although it is a contradiction in terms, we often speak of a distribution having two or more modes. This means the distribution has two or more values or ranges of values that are common. The decision to characterize a distribution as bimodal or multimodal is subjective.

Minimum and Maximum The minimum and maximum of a sample give an idea of its location, though certainly not of its central tendency. They are often ignored in statistics books, which concentrate on statistics that are used in hypothesis testing, such as the mean. Nevertheless, they are important because they help to visualize the data, outlining its location and helping to identify outliers, which are sometimes the most important elements of the data because they are most in need of explanation.

4.1.2 Spread

The following measures describe how much the data spreads out around its central value.

Range The range is simply the difference between the maximum and the minimum, and consequently indicates how spread out the data are.

Variance Variance is an important measure of how spread out a sample is around its center of mass—the mean. The higher the variance, the more spread out the sample. Mathematically, the variance is sum of the squared distances from the mean of each number in the sample. For example, the mean of (1 2 3 4 5) is 3 and the variance is 10, while the mean of (0 1 3 5 6) is also 3 but the variance is 26, reflecting the greater distance of the data from the mean.

Standard Deviation This statistic is simply the square root of the variance. The standard deviation is useful because its units are the same as the data's. For example, if the data are distances in meters, the mean will also be in meters, while the variance will

be in meters-squared. The standard deviation will be in meters, and can be usefully visualized as the “average” deviation of the data from the mean. If the data are normally distributed, 96% of the values will fall within a range of plus or minus 2 standard deviations from the mean.

Interquartile Range The interquartile range is a robust alternative to the variance and standard deviation; that is, it measures how spread out the data are, but it is less sensitive to one or two values that are wildly far away. The interquartile range is found by dividing a sorted distribution into four contiguous parts, each containing the same number of individuals. Each part is called a quartile. (The quartiles are the $1/4$, $1/2$ and $3/4$ quantiles.) The difference between the highest value in the third quartile and the lowest value in the second quartile is the interquartile range. For example, if the sorted distribution is (1 1 2 3 3 5 5 5 6 6 100), then the quartiles are (1 1 2), (3 3 5), (5 5 5), (6 6 100), and the interquartile range is $5 - 3 = 2$.

4.1.3 Covariance

Pearson’s Correlation This statistic describes the degree of association between two samples. Alternatively, it describes how well the value of the dependent, Y , variable can be predicted from the score of the independent, X , variable. This relationship does not imply causality; it simply indicates that a high value of Y is likely to co-occur with a high value of X . The idea of co-occurrence means that the data must have some kind of pairing: (x_i, y_i) . This statistic is considered by some to be overused and abused, but generally it is a good indicator of an effect or relationship between two variables.

Cross-Correlation Given two samples, X and Y , we can generalize the notion of correlation to ask what is the correlation between X and Y when Y is shifted or translated by some amount. This idea often arises when analyzing data collected over time, where you think that there will be a strong correlation between, say, daily barometric pressures in Buffalo and daily barometric pressure in Boston two days later. The time-shift is often called a “lag,” and is an additional argument to the cross-correlation function.

Autocorrelation Autocorrelation is simply the cross-correlation of a variable with itself. Obviously, this is uninteresting unless the lag is non-zero. This is a way of measuring cyclic patterns in the sample.

4.2 Test Statistics

Many different kinds of questions can be asked of our data, and so there are many kinds of test statistics. Some questions ask simply whether two groups of numbers are different, specifically whether their central tendencies (such as their means) are different. Others ask whether there are differences among many groups of data. Still others fit models, such as a straight line or a curve, to the data and ask whether it fits well enough or just fits by chance. This chapter cannot explain all of these statistical tests in detail, but there are many excellent statistics texts available. Two good general texts are *Statistical Reasoning*

in *Psychology and Education* by Edward W. Minium [8] and *Probability and Statistics for Engineering and the Sciences*, by Jay L. Devore [4].

Test statistics fall into several categories: parametric, non-parametric, and bootstrap, each of which is most applicable for certain specific types of problems. Parametric statistics make assumptions about the shape of the population from which the data was drawn; typically, it assumes the population is normal, but with unknown mean and variance. (These parameters are estimated from the sample, hence the term “parametric” statistics.) Parametric statistics are the best available, but only when their assumptions are true. If the assumptions are false, it is hard to determine the degree of error introduced into the analysis. Non-parametric statistics avoid making distributional assumptions, and so are more widely applicable. However, they are less powerful than their parametric competitors, which means that on the home turf of the parametric statistics (normal distributions), the parametric statistics are better. Non-parametric statistics typically work by ignoring the absolute magnitudes of the data and instead using only their relative magnitudes; essentially, they sort the data and map the original numbers onto $1..N$, because the behavior of the numbers $1..N$ can be analyzed and tables precomputed. James Bradley has an excellent book on non-parametric statistics [2]. Bootstrap statistics also avoid making distributional assumptions, but they retain the absolute magnitudes of the data. They substitute computation for analysis and tables, essentially by computing tables as necessary. Bootstrap statistics are even more widely applicable than non-parametric statistics because there is no reliance on precomputed tables. They are, however, very computation-intensive. A much more detailed discussion may be found in Cohen’s primer [3].

This section will briefly review the test statistics that are implemented in CLASP and their uses.

4.2.1 Confidence Intervals

Descriptive statistics of the sample are often used to estimate the equivalent statistic on the population. For example, one can estimate the mean of the population by the mean of the sample. Of course, the two will never be exactly equal, which gives rise to the question: how close is the estimate? One way of answering this is by constructing “confidence intervals.” A confidence interval is a range, a to b , centered around the sample statistic. It also has an associated probability, say 90 percent. What does this probability mean? Its meaning may be best understood with an analogy, where the size of the confidence interval $b - a$ corresponds, say, to the size of a horseshoe in a game of horseshoes. A bigger horseshoe has a higher probability of landing around the stake, which corresponds to the parameter we are estimating. Equivalently, if we were to construct 90 percent confidence intervals for many samples, 90 percent of the confidence intervals would contain the desired parameter.

CLASP has two ways of computing confidence intervals, both for means, where the difference depends on the population the mean is from. “Confidence Interval Using Z Statistic” is for large samples, where the Central Limit Theorem assures us that the sampling distribution of the mean is normal. “Confidence Interval Using T Statistic” is for small samples, where the sampling distribution of the mean is a t -distribution.

Bootstrapping is a general technique for deriving confidence intervals for any statistic on any sample, and so confidence intervals for other statistics can easily be obtained. These will be implemented in future releases of CLASP.

4.2.2 t Tests

To test whether two means are significantly different, use a t test. The t test is reasonably robust, even for small samples, provided the underlying assumption of normal populations is not badly violated.

CLASP supports three versions of t tests. The *one-sample* t test is used to test the hypothesis that a sample mean has a specific value. For example, imagine you are a burglar, driving through an unfamiliar part of town, trying to judge whether the houses are worth robbing. You decide to base your decision on the price of the cars parked on the street. One short cul de sac appears promising: it has just five parked cars, but they all cost a lot. In fact the mean price of the cars is \$20,270 and the standard deviation of this sample is \$5811. You know that the mean cost of cars in town is \$12,000, and you want to know whether the cars in the cul de sac are significantly dearer than the “average” car in town. The test statistic is:

$$t = \frac{20270 - 12000}{5811/\sqrt{5}} = 3.18$$

The probability of attaining this result by chance, under the null hypothesis that all cars cost on average \$12,000, is less than .01.

To compare two means, use a two-sample t test. For example, to test the hypothesis that one search algorithm expands significantly more nodes than another, on average, simply enter the sample for each algorithm in a separate column, and CLASP will calculate the t statistic:

$$t = \frac{\bar{x}_1 - \bar{x}_2}{\hat{\sigma}_{\bar{x}_1 - \bar{x}_2}}$$

CLASP will also return a p-value for the t statistic. If the p-value is small (conventionally, .05 is considered small and .01 very small), the test indicates that the two samples are significantly different.

A special case of a two-sample test is a *paired-sample t test*. To continue the previous example, imagine that the two search algorithms were tested on the same problems, so for each problem you have a pair of scores, one for each algorithm. Under the null hypothesis that the algorithms are equal, the expected value of the difference of scores is zero. The paired-sample t test computes the average difference (summing the differences and dividing by the number of problems) and then runs a one-sample test of the hypothesis that the average difference is zero.

Randomization versions of the one-sample, two-sample and paired-sample t tests are available in CLASP. These tests substitute massive amounts of computation for parametric assumptions about the sampling distribution of the mean. Thus, they are the preferred alternative when one suspects that one’s data are not drawn from a normal populations. Returning to the burglary example, above, we know that the distribution of automobile prices is not normal; rather, it is skewed to the right (high costs). For such data, or when you simply do not know the form of the underlying population distribution, use the randomization versions of t tests.

4.2.3 Analysis of Variance

Analysis of Variance (ANOVA) may initially be thought of as a generalization of the t-test to multiple groups. Rather than test whether the means of two groups are equal, it tests whether the means of k groups are equal. Analysis of variance is so named because it

identifies the amount of variation in a sample due to the effect of some treatment, where each different treatment results in a different group.

There are two forms of ANOVA implemented in CLASP, “One Way” and “Two Way.” One-way ANOVA is used to examine the effect of a single treatment on several samples, where the treatment is administered at several discrete levels. For example, in the Phoenix firefighting system, we once did an experiment with three levels of wind speed—low, medium and high—and performed a one-way ANOVA to test whether wind speed had an effect on performance. Two-way ANOVA is used when considering the effects of multiple treatments and their interactions. For example, in the Phoenix firefighting system, we might have run an experiment with three levels of wind speed and two kinds of firefighting plans—called multiple-shell and model-based. We could use a two-way ANOVA to test whether the performance of the plans depends on wind speed and which plan is better.

Each analysis of variance is calculated in a similar fashion and the results reported in an ANOVA table. Values used in calculating the ANOVA table are based on the sum of squares of the deviation scores, SS , and the degrees of freedom, df , associated with that sum of squares. The general relationship is in the form

$$s^2 = \frac{SS}{df}$$

The types of variation involved for one-way ANOVA are those of:

1. Variability of the values about the grand mean (the mean of all scores). Given by the deviation: $(X - \bar{X})$, where \bar{X} is the grand mean. This is known as the *total* variance, symbolized as s_T^2 . This value is useful only in the calculating of the following items.
2. Variability of scores about their subgroup sample means. Given by the deviation: $(X - \bar{X})$, where \bar{X} is the mean of the subgroup which contains X . This is known as the *within groups* variance, denoted by s_W^2 .
3. Variability of subgroup sample means about the grand mean of all scores. Given by the deviation: $(\bar{X} - \bar{X})$. This is the *among groups* variance, denoted by, s_A^2 .

The *within groups* variance indicates the amount of variation inherent in the subgroup, all of the members of the group receive the same treatment so none of the variation can be due to the treatment. The *among groups* variance contains two components, the variance inherent in the population and the variance due to the treatment. By examining these components, the effect of a treatment on a population can be determined.

Calculations for these values are described by the formulas in Figure 4.1, where X_i are the scores in subgroup i , \bar{X}_i is the mean of subgroup i , n_i is the number of observations in the i th subgroup and k is the number of subgroups.

Using these values, the F ratio is calculated. The F ratio indicates the probability that the variation in the samples could be accounted for by chance. Calculation of the F ratio uses the following formula:

$$F = \frac{s_A^2}{s_W^2}$$

For two-way ANOVA, there is an additional interaction between different treatments, sometimes called factors. These are separated into *Row* and *Column* treatments. In our second example, above, we had two treatments, wind speed and plan, and so one becomes

$$\begin{aligned}
SS_W &= \sum (X_i - \bar{X}_i)^2 + \sum (X_j - \bar{X}_j)^2 + \dots \\
SS_A &= \sum_i^k (\bar{X}_i - \bar{X})^2 \\
SS_T &= \sum (X - \bar{X})^2 \\
df_W &= \sum (n_i - 1) \\
df_A &= k - 1 \\
df_T &= \sum n_i - 1 \\
s_W^2 &= \frac{SS_W}{df_W} \\
s_A^2 &= \frac{SS_A}{df_A}
\end{aligned}$$

Figure 4.1: One Way ANOVA Calculations

$$\begin{aligned}
SS_C &= \sum (\bar{X}_{C_i} - \bar{X})^2 \\
SS_R &= \sum (\bar{X}_{R_i} - \bar{X})^2 \\
SS_{WC} &= \sum (X - \bar{X}_{cell})^2 \\
SS_{R \times C} &= SS_T - SS_C - SS_R - SS_{WC} \\
df_{WC} &= \sum^{all \ cells} (n_{WC} - 1) \\
df_C &= (C - 1) \\
df_R &= (R - 1) \\
df_{R \times C} &= (C - 1)(R - 1) \\
df_T &= (R)(C)(n_{WC}) - 1 \\
s_C^2 &= \frac{SS_C}{df_C} \\
s_R^2 &= \frac{SS_R}{df_R} \\
s_{WC}^2 &= \frac{SS_{WC}}{df_{WC}} \\
s_{R \times C}^2 &= \frac{SS_{R \times C}}{df_{R \times C}}
\end{aligned}$$

Figure 4.2: Two Way ANOVA Calculations

the row (say, wind speed, giving three rows: low, medium, and high) and the other becomes the column (plan, giving two columns: multiple-shell and model-based), and each cell of the table is a particular combination of the two. Having two treatments means there may be an additional effect called an interaction effect, where the condition of, say, low wind can combine with condition of, say, the multiple-shell plan, to produce a different performance. This interaction effect, also called a *Row by Column* effect, is calculated with the formulas in Figure 4.2, where the four variance estimates correspond to:

s_{WC}^2 *within cells estimate*, derived from the individual cell variation. This measures the inherent variation in an subgroup free from the effect of any treatment.

s_C^2 *column estimate*, derived from the differences from the column means. If there is an effect from the column treatments this value will tend to be larger than s_{WC}^2 .

s_R^2 *row estimate*, as above but for row treatments.

$s_{R \times C}^2$ *interaction estimate*, derived from the discrepancy between the means of several cells. This value will tend to be larger than s_{WC}^2 if there is an interaction effect between the *Row* and *Column* treatments.

There are three F ratios calculated—column, row, and row by column—in the form:

$$F = \frac{s_C^2}{s_{WC}^2}$$

$$F = \frac{s_R^2}{s_{WC}^2}$$

$$F = \frac{s_{R \times C}^2}{s_{WC}^2}$$

A significant value of F indicates that the hypothesis that the means of the subgroups are equal, that is that there is no variance due to the treatment, should be rejected. Two way ANOVA requires that each of the cells have the same number of observations in order to calculate a meaningful statistic, because unequal group sizes make the statistic unstable.

4.2.4 Linear Regression

Many statistical tests can be thought of as fitting a particular kind of model to the data then measuring how plausible it is to have drawn the data from that model. For t-tests and ANOVA, those models are simple Gaussian distributions. A slightly more complex model, but still very simple, is a linear model. Essentially, we think the data points are drawn from a line with “fuzz” around it—Gaussian fuzz with constant variance:

$$y_i = mx_i + b + \epsilon$$

In this equation, m is the slope of the line, b is its y -intercept, and ϵ represents the Gaussian fuzz around the line. Linear regression is a way to find m and b and also to estimate how big the ϵ 's are and how plausible it is that the data is generated by this linear model.

A plot of the regression line with a scatter plot of the individual data points provides a picture of the relationship between the variables and also gives a good indication of the correlation between the variables. It is important to remember that there may be

relationships between the variables which are non-linear, such as a logarithmic or quadratic relationship. Don't use linear regression if the relationship between the variables is clearly non-linear.

Multi-variate linear regression performs the same calculations with multiple independent variables, using the model

$$Y = a + b_1X_1 + b_2X_2 + \cdots + b_nX_n + \epsilon$$

where the b_i are the *coefficients of regression* for the corresponding X , a is the intercept of the regression line, and ϵ is the error, or residual. Each b has a corresponding t statistic indicating the probability of whether or not the relationship could occur by chance. A significant b indicates that a portion of the variance of the dependent variable is due to the influence of the associated independent variable. The amount of the influence is measured as a percentage of the total variance of the dependent variable, which is calculated as the square of the correlation between Y and X_i . The percentage of the variance due to all of the independent variables is given by R^2 , with the F statistic for the regression calculated from this value. The report for this procedure includes the correlation coefficients between the independent variables.

4.2.5 Chi Square

Chi square is a useful statistic when we are counting discrete things, such as the number of birds of each of several species, and we are interested in whether the observed proportions (such as ten percent big black birds and ninety percent little brown birds) is likely to have been drawn from some theoretical population. This is done by subtracting the expected number of events (for example, the expected number of little brown birds calculated from our theoretical population and the number of birds we saw) from the actual number of events, squaring the difference, and summing over each category of event. (There is also a normalizing factor in each term.) By assuming that deviations are normally distributed, the probability of drawing the sample from the theoretical population can be calculated. Hence, chi square is a parametric statistic. With a slightly different calculation, we can compare two samples to test whether they are the same, in the sense of having roughly the same distribution.

An important special case of the chi square is to test for independence. Again, we are counting events, and the events fall into categories, but the categories fall along two different dimensions. The question we are asking is whether the two dimensions are independent: the distribution of events along one dimension is the same for each category of events along the other, and vice versa. For example, suppose we are tossing a nickel and a penny, and we are counting the outcomes. There are four possible events, falling into categories along two dimensions, namely the identity of the coin. The following figure is a 2x2 "contingency table" to represent the data, where each cell is the observed frequency for the combination of the two categories.

		penny	
		H	T
nickel	H	a	b
	T	c	d

We can use “Chi square” to analyze the contingency table and test whether the two coins are independent. (We would be surprised if they weren’t.) Larger contingency tables can be handled by CLASP.

When the sample is small, the deviations can’t be normally distributed, since the data are discrete integers. In that case, the user can choose *Yate’s correction*, which adjusts for the discontinuity.

Chapter 5

Data Manipulation

CLASP data is organized into variables and datasets. A variable is an ordered column of numbers which have some semantic interpretation, such as “Number of Ships Allocated” in a transportation simulator. The values might represent different measurements over time or measurements during different trials in an experiment. A dataset is a collection of variables, all of which have the same number of values. Usually a dataset will represent the variables collected during a single experiment or set of experiments. In CLASP, all variables belong to a dataset. Datasets are analogous to tables of data, where variables are columns, and a row is a set of the i th value in each column. For example, in an experiment where thirty trials are run and, during each trial, “Runtime,” “Num-Nodes” and “Goal-Found?” were measured, the resulting dataset would have thirty rows, each corresponding to one of the trials, and three Variables, “Runtime,” “Num-Nodes” and “Goal-Found?” Most statistical operations involve performing some operation on the values of one or more variables. For example, the mean operates on one variable and the t-test compares two variables to test if they have statistically different means. CLASP offers a variety of ways of accessing these values for statistical manipulation.

5.1 Function Application

Most data manipulation consists of transforming one or more variables by mapping a function across them. You may think of this as equivalent to the Common Lisp ‘map’ function, in which each variable is treated as a sequence. By typing in an expression, you can transform your data in arbitrary ways. Certain standard operations, such as log transforms, are built in.

5.1.1 Formula Expressions

Transformations of a variable are accomplished by choosing “User Defined” from the transform menu. The user is prompted for a dataset and an expression. The formula can be any valid Common Lisp expression and may contain variable names. The formula is mapped over all the variables it contains and applied to each successive value. The result is a new variable which is added to the dataset. For example, the following procedure would produce the ratio of “Nodes Searched” to “Runtime” in the dataset “My Data.”

1. Click on Transform Menu
2. Click on User Defined

3. CLASP prompts for a dataset, Click on My Data
4. CLASP prompts for a formula, type: (/ Nodes-Searched Runtime)

A new variable will be created called /-Nodes-Searched-Runtime.

