

VarScreen

A program for screening predictor and/or target variables that will be employed in model building

Version 1.82, containing:

- *Univariate mutual information, with p-values compensated for selection bias, and probability of the best in-sample selection underperforming others out-of-sample*
- *Bivariate mutual information and uncertainty reduction, with p-values compensated for selection bias*
- *Optimal predictor sets defined by Relevance minus Redundancy, including solo and group p-values*
- *Hidden Markov models chosen according to their multivariate correlation with a target, including p-values compensated for selection bias*
- *Detecting change in the mean of a time series (such as deterioration of performance of a market trading system), compensated for multiple series as well as multiple comparisons across time*
- *Ensemble FREL (**F**eature **W**eighting as **R**egularized **E**nergy-based **L**earning) for high-dimensionality, small-sample applications*
- *FSCA (**F**orward **S**election **C**omponent **A**nalysis) for forward and optional backward refinement of maximum-variance-capture components from a subset of a large group of variables*
- *LFS (**L**ocal **F**eature **S**election) for identifying predictors that are optimal in localized areas of the feature space but may not be globally optimal. Such predictors can be effectively used by nonlinear models but are neglected by many other feature selection algorithms.*

Table of Contents

| | |
|---|----|
| About the VarScreen Program | 1 |
| About CUDA Processing..... | 3 |
| Reading a Dataset | 4 |
| Univariate Mutual Information | 5 |
| Bivariate Mutual Information / Uncertainty Reduction | 14 |
| Predictors having Max Relevance, Min Redundancy | 24 |
| Hidden Markov Models with Target Correlation | 31 |
| Stationarity Test for Break in Mean | 41 |
| FREL: Feature Weighting as Regularized Energy-Based Learning..... | 48 |
| FSCA: Forward Selection Component Analysis | 58 |
| LFS: Local Feature Selection | 65 |
| Appendix: Version Updates | 75 |

About the VarScreen Program

VarScreen contains in one easy-to-use program a variety of software tools useful for the developer of predictive models. These tools screen and evaluate candidates for predictors and targets. More on this later. But first, we need to issue a vitally important disclaimer:

This program is an experimental work in progress. It is provided free of charge to interested users for educational purposes only. In all likelihood this program contains errors and omissions. If you use this program for a purpose in which loss is possible, then you are fully responsible for any and all losses associated with use of this program. The developer of this program disclaims all responsibility for losses which the user may incur.

Okay, enough of that. You've been warned. The *VarScreen* program is being developed with two major goals in mind:

- 1) The program should be exceptionally easy to learn and use. Results should be obtainable with no more than a few intuitive mouse clicks and key presses. Detailed study of an exhaustive manual should not be required.
- 2) The software should provide cutting-edge statistical information, employing tests and algorithms not readily available in any standard analysis software.

I believe that these goals have been and will continue to be obtained.

Finally, understand that *VarScreen* is a work in progress. New screening algorithms will likely be added on a regular basis. Stay tuned. Updates will be reported on the author's website: **TimothyMasters.info**.

Features of the program

In keeping with the goals of simplicity plus mathematical sophistication, the following items are noteworthy:

- Most operations involve just two quick steps: read the data and select the test to be performed. Program-supplied defaults are often satisfactory, and adjusting them is easy. The next section will describe reading the data, and subsequent sections will describe the tests that can be performed.
- The program is fully multi-threaded, enabling it to take maximum advantage of modern multiple-core processors. As of this writing, many over-the-counter computers contain a CPU with six cores, each of which is hyperthreaded to perform two sets of operation streams simultaneously. *VarScreen* keeps all twelve of these threads busy as much as possible, which tremendously speeds operation compared to single-threaded programs.
- The most massively compute-intensive algorithms make use of any CUDA-enabled nVidia card in the user's computer. These widely available video cards (standard hardware on many computers) turn an ordinary desktop computer into a super-computer, accelerating computations by several orders of magnitude. Enormously complex algorithms that would require days of compute time on an ordinary computer with ordinary software can execute in several minutes using the *VarScreen* program on a computer with a modern nVidia display card.
- Rather than printing results on the screen, the program writes a log file called VARSCREEN.LOG. This way a 'permanent' copy of all results is available for optional printing and archiving.

About CUDA Processing

CUDA stands for Compute Unified Device Architecture. It is the interface system by which nVidia makes the massive parallel processing hardware of its video display cards available to applications. The power of this hardware is breathtaking; the GTX Titan video card contains nearly 3000 processors that can execute programs simultaneously. *VarScreen* makes use of this capability for especially time-consuming tasks.

There is an annoying quirk, however, which users of *VarScreen* should be aware of. Microsoft, in its infinite wisdom, forbids any Windows program from executing a CUDA application for longer than two seconds. Moreover, Windows makes it almost impossible for most users to increase or disable this limit; doing so involves tampering with the Registry, a frightening endeavor. Unfortunately, some large problems can require far more than two seconds of CUDA time.

In order to get around this issue, *VarScreen* breaks up large tasks into multiple small tasks. Each such task is called a *Launch*. An ugly tradeoff is involved in this breakup. Each launch incurs a significant overhead, so one should minimize the number of launches. On the other hand, increasing the workload of each launch increases the probability that the deadly two-second limit will be reached, with the result that Windows terminates the program, and somewhere, behind some closed door, a Microsoft programmer snickers. Due to the large variety of CUDA hardware available, it is not practical to predict in advance how long a launch will tie up the CUDA processing, so one must be conservative.

The reason I am making such an issue of this is to allow the user of *VarScreen* to understand a bit of output written to the screen. Whenever a large test involving CUDA computations is running, a progress bar is displayed. This bar also includes text similar to the following:

Max CUDA time = 23 ms in 2 launches

What this means is that each task had to be broken up into two launches, and the maximum CUDA processing time for those two launches was 23 milliseconds. There is one reason why this may be important to the user: if the time approaches 2000 (two seconds) you are near crashing (a brief black screen followed by a message that the video card has been reset). I would be grateful if you contacted me at my email: tim@TimothyMasters.info and reported this so I can continue to tweak the program.

Reading a Dataset

VarScreen reads data files that are in a common data format: the first record names the fields, and each subsequent record is a single case. For example, the first few lines of a dataset might look like this:

```
X1 X2 X3 Y
3.14 0.21 -5.33 4.01
-1.02 -0.45 2.12 -7.02
...
```

Variable names may be at most 15 characters long. Spaces, commas, and tabs may be used as delimiters. One implication of this fact is that variable names must not contain spaces. In place of a space, the underscore character (`_`) may be used. Numeric values must be strictly numeric; scientific notation (i.e. `3.14e-9`) is illegal in the current version of the program. If users scream loudly enough, this feature may be added later.

Files exported from Microsoft Excel as comma-delimited (.CSV) files are generally readable by *VarScreen*, although if dates with slashes appear, or other text fields appear, trouble may be encountered. (Text variables or otherwise non-numeric fields will typically be assigned the value 0.0.) If exporting from Excel, also beware of column headers that contain spaces. CSV files strictly use commas as delimiters, so spaces in column names are legal in Excel, but since *VarScreen* treats spaces as delimiters, the single variable name in Excel will be mistakenly treated as two or more variables in *VarScreen* if the name contains spaces.

Missing data is not allowed; every data record must have a numeric value present for every field. Note that if a file exported from Excel contains missing data, this will be represented in the file as contiguous commas, which will cause problems for *VarScreen*.

After the file is read, the log file `VARSCREEN.LOG` will contain a table of the mean and standard deviation of every variable in the file. Users should get in the habit of skimming this table as a quick sanity check of the validity of the data; a wild value in the table may indicate an unexpected flaw in the data file.

One additional variable is computed: `_SEQNUM_`. For each case this is the sequence number of the case within the dataset. The value of `_SEQNUM_` is 1 for the first case, 2 for the second, and so forth. One interesting use for this variable arises when the data is a time series. A relationship such as mutual information between `_SEQNUM_` and a variable indicates that the variable is probably nonstationary.

Univariate Mutual Information

The *Univariate Mutual Information* test computes the mutual information between a specified target variable and each of a specified set of predictor candidates. The predictors are then listed in the VARSCREEN.LOG file in descending order of mutual information. Along with each candidate, a specialized probability described later, as well as the *Solo pval* and *Unbiased pval*, are printed if Monte-Carlo replications are requested.

The *Solo pval* is the probability that a candidate that has a strictly random (no predictive power) relationship with the target could have, by sheer good luck, had a mutual information at least as high as that obtained. If this quantity is not small, the developer should strongly suspect that the candidate is worthless for predicting the target. Of course, this logic is, in a sense, accepting a null hypothesis, which is well known to be a dangerous practice. However, if a reasonable number of cases are present and a reasonable number of Monte-Carlo replications have been done, this test is powerful enough that failure to achieve a small p-value can be interpreted as the candidate having little or no predictive power.