Often a categorical variable will need to be converted to a numeric variable to allow certain statistics to be applied to it. For instance a simulator that reports wind-speed as “low”, “medium” or “high”, but for analysis, the user prefers the encoding to be 1, 2 and 3. Alternatively, categorical variables which are encoded as integers might need to be converted to a nominal encoding for ease of interpretation. For example, a variable which encodes the completion status of an algorithm as 1 or 0 might be converted to “success” or “failure” for ease of interpretation. This is called *recoding*. Formulas may include the *recode* operator. The form of this expression is:

(recode Variable old-values new-values)

Recoding will translate all the elements of variable in old-values to the equivalent new-value. For instance, if a variable is valued

(A D A C B A B D C)

and the expression

(recode Variable '(A B C D) '(1 2 3 4))

is used in a “Transform” statement, the new variable will have values

(1 4 1 3 2 1 2 4 3)

This kind of transformation can be valuable for running statistics that require numeric data when the original data is symbolic.

5.1.2 Smoothing

Time series and other series plots are sometimes difficult to read because of rapid fluctuations. For example, the time series of mean daily temperature bounces all over the place and can make it difficult to see slower-moving seasonal components. Spectral analysis is sometimes used to decompose series into frequency components, but yields a plot in the frequency domain, not the time domain. Often, it suffices to simply remove the high-frequency component from a time series by applying *smoothing* operators. The basic idea of smoothing is simple: every value in the series is replaced by an average of the value and its neighbors. (Endpoints are handled separately; for example, they might be used as is, or averaged with their neighbor.) The average need not be the arithmetic average, in fact, it is often better to use the median instead of the mean. In either case, values that are much bigger or smaller than their neighbors are “brought into line.” The most common mean and median smoothing techniques involve replacing the *i*th value in the series with the mean or median of the *i*-1, *i*, and *i*+1 values. These are called “3” smooths because they involve three values, a *window* of size 3.

Median smooths handle outlier values differently than mean smooths. Whereas median smooths ignore outliers, leaving only small values, mean smooths average outliers and neighboring values, giving the impression that the neighbors have higher values than they actually do. On the other hand, median smooths create *mesas* or short sequences of identical values.

Mesas can be handled by *resmoothing* that is, smoothing the smoothed series again. One can resmooth a sequence many times, with various sizes of windows, and various kinds of averaging—not only median and mean smoothing, but also weighted smoothing in which each value in a window is multiplied by a weight. One can create long and complicated smoothing plans. For example, 3R2H means “smooth the sequence with 3,median smooths repeatedly until the appearance of the series doesn’t change with successive smooths (that’s what the R means), then apply a 2,median smooth, and then a hanning operation (that’s what H means). A hanning operation multiplies the three values in a window by .25, .5 and .25, respectively, and sums the results. It is used to “clean up” a smooth by removing any leftover spikes.

Many smoothing plans have been developed for different kinds of data, but one that works pretty well in general is 4253H. This is four median smooths, with window sizes 4, 2, 5 and 3, followed by a Hanning operation. Currently, CLASP offers four median smoothing operators with window sizes 2, 3, 4, and 5, respectively; and the general smoothing plan 4253H.

5.2 Partition Clauses

A common manipulation of data is to reorganize it, say by grouping it or looking at a fraction of it. Most of these kinds of manipulation can be accomplished by *partitioning*, which essentially means mapping over the data, grouping it, and doing something with each group. (You might discard some groups in order to look more closely at others.)

It is often necessary to apply a command to a subset of the rows in a dataset. Sometimes this is because the statistics being used are only meaningful for that subset or because different operations must be done to different parts of the dataset. For instance, suppose you collect data on the time it takes an algorithm to solve a problem; it is meaningless to try and compute the mean solution time in cases where an algorithm couldn’t find a solution. Instead, you can select from the dataset just those cases where the algorithm succeeded and then take the mean of the runtimes in that selection. In other situations, it might be desirable to break the dataset up into a number of subsets by using one or more key variables. For every row in the dataset, the value of the key variable(s) at that row determine which new partition the row will be assigned to. This is often useful when a factorial analysis of data is being done.

In CLASP, the “Partition” command does both partitioning and selection. The form of the command is

```
partition dataset partition-clause
```

Partition clauses are in prefix notation (operator first, arguments last) and must be enclosed in parentheses. They can be nested. The operators for partition clauses are given below. A typical partition clause might look like the following:

```
(.and. (. == . success t)
      (.or. (. <= . nodes-searched 5)
            (. >><< . tree-depth 50 101)))
```

This clause would select those rows of the dataset where the “success” variable was true, and either the “tree-depth” was less than or equal to 5 or the “nodes-searched” was between 50 and 101 exclusive.

5.2.1 Partition Clause Logic Operators

The following operations can be used to combine variable comparison operations in a selection clause.

<code>.and.</code>	<code>arg1 &rest args</code> Selects rows matching every argument's selection criteria.
<code>.or.</code>	<code>arg1 &rest args</code> Selects rows matching at least one argument's selection criteria
<code>.not.</code>	<code>arg</code> Selects rows not matching argument's selection criteria.

5.2.2 Partition Clause Comparison Operators

The following operators compare the variable's value to a constant and selects those rows where the condition is true:

<code>.==.</code>	<code>variable-name value</code>
<code>.<=.</code>	<code>variable-name value</code>
<code>.>=.</code>	<code>variable-name value</code>
<code>./=.</code>	<code>variable-name value</code>
<code>.>>.</code>	<code>variable-name value</code>
<code>.<<.</code>	<code>variable-name value</code>

The following operators select rows where the variable's value is between the two given constants—either exclusively or inclusively.

<code>.=><=.</code>	<code>variable-name value-1 value-2</code>
<code>.=><<.</code>	<code>variable-name value-1 value-2</code>
<code>.>><=.</code>	<code>variable-name value-1 value-2</code>
<code>.>><<.</code>	<code>variable-name value-1 value-2</code>
<code>.<<>>.</code>	<code>variable-name value-1 value-2</code>
<code>.<<=>.</code>	<code>variable-name value-1 value-2</code>
<code>.<=>>.</code>	<code>variable-name value-1 value-2</code>
<code>.<==>.</code>	<code>variable-name value-1 value-2</code>

The following operators select rows where the variable's value is an extreme of some sort:

<code>.min.</code>	<code>variable-name</code> Selects rows where variable is the minimum value in table.
<code>.max.</code>	<code>variable-name</code> Selects rows where variable is the maximum value in table.
<code>.pred.</code>	<code>variable-name value</code> Selects rows where variable has greatest value less than value.
<code>.succ.</code>	<code>variable-name value</code> Selects rows where variable has smallest value greater than value.
<code>.floor.</code>	<code>variable-name value</code> Selects rows where variable has greatest value less than or equal to value.

`.ceiling. variable-name value`
 Selects rows where `variable` has smallest value greater than or equal to `value`.

5.2.3 Exhaustive Partitioning Operator

The exhaustive partition operator is “.on.”. A .on. clause has the form

`(.on. key-1 ...key-n)`

and the dataset is partitioned on the keys. It is not advisable to partition a dataset on a key that has many different values, since you'll get an unmanageable number of partitions. In particular, continuous variables make very bad keys, since it is likely that almost every row will end up in its own partition.

This operator makes several partitions from a dataset—one for each distinct value of the key variables. Each resulting dataset will contain all the rows of the original dataset which share the same key variable values. Here is an example. We start with the following table:

My Data			
Key-1	Key-2	Time	Temp
1	A	10	52.1
1	A	12	53.9
2	A	9	52.6
2	B	7	48.2
1	C	10	50.2
2	B	13	47.4
1	C	12	51.5
2	A	6	55.1
2	C	11	46.7

Partitioning the dataset `(.on. Key-1 Key-2)` would result in the following 5 datasets

My Data where Key-1 = 1 and Key-2 = A

Key-1	Key-2	Time	Temp
1	A	10	52.1
1	A	12	53.9

My Data where Key-1 = 1 and Key-2 = C

Key-1	Key-2	Time	Temp
1	C	10	50.2
1	C	12	51.5

My Data where Key-1 = 2 and Key-2 = A

Key-1	Key-2	Time	Temp
2	A	9	52.6
2	A	6	55.1

My Data where Key-1 = 2 and Key-2 = B

Key-1	Key-2	Time	Temp
2	B	7	48.2
2	B	13	47.4

My Data where Key-1 = 1 and Key-2 = A

Key-1	Key-2	Time	Temp
2	C	11	46.7

Chapter 6

Functions

CLASP is designed so that the functions that actually compute statistics on data are separated from the functions that manipulate data in the CLASP database and also separated from the graphical interface (CLIM) functions. The purpose of this modularity is to allow users to call statistics functions directly on their own data, from their Lisp Read-Eval-Print loops (REPLs) or even from their own Lisp programs, without having to use the graphical interface or the database. Furthermore, it allows users to manipulate their data using the database functions, calling them from their own REPLs or their own Lisp programs. By using CLIP and CLASP, data can be collected and analyzed completely under program control. (Note that we have called this section “Functions” and referred to its subject as functions, but a few of them are macros, CLOS classes or methods, and other things written in Common Lisp.)

Having divorced the statistical functions from the database functions, we had to decide how data was to be presented to the statistical functions. We decided to use sequences, which have a long history of usefulness in Lisp programs. Some of the simpler functions, such as ‘mean,’ take sequence keyword arguments, such as `:start` and `:end` to indicate a subsequence, and `:key` to map sequence elements to actual numbers. These all have their standard Common Lisp syntax and semantics. Unless otherwise specified, you should assume that all of these functions take sequences of numbers.

The documentation for specific functions is pulled from their online documentation strings. Those documentation strings occasionally refer you to this manual, which will seem odd since you are reading the manual. However, this approach helps to keep the manual up to date.

6.1 Descriptive Statistic Functions

The functions in this section compute statistics that describe a sample, without making assumptions about the population it was drawn from. All of these functions take sequences of numbers.

The ‘data-length’ function is usually used to compute the n used in statistical computations, such as the mean. It is described here because it is so simple and doesn’t fit anywhere very well.

`data-length` *data* *key* *start* *end* *key* [Function]

Returns the number of data values in ‘data.’ Essentially, this is the Common Lisp

'length' function, except it handles sequences where there is a 'start' or 'end' parameter. The 'key' parameter is ignored.

6.1.1 Location

The functions in this subsection describe the "location" of the sample, usually the "center" of the sample, for some reasonable definition of the center.

mean data *rest standard-args* *key start end key* [Function]

Returns the arithmetic mean of 'data,' which should be a sequence. Signals 'no-data' if there is no data.

The following statistics are often ignored in statistics books, which concentrate on statistics that are used in hypothesis testing, such as the mean. Nevertheless, they are important because they outline the location of the data and help to identify outliers, which are sometimes the most important elements of the data because they are most in need of explanation.

minimum data *rest standard-args* *key start end key* [Function]

Returns the element of the sequence 'data' whose 'key' is minimum. Signals 'no-data' if there is no data. If there is only one element in the data sequence, that element will be returned, regardless of whether it is valid (a number).

maximum data *rest standard-args* *key start end key* [Function]

Returns the element of the sequence 'data' whose 'key' is maximum. Signals 'no-data' if there is no data. If there is only one element in the data sequence, that element will be returned, regardless of whether it is valid (a number).

The 'quantile' function is a generalization of medians, quartiles, and percentiles. Even more than the minimum and maximum, it helps to picture the data, identifying central values and outliers.

quantile data q *rest standard-args* *key start end key* [Function]

Returns the element which is the q'th percentile of the data when accessed by 'key.' That is, it returns the element such that 'q' of the data is smaller than it and 1-'q' is above it, where 'q' is a number between zero and one, inclusive. For example, if 'q' is .5, this returns the median; if 'q' is 0, this returns the minimum (although the 'minimum' function is more efficient). If there is no 'key,' this function returns a single number, which may be the mean of two numbers in 'data.' If there is a 'key,' this function also returns the two data elements whose keys bracket the desired quantile. If the quantile happens to be exact, these data elements may be eq.

median data *rest standard-args* *key start end key* [Function]

Returns the median of the subsequence of 'data' from 'start' to 'end', using 'key'. Works by calling 'quantile,' and that function has more complete documentation.

Essentially, the 'trimmed-mean' gets rid of the outliers and takes the mean of the central values. It is a robust statistic that still provides good estimates of the center of different kinds of population distributions, such as uniform, double-exponential, and normal.

trimmed-mean data percentage *rest standard-args* *key start end key* [Function]

Returns a trimmed mean of ‘data.’ A trimmed mean is an ordinary, arithmetic mean of the data, except that an outlying percentage has been discarded. For example, suppose there are ten elements in ‘data,’ and ‘percentage’ is 0.1: the result would be the mean of the middle eight elements, having discarded the biggest and smallest elements. If ‘percentage’ doesn’t result in a whole number of elements being discarded, then a fraction of the remaining biggest and smallest is discarded. For example, suppose ‘data’ is ‘(1 2 3 4 5)’ and ‘percentage’ is 0.25: the result is $(.75(2) + 3 + .75(4)) / (.75 + 1 + .75)$ or 3. By convention, the 0.5 trimmed mean is the median, which is always returned as a number.

The mode of a distribution may not be at the “center” at all, but is certainly a representative value.

`mode data` *rest standard-args* *key start end key* [Function]

Returns the most frequent element of ‘data,’ which should be a sequence. The algorithm involves sorting, and so the data must be numbers or the ‘key’ function must produce numbers. Consider ‘sxhash’ if no better function is available. Also returns the number of occurrences of the mode. If there is more than one mode, this returns the first mode, as determined by the sorting of the numbers.

The following is a generalization of the “mode” function. If the ‘k’ parameter is quite large, this function can count the number of occurrences of each data value. Of course, that computation could take a long time, as the complexity is roughly $O(kn)$.

`multiple-modes data k` *rest standard-args* *key start end key* [Function]

Returns the ‘k’ most frequent elements of ‘data,’ which should be a sequence. The algorithm involves sorting, and so the data must be numbers or the ‘key’ function must produce numbers. Consider ‘#’sxhash if no better function is available. Also returns the number of occurrences of each mode. The value is an association list of modes and their counts. This function is a little more computationally expensive than ‘mode,’ so only use it if you really need multiple modes.

6.1.2 Spread

The functions in this subsection describe how much a sample spreads out around its center.

The ‘sum-of-squares’ function is rarely useful in itself, but is a common building block for statistical computations such as the variance or the analysis of variance. It is given by the following formula:

$$\sum_{i=1}^n (x_i - \bar{x})^2$$

The variance is this divided by $n - 1$ and the standard deviation is the square root of the variance.

`sum-of-squares data` *rest standard-args* *key start end key* [Function]

Returns the sum of squared distances from the mean of ‘data’. Signals ‘no-data’ if there is no data.

`variance data` *rest standard-args* *key start end key* [Function]

Returns the variance of ‘data,’ that is, the ‘sum-of-squares’ divided by $n-1$. Signals ‘no-data’ if there is no data. Signals ‘insufficient-data’ if there is only one datum.

`standard-deviation` *data* *rest standard-args* *key start end key* [Function]

Returns the standard deviation of ‘data,’ which is just the square root of the variance. Signals ‘no-data’ if there is no data. Signals ‘insufficient-data’ if there is only one datum.

`range` *data* *rest standard-args* *key start end key* [Function]

Returns the range of the sequence ‘data.’ Signals ‘no-data’ if there is no data. The range is given by max - min.

`interquartile-range` *data* *rest standard-args* [Function]

The interquartile range is similar to the variance of a sample because both are statistics that measure out “spread out” a sample is. The interquartile range is the difference between the 3/4 quantile (the upper quartile) and the 1/4 quantile (the lower quartile).

6.1.3 Correlation

When several measurements are taken of a system, so that we get a sample of points in two-, three-, or n -dimensions, we can not only describe the collection of values in each dimension, but also any relationships between them. One simple measure is their covariance or correlation, which is simply the normalized equivalent of covariance. (Covariance is rarely reported, since it depends on the units of measure.) Essentially, correlation measures the *linear* association between the values in one dimension and those in another; that is, the extent to which, say, high x values coincide with high y values.

The first question when we start to deal with multi-dimensional samples is representation. The two obvious possibilities are:

1. a sequence of points:

$$((x_0 \ y_0) (x_1 \ y_1) (x_2 \ y_2) \dots (x_n \ y_n))$$

2. a collection of sequences:

$$\begin{array}{l} (x_0 \ x_1 \ x_2 \ \dots \ x_n) \\ (y_0 \ y_1 \ y_2 \ \dots \ y_n) \end{array}$$

We chose to use the latter representation because the database is organized that way—around *variables*, which are sequences of data values of a particular kind or dimension.

Note that these representations are easily converted from one to the other:

```
(setq pts '((0 5) (1 6) (2 8) (3 9) (4 10)))

(values (mapcar #'first pts) (mapcar #'second pts))
=> (0 1 2 3 4)
    (5 6 8 9 10)

(mapcar #'list (multiple-value-list *))
=> ((0 5) (1 6) (2 8) (3 9) (4 10))
```


If you want to avoid the consing that this involves, the function ‘map-into’ can let you use and reuse an existing array.

Because the ‘covariance’ and ‘correlation’ functions take several sequences, we allowed them each to take ‘start’ and ‘end’ parameters. However, there is no ‘key’ parameter.

Mathematically, the covariance is

$$\sum (x_i - \bar{x})(y_i - \bar{y})$$

which is algebraically equivalent to

$$\sum x_i y_i - (\sum x_i)(\sum y_i)/n$$

but the latter is more efficient to calculate.

covariance *sample1 sample2* *ℰrest args ℰkey start1 end1 start2 end2* [Function]

Computes the covariance of two samples, which should be equal-length sequences of numbers. Covariance is the inner product of differences between sample elements and their sample means. For more information, see the manual.

correlation *sample1 sample2* *ℰrest args ℰkey start1 end1 start2 end2* [Function]

Computes the correlation coefficient of two samples, which should be equal-length sequences of numbers.

The ‘inner-product’ function was written to be called by the ‘covariance’ and ‘correlation’ functions, but may be useful in its own right.

inner-product *sample1 sample2* *ℰkey start1 end1 start2 end2* [Function]

Returns the inner product of the two samples, which should be sequences of numbers. The inner product, also called the dot product or vector product, is the sum of the pairwise multiplication of the numbers. Stops when either sample runs out; it doesn’t check that they have the same length.

Finally, the computation for correlation doesn’t really require all the data; it can get by with just a few “summary” statistics. If you have those summary statistics, you may prefer to call the following function directly. Keeping the summary statistics can be done incrementally in a constant amount of space, which can be very useful if you are collecting data during the run of a program but don’t want to keep every value. Instead you can just keep a few running sums. There are a number of other statistical functions that work on summary statistics, and the same comments apply to them.

correlation-from-summaries *n x x2 y y2 xy* [Function]

Computes the correlation of two variables given summary statistics of the variables. All of these arguments are summed over the variable: ‘x’ is the sum of the x’s, ‘x2’ is the sum of the squares of the x’s, and ‘xy’ is the sum of the cross-products, which is also known as the inner product of the variables x and y. Of course, ‘n’ is the number of data values in each variable.

See section 4.1 for a description of cross-correlation and autocorrelation

cross-correlation *sequence1 sequence2 max-lag* *ℰoptional min-lag* [Function]

Returns a list of the correlation coefficients for all lags from ‘min-lag’ to ‘max-lag,’ inclusive, where the ‘i’th list element is the correlation of the first (length-of-sequence1 - i) elements of sequence1 with with the last i elements of sequence2. Both sequences should be sequences of numbers and of equal length.

`autocorrelation` *sample max-lag* *Optional min-lag* [Function]

Autocorrelation is merely a cross-correlation between a sample and itself. This function returns a list of correlations, where the i’th element is the correlation of the sample with the sample starting at ‘i.’

6.1.4 Summary

The following function just collects many, probably too many, of the descriptive statistic functions. It’s an expensive but convenient way to say a lot about a sample in just ten numbers.

`statistical-summary` *data* *rest standard-args* *key start end key* [Function]

Compute the length, minimum, maximum, range, median, mode, mean, variance, standard deviation, and interquartile-range of ‘sequence’ from ‘start’ to ‘end’, accessed by ‘key’.

6.2 Test Statistic Functions

This section describes functions to test hypotheses and such. The reasoning behind them is somewhat complicated at times, and users should check that they understand the tests and their proper uses by referring to any good introductory statistics book, such as the ones in the references. Most of the tests are described in section 4.2, and will not be further described here.