The problem with the *Solo pval* is that if more than one candidate is tested (the usual situation!), then there is a large probability that some truly worthless candidate will be lucky enough to achieve a high level of mutual information, and hence achieve a very small *Solo pval*. In fact, if all candidates are worthless, the *Solo pvals* will follow a uniform distribution, frequently obtaining small values by random chance. This situation can be remedied by conducting a more advanced test which accounts for this *selection bias*. The *Unbiased pval* for the best performer in the candidate set is the probability that this best performer could have attained its exalted level of performance by sheer luck if all candidates were truly worthless.

The *Unbiased pval* is printed for all candidates, not just the best. For those other, lesser candidates, the *Unbiased pval* is an upper bound (a conservative measure) for the true unbiased p-value of the candidate. Thus, a very small *Unbiased pval* for any candidate is a strong indication that the candidate has true predictive power. Unfortunately, unlike the *Solo pval*, large values of the *Unbiased pval* are not necessarily evidence that the candidate is worthless. Large values, especially near the bottom of the sorted list, may be due to over-estimation of the true p-value. The author is not aware of any algorithm for computing correct unbiased p-values for any candidate other than the best. However, because this measure is conservative, it does have great utility in selecting promising predictors.

The user must be aware of a vital caveat to this procedure: The *Solo pval* and *Unbiased pval* computations fall apart if there is significant serial correlation (or any other dependency) among both the target variable as well as one or more of the predictor candidates. In most practical applications, the predictor candidates are hopelessly dependent, so the key is the target variable. If it has anything beyond tiny dependency (typically serial correlation), the test will become anti-conservative: the computed p-values will be smaller than the correct values. This is dangerous. *VarScreen* contains an option that somewhat helps in this situation, but it is not a complete cure.

The final column printed is inspired by a research report titled “The Probability of Backtest Overfitting” by David Bailey, Jonathan Borwein, Marcos Lopez de Prado, and Jim Zhu. Like the permutation test, it assumes that there is no significant serial correlation among both the target variable and one or more predictor candidates, although it tends to be fairly robust in this regard. I heavily modified their clever algorithm to apply to mutual information.

When one examines a pool of candidates and selects a predictor based on its having the maximum value of some criterion such as mutual information, one hopes that this superiority will carry over to data not yet seen (out-of-sample or OOS data). In particular, consider the (unknown at test time) median OOS performance of all predictor candidates. At a minimum, one would hope that the OOS performance of the candidate selected based on its having maximum in-sample performance would exceed the median OOS performance of all candidates. If not, the selection process is useless; no superiority is obtained by choosing the best in-sample performer.

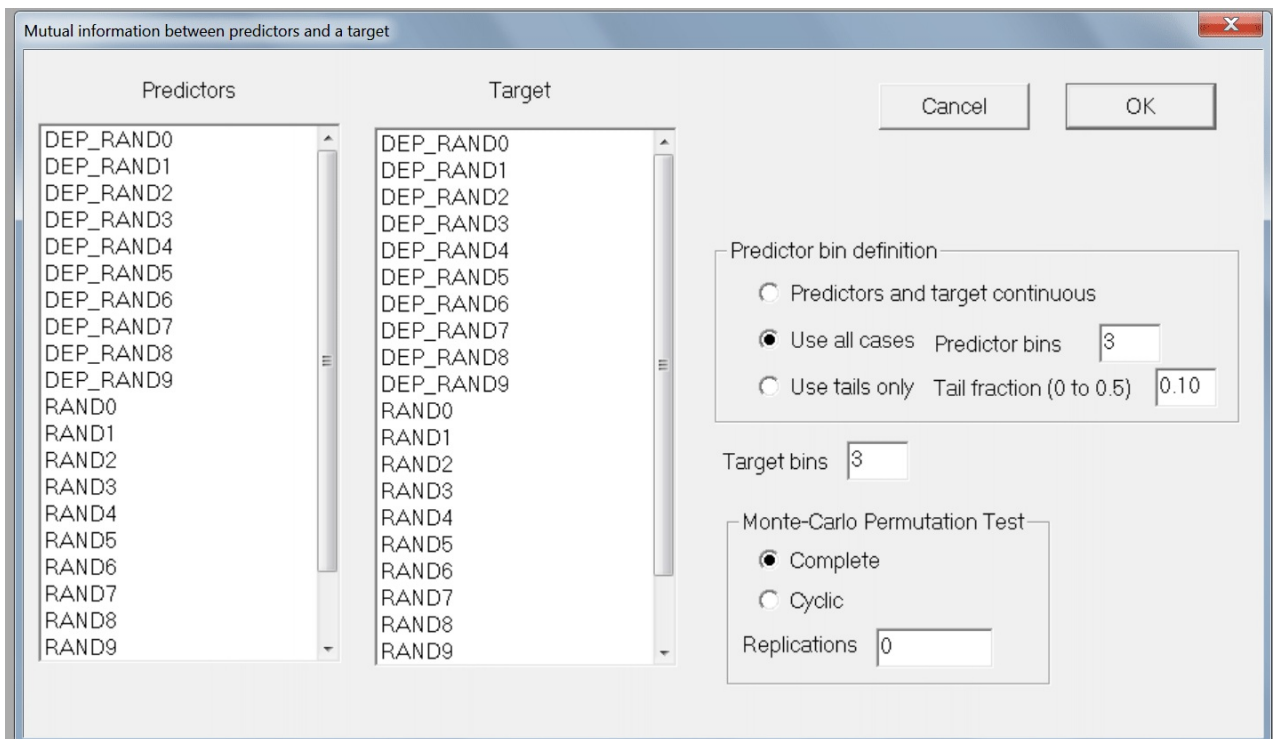
The rightmost figure printed for the first row (the best candidate) is the estimated probability that the OOS mutual information of this selected candidate will be less than or equal to the median OOS performance of all of the other candidates. Obviously, we want this probability to be small.

The figures printed for subsequent rows are the equivalent probabilities for lower rank orders. For example, the figure for the second row is the probability that the second best in-sample candidate will have OOS performance less than or equal to the median. This is subtly different from the probability for the particular candidate that was selected; it’s a more theoretical figure. Nonetheless, equating the two should not be unreasonable.

Ideally, one would see low probabilities near the top (the best in-sample candidates outperform OOS) and high probabilities near the bottom (the worst underperform). A large quantity of worthless candidates will make the distribution more random.

Specifying the Test Parameters

When the user clicks *Tests / Univariate Mutual Information*, a dialog similar to that shown below will appear. The various parameters are described on the following page.



The leftmost column is used to specify the set of predictor candidates. Multiple candidates can be selected by dragging the mouse cursor across a block, or by clicking the first candidate in a block, holding the Shift key, and clicking the last candidate in the block. Individual candidates can be toggled on and off by holding the Ctrl key while clicking on the variable.

The *Target* column is used to select a single target variable.

Three methods for computing mutual information are available, and the method to use is chosen by selecting one of the three buttons in the *Predictor bin definition* block:

Predictors and target continuous uses the Darbellay-Vajda algorithm (fully described in “Assessing and Improving Prediction and Classification” by Timothy Masters) to compute continuous mutual information. This method is appropriate (and almost always the preferred approach) when all variables are continuous or nearly so. It’s main disadvantage is that it is much slower to compute than the bin methods. Also, candidates that have tiny mutual information with the target will have their computed mutual information reduced to exactly zero by the algorithm. This will produce a sudden discontinuity in p-values, which may appear unusual but which in fact is perfectly reasonable.

Use all cases partitions each predictor into bins that are as equal in size as possible. The user must specify the number of bins to employ, and unless the dataset is huge the default of three bins is frequently appropriate.

Use tails only computes mutual information based on only the maximum and minimum collection of values of each predictor. The *tail fraction* specified by the user is the fraction of cases in each tail. So, for example, the default *tail fraction* of 0.1 would use the cases having the smallest ten percent and the largest ten percent of predictor values. The 80 percent of cases having intermediate values of the predictor candidate would be completely ignored in the mutual information calculation. This method is especially useful in high-noise situations, such as prediction of financial markets. The probability that superior mutual information will hold up out-of-sample cannot be computed when this option is selected.

Target bins must be specified if *Use all cases* or *Use tails only* is chosen. This is the number of approximately equal-size bins into which the target variable is distributed. The default value of 3 is appropriate for a wide variety of applications. This field is ignored if the *Predictors and target continuous* option is selected.

Replications defaults to zero, in which case no Monte-Carlo Permutation Test is performed. However, it is usually best to set this to at least 100, and perhaps as much as 1000, so that solo and unbiased p-values will be computed. Note that the minimum possible p-value is the reciprocal of the number of permutations. So, for example, if the user specifies 100 permutations, the minimum p-value that can appear is 0.01. Run time of this test is linearly related to the number of permutations.