6.2.1 t Tests

The functions to compute various t-tests are the following. We don’t implement a z-test because it is just a special case of the t-test in which the samples are quite large. You can use the ‘t-test’ function on large samples without any significant loss of accuracy or efficiency.

`t-test-one-sample` *data tails* *Optional (h0-mean 0)* *rest standard-args* *key* [Function]
start end key

Returns the t-statistic for the mean of the data, which should be a sequence of numbers. Let D be the sample mean. The null hypothesis is that D equals the ‘H0-mean.’ The alternative hypothesis is specified by ‘tails’: ‘:both’ means $D \neq H0\text{-mean}$, ‘:positive’ means $D > H0\text{-mean}$, and ‘:negative’ means $D < H0\text{-mean}$. The function also returns the significance, the standard error, and the degrees of freedom. Signals ‘zero-variance’ if that condition occurs. Signals ‘insufficient-data’ unless there are at least two elements in the sample.

`t-test` *sample-1 sample-2 tails* [Function]

Returns the t-statistic for the difference in the means of two samples, which should each be a sequence of numbers. Let $D = \text{mean1} - \text{mean2}$. The null hypothesis is that $D = 0$. The alternative hypothesis is specified by ‘tails’: ‘:both’ means $D \neq 0$,

‘:positive’ means $D > 0$, and ‘:negative’ means $D < 0$. Unless you’re using :both tails, be careful what order the two samples are in: it matters! The function also returns the significance, the standard error, and the degrees of freedom. Signals ‘standard-error-is-zero’ if that condition occurs. Signals ‘insufficient-data’ unless there are at least two elements in each sample.

t-test-matched *sample1 sample2 tails* [Function]

Returns the t-statistic for two matched samples, which should be equal-length sequences of numbers. Let $D = \text{mean1} - \text{mean2}$. The null hypothesis is that $D = 0$. The alternative hypothesis is specified by ‘tails’: ‘:both’ means $D \neq 0$, ‘:positive’ means $D > 0$, and ‘:negative’ means $D < 0$. Unless you’re using :both tails, be careful what order the two samples are in: it matters! The function also returns the significance, the standard error, and the degrees of freedom. Signals ‘standard-error-is-zero’ if that condition occurs. Signals ‘insufficient-data’ unless there are at least two elements in each sample.

The d-test is a competitor to the t-test. It uses bootstrapping to estimate the sampling distribution of d under the null hypothesis. Thus, it avoids the distributional assumptions of the t-test at the price of a great deal of computation. The t-test will take less than a second of computation, while the d-test may take several seconds. For example, on two samples of size 20 on a TI Explorer II, the t-test took 11.7 milliseconds, while the d-test took 2.9 seconds. Thus, the relative cost of the d-test is very high, but its absolute cost is really quite small.

d-test *sample-1 sample-2 tails* *&optional (bootstrap-times 1000)* [Function]

Two-sample test for difference in means. Competes with the unmatched, two-sample t-test. Each sample should be a sequence of numbers. We calculate the mean of ‘sample-1’ minus the mean of ‘sample-2’; call that D . Under the null hypothesis, D is zero. There are three possible alternative hypotheses: D is positive, D is negative, and D is either, and they are selected by the ‘tails’ parameter, which must be :positive, :negative, or :both, respectively. We count the number of chance occurrences of D in the desired rejection region, and return the estimated probability.

6.2.2 Confidence Intervals

Computing a confidence interval on a parameter depends on the parameter and its theoretical distribution. We have implemented just three at this time.

confidence-interval-z *data confidence* [Function]

Suppose you have a sample of 50 numbers and you want to compute a 90 percent confidence interval on the population mean. This function is the one to use. Note that it makes the assumption that the sampling distribution is normal, so it’s inappropriate for small sample sizes. Use confidence-interval-t instead. It returns three values: the mean and the lower and upper bound of the confidence interval. True, only two numbers are necessary, but the confidence intervals of other statistics may be asymmetrical and these values would be consistent with those confidence intervals. This function handles 90, 95 and 99 percent confidence intervals as special cases, so those will be quite fast. ‘Sample’ should be a sequence of numbers. ‘Confidence’ should be a number between 0 and 1, exclusive.

`confidence-interval-t` *data confidence* [Function]

Suppose you have a sample of 10 numbers and you want to compute a 90 percent confidence interval on the population mean. This function is the one to use. This function uses the t-distribution, and so it is appropriate for small sample sizes. It can also be used for large sample sizes, but the function ‘confidence-interval-z’ may be computationally faster. It returns three values: the mean and the lower and upper bound of the confidence interval. True, only two numbers are necessary, but the confidence intervals of other statistics may be asymmetrical and these values would be consistent with those confidence intervals. ‘Sample’ should be a sequence of numbers. ‘Confidence’ should be a number between 0 and 1, exclusive.

`confidence-interval-proportion` *x n confidence* [Function]

Suppose we have a sample of ‘n’ things and ‘x’ of them are “successes.” We can estimate the population proportion of successes as x/n ; call it ‘p-hat.’ This function computes the estimate and a confidence interval on it. This function is not appropriate for small samples with p-hat far from $1/2$: ‘x’ should be at least 5, and so should ‘n-‘x.’ This function returns three values: p-hat, and the lower and upper bounds of the confidence interval. ‘Confidence’ should be a number between 0 and 1, exclusive.

The following functions compute the very same values, but don’t require all of the data, just summaries of it. There is no need for a summary version of ‘confidence-interval-proportion,’ since it is already summarized.

`confidence-interval-z-summaries` *mean standard-error confidence* [Function]

This function is just like ‘confidence-interval-z,’ except that instead of its arguments being the actual data, it takes the following summary statistics: ‘mean,’ a point estimator of the mean of some normally distributed population; and the ‘standard-error’ of the estimator, that is, the estimated standard deviation of the normal population. ‘Confidence’ should be a number between 0 and 1, exclusive.

`confidence-interval-t-summaries` *mean dof standard-error confidence* [Function]

This function is just like ‘confidence-interval-t,’ except that instead of its arguments being the actual data, it takes the following summary statistics: ‘mean,’ which is the estimator of some t-distributed parameter; ‘dof,’ which is the number of degrees of freedom in estimating the mean; and the ‘standard-error’ of the estimator. In general, ‘mean’ is a point estimator of the mean of a t-distribution, which may be the slope parameter of a regression, the difference between two means, or other practical t-distributions. ‘Confidence’ should be a number between 0 and 1, exclusive.

6.2.3 Analysis of Variance

The analysis of variance is something like a t-test on many samples simultaneously. In any event, to pass data to these functions, we need a way to specify many groups of numbers. We have implemented two ways:

1. Structural Groups

In this representation, each group of data (numbers), is its own sequence. Thus, the one-way analysis of variance using structural grouping takes a sequence of sequences. For example, the following is three groups of four numbers each.

```
data => ((1 6 2 5) (8 3 6 5) (8 9 3 5))
```

2. Keyed Groups (Variables)

In this representation, each data value (number) is associated with a key or keys, and data that share the same keys are in the same group. Functions that take data represented this way have two or three arguments, one of which is called 'dv' and contains *all* the data as a flat sequence of numbers. The other arguments are sequences of keys. Keys can be any Lisp object; they are compared with 'eql.' These functions have `-variables` appended to their names because their representation is the same as that of variables in the database. For example, to specify the same three groups of four numbers, the input would be as follows:

```
iv => (a a a a b b b b c c c c)
dv => (1 6 2 5 8 3 6 5 8 9 3 5)
```

Of course, it doesn't matter what identifiers are used to distinguish groups.

The one-way anova function needs its input to be sorted in key order, so that the groups are contiguous, as in our example. The two-way anova function can have its inputs in any order, because it does a lot more processing.

The analysis of variance functions return an ANOVA table as their first value. For the one-way ANOVA, it looks like this:

```
((df-group ss-group ms-group f p)
 (df-error ss-error ms-error)
 (df-total ss-total))
```

For a two-way ANOVA, there is a larger ANOVA table because there are two groupings (a row effect and a column effect) and the interaction of the two. The table looks like this:

```
((df-interaction SS-interaction MS-interaction F-interaction p-interaction)
 (df-row          SS-row          MS-row          F-row          p-row)
 (df-column      SS-column      MS-column      F-column      p-column)
 (df-error       SS-error       MS-error)
 (df-total       SS-total       MS-total))
```

A useful function to print these out is 'print-anova-table'; it can tell which kind of table it is by the length of the list.

`anova-one-way-variables` *iv dv* *Optional (scheffe-tests-p t)* [Function]

Performs a one-way analysis of variance (ANOVA) on the input data, which should be two equal-length sequences: 'iv' is the independent variable, represented as a sequence of categories or group identifiers, and 'dv' is the dependent variable, represented as a sequence of numbers. The 'iv' variable must be "sorted," meaning that AAABBBBBCCCCDDDD is okay but ABCDABCDABDCDC is not, where A, B, C and D are group identifiers. Furthermore, each group should consist of at least 2 elements. The significance of the result indicates that the group means are not all equal; that is, at least two of the groups have significantly different means. If there were only two groups, this would be semantically equivalent to an unmatched, two-tailed t-test, so you can think of the one-way ANOVA as a multi-group, two-tailed

t-test. This function returns four values: 1. an ANOVA table; 2. a list a group means; 3. either a Scheffe table or nil depending on ‘scheffe-tests-p’; and 4. an alternate value for SST. The fourth value is only interesting if you think there are numerical accuracy problems; it should be approximately equal to the SST value in the ANOVA table. This function differs from ‘anova-one-way-groups’ only in its input representation. See the manual for more information.

`anova-one-way-groups data` *Optional (scheffe-tests-p t)* [Function]

Performs a one-way analysis of variance (ANOVA) on the ‘data,’ which should be a sequence of sequences, where each interior sequence is the data for a particular group. Furthermore, each sequence should consist entirely of numbers, and each should have at least 2 elements. The significance of the result indicates that the group means are not all equal; that is, at least two of the groups have significantly different means. If there were only two groups, this would be semantically equivalent to an unmatched, two-tailed t-test, so you can think of the one-way ANOVA as a multi-group, two-tailed t-test. This function returns four values: 1. an ANOVA table; 2. a list a group means; 3. either a Scheffe table or nil depending on ‘scheffe-tests-p’; and 4. an alternate value for SST. The fourth value is only interesting if you think there are numerical accuracy problems; it should be approximately equal to the SST value in the ANOVA table. This function differs from ‘anova-one-way-variables’ only in its input representation. See the manual for more information.

`print-anova-table anova-table` *Optional (stream *standard-output*)* [Function]

Prints ‘anova-table’ on ‘stream.’

Optionally, the ANOVA functions can also return scheffe-tables, which represent all pairwise comparisons between groups. A scheffe-table is an upper-triangular table represented with list structure. For example, with four groups, it would look like this:

```
((F-01 p-01) (F-02 p-02) (F-03 p-03))
 ((F-12 p-12) (F-13 p-13))
 ((F-22 p-23)))
```

This table makes more sense when formatted as follows:

	mean-1	mean-2	mean-3
mean-0:	F-01 p-01	F-02 p-02	F-03 p-03
mean-1:		F-12 p-12	F-13 p-13
mean-2:		F-22 p-22	F-23 p-23

The scheffe-tables can be computed and printed by the following functions.

`scheffe-tests group-means group-sizes ms-error df-error` [Function]

Performs all pairwise comparisons between group means, testing for significance using Scheffe’s F-test. Returns an upper-triangular table in a format described in the manual. Also see the function ‘print-scheffe-table.’ ‘Group-means’ and ‘group-sizes’ should be sequences. The arguments ‘ms-error’ and ‘df-error’ are the mean square error within groups and its degrees of freedom, both of which are computed by the analysis of variance. An ANOVA test should always be run first, to see if there are any significant differences.

`print-scheffe-table` *scheffe-table* *&optional group-means (stream* [Function]
**standard-output*)*

Prints ‘scheffe-table’ on ‘stream.’ If the original one-way anova data had N groups, the Scheffe table prints as an n-1 x n-1 upper-triangular table. If ‘group-means’ is given, it should be a list of the group means, which will be printed along with the table.

Two factor analysis of variance is computed by the following functions.

`anova-two-way-groups` *data-array* [Function]

Calculates the analysis of variance when there are two factors that may affect the dependent variable. Because the input is represented as an array, we can refer to these two factors as the row-effect and the column effect. Unlike the one-way ANOVA, there are mathematical difficulties with the two-way ANOVA if there are unequal cell sizes; therefore, we require all cells to be the same size, and so the input is a three-dimensional array. The result of the analysis is an anova-table, as described in the manual. This function differs from ‘anova-two-way-variables’ only in its input representation. See the manual for further discussion of analysis of variance.

`anova-two-way-variables` *dv iv1 iv2* [Function]

Calculates the analysis of variance when there are two factors that may affect the dependent variable, specifically ‘iv1’ and ‘iv2.’ Unlike the one-way ANOVA, there are mathematical difficulties with the two-way ANOVA if there are unequal cell sizes; therefore, we require all cells to be the same size; that is, the same number of values (of the dependent variable) for each combination of the independent factors. The result of the analysis is an anova-table, as described in the manual. This function differs from ‘anova-two-way-groups’ only in its input representation. See the manual for further discussion of analysis of variance. If you use ‘print-anova-table,’ the row effect is ‘iv1’ and the column effect is ‘iv2.’

6.2.4 Linear Regression

CLASP implements six different linear regression functions, the result of two kinds of input representations (data versus summaries) and three quantities of output (minimal, brief, and verbose). The verbose form returns an anova table in the following format:

```
((df-regression SS-regression MS-regression F p)
 (df-error      SS-error      MS-error)
 (df-total      SS-total))
```

`linear-regression-minimal` *dv iv* [Function]

Calculates the slope and intercept of the regression line. This function takes two equal-length sequences of raw data. Note that the dependent variable, as always, comes first in the argument list. You should first look at your data with a scatter plot to see if a linear model is plausible. See the manual for a fuller explanation of linear regression statistics.

`linear-regression-minimal-summaries` *n x y x2 y2 xy* [Function]

Calculates the slope and intercept of the regression line. This function differs from ‘linear-regression-minimal’ in that it takes summary statistics: ‘x’ and ‘y’ are the sums of the independent variable and dependent variables, respectively; ‘x2’ and ‘y2’ are the sums of the squares of the independent variable and dependent variables, respectively; and ‘xy’ is the sum of the products of the independent and dependent variables. You should first look at your data with a scatter plot to see if a linear model is plausible. See the manual for a fuller explanation of linear regression statistics.

`linear-regression-brief` *dv iv* [Function]

Calculates the main statistics of a linear regression: the slope and intercept of the line, the coefficient of determination, also known as r-square, the standard error of the slope, and the p-value for the regression. This function takes two equal-length sequences of raw data. Note that the dependent variable, as always, comes first in the argument list. You should first look at your data with a scatter plot to see if a linear model is plausible. See the manual for a fuller explanation of linear regression statistics.

`linear-regression-brief-summaries` *n x y x2 y2 xy* [Function]

Calculates the main statistics of a linear regression: the slope and intercept of the line, the coefficient of determination, also known as r-square, the standard error of the slope, and the p-value for the regression. This function differs from ‘linear-regression-brief’ in that it takes summary variables: ‘x’ and ‘y’ are the sums of the independent variable and dependent variables, respectively; ‘x2’ and ‘y2’ are the sums of the squares of the independent variable and dependent variables, respectively; and ‘xy’ is the sum of the products of the independent and dependent variables. You should first look at your data with a scatter plot to see if a linear model is plausible. See the manual for a fuller explanation of linear regression statistics.

`linear-regression-verbose` *dv iv* [Function]

Calculates almost every statistic of a linear regression: the slope and intercept of the line, the standard error on each, the correlation coefficient, the coefficient of determination, also known as r-square, and an ANOVA table as described in the manual. This function takes two equal-length sequences of raw data. Note that the dependent variable, as always, comes first in the argument list. If you don’t need all this information, consider using the “-brief,” or “-minimal” functions, which do less computation. You should first look at your data with a scatter plot to see if a linear model is plausible. See the manual for a fuller explanation of linear regression statistics.

`linear-regression-verbose-summaries` *n x y x2 y2 xy* [Function]

Calculates almost every statistic of a linear regression: the slope and intercept of the line, the standard error on each, the correlation coefficient, the coefficient of determination, also known as r-square, and an ANOVA table as described in the manual. If you don’t need all this information, consider using the “-brief” or “-minimal” functions, which do less computation. This function differs from ‘linear-regression-verbose’ in that it takes summary variables: ‘x’ and ‘y’ are the sums of the independent variable and dependent variables, respectively; ‘x2’ and ‘y2’

are the sums of the squares of the independent variable and dependent variables, respectively; and 'xy' is the sum of the products of the independent and dependent variables. You should first look at your data with a scatter plot to see if a linear model is plausible. See the manual for a fuller explanation of linear regression statistics.

The verbose functions do not currently return the standard error of the intercept. This will be fixed in the next release.

6.2.5 Contingency Table Analysis

Often, data collection will simply observe the occurrence of discrete events, rather than measuring some dependent variable. The data, then, consists of counts: n events of type A, m events of type B, and so on. Data tables of counts are called *contingency tables*. Hypotheses can then be tested about the distribution of events (their relative frequency) and whether one kind of event is independent of another. The following functions test such hypotheses.

There are generally two kinds of functions, depending on whether the input representation is a contingency table or sequences of events. In the latter case, the functions construct a contingency table by counting events and then call the other functions.

By convention, the cells in a 2x2 contingency table are labeled as follows:

a	b
c	d

This is the semantics for the arguments of 'chi-square-2x2-counts.' By symmetry, it doesn't matter if you swap the 'b' and 'c' arguments.

`chi-square-2x2-counts a b c d` *Optional (yates t)* [Function]

Runs a chi-square test for association on a simple 2 x 2 table. If 'yates' is nil, the correction for continuity is not done; default is t. Returns the chi-square statistic and the significance of the value.

`chi-square-2x2 v1 v2` [Function]

Performs a chi-square test for independence of the two variables, 'v1' and 'v2.' These should be categorical variables with only two values; the function will construct a 2x2 contingency table by counting the number of occurrences of each combination of the variables. See the manual for more details.

`chi-square-rc-counts contingency-table` [Function]

Calculates the chi-square statistic and corresponding p-value for the given contingency table. The result says whether the row factor is independent of the column factor. Does not apply Yate's correction.

`chi-square-rc v1 v2` [Function]

Performs a chi-square test for independence of the two variables, 'v1' and 'v2.' These should be categorical variables; the function will construct a contingency table by counting the number of occurrences of each combination of the variables. See the manual for more details.

6.2.6 Smoothing

The following functions “smooth” data in a non-parametric way, by replacing elements with a function of themselves and their neighbors. They aren’t statistics functions in the usual sense, but they operate on sequences and are useful in processing data before computing statistics. For example, smoothing a time series before doing a cross-correlation can produce a clearer result.

`smooth-median-2 data` *[Function]*

Smooths ‘data’ by replacing each element with the median of it and its neighbor on the left. A median of two elements is the same as their mean. The end is handled by duplicating the end element. This function is not destructive; it returns a list the same length as ‘data,’ which should be a sequence of numbers.

`smooth-median-3 data` *[Function]*

Smooths ‘data’ by replacing each element with the median of it and its two neighbors. The ends are handled by duplicating the end elements. This function is not destructive; it returns a list the same length as ‘data,’ which should be a sequence of numbers.

`smooth-median-4 data` *[Function]*

Smooths ‘data’ by replacing each element with the median of it, its left neighbor, and its two right neighbors. The ends are handled by duplicating the end elements. This function is not destructive; it returns a list the same length as ‘data,’ which should be a sequence of numbers.

`smooth-median-5 data` *[Function]*

Smooths ‘data’ by replacing each element with the median of it, its two left neighbors and its two right neighbors. The ends are handled by duplicating the end elements. This function is not destructive; it returns a list the same length as ‘data,’ which should be a sequence of numbers.

`smooth-hanning data` *[Function]*

Smooths ‘data’ by replacing each element with the weighted mean of it and its two neighbors. The weights are 1/2 for itself and 1/4 for each neighbor. The ends are handled by duplicating the end elements. This function is not destructive; it returns a list the same length as ‘data,’ which should be a sequence of numbers.

`smooth-4253h data` *[Function]*

Smooths ‘data’ by successive smoothing: 4,median; then 2,median; then 5,median; then 3,median; then hanning. The ends are handled by duplicating the end elements. This function is not destructive; it returns a list the same length as ‘data,’ which should be a sequence of numbers.

6.3 Data Manipulation Functions

This section describes the representation and functions for CLASP’s internal data manipulation code. These functions are called by the graphical interface when data is loaded, manipulated, transformed, saved, and so forth. The functions in this section could allow

you to manipulate your own data under program control, or they could be the basis for new manipulation functions.

The basic class definitions and methods are the following:

data	[Class]
'data' is the basic data class for clasp. Dataset, variable etc. are all subclasses of data.	
name object	[Method]
data 'name' returns the name of a data object.	
name object	[Method]
data 'description' returns the description of a data object.	
name object	[Method]
data 'id' returns the internal id of a data object.	
dataset	[Class]
The 'dataset' is the top-level user class for data representation in CLASP. A dataset is a collection of variables.	
variable	[Class]
'variable's are ordered collections of values. Valid data types for variables are numbers (integers, rationals and floats), symbols and strings. Symbols and strings can be used to represent for categorical data).	

6.3.1 Dataset Functions

With the following functions, we start creating and manipulating data, variables, and datasets. These are the fairly low-level functions for creating, accessing and destroying the data structures.

variable-value <i>the-variable</i> <i>&key</i>	[Function]
'variable-value' returns the value of a variable, expressed as a list.	
make-dataset <i>&key</i> (<i>name nil</i>) (<i>description nil</i>)	[Function]
Creates a new dataset and fills the slots with 'name', 'data', 'variable-list' and 'description'. 'data' is an rtm-table.	
make-dataset-from-rows <i>name data column-names</i> <i>&optional description</i>	[Function]
Creates a new data set and fills the slots with 'name', 'data' and 'description'. 'data' is a list of lists where each interior list represents one row of the data.	
make-dataset-from-columns <i>name data column-names</i> <i>&optional description</i>	[Function]
Creates a new data set and fills the slots with 'name', 'data' and 'description'. 'data' is a list of lists where each interior list represents one column of the data.	
rename-dataset <i>dataset new-name</i>	[Function]
rename-dataset changes the name of 'dataset' to 'new-name', insuring the new name is unique. 'dataset' can either be a dataset or the name of a dataset, and 'new-name' must be a string.	
delete-dataset <i>dataset</i>	[Function]

Destroys a dataset and removes it from the dataset lists

get-dataset *the-dataset* *&optional active-datasets-only-p* [Function]

Given a name or id, finds the dataset. If ‘active-datasets-only-p’ is nil, will search through all current clasp datasets, otherwise, only the active datasets are searched.

activate-dataset *dataset* [Function]

Takes a dataset on the available list **clasp-datasets** and adds it to the active list **clasp-active-datasets**

deactivate-dataset *dataset* [Function]

Takes a dataset off of the **clasp-active-datasets** list.

make-variable *&key name type attribute value description dataset rtm-table* [Function]

make-variable creates a new variable. ‘name’ is a string. If ‘type’ is one of :number, :symbol or :string, then that will be used as the type, otherwise the type will be inferred from (type-of ‘type’). ‘value’ is the data the variable should be initialized with. ‘description’ is a description of the variable. ‘dataset’ is the dataset to which the variable belongs.

add-variable-to-dataset *dataset data variable-name &optional (description nil)* [Function]

add-variable-to-dataset adds the ‘data’, which must be a list of values, to ‘dataset’ using ‘variable-name’ as the name of the new variable and ‘description’ as a description.

add-variables-to-dataset *dataset data-list variable-names &optional (descriptions nil)* [Function]

add-variables-to-dataset adds a number of variables to ‘dataset’. For each variable, *i*, its value (which must be a list) is the *i*th element of ‘data-list’, its name is the *i*th element of ‘variable-names’, and its description is the *i*th element of descriptions.

rename-variable *variable new-name* [Function]

rename-variable changes the name of ‘variable’ to ‘new-name’, insure the new name is unique. ‘variable’ can either be a variable or the name of a variable, and ‘new-name’ must be a string.

delete-variable *variable* [Function]

Destroys a variable. ‘variable’ can either be a variable or the name of a variable.

get-variable *the-variable &optional dataset active-datasets-only-p* [Function]

Given a name or id, finds the variable. If ‘dataset’ is specified, that is the dataset which will be searched for ‘variable’. Otherwise, a search will be made across datasets. If this is the case, then ‘active-datasets-only-p’ will determine whether all datasets currently loaded are searched, or just the active ones. Note, since variable names must only be unique within a dataset, not across datasets, care must be taken to request the correct variable. It is advisable to call get-variable with the dataset optional variable set.

add-row-to-dataset *dataset data variables* [Function]

This function will add ‘data’ to the dataset ‘dataset’ as a new row. ‘variables’ is a list of variables and specifies what order the new row is in with respect to the

variables in the dataset. Therefore, ‘data’ must match ‘variables’ in both length and the types of its values. Note: Whenever this function is called, `clear-cache` should be called on the same dataset.

`add-rows-to-dataset` *dataset data-list variables* [Function]

This function will add ‘data-list’ to the dataset ‘dataset’. ‘data-list’ is a list of lists representing rows of data. The rows must be of equal length. ‘variables’ is a list of variables and specifies what order the new rows are in with respect to the variables in the dataset. Therefore, each element of ‘data-list’ must match variables in both length and the types of its elements.

6.3.2 Input/Output

The following functions do input and output of data, so that they can be saved between sessions and processed by other programs. CLASP supports two ways of representing datasets in files: a special CLASP format and a general non-CLASP format that should be readable by most other statistics programs.

The CLASP file format is read by ‘load-dataset’ and saved by ‘save-dataset.’ The loading function takes a filename as its only argument, while the saving function takes a dataset and a filename. This format is also produced by CLIP, because we want CLIP and CLASP to have a seamless interface. The CLASP format is as follows:*

- The first line of the file is a string (surrounded by double-quotes) containing the name of the dataset.
- The following n lines of the file each have a single string with the name of one of the n variables in the dataset.
- At any point after the variable names, comments may be inserted by starting a line with a semicolon (;).
- The rest of the lines of the file each contain a Common Lisp list representing one row of data. The order of the elements must match the order of the variable names listed at the top of the file.

The non-CLASP file format is read by ‘import-dataset’ and saved by ‘export-dataset.’ They take extra arguments because the file representation allows user-specified separator characters between data values. Some statistics packages use commas between values, others use tabs or spaces. The separator character is specified by `:separator`, defaulting to a comma (`#\,`). Furthermore, the file representation allows, optionally, the names of the variables to appear in the file. You can allow them to appear using `:include-labels-p`, which defaults to `nil`. The non-CLASP format is as follows:

- The first line is, optionally, a line of variable names, expressed as symbols, and separated by separator characters.
- The rest of the file is lines of data values, separated by separator characters.

*For examples of CLASP format, see Section 2.4 and Appendix B.

Both ‘load-dataset’ and ‘import-dataset’ will signal an error if ‘filename’ does not specify an extant file. No other error checking is currently done, and so, for example, if some of the data rows have too few elements, as in the case of missing data, the dataset will be incorrect. In addition, the datatype of each variable will be inferred from its value in the first row of the data.

Both ‘save-dataset’ and ‘export-dataset’ will signal an error if ‘filename’ specifies an already extant file.

load-dataset *filename* [Function]

Reads a dataset in from a file in standard CLASP format. Specifically, the first line is a string containing the name of the dataset. The following *n* lines are strings containing the names of the variables in the dataset, and the following *m* lines are lists containing the rows of the dataset.

save-dataset *data filename* [Function]

Writes a dataset out to a file in CLASP format.

import-dataset *filename* *key (separator ,) include-labels-p* [Function]

Reads a dataset in from a file, where ‘filename’ is a string or pathname. If ‘filename’ doesn’t exist, an error is signaled. The data on each line in the file must be separated by ‘separator’ characters or whitespace. No error checking is done to insure that each line of the file contains the same number of data. If ‘include-labels-p’ is non-nil, the first line will be used a variable names, otherwise they will be named variable, variable-1, variable-2, etc. The filename will be used for the dataset name.

export-dataset *dataset filename* *key (separator #\,) include-labels-p* [Function]

Writes out the values in ‘dataset’ in row-major order to the file ‘filename.’ Each row in ‘filename’ will be one row of data in ‘dataset.’ The data in a row will be separated by the character ‘separator’. If include-labels-p is non-nil, the variable-names will be written out on the first line. ‘export-dataset’ will signal an error if a file of name ‘filename’ already exists.

6.3.3 Manipulation

Finally, we get to functions that transform or otherwise manipulate existing data.

create-new-column-from-function *dataset expression* [Function]

create-new-column-from-function adds a new variable to ‘dataset’, using ‘expression’ to calculate the values. ‘expression’ may be any numerical expression. If variables are included in ‘expression’, they should be referenced by name. For instance, if a dataset has two variables, “hits” and “misses”, then

```
(create-new-column-from-function dataset
      (/ hits (+ hits misses)))
```

will add a new variable with the hit-to-attempt ratio.

When create-new-column-from-function is called, the rows of ‘dataset’ are mapped over, and for each row, substitutions are made into ‘expression’ for any variable references. The value returned is added to the new variable in the appropriate row of ‘dataset’. Because the dataset is mapped over, aggregating functions (such as mean

or data-length) cannot be used. A new function, ‘recode’, has been implemented to allow for mapping of categorical variables to new values.

recode *value old-values new-values* [Function]

Maps ‘value’ from its position in ‘original-values’ to the value at the equivalent position in ‘new-values’. For instance, (recode 2 '(1 2 3) '(a b c)) would return 'b. This function can be used with create-new-column-from-function to recode a categorical variable. From the program interface, it is recommended that “recode-list” be used instead.

recode-list *list old-values new-values* [Function]

Maps each element of ‘list’ from its position in ‘old-values’ to the value at the equivalent position in ‘new-values’.

partition-dataset *old-dataset partition-clause* [Function]

Often it is desirable to operate on a portion of a dataset, or to separate out different parts of a dataset according to the values of some of the variables. When a dataset is broken up this way, it is called partitioning, and the parts of the dataset created are called partitions. partition-dataset takes a dataset and a partition-clause which describes how to split the dataset up, and returns a list of datasets, each one containing one partition. (Note, this does not destroy the original dataset.) Partition clauses are made up of boolean operators, and the partitioning operator, “.on.”. The boolean operators act as expected and are explained in detail in the reference manual. The .on. operator takes any number of variables as its operands, and separates the dataset into different partitions for each unique set of values of its operands. For example, if a dataset consisted of a boolean variable “key1”, a categorical variable “key2”, with values 'a, 'b and 'c, and a data variable “y”, and was partitioned with a partitioning clause of (.on. key1 key2), then at most six new datasets would be created, one containing all the rows of the original dataset where key1 = true and key2 = 'a, one with all the rows where key1 = false and key2 = 'a, one with all the rows where key1 = true and key2 = 'b, and so on. If a partition would be empty (i.e. if there were no rows where key1 = false and key2 = 'c), then no dataset is created for that partition, which is why AT MOST six new datasets would be created in the above example. The .on. operator can be used in conjunction with the boolean operators. When this happens, the .on. operator can be thought of as a macro which causes a set of partitioning clauses to be created, each exactly like the original, except that .on. is replaced with an .==. operator and one of the unique values of the key. If there are multiple keys in the .on., it is replaced by (.and. (.==. key1 value1) (.==. key2 value2), etc.). In the above example, a partition clause of (.and. (.==. key1 true) (.on. key2)) would cause up to three datasets to be created. Each would contain rows where key1 was true, and one would be created for key2 = 'a, one for key2 = 'b and one for key2 = 'c.

dataset-to-rows *data &key rows columns* [Function]

Takes a dataset ‘data’ and extracts tuples representing the rows of data in the dataset.

dataset-to-columns *data &key rows columns* [Function]

Takes a dataset ‘data’ and extracts tuples representing the columns of data in the dataset.

`transpose data`

[Function]

Takes a list of lists, ‘data’, and makes a list of tuples `t` where `t[i]` contains all the `i`’th elements of the lists in ‘data’.

6.4 Density and Combinatorial Functions

Given a statistic, such as the t-statistic, the final step in hypothesis testing is to compute the p-value of the statistic—that is, the probability that a value for the statistic as extreme or more so could occur by chance. In statistics books, p-values are found by looking in tables in the appendices, but our computer programs can’t do that. Instead, they must calculate the p-value.

The following functions are how CLASP calculates p-values. Some of them may be completely useless to you unless you are implementing new statistical functions. (If you do, let us know; perhaps they can be incorporated into the next release.) Some are useful if you want to graph theoretical distributions or otherwise experiment with them. Finally, there are some functions that calculate combinatorial probabilities, such as the binomial.

Most of the following functions are implemented from *Numerical Recipes in C* [10]. The code has been extensively tested and the values compared to tables of statistics or, in some cases, to figures in the *Numerical Recipes* book.

`beta z w`

[Function]

Returns the value of the Beta function, defined in terms of the complete gamma function, G , as: $G(z)G(w)/G(z+w)$. The implementation follows *Numerical Recipes in C*, section 6.1.

The complete Beta function is not used in CLASP, but may be helpful to some users.

`gamma-incomplete a x`

[Function]

This is an incomplete gamma function, what *Numerical Recipes in C* calls “gammmp.” This function also returns, as the second value, $g(a,x)$. See the manual for more information.

The incomplete gamma function, notated $p(a,x)$ is mathematically defined as

$$p(a, x) = g(a, x)/G(a)$$

where $g(a, x)$ is what Mathematicatm calls the incomplete gamma function and is mathematically defined as

$$\int_0^x t^{a-1} e^{-t} dt$$

and G is the complete gamma function, defined as

$$\int_0^\infty t^{a-1} e^{-t} dt$$

Yes, it’s very confusing; apparently, there is no standard naming scheme for these functions. Our implementation follows *Numerical Recipes in C*, section 6.2 and is the function that they call “gammmp.” We did not implement the function they call “gammq,” since it is just 1.0-gammmp.

error-function x

[Function]

Computes the error function, which is typically used to compute areas under the Gaussian probability distribution. See the manual for more information. Also see the function ‘gaussian-cdf.’ This implementation follows *Numerical Recipes in C*, section 6.2

The error function is not well described in *Numerical Recipes* or other books that we checked, and it mystified us for some time. Once we understood it, it made perfect sense, so we offer the following explanation.

The error function is defined as

$$\operatorname{erf}(x) = \frac{2}{\sqrt{\pi}} \int_0^x e^{-t^2} dt$$

The error function is related to the integral of the Gaussian probability density function (pdf), as we can see in the following proof. First, define $A(x)$ to be the area under the Gaussian pdf and less than x in absolute value. (This is the complement of a two-tailed test of the significance of x .)

$$A(x) = \int_{-x}^x \frac{1}{\sqrt{2\pi}} e^{-t^2/2} dt$$

Because the Gaussian is symmetrical, if $x > 0$, we can just integrate from 0 to x and double it:

$$A(x) = 2 \int_0^x \frac{1}{\sqrt{2\pi}} e^{-t^2/2} dt$$

Now we can show that A is closely related to erf . The proof works by a change of variable, substituting $t = u\sqrt{2}$. All the rest is just algebra.

$$\begin{aligned} A(x) &= 2 \int_0^x \frac{1}{\sqrt{2\pi}} e^{-t^2/2} dt \\ &= 2 \int_0^{x/\sqrt{2}} \frac{1}{\sqrt{2\pi}} e^{-(u\sqrt{2})^2/2} du\sqrt{2} \\ &= 2 \int_0^{x/\sqrt{2}} \frac{1}{\sqrt{\pi}} e^{-u^2} du \\ &= \frac{2}{\sqrt{\pi}} \int_0^{x/\sqrt{2}} e^{-u^2} du \\ &= \operatorname{erf}(x/\sqrt{2}) \end{aligned}$$

That change of variable is the key to understanding erf : erf is the integral (from zero outwards) of a function that is shaped exactly like the Gaussian, except that the axis has been “stretched” by a factor of $\sqrt{2}$. So, to find $A(x)$, we just compute $\operatorname{erf}(x/\sqrt{2})$.

Note, by defining erf to be the integral from 0 to x , we have $\operatorname{erf}(-x) = -\operatorname{erf}(x)$. With this in mind, we can compute other integrals of the Gaussian using erf . Define $F(x)$ to be the cumulative distribution function of the Gaussian. Then

$$F(x) = \frac{1}{2}(1 + \operatorname{erf}(x/\sqrt{2}))$$

Indeed, erf can be used to compute the cdf of Gaussians with mean m and standard deviation s by standardizing and dividing by the square root of 2.

`gaussian-cdf` *x* *ℰoptional (mean 0.0) (sd 1.0)* [Function]

Computes the cumulative distribution function for a Gaussian random variable (defaults: mean=0.0, s.d.=1.0) evaluated at ‘x.’ The result is the probability of getting a random number less than or equal to ‘x,’ from the given Gaussian distribution.

`error-function-complement` *x* [Function]

This function computes the complement of the error function, “erfc(x),” defined as $1 - \text{erf}(x)$. See the documentation for ‘error-function’ for a more complete definition and description. Essentially, this function on $z/\sqrt{2}$ returns the two-tailed significance of z in a standard Gaussian distribution. This function implements the function that *Numerical Recipes in C* calls `erfcc`, see section 6.3; that is, it’s the one using the Chebyshev approximation, since that is the one they call from their statistical functions. It is quick to compute and has fractional error everywhere less than 1.2×10^{-7} .

`gaussian-significance` *x tails ℰoptional mean sd* [Function]

Computes the significance of ‘x’ in a Gaussian distribution with mean=‘mean’ (default 0.0) and standard deviation=‘sd’ (default 1.0); that is, it returns the area which farther from the mean than ‘x’ is. The null hypothesis is roughly that ‘x’ is zero; you must specify your alternative hypothesis (H1) via the ‘tails’ parameter, which must be `:both`, `:positive` or `:negative`. The first corresponds to a two-tailed test: H1 is that ‘x’ is not zero, but you are not specifying a direction. If the parameter is `:positive`, H1 is that ‘x’ is positive, and similarly for `:negative`.

Departing from the Gaussian distribution, we can compute the significance of other kinds of statistics.

`poisson-cdf` *k x* [Function]

Computes the cumulative distribution function for a Poisson random variable with mean ‘x’ evaluated at ‘k.’ The result is the probability that the number of Poisson random events occurring will be between 0 and k-1 inclusive, if the expected number is ‘x.’ The argument ‘k’ should be an integer, while ‘x’ should be a float. The implementation follows *Numerical Recipes in C*, section 6.2

`chi-square-significance` *x dof* [Function]

Computes the complement of the cumulative distribution function for a Chi-square random variable with ‘dof’ degrees of freedom evaluated at ‘x.’ The result is the probability that the observed chi-square for a correct model should be greater than ‘x.’ The implementation follows *Numerical Recipes in C*, section 6.2. Small values suggest that the null hypothesis should be rejected; in other words, this computes the significance of ‘x.’

The incomplete “beta” function is not particularly useful by itself, but is useful in defining the cumulative distributions for Student’s t and the F distribution. It is mathematically defined as

$$I_x(a, b) = \frac{1}{B(a, b)} \int_0^x t^{a-1} (1-t)^{b-1} dt$$

where $B(a, b)$ is the complete “beta” function.

`beta-incomplete` *a b x* [Function]

This function is useful in defining the cumulative distributions for Student's t and the F distribution. All arguments must be floating-point numbers; 'a' and 'b' must be positive and 'x' must be between 0.0 and 1.0, inclusive.

`students-t-significance` *t-statistic dof tails* [Function]

Student's distribution is much like the Gaussian distribution except with heavier tails, depending on the number of degrees of freedom, 'dof.' As 'dof' goes to infinity, Student's distribution approaches the Gaussian. This function computes the significance of 't-statistic.' Values range from 0.0 to 1.0: small values suggest that the null hypothesis—that 't-statistic' is drawn from a t distribution—should be rejected. The 't-statistic' parameter should be a float, while 'dof' should be an integer. The null hypothesis is roughly that 't-statistic' is zero; you must specify your alternative hypothesis (H1) via the 'tails' parameter, which must be `:both`, `:positive` or `:negative`. The first corresponds to a two-tailed test: H1 is that 't-statistic' is not zero, but you are not specifying a direction. If the parameter is `:positive`, H1 is that 't-statistic' is positive, and similarly for `:negative`. This implementation follows *Numerical Recipes in C*, section 6.3.