The user must choose either *Complete* or *Cyclic* permutations. If the user is confident that there is no dependency as described earlier, then *Complete* should be used; it is the traditional approach which does a complete random shuffle for each permutation. However, if there is dependency, this type of shuffling will produce underestimation of *p-values*, a very dangerous situation. If the dependency is serial (the data is a time series and the dependency is among samples close in time) then a slight improvement in the situation can be obtained by using *Cyclic* permutation. In this type of shuffle, the time order of the target is kept intact except at the ends by rotating the targets with end-point wraparound. Shuffling this way preserves most of the serial dependency in the permuted targets, which makes the algorithm more accurate. The *p-values* computed this way will generally be larger than those computed with complete shuffling, and hence less likely to lead to false rejection of the null hypothesis of no predictive power. But be warned that the cure is far from complete; computed *p-values* will still underestimate the true values, just not as badly.

Note that in most cases it is legitimate to use *Cyclic* permutation instead of *Complete* when there is no dependency. However, if the dataset is small, *Cyclic* permutation will limit the number of unique permutations and hence increase the random error inherent in the process. As long as the dataset is large, some users may prefer to use *Cyclic* permutation even if it is assumed that there is no serial dependency; in case there really is hidden serial dependency, this is a cheap insurance policy. Still, the best practice is to make sure that the data does not contain dependency and then use *Complete* permutation. Relying on *Cyclic* permutation to take care of dependency problems is living dangerously. And if the dataset contains fewer than 1000 or so cases, use of *Cyclic* permutation is not recommended unless it is necessary to handle dependency.

Examples of Univariate Mutual Information

This section demonstrates three situations, all using synthetic data to clarify the presentation. The variables in the dataset are as follows:

RAND0 - RAND9 are independent (within themselves and with each other) random time series.

DEP_RAND0 - DEP_RAND9 are derived from *RAND0 - RAND9* by introducing strong serial correlation up to a lag of nine observations. They are independent of one another.

SUM12 = RAND1 + RAND2

SUM34 = RAND3 + RAND4

SUM1234 = SUM12 + SUM34

The first test run attempts to predict *SUM1234* from *RAND0 - RAND9*, *SUM12*, and *SUM34*. The output looks like this:

```
*****
*
* Computing univariate mutual information (one predictor, one target) *
* 12 predictor candidates *
* 5 predictor bins *
* 5 target bins *
* 10000 replications of complete Monte-Carlo Permutation Test *
*
*****
```

The bounds that define bins are now shown

Target bounds are based on the entire dataset...

```
-0.97362      -0.27795      0.31417      1.00879
```

Variable Bounds...

| | | | | |
|-------|----------|----------|---------|---------|
| RAND0 | -0.59427 | -0.18805 | 0.20723 | 0.60549 |
| RAND1 | -0.58905 | -0.18795 | 0.22570 | 0.62047 |
| RAND2 | -0.59430 | -0.18090 | 0.21697 | 0.61045 |
| RAND3 | -0.62008 | -0.20843 | 0.19894 | 0.59159 |
| RAND4 | -0.59696 | -0.18753 | 0.21087 | 0.61077 |
| RAND5 | -0.59819 | -0.21468 | 0.18130 | 0.56676 |
| RAND6 | -0.61150 | -0.21273 | 0.19102 | 0.59680 |
| RAND7 | -0.61383 | -0.22039 | 0.18521 | 0.58843 |
| RAND8 | -0.59055 | -0.19032 | 0.20591 | 0.59859 |

| | | | | |
|-------|----------|----------|---------|---------|
| RAND9 | -0.60422 | -0.19932 | 0.20315 | 0.58792 |
| SUM12 | -0.67798 | -0.17129 | 0.22588 | 0.74242 |
| SUM34 | -0.73810 | -0.21209 | 0.21164 | 0.74363 |

The marginal distributions are now shown.

If the data is continuous, the marginals will be nearly equal.

Widely unequal marginals indicate potentially problematic ties.

Target marginals are based on the entire dataset...

0.19987 0.20003 0.20003 0.20003 0.20003

Variable Marginal...

| | | | | | |
|-------|---------|---------|---------|---------|---------|
| RAND0 | 0.19987 | 0.20003 | 0.20003 | 0.20003 | 0.20003 |
| RAND1 | 0.19987 | 0.20003 | 0.20003 | 0.20003 | 0.20003 |
| RAND2 | 0.19987 | 0.20003 | 0.20003 | 0.20003 | 0.20003 |
| RAND3 | 0.19987 | 0.20003 | 0.20003 | 0.20003 | 0.20003 |
| RAND4 | 0.19987 | 0.20003 | 0.20003 | 0.20003 | 0.20003 |
| RAND5 | 0.19987 | 0.20003 | 0.20003 | 0.20003 | 0.20003 |
| RAND6 | 0.19987 | 0.20003 | 0.20003 | 0.20003 | 0.20003 |
| RAND7 | 0.19987 | 0.20003 | 0.20003 | 0.20003 | 0.20003 |
| RAND8 | 0.19987 | 0.20003 | 0.20003 | 0.20003 | 0.20003 |
| RAND9 | 0.19987 | 0.20003 | 0.20003 | 0.20003 | 0.20003 |
| SUM12 | 0.19987 | 0.20003 | 0.20003 | 0.20003 | 0.20003 |
| SUM34 | 0.19987 | 0.20003 | 0.20003 | 0.20003 | 0.20003 |

-----> Mutual Information with SUM1234 <-----

| Variable | MI | Solo pval | Unbiased pval | P(<=median) |
|----------|--------|-----------|---------------|-------------|
| SUM34 | 0.2877 | 0.0001 | 0.0001 | 0.0000 |
| SUM12 | 0.2610 | 0.0001 | 0.0001 | 0.0000 |
| RAND3 | 0.1307 | 0.0001 | 0.0001 | 0.0000 |
| RAND4 | 0.1263 | 0.0001 | 0.0001 | 0.0000 |
| RAND1 | 0.1129 | 0.0001 | 0.0001 | 0.0000 |
| RAND2 | 0.1085 | 0.0001 | 0.0001 | 0.0000 |
| RAND8 | 0.0015 | 0.2994 | 0.9828 | 1.0000 |
| RAND5 | 0.0014 | 0.3673 | 0.9950 | 1.0000 |
| RAND6 | 0.0012 | 0.5303 | 1.0000 | 1.0000 |
| RAND7 | 0.0010 | 0.7384 | 1.0000 | 1.0000 |
| RAND0 | 0.0008 | 0.8332 | 1.0000 | 1.0000 |
| RAND9 | 0.0006 | 0.9605 | 1.0000 | 1.0000 |

The bounds that define the target and predictor bins are shown, along with the marginal probabilities. If any marginal is far from being equal, that variable has significant ties and the situation should be investigated.

As expected, the best predictors of SUM1234 are SUM12 and SUM34. RAND1 - RAND4 are the next best. All other predictors are obviously worthless. Note how dramatically the unbiased p-value delineates the break.

The next example shows what happens when worthless and serially correlated predictors are tested with a serially correlated target. We use DEP_RANDOM1 - DEP_RANDOM9 to predict DEP_RANDOM0, a situation which should demonstrate no predictive power whatsoever. The mutual information table is as follows:

```
-----> Mutual Information with DEP_RANDOM0 <-----
```

| Variable | MI | Solo pval | Unbiased pval | P(<=median) |
|-------------|--------|-----------|---------------|-------------|
| DEP_RANDOM2 | 0.0044 | 0.0001 | 0.0002 | 0.6944 |
| DEP_RANDOM4 | 0.0030 | 0.0018 | 0.0175 | 0.6190 |
| DEP_RANDOM3 | 0.0025 | 0.0110 | 0.0881 | 0.6270 |
| DEP_RANDOM6 | 0.0023 | 0.0249 | 0.2004 | 0.5516 |
| DEP_RANDOM9 | 0.0023 | 0.0242 | 0.2062 | 0.5397 |
| DEP_RANDOM8 | 0.0023 | 0.0287 | 0.2284 | 0.5079 |
| DEP_RANDOM1 | 0.0022 | 0.0317 | 0.2494 | 0.4960 |
| DEP_RANDOM5 | 0.0019 | 0.0883 | 0.5509 | 0.4325 |
| DEP_RANDOM7 | 0.0008 | 0.8682 | 1.0000 | 0.5317 |

The mutual information figures are all tiny, yet the p-values show extreme significance. The careless user would surely be fooled by this, because not only are the solo p-values mostly small, but even the unbiased p-value has been fooled for one or two of the candidates.

It should be emphasized that this phenomenon is not an artifact of just the Monte-Carlo Permutation Test. This is a universal phenomenon, which is why Statistics 101 courses always emphasize the importance of independent observations. The simple explanation of why this occurs is that any sort of dependence reduces the effective degrees of freedom of the test. The testing procedure looks at the number of cases and proceeds accordingly, but the dependence in the data increases the variance of the test statistic beyond what would be expected from a sample of the given size. Thus we are more likely to falsely reject the null hypothesis.