`f-significance` *f-statistic numerator-dof denominator-dof* *optional one-tailed-p* [Function]

This function occurs in the statistical test of whether two observed samples have the same variance. A certain statistic, F, essentially the ratio of the observed dispersion of the first sample to that of the second one, is calculated. This function computes the tail areas of the null hypothesis: that the variances of the numerator and denominator are equal. It can be used for either a one-tailed or two-tailed test. The default is two-tailed, but one-tailed can be computed by setting the optional argument 'one-tailed-p' to true. For a two-tailed test, this function computes the probability that F would be as different from 1.0 (larger or smaller) as it is, if the null hypothesis is true. For a one-tailed test, this function computes the probability that F would be as LARGE as it is if the first sample's underlying distribution actually has SMALLER variance than the second's, where 'numerator-dof' and 'denominator-dof' is the number of degrees of freedom in the numerator sample and the denominator sample. In other words, this computes the significance level at which the hypothesis "the numerator sample has smaller variance than the denominator sample" can be rejected. A small numerical value implies a very significant rejection. The 'f-statistic' must be a non-negative floating-point number. The degrees of freedom arguments must be positive integers. The 'one-tailed-p' argument is treated as a boolean. This implementation follows *Numerical Recipes in C*, section 6.3 and the 'ftest' function in section 13.4. Some of the documentation is also drawn from the section 6.3, since I couldn't improve on their explanation.

Some useful combinatorial functions are the following. The 'gamma-ln' function is the natural logarithm of the gamma function, which is the continuous form of the factorial function.

`gamma-ln` *x* [Function]

Returns the natural logarithm of the Gamma function evaluated at 'x.' Mathematically, the Gamma function is defined to be the integral from 0 to Infinity of $t^x \exp(-t) dt$. The implementation is copied, with extensions for the reflection formula, from *Numerical Recipes in C*, section 6.1. The argument 'x' must be positive.

Full accuracy is obtained for $x > 1$. For $x < 1$, the reflection formula is used. The computation is done using double-floats, and the result is a double-float.

factorial-exact n [Function]

Returns the factorial of ‘ n ,’ which should be an integer. The result will returned as an integer or bignum. This implementation is exact, but is more computationally expensive than ‘factorial,’ which is to be preferred.

factorial n [Function]

Returns the factorial of ‘ n ,’ which should be a non-negative integer. The result will returned as a floating-point number, single-float if possible, otherwise double-float. If it is returned as a double-float, it won’t necessarily be integral, since the actual computation is $(\exp(\gamma - \ln(1+n)))$. Implementation is loosely based on *Numerical Recipes in C*, section 6.1. On the TI Explorer, the largest argument that won’t cause a floating overflow is 170.

factorial-ln n [Function]

Returns the natural logarithm of $n!$; ‘ n ’ should be an integer. The result will be a single-precision, floating point number. The implementation follows *Numerical Recipes in C*, section 6.1

binomial-coefficient $n k$ [Function]

Returns the binomial coefficient, ‘ n ’ choose ‘ k ,’ as an integer. The result may not be exactly correct, since the computation is done with logarithms. The result is rounded to an integer. The implementation follows *Numerical Recipes in C*, section 6.1

binomial-coefficient-exact $n k$ [Function]

This is an exact but computationally intensive form of the preferred function, ‘binomial-coefficient.’

binomial-probability $p n k$ [Function]

Returns the probability of ‘ k ’ successes in ‘ n ’ trials, where at each trial the probability of success is ‘ p .’ This function uses floating-point approximations, and so is computationally efficient but not necessarily exact.

binomial-probability-exact $p n k$ [Function]

This is an exact but computationally intensive form of the preferred function, ‘binomial-probability.’

binomial-cdf $p n k$ [Function]

Suppose an event occurs with probability ‘ p ’ per trial. This function computes the probability of ‘ k ’ or more events occurring in ‘ n ’ trials. Note that this is the complement of the usual definition of cdf. This function approximates the actual computation using the incomplete beta function, but is preferable for large ‘ n ’ (greater than a dozen or so) because it avoids summing many tiny floating-point numbers. The implementation follows *Numerical Recipes in C*, section 6.3.

binomial-cdf-exact $p n k$ [Function]

This is an exact but computationally intensive form of the preferred function, ‘binomial-cdf.’

Distribution	Slots	Legal value
Uniform-distribution	Minimum	Real number
	Maximum	Real number
Normal-distribution	Mean	Real number
	Standard-deviation	Non-negative real number
Binomial-distribution	P	Real number from 0 to 1
	N	Positive integer
Poisson-distribution	Mean	Real number
Gamma-distribution	IA	Positive real number

Table 6.1: Distribution classes and parameters

The following function uses binary search to find a critical value of a statistic, given a distribution function for the statistic.

```
find-critical-value p-function p-value Optional (x-tolerance 1.0e-5) [Function]
(y-tolerance 1.0e-5)
```

Returns the critical value of some statistic. The function ‘p-function’ should be a unary function mapping statistics—x values—to their significance—p values. The function will find the value of x such that the p-value is ‘p-value.’ The function works by binary search. A secant method might be better, but this seems to be acceptably fast. Only positive values of x are considered, and ‘p-function’ should be monotonically decreasing from its value at x=0. The binary search ends when either the function value is within ‘y-tolerance’ of ‘p-value’ or the size of the search region shrinks to less than ‘x-tolerance.’

6.5 Sampling

CLASP provides random sampling from uniform, normal, binomial, poisson and gamma distributions. Distributions are CLOS classes. The distribution classes are defined as follows:

Sampling from a distribution is a two step process. The first step is to create the distribution, and the second step is to draw a sample from the distribution. To create a distribution, use `make-instance`. To draw a sample from a distribution use either `sample-to-list` or `sample-to-dataset`.

The following example shows how to create a standard-normal distribution and then make a dataset from two samples of length 50 drawn from the distribution.

```
;;; create a standard normal distribution
(setf dist (make-instance normal-distribution :mean 0
                                           :standard-deviation 1))

;;; draw two 50 element samples from the distribution and create a
;;; dataset from the result
(sample-to-dataset dist 50 2)
```

```
sample-to-list ((dist distribution) size-of-sample) [Method]
```

Gets a sample from the distribution 'dist' of size 'size-of-sample'

`sample-to-dataset` (*(dist distribution) size-of-samples* &optional *number-of-samples dataset-name variable-names*) [Method]

Gets a sample from the distribution 'dist' of size 'size-of-sample' and creates a dataset out of it. If 'number-of-samples' is > 1 then multiple samples are taken from the distribution and each one forms one column of the dataset.

6.6 Graphics

All graphics in CLASP are obtained via the SCIGRAPH graphing system from BBN. Please refer to the documentation available with that system to use graphical functions in CLASP. SCIGRAPH is available via anonymous ftp from *cambridge.apple.com* in the directory */pub/clim/clim-1-and-2/scigraph*.

Appendix A

Computational Methods

Essentially, this chapter acts as extended documentation for some of the more complicated statistical functions. Statistical functions, by their very nature, take collections of numbers, do some calculations that are often mysterious and unintuitive, and produce some other numbers that are difficult to check by hand. Consequently, these functions are hard for programmers to debug and hard for users to trust. This chapter will try to explain the calculations for some of the more arcane functions.

A.1 Quantiles

Quantiles seem like they ought to be simple, since they are just a generalization of the median, which splits the data in half. All we need to do is split the data into q and $1-q$. But there are other valid ways of viewing the quantile, which lead to other ways of dividing the data. As Freund says in his discussion of quartiles, “there is ample room for arbitrariness” [5, p. 52].

x_0	x_1	x_2	x_3	x_4
2	3	5	8	13

Consider the preceding five-element sample, which we will use throughout our discussion of quantiles. Clearly, x_2 is the median. What is x_1 ? There are three reasonable answers:

- **First Quartile** x_1 splits the sample so that three times as many elements are above it as below it; splitting the sample in a ratio of 1:3 is one definition of the first quartile.
- **First Quintile** One fifth of the sample is less than x_1 , and so it is the first quintile.
- **Second Quintile** Two fifths of the sample is less than or equal to x_1 , and so it is the second quintile.

There are two kinds of definitions here, resulting in different computational methods. The first is what we will call the “bisection” method because we look at how the data are split, and the second is the “cumulative” method, because we refer to the cumulative distribution function of the data. Of course, in the cumulative method, we also have to think about whether we want “ $<$ ” or “ \leq .” These methods will be discussed in turn. Only the bisection method is implemented in CLASP, so some readers may want to read only subsection A.1.1; readers who want to understand why we chose the bisection method over the cumulative method will want to read this whole section.

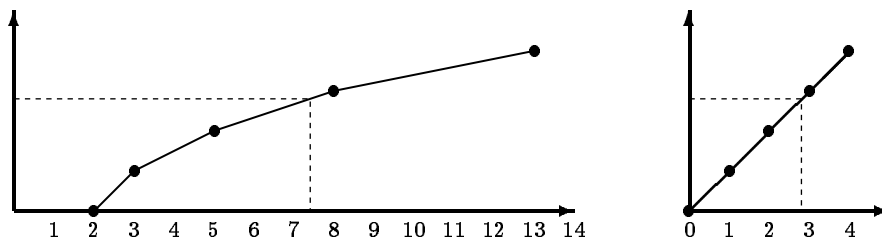


Figure A.1: In these graphs, the vertical axis is q , the quantile, and the horizontal axis is either x (at left) or i (at right). The bisection method maps q directly onto $i \in [0, n-1]$, as shown at right. The graph at the left shows how q maps onto the data values, x_i .

In either case, the way to think about the quantile is as a function from the quantile $q \in [0, 1]$ to the indices of the data, $i \in [0, n-1]$, where n is the number of sample items. This makes sense because the actual data values have little influence on the quantile; their relative ranks matter more. Their magnitudes will matter when the quantile falls between two values and we interpolate.

A.1.1 Bisection

The bisection method is very easy, because it maps q directly onto i in the obvious manner, $i = q(n-1)$, as shown in Figure A.1. The median conforms easily to this definition: $0.5(5-1) = 2$, and x_2 is the median of our example. Similarly, $0.25(5-1) = 1$ and x_1 is the first quartile. Using the bisection method, $q = 0.0$ is the smallest element in the sample and $q = 1.0$ is the largest.

What about when i is not an integer? For example, suppose $q = 0.7$, so the element “index” is 2.8, as illustrated in Figure A.1. Letting $i = \lfloor q(n-1) \rfloor$, we can linearly interpolate between the two bracketing values, x_i and x_{i+1} , which have exact quantiles of $q_i = i/(n-1)$ and $q_{i+1} = (i+1)/(n-1)$. The linear interpolation is a weighted sum of x_i and x_{i+1} where the weights are such that q is the same weighted sum of q_i and q_{i+1} . Let the weights be λ and $(1-\lambda)$ for elements $i+1$ and i , respectively.

$$\begin{aligned} q &= \lambda q_{i+1} + (1-\lambda)q_i \\ q &= \lambda \left(\frac{i+1}{n-1} \right) + (1-\lambda) \left(\frac{i}{n-1} \right) \\ q(n-1) &= \lambda(i+1) + (1-\lambda)i \\ q(n-1) &= \lambda + i \\ \lambda &= q(n-1) - i \\ \lambda &= q(n-1) - \lfloor q(n-1) \rfloor \end{aligned}$$

In other words, λ is the fractional part left over after we compute the floor of $q(n-1)$, so we get it for free in our calculation of i . We then compute:

$$x_{q(n-1)} = \lambda x_{i+1} + (1-\lambda)x_i$$

In our example with $q = 0.7$, $\lambda = 2.8 - \lfloor 2.8 \rfloor$, and $x_{2.8} = (0.2)5 + (0.8)8 = 7.4$.

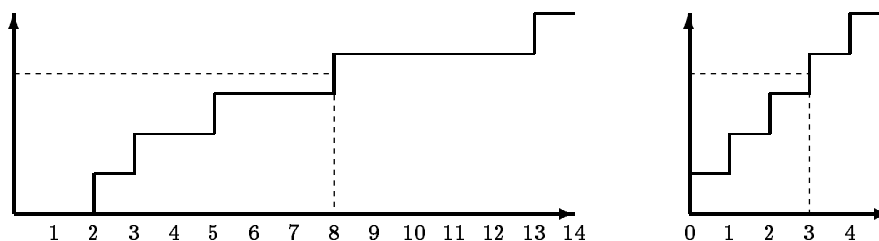


Figure A.2: Inverting the discrete CDF for our sample. In both graphs, each vertical step is $1/n$, since each element is that fraction of the whole. On the left, we map directly to the x_i , while on the right we just map to i , which has a very easy structure, since each horizontal step is of unit length.

A.1.2 Cumulative

The cumulative distribution function (CDF) of a random variable \mathbf{X} is defined as $F(x) = \Pr(\mathbf{X} \leq x)$. As such, the CDF maps from the domain of the random variable to the interval $[0, 1]$. The quantile function is simply the inverse of the CDF! For example, we know that if \mathbf{X} has a standard normal distribution, $\Pr(\mathbf{X} < -1.96) = 0.025$, so -1.96 is the $q = 0.025$ quantile of the standard normal.

Figure A.2 shows two cumulative distribution functions for our sample. They are step functions because our sample is finite, and each element adds $1/n$ to the running total. The function on the left has the x_i as its domain, while the function on the right has i as its domain. The idea is to invert this function, so start at q on the vertical axis, go right until we hit the function, and then go down until we hit the horizontal axis, just as we did in Figure A.1. In the graph on the right, we hit i such that:

$$\begin{aligned} i/n &< q < (i+1)/n \\ i &< nq < (i+1) \end{aligned}$$

So, $i = \lfloor qn \rfloor$. Compare this with the bisection method, which gives $i = \lfloor q(n-1) \rfloor$; the algorithmic difference is obvious, but the semantic difference is subtle.

Note that the cumulative method is the inverse of the bisection method with respect to non-integer values of i : with bisection, we knew what to do when $q(n-1)$ was an integer, and we had to resolve what to do when it is not, while with the cumulative method, we have not resolved what to do when qn is an integer. When qn is an integer, the horizontal line hits the top of one vertical segment and the bottom of the next vertical segment. To which vertical segment does q belong? This question is identical to whether we define quantiles in terms of “ $<$ ” or “ \leq .” In other words, is $i/n < q$ or is $i/n \leq q$? This is a fairly arbitrary decision, but we can reduce its effects by interpolating between values.

Suppose we want to interpolate with the cumulative method. Remember that, currently, the cumulative method always returns one of the elements of the sample. Indeed, this is a major flaw of the cumulative method, because the cumulative method cannot yet compute the median when n is even. In our example, $n = 5$, so it computes the median as x_i where $i = \lfloor qn \rfloor = 2$, which is correct. But if $n = 4$, the median is defined to be the average of the middle two elements and is not equal to any element of the sample.

Linear interpolation in the cumulative method can be visualized by connecting the vertical lines in our cumulative distribution function (see Figure A.3). Unfortunately, because

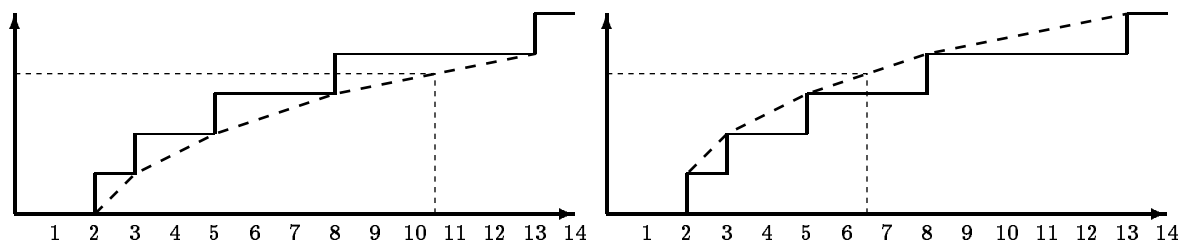


Figure A.3: The ways of interpolating in the cumulative method: either connecting the bottoms of the vertical segments (left) or the tops (right).

of the “ $<$ ” versus “ \leq ” problem, it can be done in two different ways, either connecting the bottoms of the vertical lines (“ $<$ ”) or the tops (“ \leq ”). Once we find the right interval, the interpolation is done the same as in the bisection method, computing the weight λ from the fractional part of $i = qn$. In our $q = 0.7$ example, the interpolated value is either 10.5 (halfway between the 0.8 quantile, $x_3 = 8$, and the 1.0 quantile, $x_4 = 13$) or 6.5 (halfway between the 0.8 quantile, $x_2 = 5$, and the 0.8 quantile, $x_3 = 8$), as is shown in Figure A.3. Contrast these values with the 7.4 obtained with the bisection method.

Note in figure A.3 that when we connect the bottoms of the segments, the last vertical segment (over $x_4 = 13$) belongs to the CDF, while when we connect the tops, the first vertical segment (over $x_0 = 2$) belongs to the CDF. This produces an essential asymmetry in each curve. Most elements influence two n ths of the curve (the interval on either side), but the extremes are different. One extreme influences only one interval, while the other extreme influences one interval and completely determines another. For example, when connecting the bottoms of the verticals, it seems that $x_4 = 13$ has more than its fair share of influence. One solution to this asymmetry is to count the quantile from the nearer end. This amounts to connecting the tops of the segments for $q < 0.5$ and the bottoms of the segments for $q > 0.5$ or vice versa. Thus either both extreme elements get undue influence (relative to middle elements) or neither does.

A.1.3 Discussion

We have presented two ways of viewing the quantile. One views the quantile as bisecting a line from 0 to $n-1$, with data elements sitting on the line at the integer positions. The other views the quantile with reference to the CDF and leads to unending problems all stemming from how to treat the vertical steps in that CDF.

One problem that still remains with the cumulative method is how to find the median. In our example, element $x_2 = 5$ is either at quantile $q = .4$ or $q = .6$, and so it cannot be the median. A solution can certainly be found, but it doesn’t fall easily out of the cumulative method.

For these reasons, we have chosen to implement the bisection method as simpler and more elegant than any version of the cumulative method. Figure A.4 shows all three functions in one large graph; you can see how the three methods will disagree on every quantile (except the extremes), and sometimes even interpolate in different intervals. Nevertheless, the disagreements will be fairly small and will shrink as the samples get larger.

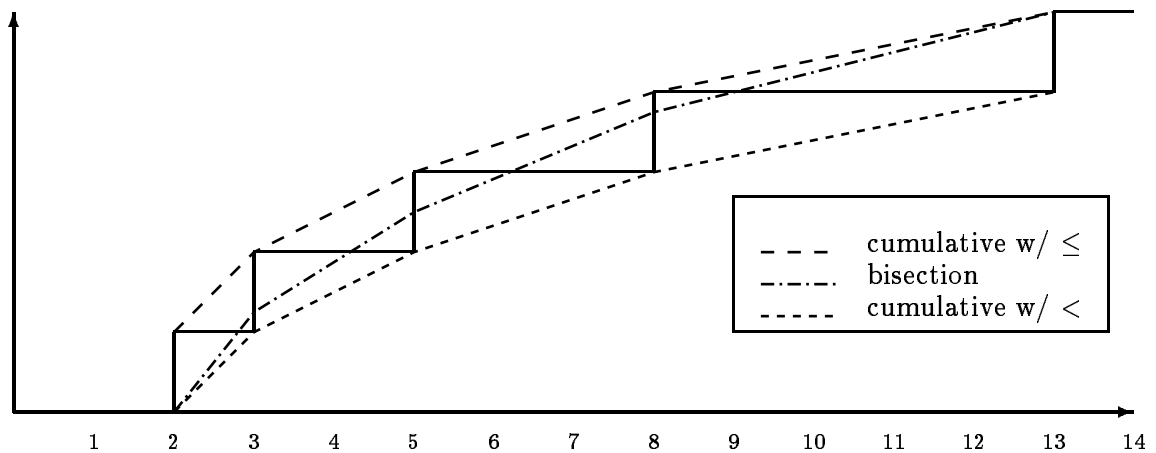


Figure A.4: A graphical comparison of the bisection method (middle line) the cumulative method with “<” (bottom line) and the cumulative method with “≤” (top line).

A.2 Analysis of Variance

The heart of the analysis of variance is, of course, calculating different kinds of variance in the data. Calculations of sums-of-squares (variance is just a sum-of-squares divided by the number of degrees of freedom) can generally be done in two different ways:

$$\sum_{i=1}^n (x_i - \bar{x})^2 = \sum_{i=1}^n x_i^2 - 1/n \left(\sum_{i=1}^n x_i \right)^2$$

These expressions are algebraically equivalent, as can easily be shown. We will do the proof carefully once, so that when we talk about different kinds of variance below, we can give sketches of the algebra, skipping steps that are similar to this proof. Here’s the careful proof.

$$\sum (x_i - \bar{x})^2$$

We dropped the indices of summation, since they will always be i from 1 to n . The next step is to expand the squaring inside the summation.

$$\sum (x_i^2 - 2x_i\bar{x} + \bar{x}^2)$$

We now split this into three different sums.

$$\sum x_i^2 - \sum 2x_i\bar{x} + \sum \bar{x}^2$$

Factor out constant quantities.

$$\sum x_i^2 - 2\bar{x} \sum x_i + \bar{x}^2 \sum$$

The middle sum is just the total of the numbers, call it t , and the last sum is the number of numbers, n

$$\sum x_i^2 - 2\bar{x}(t) + \bar{x}^2(n)$$

Note that $\bar{x} = t/n$, so let's eliminate \bar{x} .

$$\sum x_i^2 - 2t/n(t) + (t/n)^2(n)$$

Combine factors in the last two terms.

$$\sum x_i^2 - 2t^2/n + t^2/n$$

Combine the last two terms.

$$\sum x_i^2 - t^2/n$$

This last expression is just the right side of the equation, once we substitute the definition of t back in.

$$\sum x_i^2 - 1/n(\sum x_i)^2$$

Which formula should we use in our implementation? There are several issues in answering that question:

Efficiency The first expression requires computing the mean, which is n additions and a division. Then it requires n subtractions, squarings and additions. The second expression requires n squarings and additions to compute the first sum, n additions, a squaring and a division for the second, and a final subtraction. We can compare them as follows:

	Adds	Subs	Squares	Divides
$\sum (x_i - \bar{x})^2$	2n	n	n	1
$\sum x_i^2 - t^2/n$	2n	1	n+1	1

Basically, there's only a hair's breadth difference between the two. On most machines, squaring is more expensive than even quite a few subtractions, but the efficiency issue will clearly be our least important.

Data Representation The first expression clearly requires two passes over the data—once to compute the mean and the second time to compute the variance. The second expression requires only one pass over the data, because we can compute the total while we're summing the x_i^2 ; we can even count the data to compute n if that's necessary. For some data representations, this issue will dominate.

In CLASP, there are two ANOVA functions, one in which the values are structurally grouped (sequences of sequences) and another in which each value is paired with a key that tells what group it's in; for one-way ANOVA, these functions are 'anova-one-way-groups' and 'anova-one-way-variables' respectively. With the structural representation, it's easy to do several passes over a group. The other representation requires

figuring out the groups, which is not too hard, but may tip the scales towards the one-pass expression.

Numerical Accuracy If a data value lies quite close to the mean, the difference will be very small, and squaring it could produce a numerical underflow. At the very least, the first expression (subtract and square) will tend to cancel the significant digits and square the uncertain digits, resulting in numerical instability. The second expression will preserve the significant digits for as long as possible, but because it is adding up squared data values, it runs the risk of numerical overflow. (Numerical overflow and underflow will only happen if the data are floating point; if they are exact (integer or rational), the exactness will be preserved by the CLASP functions.) Numerical accuracy suggests that the second expression is better.

Since there are two ways of computing the variance, couldn't we compute *both* and use them to improve our numerical accuracy? For example, we could signal an error if they are very different, or we could simply take their average as the "true" variance. We could do this, but it seems excessive for most purposes. However, we will shortly see a way to exploit this redundancy.

The one-way ANOVA computes three sums of squares, called the sums of squares total, group, and error, or SST, SSG and SSE, respectively. By definition, they are the following:

$$\begin{aligned} SST &= \sum_{i=1}^I \sum_{j=1}^{J_i} (x_{ij} - M)^2 \\ SSG &= \sum_{i=1}^I \sum_{j=1}^{J_i} (M_i - M)^2 \\ SSE &= \sum_{i=1}^I \sum_{j=1}^{J_i} (x_{ij} - M_i)^2 \end{aligned}$$

The notation is a little complicated, but not too bad. There are I groups, and they all might be different sizes, denoted J_i . Thus, each of these summations has one term for each data value, x_{ij} . The first sum of squares, SST, subtracts the data values from the grand mean, denoted M . The second sum of squares, SSG, subtracts the mean of the i th group, denoted M_i , from the grand mean. The third sum of squares, SSE, subtracts the data values from the group means.

There is a fundamental identity, called the ANOVA identity, that states:

$$SST = SSG + SSE$$

This identity doesn't look obvious from the formulas above; however, it becomes more obvious when the computing formulas are derived. We will denote the grand total by T and the number of data values by N , so the grand mean is just $M = T/N$. Similarly, we will denote the group totals by T_i , so the group means are $M_i = T_i/J_i$.

$$\begin{aligned} SST &= \sum \sum (x_{ij} - M)^2 \\ &= \sum \sum x_{ij}^2 - 2M \sum \sum x_{ij} + \sum \sum M^2 \end{aligned}$$

$$\begin{aligned}
SSG &= \sum \sum x_{ij}^2 - 2MT + M^2 N \\
&= \sum \sum x_{ij}^2 - T^2/N \\
&= \sum \sum (M_i - M)^2 \\
&= \sum J_i (M_i - M)^2 \\
&= \sum J_i M_i^2 - 2J_i M_i M + M^2 \\
&= \sum J_i M_i^2 - 2M \sum J_i M_i + M^2 \sum J_i \\
&= \sum J_i (T_i/J_i)^2 - 2MT + M^2 N \\
&= \sum T_i^2/J_i - T^2/N \\
SSE &= \sum \sum (x_{ij} - M_i)^2 \\
&= \sum \sum x_{ij}^2 - 2 \sum \sum x_{ij} M_i + \sum \sum M_i^2 \\
&= \sum \sum x_{ij}^2 - 2 \sum M_i T_i + \sum J_i M_i^2 \\
&= \sum \sum x_{ij}^2 - 2 \sum T_i^2/J_i + \sum J_i (T_i/J_i)^2 \\
&= \sum \sum x_{ij}^2 - \sum T_i^2/J_i
\end{aligned}$$

$$\begin{array}{ccc}
\begin{array}{l}
SST \\
\sum \sum (x_{ij} - M)^2 \\
\sum \sum x_{ij}^2 - 2M \sum \sum x_{ij} + \sum \sum M^2 \\
\sum \sum x_{ij}^2 - 2MT + M^2 N \\
\sum \sum x_{ij}^2 - T^2/N
\end{array} &
\begin{array}{l}
SSG \\
\sum \sum (M_i - M)^2 \\
\sum J_i (M_i - M)^2 \\
\sum J_i M_i^2 - 2J_i M_i M + M^2 \\
\sum J_i M_i^2 - 2M \sum J_i M_i + M^2 \sum J_i \\
\sum J_i (T_i/J_i)^2 - 2MT + M^2 N \\
\sum T_i^2/J_i - T^2/N
\end{array} &
\begin{array}{l}
SSE \\
\sum \sum (x_{ij} - M_i)^2 \\
\sum \sum x_{ij}^2 - 2 \sum \sum x_{ij} M_i + \sum \sum M_i^2 \\
\sum \sum x_{ij}^2 - 2 \sum M_i T_i + \sum J_i M_i^2 \\
\sum \sum x_{ij}^2 - 2 \sum T_i^2/J_i + \sum J_i (T_i/J_i)^2 \\
\sum \sum x_{ij}^2 - \sum T_i^2/J_i
\end{array}
\end{array}$$

We see a lot of recurring expressions here, so we can see the relationships of these expressions better by making the following substitutions:

$$\begin{aligned}
A &\Rightarrow \sum \sum x_{ij}^2 \\
B &\Rightarrow T^2/N \\
C &\Rightarrow \sum T_i^2/J_i
\end{aligned}$$

With those substitutions, we get the following formulas for the sums of squares:

$$\begin{aligned}
SST &= A - B \\
SSG &= C - B \\
SSE &= A - C
\end{aligned}$$

It's now obvious why the ANOVA identity is true.

What's the point of all this? First, the code for calculating the sums of squares will calculate A, B, and C, and that code wouldn't make any sense to a reader without this explanation. Second, we can calculate any one of these sums of squares by the two-pass method and compare it to the one-pass method as a test for numerical accuracy. The largest

sum is SST and so it is the one most liable to produce floating overflow, especially with these calculating formulas, since A will be bigger than SST. SST also computes with the largest differences (subtracting data from the grand mean rather than a group mean), and so is the least likely to suffer from numerical inaccuracy by the two-pass method. Therefore, we will calculate SST by the two-pass method as well as by $A - B$. Rather than enforce some arbitrary precision, we will simply report both values for SST to the users, who can judge the numerical accuracy for themselves.

A.3 Linear Regression

Linear regression fits a linear model to the data using a least squares fit. If the linear model has more than one independent (x) variable, we will call it multiple linear regression. Multiple linear regression is implemented by matrix methods, which are covered in the next section. This section is devoted to the special case of linear regression in which we fit a line to bivariate data. Thus, we only have two variables, x and y . We will give the computational formulas we implemented. The main references we used are Devore [4], Freund [5], and Bohrnstedt and Knoke [1].

The linear regression statistics can be calculated from just six summary statistics of the original data. Since authors use different notation for the regression statistics and the summary statistics, we have also decided to invent our own notation. This will be useful later, when we attempt to resolve discrepancies between the formulas given by different authors. The six summary statistics are the following. The second column is the symbol used in the Common Lisp code, the third column is the mathematical notation we are inventing. These are all easy summary statistics, computable in one pass over the data, which makes them useful for bootstrapping.

	CL symbol	math symbol		definitions
1	n	n		the number of data points, (x_i, y_i)
2	x	x^+	$\sum_{i=1}^n x_i$	sum of the x's
3	y	y^+	$\sum_{i=1}^n y_i$	sum of the y's
4	x2	$(x^2)^+$	$\sum_{i=1}^n (x_i)^2$	sum after squaring
5	y2	$(y^2)^+$	$\sum_{i=1}^n (y_i)^2$	sum after squaring
6	xy	$(xy)^+$	$\sum_{i=1}^n x_i y_i$	AKA dot product of \vec{x} and \vec{y}

CLASP implements linear regression functions that return three amounts of information: minimal, brief, and verbose. The functions take two different kinds of argument lists: two sequences of data, or the six summary statistics. Thus, there are six linear regression functions, depending on how you want to supply your arguments and how much computation you want done. All of the ones that take data sequences just calculate the summary statistics and pass them on to the ones that take summary statistics, so we will only discuss the latter in this section.

From these summary statistics, we will calculate the following sums of squares. The first formula on each line is the intuitive definition in terms of variation or covariation around the means and the second formula is an easy computing formula in terms of summary statistics.

$$\begin{aligned}
 \text{SSX} &= \sum (x_i - \bar{x})^2 &= (x^2)^+ - (x^+)^2/n \\
 \text{SSY} &= \sum (y_i - \bar{y})^2 &= (y^2)^+ - (y^+)^2/n \\
 \text{SSXY} &= \sum (x_i - \bar{x})(y_i - \bar{y}) &= (xy)^+ - x^+ y^+ / n
 \end{aligned}$$

Notice that all of these formulas involve a division by n . Since division is an expensive arithmetic operation, second only to square root, we try to avoid it wherever possible. Consequently, we found it efficient to compute the following instead.

$$\begin{aligned} \text{NSSY} &= n(x^2)^+ - (x^+)^2 \\ \text{NSSY} &= n(y^2)^+ - (y^+)^2 \\ \text{NSSXY} &= n(xy)^+ - x^+y^+ \end{aligned}$$

Now we are ready to give computational formulas for the slope and intercept of the regression line. Our formula for slope matches Devore's, and is algebraically equivalent to the other authors' formulas. Indeed, it was Devore's use of multiplying through by n that inspired us to do so everywhere.

$$\begin{aligned} \text{slope} &: \frac{\text{NSSXY}}{\text{NSSX}} \\ \text{intercept} &: \frac{y^+ - (\text{slope})x^+}{n} \end{aligned}$$

The slope and intercept are the only statistics calculated by the function 'linear-regression-minimal.' That function doesn't need the 'y2' variable, $(y^2)^+$, but takes it anyhow, just for consistency with the other linear regression functions.

The next statistics all depend on the sum of squares of the residuals—the remaining error not accounted for by the regression model. Most authors denote this quantity with SSE. All define it the same way, based on the difference between y_i and the y value of the line at the corresponding x_i , denoted \hat{y}_i .

$$\text{SSE} = \sum (y_i - \hat{y}_i)^2 = \sum (y_i - [(\text{slope})x_i + \text{intercept}])^2$$

Just as with the analysis of variance, there is an identity relating the total sum of squares, usually denoted SST, to the regression sum of squares, denoted SSR, and the error sum of squares, SSE. That identity is

$$\text{SST} = \text{SSR} + \text{SSE}$$

By definition, $\text{SST} = \text{SSY}$, so we will use SSY to avoid introducing any new symbols. Because of this identity, and the fact that we already have SSY , we should compute SSR or SSE , whichever is easier, and then obtain the other by subtraction.

The authors all give different computing formulas for SSE. Translating into our notation, their computing formulas are:

$$\begin{aligned} \text{BK} & \text{SSY}(1 - r^2) \\ \text{Devore} & (y^2)^+ - (\text{intercept})y^+ - (\text{slope})(xy)^+ \\ \text{Freund} & \text{SSY} - (\text{slope})\text{SSXY} \end{aligned}$$

It's hard to believe that these are all algebraically equivalent! Still, because of the subtractions in Freund' and BK's formulas, it's easy to see that $\text{SSR} = (\text{slope})\text{SSXY}$. Therefore, our code will compute the following:

$$\begin{aligned} \text{NSSR} &= (\text{slope})\text{NSSXY} \\ \text{NSSE} &= \text{NSSY} - \text{NSSR} \end{aligned}$$

At last, we are ready to compute the coefficient of determination, denoted r^2 . We are also ready to compute the standard error of the slope, which is used in t-tests on the slope

and also in building confidence intervals on the slope. Fortunately, all the authors agree on how to compute these statistics.

$$r^2 = \frac{\text{NSSR}}{\text{NSSY}}$$

$$\text{std. err. slope} = \sqrt{\frac{\text{NSSE}}{(n-2)\text{NSSX}}}$$

Note that the degrees of freedom for the t-distribution of the slope is $n - 2$. Also note how we have taken advantage of multiplying through by n in all our sums of squares.

The next step is to perform a t-test on the slope parameter to test the hypothesis that the slope is non-zero. This t-test is equivalent to testing whether the correlation coefficient, r , is non-zero and equivalent to testing whether the F statistic (mean square regression divided by mean square error) is greater than one. The t-statistic is:

$$t = \frac{\text{slope}}{\text{std. err. slope}}$$

and is tested with $n - 2$ degrees of freedom. This computes a p-value which is also returned.

In summary, the “brief” linear regression functions return the following information: slope, intercept, coefficient of determination, standard error of the slope, and the p-value.

The “verbose” linear regression functions return even more statistics, but most are redundant with others. They return the correlation coefficient, r , but the user can get that by simply taking the square root of the coefficient of determination, r^2 . It also returns an ANOVA table, computing the F statistic and p-value rather than the t-statistic and p-value that the “brief” functions do. The ANOVA table looks like this:

Source of Variation	degrees of freedom	Sum of squares	Mean Square	F	p-value
Regression	1	SSR	MSR = SSR/1	MSR/MSE	p
Error	$n - 2$	SSE	MSE = SSE/($n - 2$)		
Total	$n - 1$	SSY			

In addition to all this, the “verbose” functions return the standard error of the intercept. Unfortunately, there are several conflicting formulas for the standard error of the intercept, which we have not yet resolved as of this release. For now, we return nil, deeming less information better than bad information.

A.4 Regression Analysis: Matrix Method*

To implement ‘linear regression’ in CLASP the decision was made to use matrix algebra methods for the calculation of the statistics. The basis of the method is in the ability to calculate the sum of squares and cross products for the independent and dependent variables through matrix multiplication, then using these values to calculate the regression coefficients, correlation coefficients, and percentage variance of the dependent variable, Y . The advantage of using matrix algebra over direct raw score calculation is not obvious when considering an example with only two independent variables, however the complexity of the

*This discussion is a condensation of chapters 2 and 3 from Pedhazur’s book, *Multiple Regression in Behavioral Research* [9].

direct calculations becomes apparent as the number of independent variables increase. The following formulae illustrate the calculations involved using the direct raw score method, for two independent, X , variables, with $\sum x_i$ the sum of deviation scores for the i 'th X variable, and $\sum x_i y$ the sum of the cross products of the deviation scores of Y and the i 'th X variable.

$$b_1 = \frac{(\sum x_2^2)(\sum x_1 y) - (\sum x_1 x_2)(\sum x_2 y)}{(\sum x_1^2)(\sum x_2^2) - (\sum x_1 x_2)^2}$$

$$b_2 = \frac{(\sum x_1^2)(\sum x_2 y) - (\sum x_1 x_2)(\sum x_1 y)}{(\sum x_1^2)(\sum x_2^2) - (\sum x_1 x_2)^2}$$

As can be seen these equations expand with the introduction of each new independent variable, making them very unwieldy. In order to perform the calculations with matrices we first introduce a new variable, X_0 , a unit vector (a vector of ones), of the same size as the vectors X_i , with this the matrix \mathbf{Z} is built, with the columns being defined as $X_0 \cdots X_i$ and Y . Using the equation

$$Y = a + b_1 X_1 + b_2 X_2 + \cdots + b_k X_k + e$$

in matrix form, where \mathbf{Y} is a $N \times 1$ column vector, \mathbf{X} is a $N \times k$ matrix with each column corresponding to X_i , where X_0 is the unit vector described above, a is the first element of the $1 \times k$ row vector of regression coefficients, \mathbf{b} , and e is the $1 \times k$ column vector of the error terms of the equation, given by the form $Y = \hat{Y} + e$, where \hat{Y} is the predicted value of Y from the regression equation, produces the following equation:

$$\begin{bmatrix} Y \\ Y_1 \\ Y_2 \\ Y_3 \\ \vdots \\ Y_N \end{bmatrix} = a + \begin{bmatrix} b \\ b_1 \\ b_2 \\ \vdots \\ b_N \end{bmatrix} \begin{bmatrix} X \\ X_{10} & X_{11} & X_{12} & \cdots & X_{1k} \\ X_{20} & X_{21} & X_{22} & \cdots & X_{2k} \\ X_{30} & X_{31} & X_{32} & \cdots & X_{3k} \\ \vdots & \vdots & \vdots & \vdots & \vdots \\ X_{N0} & X_{N1} & X_{N2} & \cdots & X_{Nk} \end{bmatrix} + \begin{bmatrix} e \\ e_1 \\ e_2 \\ e_3 \\ \vdots \\ e_N \end{bmatrix}$$

To calculate the values of \mathbf{b} and e we use the sum of squares and cross product matrix produced by multiplying the transpose of the matrix \mathbf{Z} , \mathbf{Z}' , by \mathbf{Z} , $\mathbf{Z}'\mathbf{Z}$. As can be seen the result of this multiplication contains the summed squares of the individual variables along the diagonal, with $\sum X_0^2 = \sum X_0$, and the cross products along the off diagonals.