Observe that in this 'no predictive power' case, despite the serial correlation, the probabilities in the final column are distributed around 0.5, which would be expected when none of the candidates has predictive power. This is because the best in-sample candidate is random, and hence its associated out-of-sample performance has about a 50-50 chance of lying above or below the median. This is the pattern usually seen when all candidates are worthless.

The final example shows how the cyclic modification of the Monte-Carlo Permutation Test can at least partially remedy the situation. We repeat the same test as that just shown, except that instead of using *Complete* permutation we use *Cyclic* permutation. The results are shown below:

```
-----> Mutual Information with DEP_RANDOM <-----
```

| Variable | MI | Solo pval | Unbiased pval | P(<=median) |
|-------------|--------|-----------|---------------|-------------|
| DEP_RANDOM2 | 0.0044 | 0.0513 | 0.3529 | 0.6944 |
| DEP_RANDOM4 | 0.0030 | 0.2408 | 0.9316 | 0.6190 |
| DEP_RANDOM3 | 0.0025 | 0.3976 | 0.9918 | 0.6270 |
| DEP_RANDOM6 | 0.0023 | 0.5007 | 0.9976 | 0.5516 |
| DEP_RANDOM9 | 0.0023 | 0.5237 | 0.9982 | 0.5397 |
| DEP_RANDOM8 | 0.0023 | 0.4719 | 0.9988 | 0.5079 |
| DEP_RANDOM1 | 0.0022 | 0.5344 | 0.9990 | 0.4960 |
| DEP_RANDOM5 | 0.0019 | 0.6643 | 1.0000 | 0.4325 |
| DEP_RANDOM7 | 0.0008 | 0.9920 | 1.0000 | 0.5317 |

Now observe that even the largest random relationship is not significant at the 0.05 level on a solo basis, and the unbiased p-value is far from significant.

14 Bivariate Mutual Information / Uncertainty Reduction

Bivariate Mutual Information / Uncertainty Reduction

Sometimes a single variable acting alone has little or no predictive power, but in conjunction with another it becomes useful. The classic example is the height and weight of an individual, predicting coronary health. Either predictor alone has relatively little predictive power, but the two taken together can have great power.

Also, sometimes we have several equally useful candidates for the target variable, and we are not sure which will be most predictable. One example of this situation is when the application is predicting future movement of a financial market with the goal of taking a position and then hopefully closing the position with a profit. Should we employ a tight stop to discourage severe losses? Or should we use a loose stop to avoid being closed out by random noise? We might test multiple targets corresponding to various degrees of stop positioning, and then determine which of the competitors is most predictable.

The *Bivariate Mutual Information* test handles both of these situations. It computes the mutual information or uncertainty reduction between each of one or more specified target variables and each possible pair of predictors taken from a specified set of predictor candidates. The predictor pairs and associated targets are then listed in the VARSCREEN.LOG file in descending order of mutual information. Along with each such set, the *Solo pval* and *Unbiased pval* are printed if Monte-Carlo replications are requested.

The *Solo pval* is the probability that a pair of candidates that has a strictly random (no predictive power) relationship with the target could have, by sheer good luck, had a relationship at least as high as that obtained. If this quantity is not small, the developer should strongly suspect that the candidate is worthless for predicting the target. Of course, this logic is, in a sense, accepting a null hypothesis, which is well known to be a dangerous practice. However, if a reasonable number of cases are present and a reasonable number of Monte-Carlo replications have been done, this test is powerful enough that failure to achieve a small p-value can be interpreted as the candidate having little or no predictive power.

The problem with the *Solo pval* is that if more than one candidate set (a set being two predictors and a target) is tested (the usual situation!), then there is a large probability that some truly worthless candidate set will be lucky enough to achieve a high level of the relationship criterion, and hence achieve a very small *Solo pval*. In fact, if all candidate sets are worthless, the *Solo pvals* will follow a uniform distribution, frequently

obtaining small values by random chance. This situation can be remedied by conducting a more advanced test which accounts for this *selection bias*. The *Unbiased pval* for the best performing candidate set is the probability that this best performer could have attained its exalted level of performance by sheer luck if all candidate sets were truly worthless.

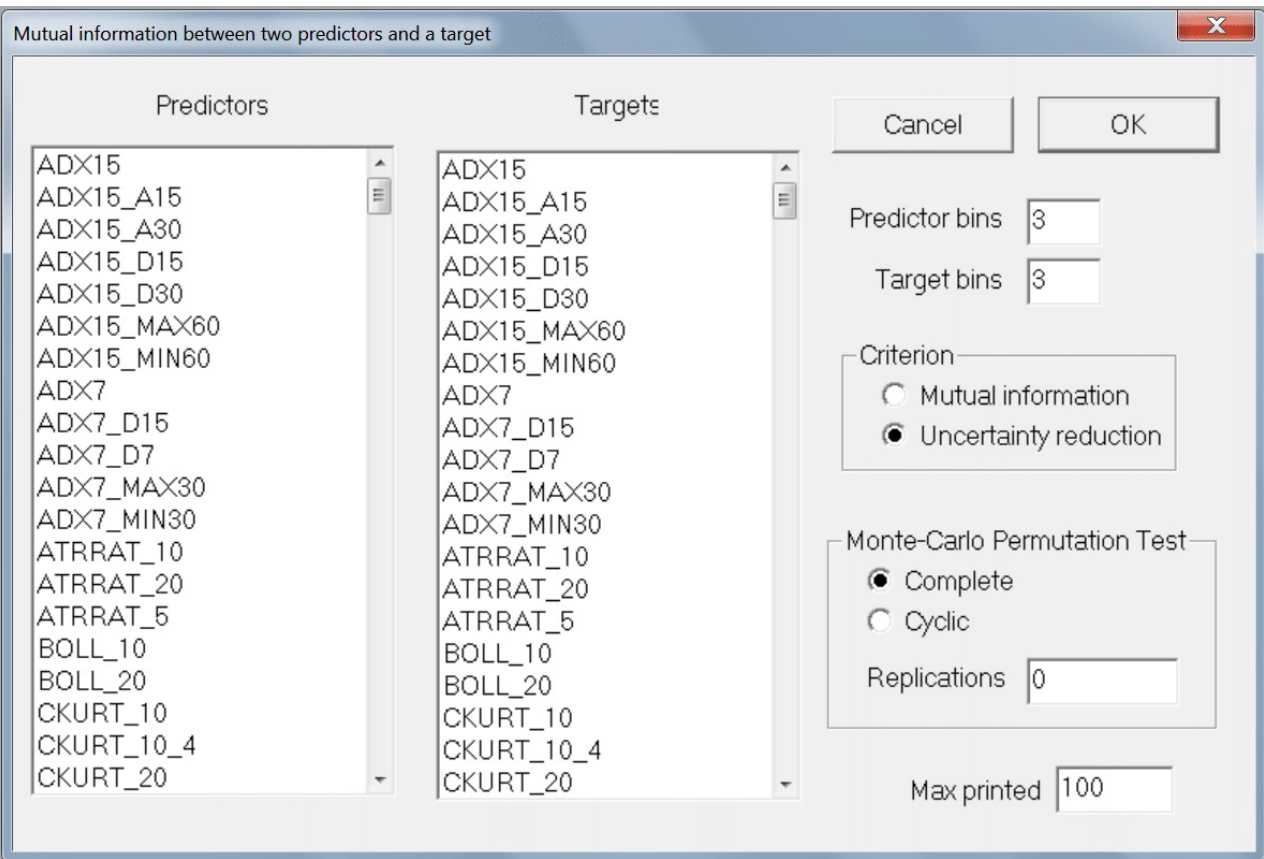
The *Unbiased pval* is printed for all candidate sets, not just the best. For those other, lesser candidates, the *Unbiased pval* is an upper bound (a conservative measure) for the true unbiased p-value of the candidate set. Thus, a very small *Unbiased pval* for any candidate set is a strong indication that the pair of predictors has true predictive power for the target. Unfortunately, unlike the *Solo pval*, large values of the *Unbiased pval* are not necessarily evidence that the candidate set is worthless. Large values, especially near the bottom of the sorted list, may be due to over-estimation of the true p-value. The author is not aware of any algorithm for computing correct unbiased p-values for any candidate set other than the best. However, because this measure is conservative, it does have great utility in selecting promising predictors.

The user must be aware of a vital caveat to this procedure: The *Solo pval* and *Unbiased pval* computations fall apart if there is significant serial correlation (or any other dependency) among one or more target variables as well as one or more of the predictor candidates. In most practical applications, the predictor candidates are hopelessly dependent, so the key is the target variable. If it has anything beyond tiny dependency (typically serial correlation), the test will become anti-conservative: the computed p-values will be smaller than the correct values. This is dangerous. *VarScreen* contains an option that somewhat helps in this situation, but it is not a complete cure.

16 Bivariate Mutual Information / Uncertainty Reduction

Specifying the Test Parameters

When the user clicks *Tests / Bivariate Mutual Information*, a dialog similar to that shown below will appear. The various parameters are described after the dialog.



The leftmost column is used to specify the set of predictor candidates. Multiple candidates can be selected by dragging the mouse cursor across a block, or by clicking the first candidate in a block, holding the Shift key, and clicking the last candidate in the block. Individual candidates can be toggled on and off by holding the Ctrl key while clicking on the variable.

The *Target* column is used to select one or more target variables, with multiple selections obtained as described for predictors.

The predictors and the targets are partitioned into bins that are as equal in size as possible. The user must specify the number of bins to employ for each, and unless the dataset is huge the default of three bins is frequently appropriate.

There can be an annoying problem when using mutual information as a measure of relationship when more than one target is in competition. Mutual information is highly related to the entropy of the predictor and target. If there is only one target in play, the mutual information between it and each predictor candidate will have the same rank order as the uncertainty reduction. But if there are several targets in competition and they have widely disparate entropies, then mutual information is not a good measure of their relationship because the target entropies can confound the rank ordering.

What you are really interested in is the degree to which uncertainty about a target is reduced by having knowledge of a predictor. It can be thought of as their mutual information divided by the entropy of the target. Equivalently, it is the fraction of the target's entropy which is mutual information. For example, if they have zero mutual information, there will be zero uncertainty reduction (about the target) by knowing the predictor. At the other extreme, if their mutual information equals the target entropy, then knowing the predictor will provide perfect (1.0) uncertainty reduction regarding the target.

Thus, a target with high entropy will need high mutual information in order to have a high relationship score. For this reason, *uncertainty reduction* is the default for this test. Much more detail on this important concept can be found in "Assessing and Improving Prediction and Classification" by Timothy Masters.

Replications defaults to zero, in which case no Monte-Carlo Permutation Test is performed. However, it is usually best to set this to at least 100, and perhaps as much as 1000, so that solo and unbiased p-values will be computed. Note that the minimum possible p-value is the reciprocal of the number of permutations. So, for example, if the user specifies 100 permutations, the minimum p-value that can appear is 0.01. Run time of this test is linearly related to the number of permutations.

The user must choose either *Complete* or *Cyclic* permutations. If the user is confident that there is no dependency as described earlier, then *Complete* should be used; it is the traditional approach which does a complete random shuffle for each permutation. However, if there is dependency, this type of shuffling will produce underestimation of *p-values*, a very dangerous situation. If the dependency is serial (the data is a time series and the dependency is among samples close in time) then a slight improvement in the situation can be obtained by using *Cyclic* permutation. In this type of shuffle, the time

18 Bivariate Mutual Information / Uncertainty Reduction

order of the target is kept intact except at the ends by rotating the targets with end-point wraparound. Shuffling this way preserves most of the serial dependency in the permuted targets, which makes the algorithm more accurate. The *p-values* computed this way will generally be larger than those computed with complete shuffling, and hence less likely to lead to false rejection of the null hypothesis of no predictive power. But be warned that the cure is far from complete; computed *p-values* will still underestimate the true values, just not as badly.

Note that in most cases it is legitimate to use *Cyclic* permutation instead of *Complete* when there is no dependency. However, if the dataset is small, *Cyclic* permutation will limit the number of unique permutations and hence increase the random error inherent in the process. As long as the dataset is large, some users may prefer to use *Cyclic* permutation even if it is assumed that there is no serial dependency; in case there really is hidden serial dependency, this is a cheap insurance policy. Still, the best practice is to make sure that the data does not contain dependency and then use *Complete* permutation. Relying on *Cyclic* permutation to take care of dependency problems is living dangerously. And if the dataset contains fewer than 1000 or so cases, use of *Cyclic* permutation is not recommended unless it is necessary to handle dependency.

Examples of Bivariate Mutual Information

This section demonstrates three situations, all using synthetic data to clarify the presentation. The variables in the dataset are as follows:

RAND0 - RAND9 are independent (within themselves and with each other) random time series.

DEP_RAND0 - DEP_RAND9 are derived from *RAND0 - RAND9* by introducing strong serial correlation up to a lag of nine observations. They are independent of one another.

SUM12 = RAND1 + RAND2

SUM34 = RAND3 + RAND4

SUM1234 = SUM12 + SUM34

The first test run attempts to predict *SUM1234* from *RAND0 - RAND9*, *SUM12*, and *SUM34*. Two predictors at a time will be used. The output is shown below. For bin boundaries and marginals, the predictor candidates are shown first, followed by a single blank line, and then the target candidates (just one in this example) appear.

```
*****
*
* Computing bivariate mutual information (two predictors, one target) *
*   12 predictor candidates                                           *
*    1 target candidates                                             *
*   66 predictor/target combinations to test                         *
*  100 best combinations will be printed                             *
*    5 predictor bins                                                *
*    5 target bins                                                  *
*  10000 replications of complete Monte-Carlo Permutation Test      *
*
*****
```

The bounds that define bins are now shown

| | | | | |
|-------|----------|----------|---------|---------|
| RAND0 | -0.59427 | -0.18805 | 0.20723 | 0.60549 |
| RAND1 | -0.58905 | -0.18795 | 0.22570 | 0.62047 |
| RAND2 | -0.59430 | -0.18090 | 0.21697 | 0.61045 |
| RAND3 | -0.62008 | -0.20843 | 0.19894 | 0.59159 |
| RAND4 | -0.59696 | -0.18753 | 0.21087 | 0.61077 |
| RAND5 | -0.59819 | -0.21468 | 0.18130 | 0.56676 |
| RAND6 | -0.61150 | -0.21273 | 0.19102 | 0.59680 |

20 Bivariate Mutual Information / Uncertainty Reduction

| | | | | |
|---------|----------|----------|---------|---------|
| RAND7 | -0.61383 | -0.22039 | 0.18521 | 0.58843 |
| RAND8 | -0.59055 | -0.19032 | 0.20591 | 0.59859 |
| RAND9 | -0.60422 | -0.19932 | 0.20315 | 0.58792 |
| SUM12 | -0.67798 | -0.17129 | 0.22588 | 0.74242 |
| SUM34 | -0.73810 | -0.21209 | 0.21164 | 0.74363 |
| SUM1234 | -0.97362 | -0.27795 | 0.31417 | 1.00879 |

The marginal distributions are now shown.

If the data is continuous, the marginals will be nearly equal.

Widely unequal marginals indicate potentially problematic ties.

| | | | | | |
|---------|---------|---------|---------|---------|---------|
| RAND0 | 0.19987 | 0.20003 | 0.20003 | 0.20003 | 0.20003 |
| RAND1 | 0.19987 | 0.20003 | 0.20003 | 0.20003 | 0.20003 |
| RAND2 | 0.19987 | 0.20003 | 0.20003 | 0.20003 | 0.20003 |
| RAND3 | 0.19987 | 0.20003 | 0.20003 | 0.20003 | 0.20003 |
| RAND4 | 0.19987 | 0.20003 | 0.20003 | 0.20003 | 0.20003 |
| RAND5 | 0.19987 | 0.20003 | 0.20003 | 0.20003 | 0.20003 |
| RAND6 | 0.19987 | 0.20003 | 0.20003 | 0.20003 | 0.20003 |
| RAND7 | 0.19987 | 0.20003 | 0.20003 | 0.20003 | 0.20003 |
| RAND8 | 0.19987 | 0.20003 | 0.20003 | 0.20003 | 0.20003 |
| RAND9 | 0.19987 | 0.20003 | 0.20003 | 0.20003 | 0.20003 |
| SUM12 | 0.19987 | 0.20003 | 0.20003 | 0.20003 | 0.20003 |
| SUM34 | 0.19987 | 0.20003 | 0.20003 | 0.20003 | 0.20003 |
| SUM1234 | 0.19987 | 0.20003 | 0.20003 | 0.20003 | 0.20003 |