$$\mathbf{Z}'\mathbf{Z} = \begin{bmatrix} \sum X_0 & \sum X_1 & \sum X_2 & \cdots & \sum X_k & \sum Y \\ \sum X_1 & \sum X_1^2 & \sum X_1 X_2 & \cdots & \sum X_1 X_k & \sum X_1 Y \\ \sum X_2 & \sum X_1 X_2 & \sum X_2^2 & \cdots & \sum X_2 X_k & \sum X_2 Y \\ \vdots & \vdots & \vdots & \vdots & \vdots & \vdots \\ \sum X_k & \sum X_k X_1 & \sum X_k X_2 & \cdots & \sum X_k^2 & \sum X_k Y \\ \sum Y & \sum Y X_1 & \sum Y X_2 & \cdots & \sum Y X_k & \sum Y^2 \end{bmatrix}$$

From this we get the two component matrices $\mathbf{X}'\mathbf{X}$, X transpose X , and $\mathbf{X}'\mathbf{Y}$, X transpose Y ,

$$\mathbf{X}'\mathbf{X} = \begin{bmatrix} \sum X_0 & \sum X_1 & \sum X_2 & \cdots & \sum X_k \\ \sum X_1 & \sum X_1^2 & \sum X_1 X_2 & \cdots & \sum X_1 X_k \\ \sum X_2 & \sum X_1 X_2 & \sum X_2^2 & \cdots & \sum X_2 X_k \\ \vdots & \vdots & \vdots & \vdots & \vdots \\ \sum X_k & \sum X_k X_1 & \sum X_k X_2 & \cdots & \sum X_k^2 \end{bmatrix} \quad \mathbf{X}'\mathbf{Y} = \begin{bmatrix} \sum Y \\ \sum X_1 Y \\ \sum X_2 Y \\ \vdots \\ \sum X_k Y \\ \sum Y^2 \end{bmatrix}$$

which are then used to determine the matrix \mathbf{b} , using the following equation:

$$\mathbf{b} = (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{Y}$$

where $(\mathbf{X}'\mathbf{X})^{-1}$ is the inverse of the matrix $\mathbf{X}'\mathbf{X}$. Recall that the inverse of a matrix is a matrix such that $\mathbf{X}^{-1}\mathbf{X} = \mathbf{I}$, where \mathbf{I} is the identity matrix. So multiplying by the inverse of a matrix is equivalent to division. The significance of the coefficients can be tested using a standard t test, using the formula $t = \frac{b}{s_b}$, where s_b is the error term of the coefficient. The value of t indicates whether or not the coefficient is significantly different from zero, with a critical value at the 0.05 level of 1.96 for a two tailed test. A two tailed test is used as it does not matter whether the coefficient is positive or negative. With \mathbf{b} the sum of squares of the regression and the sum of squares of the residual are calculated, with:

$$ss_{reg} = \mathbf{b}'\mathbf{X}'\mathbf{Y} - \frac{(\sum Y)^2}{N}$$

$$ss_{res} = \mathbf{e}'\mathbf{e} = \sum e^2 = \mathbf{Y}'\mathbf{Y} - \mathbf{b}'\mathbf{X}'\mathbf{Y}$$

where $\mathbf{e}'\mathbf{e}$ is the sum of squared error terms for the regression. The sum of squares of the regression indicate what portion of the total squared deviation scores, $\sum y^2$, of the dependent variable, are due to the effects of the independent variables. The sum of squares of the residual indicate the proportion which is unaccounted for by the independent variables. Recall that

$$ss_{total} = ss_{reg} + ss_{res} = \sum y^2 = \mathbf{Y}'\mathbf{Y} - \frac{(\sum Y)^2}{N}$$

using these values the significance of the regression can be tested.

The key value of interest from the regression is the F statistic of the R^2 value, which is the percentage of the variance accounted for by the regression. The value of R^2 can be calculated using either of the following two equations:

$$R^2 = \frac{\mathbf{b}'\mathbf{X}'\mathbf{Y} - \frac{(\sum Y)^2}{N}}{\mathbf{Y}'\mathbf{Y} - \frac{(\sum Y)^2}{N}} = \frac{ss_{reg}}{\sum y^2}$$

$$R^2 = 1 - \frac{\mathbf{e}'\mathbf{e}}{\mathbf{Y}'\mathbf{Y} - \frac{(\sum Y)^2}{N}} = 1 - \frac{ss_{res}}{\sum y^2}$$

where $\sum y^2$ is the sum of squared deviations for the dependent variable Y . The value of F indicates the probability that such a value could occur due to chance. F is calculated as:

$$F = \frac{ss_{reg}/df_{reg}}{ss_{res}/df_{res}} = \frac{R^2/k}{(1 - R^2)/(N - k - 1)}$$

where k is the number of independent variables and N is the number of observations. Additionally the coefficients of correlation between each X_i and Y , as well as between the X_i 's can be calculated, with the squared coefficient for each of the X_i 's indicating its contribution to the variance of Y .

It is possible, as with other statistics, to produce results which are significant, but not meaningful, as well as meaningful, but not significant. Evaluating the result of a linear regression, especially when there is interaction between the independent variables, can be a very difficult task. Often a small value of R^2 can discourage a researcher, even though

its value is significant, and, by the same token, a large value of R^2 does not immediately confer meaning to the results of the regression. In cases where the results of the analysis are not clear, or appear to be at a border line, a useful strategy can be to reduce the number of independent variables being used, looking for a change in the analysis as the interactive effects are reduced. The covariance of the independent variables, which is indicated by the correlation between the X_i 's, is a good indicator of the interaction between the independent variables. In general it is best to refer to a statistics text when evaluating the results of a linear regression.

Appendix B

Clip Examples

B.1 More Experiment Definitions Using Clip

The following three examples, taken from different domains, illustrate experiment control and data collection using CLIP. The third example includes output files (in CLASP format) for time-series clips.

B.1.1 Measuring an Airport Gate Maintenance Scheduler's Performance

This extended example of CLIP code was provided by Zachary Rubinstein. The system being studied is an AI scheduler developed by David Hildum for a domain called Airport Resource Management (ARM). The task is to schedule gate maintenance, providing and coordinating vehicles for servicing planes at terminal gates. The experiments shown here were part of a baseline study designed to measure the performance characteristics of several priority rating schemes under resource-constrained scenarios. The dependent variables are delay and fragmentation (how much shuttling between gates occurs). The independent variables are the competing heuristic rating schemes for determining the order in which resource requests are scheduled. For clarity, many of the implementation details have been omitted.

The example uses several super clips to organize the collection of similar kinds of data, such as all measures of delay. All data collection occurs *post hoc* and is written to a single summary file. Note that two experiments are defined, each of which uses the same clips. The function `delays-information`, called by the clip `delay`, provides an interesting example of the interface between clips and the system being instrumentated.

The `delay` and `fragmentation` clips are examples of *composite* clips. The values returned by the body of a composite clip are one-to-one mapped onto the components of the clip. Each component is itself a clip which accepts an argument and has a body that returns a value. This value is what is output to the data file. The number of values a composite clip produces is equal to the number of its components.

The components of a composite clip have a clip definition automatically generated for them. The default action for this automatically generated component clip is to return its single argument as a single value. For example, the `delay` clip below generates a clip for each of its components automatically. The clip designer can also override this default behavior by writing a clip for that component that does something with the argument before returning it.

```

(define-simulator arm :start-system arm:run-arm)

;;; -----
;;; CLIP Definitions

(defclip airport ()
  "Return the Currently Loaded Airport"
  ()
  (arm::get-loaded-airport))

(defclip timetable ()
  "Return the Currently Loaded timetable"
  ()
  (arm::get-loaded-timetable))

(defclip delay ()
  "Composite clip to calculate the various delay clips."
  (:components (total-delay number-of-delays average-delay
                std-dev-delay maximum-delay duration-ratio))

  (delays-information))

(defclip fragmentation ()
  "Composite fragmentation clip to calculate the various fragmentation clips."
  (:components (total-servicing-time
                total-sig-frag-time
                total-setup-time
                total-travel-time
                ratio-of-frag-to-servicing
                ratio-of-setup-to-servicing
                ratio-of-travel-to-servicing))

  (resource-fragmentation-information 'arm:baggage-truck-planner-unit))

;;; -----
;;; Experiments

(define-experiment limit-resources-numerically-over-rating-schemes
  (&key (rating-schemes *defined-rating-schemes*)
        (airport-delimiter 10)
        (airports user::*10-FTS-AIRPORTS*)
        (timetable-name :20-flights))

  "Using the various rating schemes, constrain the resources
and report the delays."
  :locals ((report-filename nil))

```

```

:variables ((airport-name in airports)
            (rating-scheme in rating-schemes))
:instrumentation (timetable
                 total-delay number-of-delays average-delay
                 std-dev-delay maximum-delay duration-ratio
                 total-servicing-time total-sig-frag-time
                 total-setup-time total-travel-time
                 ratio-of-frag-to-servicing
                 ratio-of-setup-to-servicing
                 ratio-of-travel-to-servicing)
:before-experiment (progn (arm:zack-set-default-demo-parameters)
                        (arm:execute-load-timetable
                         :LOAD-SPECIFIC-FILE timetable-name))
:before-trial (setf report-filename
                (delay-initialize-trial rating-scheme
                 airport-delimiter airport-name))
:after-trial (delay-report report-filename
              airport-delimiter rating-scheme))

;;; -----
;;; Example Invocation

#+COMMENT
(run-experiment 'dss:limit-resources-numerically-over-rating-schemes
               :args '(:rating-schemes
                      (:zack-1 :default :obo-minest-heuristic)
                      :airport-delimiter :20-ref
                      :airports (:detroit-17-bts :detroit-16-bts
                                :detroit-15-bts :detroit-14-bts
                                :detroit-13-bts :detroit-12-bts)
                      :timetable-name :20-flights)
               :output-file "hillary:reports;20-REF-FTS-DELAY-STATS.TEXT")

;;; -----
;;; Experiments

(define-experiment run-over-rating-schemes
  (&key (rating-schemes *defined-rating-schemes*)
        (airport-delimiter 10))
  "Using the various rating schemes, constrain the resources, introduce
orders, and report the delays."
  :locals ((report-stream nil))
  :variables ((rating-scheme in rating-schemes))
  :instrumentation (airport timetable
                   total-delay number-of-delays
                   average-delay std-dev-delay

```

```

maximum-delay duration-ratio
total-servicing-time total-sig-frag-time
total-setup-time total-travel-time
ratio-of-frag-to-servicing
ratio-of-setup-to-servicing
ratio-of-travel-to-servicing)
:before-trial (setf report-stream
               (delay-initialize-trial rating-scheme
               airport-delimiter))
:after-trial (write-current-experiment-data ))

;;; -----
;;; Example Invocation

#+COMMENT
(run-experiment 'dss:run-over-rating-schemes
               :ARGS '(:RATING-SCHEMES
                       ,(cons :DEFAULT
                               (cons :ZACK-1
                                      (remove :OPPORTUNISTIC
                                              *defined-rating-schemes*
                                              :TEST #'eq)))
                       :airport-delimiter 10)
               :OUTPUT-FILE "hillary:reports;HILDUM-EXP-1.TEXT")

;;; -----
;;; Delay Information

(defun delays-information ()
  (LOOP WITH local-delay-time = 0
        AND network = nil
        AND release-time = nil
        AND due-date = nil
        AND start-time = nil
        AND finish-time = nil
        AND avg-delay = 0
        AND std-dev-delay = 0
        AND max-delay = 0
        FOR order IN (find-units 'order
                                (make-paths '(scheduler * order)) :ALL)
        DO
          (setf network (order$network order))
          (setf release-time (order$release-time order))
          (setf due-date (order$due-date-time order))
          (setf start-time (task-network$start-time network))

```



```

(setf finish-time (task-network$finish-time network))
(setf local-delay-time (- finish-time due-date))
SUM (- finish-time start-time) INTO total-actual-duration
SUM (- due-date release-time) INTO total-desired-duration
WHEN (not (= local-delay-time 0))
COUNT local-delay-time INTO number-of-delays
AND
SUM local-delay-time INTO total-delay-time
AND
COLLECT local-delay-time INTO delays
FINALLY
(unless (zerop number-of-delays)
  (setf avg-delay (float (/ total-delay-time number-of-delays)))
  (setf std-dev-delay
    (LOOP FOR delay IN delays
      SUM (expt (- delay avg-delay) 2)
      INTO sum-of-squares
      FINALLY (return
        (sqrt
          (float
            (/ sum-of-squares
              (1- number-of-delays)))))))
  (setf max-delay (apply #'max delays)))
(return (values total-delay-time
  number-of-delays
  avg-delay
  std-dev-delay
  max-delay
  (float (/ total-actual-duration
    total-desired-duration))))))

```

B.1.2 Phoenix Real-Time-Knob Experiment

This section describes an experiment done in the PHOENIX fire-fighting simulation system. We give some background on the Phoenix system and on the experiment design to help the reader better interpret the experiment definition that follows. For more detail on the experiment and its results see Hart & Cohen, 1992 [6].

The Phoenix System

PHOENIX is a multi-agent planning system that fights simulated forest-fires. The simulation uses terrain, elevation, and feature data from Yellowstone National Park and a model of fire spread from the National Wildlife Coordinating Group Fireline Handbook (National Wildlife Coordinating Group 1985). The spread of fires is influenced by wind and moisture conditions, changes in elevation and ground cover, and is impeded by natural and man-made boundaries such as rivers, roads, and fireline. The Fireline Handbook also prescribes many of the characteristics of our firefighting agents, such as rates of movement and effectiveness of various firefighting techniques. For example, the rate at which bulldozers dig fireline varies with the terrain. PHOENIX is a real-time simulation environment PHOENIX agents must think and act as the fire spreads. Thus, if it takes too long to decide on a course of action, or if the environment changes while a decision is being made, a plan is likely to fail.

One PHOENIX agent, the Fireboss, coordinates the firefighting activities of all field agents, such as bulldozers and watchtowers. The Fireboss is essentially a thinking agent, using reports from field agents to form and maintain a global assessment of the world. Based on these reports (e.g., fire sightings, position updates, task progress), it selects and instantiates fire-fighting plans and directs field agents in the execution of plan subtasks.

A new fire is typically spotted by a watchtower, which reports observed fire size and location to the Fireboss. With this information, the Fireboss selects an appropriate fire-fighting plan from its plan library. Typically these plans dispatch bulldozer agents to the fire to dig fireline. An important first step in each of the three plans in the experiment described below is to decide where fireline should be dug. The Fireboss projects the spread of the fire based on prevailing weather conditions, then considers the number of available bulldozers and the proximity of natural boundaries. It projects a bounding polygon of fireline to be dug and assigns segments to bulldozers based on a periodically updated assessment of which segments will be reached by the spreading fire soonest. Because there are usually many more segments than bulldozers, each bulldozer digs multiple segments. The Fireboss assigns segments to bulldozers one at a time, then waits for each bulldozer to report that it has completed its segment before assigning another. This ensures that segment assignment incorporates the most up-to-date information about overall progress and changes in the prevailing conditions.

Once a plan is set into motion, any number of problems might arise that require the Fireboss's intervention. The types of problems and mechanisms for handling them are described in Howe & Cohen, 1990 [7], but one is of particular interest here: As bulldozers build fireline, the Fireboss compares their progress to expected progress. If their actual progress falls too far below expectations, a plan failure occurs, and (under the experiment scenario described here) a new plan is generated. The new plan uses the same bulldozers to fight the fire and exploits any fireline that has already been dug. We call this error recovery method *vreplanning*. PHOENIX is built to be an adaptable planning system that can recover from plan failures. Although it has many failure-recovery methods, replanning is the focus

of the experiment described in the next section.

Identifying the Factors that Affect Performance

We designed an experiment with two purposes. A confirmatory purpose was to test predictions that the planner's performance is sensitive to some environmental conditions but not others. In particular, we expected performance to degrade when we change a fundamental relationship between the planner and its environment—the amount of time the planner is allowed to think relative to the rate at which the environment changes—and not be sensitive to common dynamics in the environment such as weather, and particularly, wind speed. We tested two specific predictions: 1) that performance would not degrade or would degrade gracefully as wind speed increased; and 2) that the planner would not be robust to changes in the Fireboss's thinking speed due to a bottleneck problem described below. An exploratory purpose of the experiment was to identify the factors in the Fireboss architecture and PHOENIX environment that most affected the planner's behavior, leading to a causal model.

The Fireboss must select plans, instantiate them, dispatch agents and monitor their progress, and respond to plan failures as the fire burns. The rate at which the Fireboss thinks is determined by a parameter called the Real Time Knob. By adjusting the Real Time Knob we allow more or less simulation time to elapse per unit CPU time, effectively adjusting the speed at which the Fireboss thinks relative to the rate at which the environment changes.

The Fireboss services bulldozer requests for assignments, providing each bulldozer with a task directive for each new fireline segment it builds. The Fireboss can become a bottleneck when the arrival rate of bulldozer task requests is high or when its thinking speed is slowed by adjusting the Real Time Knob. This bottleneck sometimes causes the overall digging rate to fall below that required to complete the fireline polygon before the fire reaches it, which causes replanning. In the worst case, a Fireboss bottleneck can cause a thrashing effect in which plan failures occur repeatedly because the Fireboss can't assign bulldozers during replanning fast enough to keep the overall digging rate at effective levels. We designed our experiment to explore the effects of this bottleneck on system performance and to confirm our prediction that performance would vary in proportion to the manipulation of thinking speed. Because the current design of the Fireboss is not sensitive to changes in thinking speed, we expect it to take longer to fight fires and to fail more often to contain them as thinking speed slows.

In contrast, we expect PHOENIX to be able to fight fires at different wind speeds. It might take longer and sacrifice more area burned at high wind speeds, but we expect this effect to be proportional as wind speed increases and we expect PHOENIX to succeed equally often at a range of wind speeds, since it was designed to do so.

Experiment Design

We created a straightforward fire fighting scenario that controlled for many of the variables known to affect the planner's performance. In each trial, one fire of a known initial size was set at the same location (an area with no natural boundaries) at the same time (relative to the start of the simulation). Four bulldozers were used to fight it. The wind's speed and direction were set initially and not varied during the trial. Thus, in each trial, the Fireboss receives the same fire report, chooses a fire-fighting plan, and dispatches the bulldozers to implement it. A trial ends when the bulldozers have successfully surrounded the fire or after

120 hours without success. The experiment's first dependent variable then is Success, which is true if the fire is contained, and false otherwise. A second dependent variable is shutdown time (SD), the time at which the trial was stopped. For successful trials, shutdown time tells us how long it took to contain the fire. Two independent variables were wind speed (WS) and the setting of the Fireboss's Real Time Knob (RTK). A third variable, the first plan chosen by the Fireboss in a trial (FPLAN), varied randomly between trials. It was not expected to influence performance, but because it did, we treat it here as an independent variable.

WS The settings of WS in the experiment were 3, 6, and 9 kilometers per hour. As wind speed increases, fire spreads more quickly in all directions, and most quickly downwind. The Fireboss compensates for higher values of wind speed by directing bulldozers to build fireline further from the fire.

RTK The default setting of RTK for PHOENIX agents allows them to execute 1 CPU second of Lisp code for every 5 minutes that elapses in the simulation. We varied the Fireboss's RTK setting in different trials (leaving the settings for all other agents at the default). We started at a ratio of 1 simulation-minute/cpu-second, a thinking speed 5 times as fast as the default, and varied the setting over values of 1, 3, 5, 7, 9, 11, and 15 simulation-minutes/cpu-second. These values range from 5 times the normal speed at a setting of 1 down to one-third the normal speed at 15. The values of RTK reported here are rescaled. The normal thinking speed (5) has been set to RTK=1, and the other settings are relative to normal. The scaled values (in order of increasing thinking speed) are .33, .45, .56, .71, 1, 1.67, and 5. RTK was set at the start of each trial and held constant throughout.

FPLAN The Fireboss randomly selects one of three plans as its first plan in each trial. The plans differ mainly in the way they project fire spread and decide where to dig fireline. SHELL is aggressive, assuming an optimistic combination of low fire spread and fast progress on the part of bulldozers. MODEL is conservative in its expectations, assuming a high rate of spread and a lower rate of progress. The third, MBIA, generally makes an assessment intermediate with respect to the others. When replanning is necessary, the Fireboss again chooses randomly from among the same three plans. We adopted a basic factorial design, systematically varying the values of WS and RTK. Because we had not anticipated a significant effect of FPLAN, we allowed it to vary randomly.