-----> Mutual Information <-----

| Predictor 1 | Predictor 2 | Target | MI | Solo pval | Unbiased pval |
|-------------|-------------|---------|--------|-----------|---------------|
| SUM12 | SUM34 | SUM1234 | 1.0781 | 0.0001 | 0.0001 |
| RAND1 | SUM34 | SUM1234 | 0.5363 | 0.0001 | 0.0001 |
| RAND3 | SUM12 | SUM1234 | 0.5356 | 0.0001 | 0.0001 |
| RAND2 | SUM34 | SUM1234 | 0.5333 | 0.0001 | 0.0001 |
| RAND4 | SUM12 | SUM1234 | 0.5242 | 0.0001 | 0.0001 |
| RAND3 | RAND4 | SUM1234 | 0.3094 | 0.0001 | 0.0001 |
| RAND3 | SUM34 | SUM1234 | 0.2994 | 0.0001 | 0.0001 |
| RAND4 | SUM34 | SUM1234 | 0.2985 | 0.0001 | 0.0001 |
| RAND6 | SUM34 | SUM1234 | 0.2947 | 0.0001 | 0.0001 |
| RAND9 | SUM34 | SUM1234 | 0.2946 | 0.0001 | 0.0001 |
| RAND8 | SUM34 | SUM1234 | 0.2944 | 0.0001 | 0.0001 |
| RAND5 | SUM34 | SUM1234 | 0.2939 | 0.0001 | 0.0001 |
| RAND0 | SUM34 | SUM1234 | 0.2937 | 0.0001 | 0.0001 |
| RAND7 | SUM34 | SUM1234 | 0.2925 | 0.0001 | 0.0001 |
| RAND2 | RAND3 | SUM1234 | 0.2881 | 0.0001 | 0.0001 |
| RAND1 | RAND3 | SUM1234 | 0.2879 | 0.0001 | 0.0001 |
| RAND1 | RAND4 | SUM1234 | 0.2861 | 0.0001 | 0.0001 |
| RAND2 | RAND4 | SUM1234 | 0.2811 | 0.0001 | 0.0001 |
| RAND1 | RAND2 | SUM1234 | 0.2755 | 0.0001 | 0.0001 |
| RAND2 | SUM12 | SUM1234 | 0.2709 | 0.0001 | 0.0001 |
| RAND1 | SUM12 | SUM1234 | 0.2705 | 0.0001 | 0.0001 |

Bivariate Mutual Information / Uncertainty Reduction 21

| | | | | | |
|-------|-------|---------|--------|--------|--------|
| RAND5 | SUM12 | SUM1234 | 0.2697 | 0.0001 | 0.0001 |
| RAND6 | SUM12 | SUM1234 | 0.2692 | 0.0001 | 0.0001 |
| RAND0 | SUM12 | SUM1234 | 0.2673 | 0.0001 | 0.0001 |
| RAND8 | SUM12 | SUM1234 | 0.2664 | 0.0001 | 0.0001 |
| RAND7 | SUM12 | SUM1234 | 0.2661 | 0.0001 | 0.0001 |
| RAND9 | SUM12 | SUM1234 | 0.2656 | 0.0001 | 0.0001 |
| RAND3 | RAND7 | SUM1234 | 0.1371 | 0.0001 | 0.0001 |
| RAND3 | RAND5 | SUM1234 | 0.1369 | 0.0001 | 0.0001 |
| RAND3 | RAND9 | SUM1234 | 0.1363 | 0.0001 | 0.0001 |
| RAND0 | RAND3 | SUM1234 | 0.1362 | 0.0001 | 0.0001 |
| RAND3 | RAND6 | SUM1234 | 0.1361 | 0.0001 | 0.0001 |
| RAND3 | RAND8 | SUM1234 | 0.1358 | 0.0001 | 0.0001 |
| RAND4 | RAND6 | SUM1234 | 0.1344 | 0.0001 | 0.0001 |
| RAND0 | RAND4 | SUM1234 | 0.1341 | 0.0001 | 0.0001 |
| RAND4 | RAND5 | SUM1234 | 0.1328 | 0.0001 | 0.0001 |
| RAND4 | RAND9 | SUM1234 | 0.1322 | 0.0001 | 0.0001 |
| RAND4 | RAND7 | SUM1234 | 0.1321 | 0.0001 | 0.0001 |
| RAND4 | RAND8 | SUM1234 | 0.1313 | 0.0001 | 0.0001 |
| RAND1 | RAND6 | SUM1234 | 0.1207 | 0.0001 | 0.0001 |
| RAND1 | RAND5 | SUM1234 | 0.1205 | 0.0001 | 0.0001 |
| RAND1 | RAND7 | SUM1234 | 0.1191 | 0.0001 | 0.0001 |
| RAND1 | RAND9 | SUM1234 | 0.1185 | 0.0001 | 0.0001 |
| RAND1 | RAND8 | SUM1234 | 0.1183 | 0.0001 | 0.0001 |
| RAND0 | RAND1 | SUM1234 | 0.1180 | 0.0001 | 0.0001 |
| RAND2 | RAND5 | SUM1234 | 0.1162 | 0.0001 | 0.0001 |
| RAND2 | RAND8 | SUM1234 | 0.1154 | 0.0001 | 0.0001 |
| RAND2 | RAND6 | SUM1234 | 0.1153 | 0.0001 | 0.0001 |
| RAND2 | RAND7 | SUM1234 | 0.1150 | 0.0001 | 0.0001 |
| RAND2 | RAND9 | SUM1234 | 0.1144 | 0.0001 | 0.0001 |
| RAND0 | RAND2 | SUM1234 | 0.1131 | 0.0001 | 0.0001 |
| RAND6 | RAND7 | SUM1234 | 0.0091 | 0.0952 | 0.9775 |
| RAND7 | RAND8 | SUM1234 | 0.0090 | 0.1081 | 0.9905 |
| RAND0 | RAND8 | SUM1234 | 0.0088 | 0.1563 | 0.9982 |
| RAND5 | RAND9 | SUM1234 | 0.0086 | 0.1904 | 0.9994 |
| RAND0 | RAND9 | SUM1234 | 0.0084 | 0.2327 | 0.9997 |
| RAND5 | RAND6 | SUM1234 | 0.0083 | 0.2549 | 0.9998 |
| RAND0 | RAND5 | SUM1234 | 0.0080 | 0.3693 | 1.0000 |
| RAND8 | RAND9 | SUM1234 | 0.0079 | 0.3949 | 1.0000 |
| RAND0 | RAND6 | SUM1234 | 0.0074 | 0.5647 | 1.0000 |
| RAND5 | RAND8 | SUM1234 | 0.0074 | 0.5734 | 1.0000 |
| RAND7 | RAND9 | SUM1234 | 0.0074 | 0.5830 | 1.0000 |
| RAND0 | RAND7 | SUM1234 | 0.0069 | 0.7550 | 1.0000 |
| RAND6 | RAND8 | SUM1234 | 0.0065 | 0.8598 | 1.0000 |
| RAND5 | RAND7 | SUM1234 | 0.0064 | 0.8652 | 1.0000 |
| RAND6 | RAND9 | SUM1234 | 0.0058 | 0.9657 | 1.0000 |

It should be no surprise that the best pair of predictors for SUM1234 are SUM12 and SUM34. Mutual information trails off according to how many components of the sum are present. Note the sharp transition in the unbiased p-value when we reach the point of having no component present!

22 Bivariate Mutual Information / Uncertainty Reduction

The next example shows what happens when worthless and serially correlated predictors are tested with a serially correlated target. We use DEP_RANDOM1 - DEP_RANDOM9 to predict DEP_RANDOM0, a situation which should demonstrate no predictive power whatsoever. The mutual information table is as follows:

-----> Mutual Information with DEP_RANDOM0 <-----

| Predictor 1 | Predictor 2 | Target | MI | Solo pval | Unbiased pval |
|-------------|-------------|-------------|--------|-----------|---------------|
| DEP_RANDOM2 | DEP_RANDOM7 | DEP_RANDOM0 | 0.0159 | 0.0001 | 0.0001 |
| DEP_RANDOM2 | DEP_RANDOM3 | DEP_RANDOM0 | 0.0145 | 0.0001 | 0.0001 |
| DEP_RANDOM2 | DEP_RANDOM9 | DEP_RANDOM0 | 0.0138 | 0.0001 | 0.0001 |
| DEP_RANDOM2 | DEP_RANDOM6 | DEP_RANDOM0 | 0.0132 | 0.0001 | 0.0005 |
| DEP_RANDOM4 | DEP_RANDOM8 | DEP_RANDOM0 | 0.0132 | 0.0001 | 0.0005 |
| DEP_RANDOM3 | DEP_RANDOM4 | DEP_RANDOM0 | 0.0132 | 0.0001 | 0.0005 |
| DEP_RANDOM2 | DEP_RANDOM4 | DEP_RANDOM0 | 0.0132 | 0.0001 | 0.0005 |
| DEP_RANDOM5 | DEP_RANDOM7 | DEP_RANDOM0 | 0.0131 | 0.0001 | 0.0005 |
| DEP_RANDOM1 | DEP_RANDOM2 | DEP_RANDOM0 | 0.0131 | 0.0001 | 0.0005 |
| DEP_RANDOM2 | DEP_RANDOM5 | DEP_RANDOM0 | 0.0129 | 0.0001 | 0.0011 |
| DEP_RANDOM2 | DEP_RANDOM8 | DEP_RANDOM0 | 0.0129 | 0.0001 | 0.0011 |
| DEP_RANDOM4 | DEP_RANDOM9 | DEP_RANDOM0 | 0.0127 | 0.0002 | 0.0016 |
| DEP_RANDOM1 | DEP_RANDOM3 | DEP_RANDOM0 | 0.0125 | 0.0001 | 0.0020 |
| DEP_RANDOM3 | DEP_RANDOM6 | DEP_RANDOM0 | 0.0125 | 0.0001 | 0.0022 |
| DEP_RANDOM1 | DEP_RANDOM5 | DEP_RANDOM0 | 0.0123 | 0.0001 | 0.0038 |
| DEP_RANDOM3 | DEP_RANDOM5 | DEP_RANDOM0 | 0.0122 | 0.0002 | 0.0056 |
| DEP_RANDOM6 | DEP_RANDOM8 | DEP_RANDOM0 | 0.0121 | 0.0003 | 0.0074 |
| DEP_RANDOM1 | DEP_RANDOM6 | DEP_RANDOM0 | 0.0117 | 0.0010 | 0.0213 |
| DEP_RANDOM6 | DEP_RANDOM9 | DEP_RANDOM0 | 0.0115 | 0.0006 | 0.0323 |
| DEP_RANDOM4 | DEP_RANDOM6 | DEP_RANDOM0 | 0.0110 | 0.0021 | 0.0893 |
| DEP_RANDOM1 | DEP_RANDOM4 | DEP_RANDOM0 | 0.0110 | 0.0027 | 0.0904 |
| DEP_RANDOM5 | DEP_RANDOM8 | DEP_RANDOM0 | 0.0110 | 0.0032 | 0.0906 |
| DEP_RANDOM5 | DEP_RANDOM9 | DEP_RANDOM0 | 0.0108 | 0.0044 | 0.1298 |
| DEP_RANDOM7 | DEP_RANDOM9 | DEP_RANDOM0 | 0.0108 | 0.0051 | 0.1442 |
| DEP_RANDOM7 | DEP_RANDOM8 | DEP_RANDOM0 | 0.0107 | 0.0060 | 0.1584 |
| DEP_RANDOM4 | DEP_RANDOM5 | DEP_RANDOM0 | 0.0107 | 0.0063 | 0.1610 |
| DEP_RANDOM3 | DEP_RANDOM9 | DEP_RANDOM0 | 0.0107 | 0.0051 | 0.1620 |
| DEP_RANDOM1 | DEP_RANDOM9 | DEP_RANDOM0 | 0.0104 | 0.0096 | 0.2819 |
| DEP_RANDOM6 | DEP_RANDOM7 | DEP_RANDOM0 | 0.0103 | 0.0132 | 0.3179 |
| DEP_RANDOM8 | DEP_RANDOM9 | DEP_RANDOM0 | 0.0102 | 0.0147 | 0.3827 |
| DEP_RANDOM3 | DEP_RANDOM7 | DEP_RANDOM0 | 0.0101 | 0.0181 | 0.4380 |
| DEP_RANDOM5 | DEP_RANDOM6 | DEP_RANDOM0 | 0.0099 | 0.0249 | 0.5409 |
| DEP_RANDOM1 | DEP_RANDOM8 | DEP_RANDOM0 | 0.0098 | 0.0294 | 0.5901 |
| DEP_RANDOM3 | DEP_RANDOM8 | DEP_RANDOM0 | 0.0097 | 0.0347 | 0.6486 |
| DEP_RANDOM4 | DEP_RANDOM7 | DEP_RANDOM0 | 0.0087 | 0.1757 | 0.9908 |
| DEP_RANDOM1 | DEP_RANDOM7 | DEP_RANDOM0 | 0.0084 | 0.2498 | 0.9983 |

Notice how many truly worthless predictive pairs have tiny p-values, even in the unbiased case. This is a severe problem that affects *all* common statistical tests, not just Monte-Carlo Permutation Tests.

Bivariate Mutual Information / Uncertainty Reduction 23

The final example shows how the cyclic modification of the Monte-Carlo Permutation Test can at least partially remedy the situation. We repeat the same test as that just shown, except that instead of using *Complete* permutation we use *Cyclic* permutation. The results are shown below:

```
-----> Mutual Information with DEP_RANDOM0 <-----
```

| Predictor 1 | Predictor 2 | Target | MI | Solo pval | Unbiased pval |
|-------------|-------------|-------------|--------|-----------|---------------|
| DEP_RANDOM2 | DEP_RANDOM7 | DEP_RANDOM0 | 0.0159 | 0.0261 | 0.4007 |
| DEP_RANDOM2 | DEP_RANDOM3 | DEP_RANDOM0 | 0.0145 | 0.0813 | 0.8015 |
| DEP_RANDOM2 | DEP_RANDOM9 | DEP_RANDOM0 | 0.0138 | 0.1404 | 0.9240 |
| DEP_RANDOM2 | DEP_RANDOM6 | DEP_RANDOM0 | 0.0132 | 0.1968 | 0.9761 |
| DEP_RANDOM4 | DEP_RANDOM8 | DEP_RANDOM0 | 0.0132 | 0.1660 | 0.9776 |
| DEP_RANDOM3 | DEP_RANDOM4 | DEP_RANDOM0 | 0.0132 | 0.1859 | 0.9792 |
| DEP_RANDOM2 | DEP_RANDOM4 | DEP_RANDOM0 | 0.0132 | 0.1768 | 0.9804 |
| DEP_RANDOM5 | DEP_RANDOM7 | DEP_RANDOM0 | 0.0131 | 0.2354 | 0.9837 |
| DEP_RANDOM1 | DEP_RANDOM2 | DEP_RANDOM0 | 0.0131 | 0.2077 | 0.9858 |
| DEP_RANDOM2 | DEP_RANDOM5 | DEP_RANDOM0 | 0.0129 | 0.2329 | 0.9915 |
| DEP_RANDOM2 | DEP_RANDOM8 | DEP_RANDOM0 | 0.0129 | 0.2162 | 0.9925 |
| DEP_RANDOM4 | DEP_RANDOM9 | DEP_RANDOM0 | 0.0127 | 0.2594 | 0.9949 |
| DEP_RANDOM1 | DEP_RANDOM3 | DEP_RANDOM0 | 0.0125 | 0.3104 | 0.9972 |
| DEP_RANDOM3 | DEP_RANDOM6 | DEP_RANDOM0 | 0.0125 | 0.3243 | 0.9977 |
| DEP_RANDOM1 | DEP_RANDOM5 | DEP_RANDOM0 | 0.0123 | 0.3545 | 0.9978 |
| DEP_RANDOM3 | DEP_RANDOM5 | DEP_RANDOM0 | 0.0122 | 0.3621 | 0.9982 |
| DEP_RANDOM6 | DEP_RANDOM8 | DEP_RANDOM0 | 0.0121 | 0.3613 | 0.9984 |
| DEP_RANDOM1 | DEP_RANDOM6 | DEP_RANDOM0 | 0.0117 | 0.4874 | 0.9998 |
| DEP_RANDOM6 | DEP_RANDOM9 | DEP_RANDOM0 | 0.0115 | 0.5108 | 0.9998 |
| DEP_RANDOM4 | DEP_RANDOM6 | DEP_RANDOM0 | 0.0110 | 0.6064 | 1.0000 |
| DEP_RANDOM1 | DEP_RANDOM4 | DEP_RANDOM0 | 0.0110 | 0.5907 | 1.0000 |
| DEP_RANDOM5 | DEP_RANDOM8 | DEP_RANDOM0 | 0.0110 | 0.5737 | 1.0000 |
| DEP_RANDOM5 | DEP_RANDOM9 | DEP_RANDOM0 | 0.0108 | 0.6308 | 1.0000 |
| DEP_RANDOM7 | DEP_RANDOM9 | DEP_RANDOM0 | 0.0108 | 0.6902 | 1.0000 |
| DEP_RANDOM7 | DEP_RANDOM8 | DEP_RANDOM0 | 0.0107 | 0.6681 | 1.0000 |
| DEP_RANDOM4 | DEP_RANDOM5 | DEP_RANDOM0 | 0.0107 | 0.6274 | 1.0000 |
| DEP_RANDOM3 | DEP_RANDOM9 | DEP_RANDOM0 | 0.0107 | 0.6552 | 1.0000 |
| DEP_RANDOM1 | DEP_RANDOM9 | DEP_RANDOM0 | 0.0104 | 0.7349 | 1.0000 |
| DEP_RANDOM6 | DEP_RANDOM7 | DEP_RANDOM0 | 0.0103 | 0.7587 | 1.0000 |
| DEP_RANDOM8 | DEP_RANDOM9 | DEP_RANDOM0 | 0.0102 | 0.7330 | 1.0000 |
| DEP_RANDOM3 | DEP_RANDOM7 | DEP_RANDOM0 | 0.0101 | 0.7944 | 1.0000 |
| DEP_RANDOM5 | DEP_RANDOM6 | DEP_RANDOM0 | 0.0099 | 0.8103 | 1.0000 |
| DEP_RANDOM1 | DEP_RANDOM8 | DEP_RANDOM0 | 0.0098 | 0.8036 | 1.0000 |
| DEP_RANDOM3 | DEP_RANDOM8 | DEP_RANDOM0 | 0.0097 | 0.8085 | 1.0000 |
| DEP_RANDOM4 | DEP_RANDOM7 | DEP_RANDOM0 | 0.0087 | 0.9581 | 1.0000 |
| DEP_RANDOM1 | DEP_RANDOM7 | DEP_RANDOM0 | 0.0084 | 0.9731 | 1.0000 |

This time, the unbiased p-values are not fooled at all by the serial correlation, and even the solo p-values behave well.