Rtk Experiment Clips

All data collection in this experiment occurs at the end of each trial. Many of these clips invoke system specific calls to the simulator and to agents' state memories, again, as in the first example, illustrating the interface between clips and the system under study. Of particular interest here is the script invoked to recreate the same experimental scenario for each trial (see the `:script` keyword to `define-experiment` below). This script has a set of instructions to the Phoenix simulator's event-scheduler that introduce a fixed set of environmental changes that are part of the experimental control.

```
(define-experiment real-time-knob-experiment (use-exp-style)
  "Simple experiments for exercising and testing the real time knob."
```

```

:before-experiment (real-time-knob-experiment-init-before-experiment
                    use-exp-style)
:before-trial (real-time-knob-experiment-init-before-trial)
:after-trial (real-time-knob-experiment-after-trial)
:after-experiment (real-time-knob-experiment-reset-after-experiment)
:variables ((real-time-knob in '(1 3 5 7 9 11))
           (wind-speed in '(3 9 12)))
:instrumentation (number-of-bulldozers
                 plan-to-contain-fire          ; FPLAN
                 all-plans-to-contain-fire
                 fires-started
                 shutdown-time                ; SD
                 number-of-fires-contained
                 total-fire-line-built
                 r-factor
                 area-burned
                 agents-lost
                 all-agent-instrumentation
                 fireboss-instrumentation
                 bulldozer-instrumentation)

:script
((setup-starting-conditions "12:29"
  (progn
    (send (fire-system)
          :alter-environment-parameter 'wind-direction 315)
    (send (fire-system)
          :alter-environment-parameter 'wind-speed (wind-speed))))
 (start-fire "12:30"
  (send (fire-system) :start-fire 700 (point 50000 40000))))

;;; -----
;;; Utility functions used by the :before- and :after- forms to
;;; initialize and reset the experiment environment

(defun real-time-knob-experiment-init-before-experiment (use-exp-style)
  (when use-exp-style
    (setf (interaction-style t) 'experiment))
  ;; Don't allow fires to skip over fire lines.
  (send (fire-simulation) :set-spotting-scale-factor 0)
  ;; Get rid of flank attack, etc.
  (modify-knowledge-base))

(defun real-time-knob-experiment-init-before-trial ()
  (gc-immediately :silent t))

(defun real-time-knob-experiment-after-trial ()
  ;; This is now done in the 'shutdown-trial' method.

```

```

(write-current-experiment-data))

(defun real-time-knob-experiment-reset-after-experiment ()
  (setf (interaction-style t) 'normal))

;;; -----

(defmacro pct (part wh)
  '(if (zerop ,wh) 0 (* 100.0 (/ ,part ,wh))))

(defmacro /-safe (dividend divisor)
  '(if (zerop ,divisor) 0 (/ ,dividend ,divisor)))

;;; -----
;;; Instrumentation definitions...

(defclip number-of-bulldozers ()
  (:report-key "Number of bulldozers")
  (length (find-agent-by-name 'bulldozer :multiple-allowed t)))

(defclip area-burned ()
  (:report-key "Area burned (sq. km)")
  (send (real-world-firemap) :fire-state))

(defclip shutdown-time ()
  (:report-key "Shutdown time (hours)")
  (current-time))

;;; -----
;;; Fire clips

(defclip fires-started ()
  (:report-key "Fires started")
  (let ((cnt 0))
    (map-over-fires (fire) (:delete-fire-frames nil)
      (incf cnt))
    (values cnt)))

(defclip fires-contained ()
  (:report-key "Fires contained")
  (mapcan #'(lambda (fire)
    (when (eq (f:get-value* fire 'status) 'under-control)
      (list fire)))
    (f:get-values* 'actual-fire 'instance+inv)))

(defclip number-of-fires-contained ()
  (:report-key "Fires extinguished")
  (length (fires-contained)))

```

```

;;; -----
;;; Fireline clips

(defun expand-extent-in-all-directions-by (extent distance-in-meters)
  (let ((temp-point (point distance-in-meters distance-in-meters)))
    (extent (point-clip (point-difference (extent-upper-left extent)
                                          temp-point))
            (point-clip (point-sum (extent-lower-right extent) temp-point)))))

(defun length-of-built-line-in-extent (extent)
  (let ((line-length 0))
    (dofiremap (point :upper-left (extent-upper-left extent)
                     :lower-right (extent-lower-right extent))
              (do-feature-edges (edge point (real-world-firemap) :edge-type :dynamic)
                              (incf line-length (feature-edge-length edge))))
    (values line-length)))

(defclip r-factor ()
  (:report-key "R factor")
  (/safe (total-fire-line-built) (total-perimeter)))

(defclip total-perimeter ()
  ()
  (reduce #' + (fires-contained)
         :key #' (lambda (fire)
                   (fast-polyline-length
                    (fire-perimeter-polyline (fire-origin fire)
                                             *fire-perimeter-resolution*)
                    t))))

(defclip total-fire-line-built ()
  (:report-key "Fireline Built (meters)")
  (reduce #' + (fires-contained)
         :key #' (lambda (fire)
                   (length-of-built-line-in-extent
                    (expand-extent-in-all-directions-by
                     (accurate-fire-extent (fire-center-of-mass fire)
                                           (point-on-polyline-furthest-from
                                            (fire-center-of-mass fire)
                                            (fire-boundary fire)
                                            nil))
                     *fire-sector-extension*)))))

;;; -----
;;; Agent clips

(defclip agents-lost ()

```

```

()
(mapcan #'(lambda (agent)
           (when (eq (f:get-value (send agent :self-frame) 'status) :dead)
               (list (name-of agent))))
        (all-agents)))

;;; "all-agent-instrumentation" is a mapping clip that collects information
;;; from its components, all of the Phoenix agents defined in this
;;; scenario. Each of the :component clips is run for each agent found
;;; by the :mapping function. Values returned by the :component clips
;;; are written out to the data file in sequence.

(defclip all-agent-instrumentation ()
  "Records the utilization of all the agents."
  (:components (agent-overall-utilization
                agent-cognitive-utilization
                agent-message-handling-time-pct
                agent-action-selection-time-pct
                agent-error-recovery-cost
                agent-error-recovery-percentage-of-cognitive-time
                number-of-frames-on-timeline)
  :map-function (cons
                 (find-agent-by-name 'fireboss)
                 (find-agent-by-name 'bulldozer :multiple-allowed t))))

(defclip agent-overall-utilization (agent)
  (:report-key "~a overall utilization")
  (pct (task-cumulative-cpu-time agent) (current-time)))

(defclip agent-cognitive-utilization (agent)
  (:report-key "~a cognitive utilization")
  (pct (phoenix-agent-cumulative-action-execution-time agent)
       (current-time)))

(defclip agent-message-handling-time-pct (agent)
  (:report-key "~a message handling pct")
  (pct (phoenix-agent-cumulative-message-handling-time agent)
       (current-time)))

(defclip agent-action-selection-time-pct (agent)
  (:report-key "~a action selection pct")
  (pct (phoenix-agent-cumulative-next-action-selection-time agent)
       (current-time)))

(defclip agent-error-recovery-cost (agent)
  (:report-key "~a error recovery cost")

  (f:using-frame-system ((name-of agent))

```



```

;;; -----
;;; Bulldozer clips

(defclip bulldozer-instrumentation ()
  "Mapping clip mapping reflexes-executed :component over all bulldozers."
  (:components (reflexes-executed)
    :map-function (find-agent-by-name 'bulldozer :multiple-allowed t)))

(defclip reflexes-executed (agent)
  (:report-key "~a reflexes executed")
  (reduce #'+ (standard-agent-model-reflexes agent)
    :key #'reflex-execution-count))

(defclip count-of-deadly-object-in-path-messages ()
  (:enable-function (trace-message-patterns
    '(:message-type :msg-reflex :type :error
      :reason :deadly-object-in-path))
  :disable-function (untrace-message-patterns
    '(:message-type :msg-reflex :type :error
      :reason :deadly-object-in-path)))
  (message-pattern-count '(:message-type :msg-reflex :type :error
    :reason :deadly-object-in-path)))

;;; -----
;;; Plan clips record which plan(s) was used during a trial to fight the fire.

(defun the-first-fire-started ()
  (first (last (f:get-values* 'actual-fire 'instance+inv))))

(defclip plan-to-contain-fire ()
  (:report-key "Plan to Contain Fire")

  (get-primitive-action
    (tl-entry-of-plan-to-contain-fire)))

(defun tl-entry-of-plan-to-contain-fire ()
  (f:using-frame-system ((name-of (find-agent-by-name 'fireboss)))
    (let (top-level-actions)
      (dolist (possible (f:get-values* 'act-deal-with-new-fire 'instance+inv))
        (when (equal (f:framer (the-first-fire-started))
          (variable-value 'fire :action possible))
          (push possible top-level-actions))))
      (f:get-value*
        (first (sort top-level-actions
          #'>
            :key #'(lambda (x)
              (f:get-value* x 'creation-time))))
          'has-end-action))))))

```

```
(defclip all-plans-to-contain-fire ()
  (:report-key "All plans to Contain Fire")
  (f:using-frame-system ((name-of (find-agent-by-name 'fireboss)))
    (mapcar #'(lambda (act-deal-with-our-fire)
      (get-primitive-action (f:get-value* act-deal-with-our-fire
        'has-end-action)))
      (sort
        (mapcan #'(lambda (act-deal-with-new-fire)
          (when (equal (f:framer (the-first-fire-started))
            (variable-value 'fire :action
              act-deal-with-new-fire))
            (list act-deal-with-new-fire)))
          (f:get-values* 'act-deal-with-new-fire 'instance+inv))
        #'< :key #'(lambda (x) (f:get-value* x 'creation-time))))))
```

B.1.3 Example from a Transportation Planning Simulation

This example comes from a baseline experiment in bottleneck prediction at a shipping port in a transportation planning simulator called TRANSIM. On each day it predicts the occurrence of bottlenecks at a single port in the shipping network, then captures data about actual bottlenecks for later comparison. The experiment collects time-series data, a fragment of which is shown after the example. The summary output file (data collected after each trial) is also shown. Note the use of the keywords to `define-simulator` (`:schedule-function`, `:deactivate-scheduled-function`, `:seconds-per-time-unit` and `:timestamp`) that define the time-series clips.

```
;;; -----
;;; Clips for the collection of time series data over the course of a trial

(defun current-day ()
  "Return the current day of simulated time from 0."
  (/ (current-time) 24))

(defclip port-state-snapshot ()
  "Record state information for a port at the end of each day."
  (:output-file "port-state"
   :schedule (:period "1 day")
   :map-function (list (port 'port-1))
   :components (ships-en-route ships-queued ships-docked
                  expected-ship-arrivals predicted-queue-length)))

(defclip ships-en-route (port)
  "Record the number of ships en route to a port."
  ()
  (length (apply #'append (mapcar 'contents (incoming-channels port)))))

(defclip ships-queued (port)
  "Record the number of ships queued at a port."
  ()
  (length (contents (docking-queue port))))

(defclip ships-docked (port)
  "Record the number of ships docked at a port."
  ()
  (length (apply #'append (mapcar 'contents (docks port)))))

(defclip expected-ship-arrivals (port) ()
  (progn
   (update-prediction)
   (let ((prediction-units (find-units 'prediction '(port-model) :all)))
     (float (/ (round (* 10000
                       (expected-value (first prediction-units)))) 10000))))))
```

```

(defclip port-predicted-value (port) ()
  (let ((prediction-units (find-units 'prediction '(port-model) :all)))
    (generate-prediction (expected-value (first prediction-units)))))

(defclip port-actual-change (port) ()
  (let ((prediction-units (find-units 'prediction '(port-model) :all)))
    (length
     (set-difference
      (ships-previously-in-port port 0)
      (ships-previously-in-port
       port (days-in-future (first prediction-units)))))))

(defclip predicted-queue-length (port) ()
  (let ((prediction-units (find-units 'prediction '(port-model) :all)))
    (pred-queue-length (first prediction-units))))

;;; -----
;;; Clips for the collection of data at the end of a trial

(defclip prediction-score ()
  (:output-file "score"
   :components (score-for-day)
   :map-function '(2 5)))

(defclip score-for-day (day)
  ()
  (compute-score-for-day day))

;;; -----
;;; Experiment and simulator definitions

(define-simulator transsim
  :system-name "TransSim"
  :start-system (simulate nil)
  :reset-system reset-transsim-experiment
  :schedule-function schedule-function-for-clips
  :deactivate-scheduled-function transsim::reset
  :seconds-per-time-unit 3600
  :timestamp current-day)

(define-experiment test-experiment ()
  :simulator transsim
  :instrumentation (prediction-score port-state-snapshot)
  :variables ((prediction-threshold in '(0.6 0.75 0.9))
              (eta-variance-multiplier in '(0.2 0.4 0.6))
              (prediction-point in '(5)))
  :after-trial
  ;; Other options here include building CLASP datasets,

```

```

;; exporting to some database, massaging the data or some
;; combination fo these.
(write-current-experiment-data))

;;; -----
;;; Utilities

(defun reset-transsim-experiment (prediction-threshold
                                eta-variance-multiplier
                                prediction-point)
  (setf *prediction-threshold* prediction-threshold)
  (format t "~&Prediction threshold = ~a" prediction-threshold)
  (setf *eta-variance-multiplier* eta-variance-multiplier)
  (format t "~&ETA variance multiplier = ~a" eta-variance-multiplier)
  (setf *prediction-point* prediction-point)
  (format t "~&Prediction point = ~a" prediction-point)
  (initialize-simulation))

(defun schedule-function-for-clips (function time period name)
  (if period
      (transsim::schedule-recurring-event (transsim::event-actuator :external)
                                          :function function
                                          :time time
                                          :period period
                                          :type (or name :instrumentation))
      (transsim::schedule-event (transsim::event-actuator :external)
                               :function function
                               :time time
                               :type (or name :instrumentation))))

(defun rexp (&key &optional number-of-trials)
  (run-experiment 'test-experiment :output-file "ed-buffer:out"
                 :number-of-trials number-of-trials))

```

Shown below is the summary output file produced for 9 trials. This file is in CLASP format. It begins with an informative header string, followed by a series of strings, each of which is a column name in the data table. Rows of the table follow, stored as lists. Each list contains one element per column name.

```

"
*****
****
**** Experiment: Test-Experiment
**** Machine: Miles
**** TransSim version: Unknown
**** Date: 10/1/93 11:26
**** Scenario: None
**** Script-name: None

```

```
**** First trial number: 1
**** Last trial number: 9
**** Number of trials: 9
**** Max trial length: Unknown hours
*****
```

```
The key follows:"
"Trial"
"Prediction-Threshold"
"Eta-Variance-Multiplier"
"Prediction-Point"
"Score-For-Day 2"
"Score-For-Day 5"
(1 0.75 0.2 5 100 100 )
(2 0.75 0.3 5 100 100 )
(3 0.75 0.4 5 100 100 )
(4 0.8 0.2 5 100 100 )
(5 0.8 0.3 5 100 100 )
(6 0.8 0.4 5 100 100 )
(7 0.85 0.2 5 100 100 )
(8 0.85 0.3 5 100 100 )
(9 0.85 0.4 5 100 100 )
```

This is a fragment of the time-series data (from the first trial only) showing the flat file structure produced by component and time-series clip relationships.

```
"
*****
****
**** Experiment: Test-Experiment
**** Machine: Miles
**** TransSim version: Unknown
**** Date: 10/5/93 15:17
**** Scenario: None
**** Script-name: None
**** First trial number: 1
**** Last trial number: 12
**** Number of trials: 12
**** Max trial length: Unknown hours
*****
```

```
The key follows:"
"Trial"
"Timestamp"
"Prediction-Threshold"
"Eta-Variance-Multiplier"
"Prediction-Point"
"Ships-En-Route port-1"
```

"Ships-Queued port-1"
"Ships-Docked port-1"
"Expected-Ship-Arrivals port-1"
"Predicted-Queue-Length port-1"

(1 0 0.6 0.2 5 0 0 0 0.0 0)
(1 1 0.6 0.2 5 0 0 0 0.0 0)
(1 2 0.6 0.2 5 0 0 0 0.0 0)
(1 3 0.6 0.2 5 0 0 0 0.0 0)
(1 4 0.6 0.2 5 0 0 0 0.0 0)
(1 5 0.6 0.2 5 2 0 0 0.0597 0)
(1 6 0.6 0.2 5 2 0 0 0.3357 0)
(1 7 0.6 0.2 5 2 0 0 0.7643 0)
(1 8 0.6 0.2 5 2 0 0 0.9961 0)
(1 9 0.6 0.2 5 2 0 0 0.9977 0)
(1 10 0.6 0.2 5 4 0 0 0.683 0)
(1 11 0.6 0.2 5 3 0 1 0.2361 0)
(1 12 0.6 0.2 5 3 0 1 0.6389 1)
(1 13 0.6 0.2 5 3 0 1 1.4062 1)
(1 14 0.6 0.2 5 3 0 1 1.9885 1)
(1 15 0.6 0.2 5 3 0 0 2.071 1)
(1 16 0.6 0.2 5 3 0 0 1.5729 0)
(1 17 0.6 0.2 5 1 1 1 0.5053 1)
(1 18 0.6 0.2 5 1 1 1 0.8974 2)
(1 19 0.6 0.2 5 1 1 1 0.9996 2)
(1 20 0.6 0.2 5 1 1 0 1.0 1)
(1 21 0.6 0.2 5 1 0 1 1.0 1)
(1 22 0.6 0.2 5 0 1 1 0.0 1)
(1 23 0.6 0.2 5 0 1 1 0.0 1)
(1 24 0.6 0.2 5 0 1 0 0.0 0)
(1 25 0.6 0.2 5 2 0 1 0.1331 0)
(1 26 0.6 0.2 5 2 0 1 0.355 0)
(1 27 0.6 0.2 5 2 0 1 0.6276 1)
(1 28 0.6 0.2 5 2 0 0 0.946 0)
(1 29 0.6 0.2 5 4 0 0 1.2677 0)
(1 30 0.6 0.2 5 4 0 0 2.0143 1)
(1 31 0.6 0.2 5 4 0 0 2.8832 2)
(1 32 0.6 0.2 5 4 0 0 3.3096 2)
(1 33 0.6 0.2 5 4 0 0 2.4576 1)
(1 34 0.6 0.2 5 2 1 1 1.9172 3)
(1 35 0.6 0.2 5 1 2 1 0.9961 3)

References

- [1] George W. Bohrnstedt and David Knoke. *Statistics for Social Data Analysis*. F. E. Peacock Publishers, Itasca, Illinois, second edition, 1988.
- [2] James V. Bradley. *Distribution-Free Statistical Tests*. Prentice-Hall, Inc., Englewood Cliffs, NJ, 1968.
- [3] Paul R. Cohen. *Empirical Methods in Artificial Intelligence*. MIT Press, forthcoming.
- [4] Jay L. Devore. *Probability and Statistics for Engineering and the Sciences*. Brooks/Cole Publishing Company, Monterey, California, first edition, 1982.
- [5] John E. Freund. *Modern Elementary Statistics*. Prentice-Hall, Englewood Cliffs, NJ, seventh edition, 1988.
- [6] David M. Hart and Paul R. Cohen. Predicting and explaining success and task duration in the phoenix planner. In *Proceedings of the First International Conference on AI Planning Systems*, pages 106–115. Morgan Kaufmann, 1992.
- [7] Adele E. Howe and Paul R. Cohen. Responding to environmental change. In *DARPA Workshop on Innovative Approaches to Planning, Scheduling and Control*, pages 85–92. Morgan Kaufmann, 1990.
- [8] Edward W. Minium. *Statistical Reasoning in Psychology and Education*. John Wiley and Sons, 1978.
- [9] Elazar J. Pedhazur. *Multiple Regression in Behavioral Research*. Holt, Rinehart and Winston, 1973.
- [10] William H. Press, Brian P. Flannery, Saul A. Teukolsky, and William T. Vetterling. *Numerical Recipes in C: The Art of Scientific Computing*. Cambridge University Press, 1988.