24 Predictors having Max Relevance, Min Redundancy

Predictors having Max Relevance, Min Redundancy

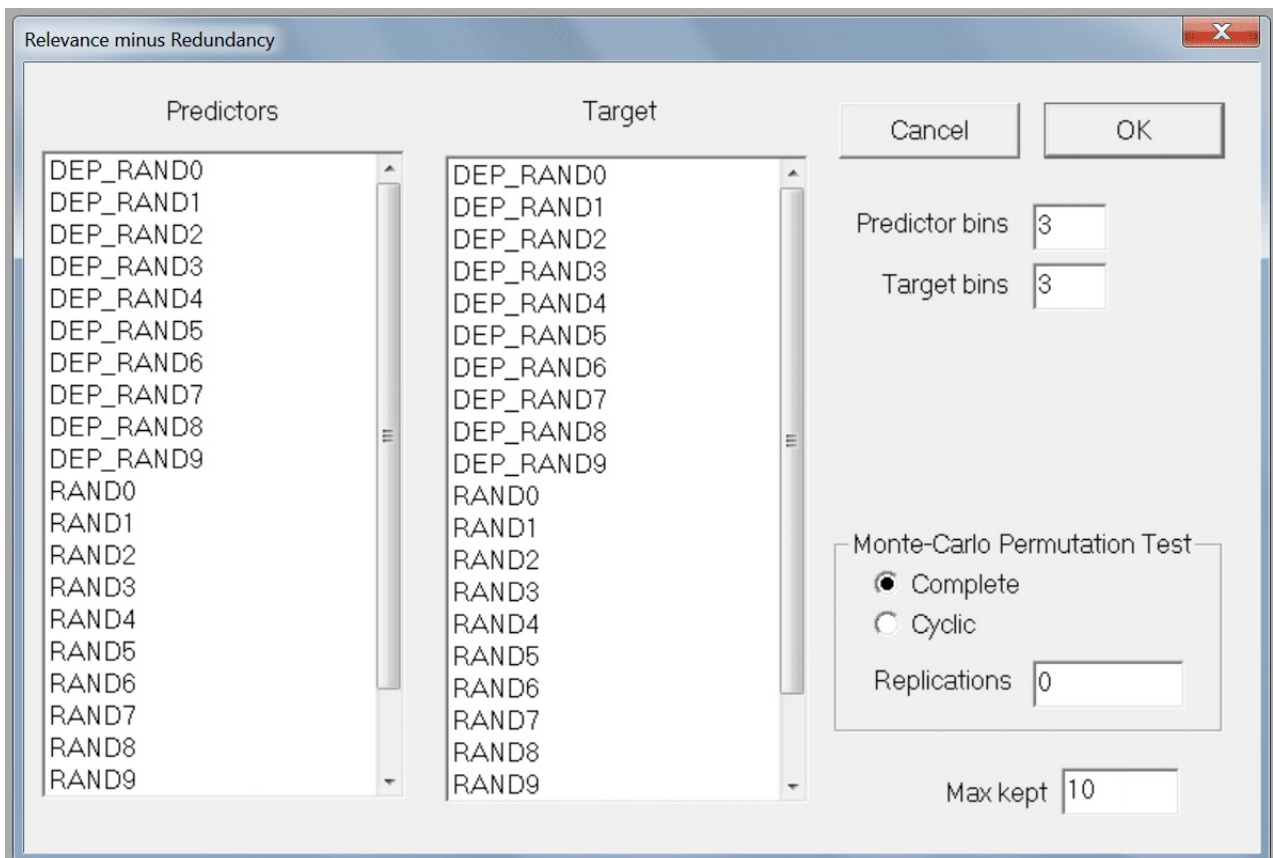
Selection of predictors by examining individual or even pairwise performance is useful for quickly identifying the most promising candidates. However, this simplistic approach suffers from redundancy. If two predictor candidates are nicely related to a target, chances are good that they are also closely related to each other; they may provide similar if not identical predictive information. Thus, if one examines a large number of candidates and chooses a subset of predictors that are all good at predicting the target, this subset will in most cases be unnecessarily large; many of them will provide nearly or exactly the same predictive information as other candidates in the subset. A much more efficient approach to selecting a good subset of predictor candidates would be to consider not only the relevance of the members at predicting the target, but also their redundancy with other members of the subset.

Peng, Long and Ding (2005) provide a fabulous algorithm for handling this redundancy problem in their paper “Feature Selection Based on Mutual Information: Criteria of Max-Dependency, Max-Relevance, and Min Redundancy”. An intuitive summary of the algorithm, along with C++ code, appears in my book “Assessing and Improving Prediction and Classification,” so details will be omitted here. However, it must be stressed that this algorithm has a powerful optimality property: suppose one were to consider the mutual information between a set of predictors (taken as a group) and a target. This is called *joint dependency*. A reasonable method for choosing an optimal subset of predictors is to use forward stepwise selection to maximize the joint dependency of the subset with the target. Unfortunately, this quantity is difficult if not impossible to compute in practical applications. But the Pen, Long, and Ding algorithm is an elegant work-around that produces the same subset of predictors as stepwise selection based on maximizing joint dependency, but it does so in a computationally feasible way.

At each step, the algorithm considers the relevance of a candidate for predicting the target, as well as the redundancy of the candidate with predictors already in the chosen subset. These quantities are subtracted to provide a selection criterion. The candidate with the maximum relevance-minus-redundancy criterion is chosen.

Specifying the Test Parameters

When the user clicks *Tests / Relevance minus Redundancy*, a dialog similar to that shown will appear. The various parameters are described below.



The leftmost column is used to specify the set of predictor candidates. Multiple candidates can be selected by dragging the mouse cursor across a block, or by clicking the first candidate in a block, holding the Shift key, and clicking the last candidate in the block. Individual candidates can be toggled on and off by holding the Ctrl key while clicking on the variable.

The *Target* column is used to select the target variable.

The predictors and the target are partitioned into bins that are as equal in size as possible. The user must specify the number of bins to employ for each, and unless the dataset is huge the default of three bins is frequently appropriate.

26 Predictors having Max Relevance, Min Redundancy

Replications defaults to zero, in which case no Monte-Carlo Permutation Test is performed. However, it is usually best to set this to at least 100, and perhaps as much as 1000, so that solo and group p-values will be computed. Note that the minimum possible p-value is the reciprocal of the number of permutations. So, for example, if the user specifies 100 permutations, the minimum p-value that can appear is 0.01. Run time of this test is linearly related to the number of permutations.

The user must choose either *Complete* or *Cyclic* permutations. If the user is confident that there is no dependency as described earlier in this document, then *Complete* should be used; it is the traditional approach which does a complete random shuffle for each permutation. However, if there is dependency, this type of shuffling will produce underestimation of *p-values*, a very dangerous situation. If the dependency is serial (the data is a time series and the dependency is among samples close in time) then a considerable improvement in the situation can be obtained by using *Cyclic* permutation. In this type of shuffle, the time order of the target is kept intact except at the ends by rotating the target with end-point wraparound. Shuffling this way preserves most of the serial dependency in the permuted target, which makes the algorithm more accurate. The *p-values* computed this way will generally be larger than those computed with complete shuffling, and hence less likely to lead to false rejection of the null hypothesis of no predictive power. But be warned that the cure is far from complete; computed *p-values* will still underestimate the true values, just not as badly.

Note that in most cases it is legitimate to use *Cyclic* permutation instead of *Complete* when there is no dependency. However, if the dataset is small, *Cyclic* permutation will limit the number of unique permutations and hence increase the random error inherent in the process. As long as the dataset is large, some users may prefer to use *Cyclic* permutation even if it is assumed that there is no serial dependency; in case there really is hidden serial dependency, this is a cheap insurance policy. Still, the best practice is to make sure that the data does not contain dependency and then use *Complete* permutation. Relying on *Cyclic* permutation to take care of dependency problems is living dangerously. And if the dataset contains fewer than 1000 or so cases, use of *Cyclic* permutation is not recommended unless it is necessary to handle dependency.

Max kept is the maximum size of the selected subset. Execution time is approximately linearly related to this quantity, so it should be kept as small as possible if run time is critical.

Note that this algorithm employs CUDA processing if available. However, unless there are many hundreds of predictor candidates, its overhead may actually slow execution.

An Example of Relevance Minus Redundancy

This section demonstrates a revealing example of the algorithm using synthetic data to clarify the presentation. The variables in the dataset are as follows:

RAND0 - RAND9 are independent (within themselves and with each other) random time series.

$$SUM12 = RAND1 + RAND2$$

$$SUM34 = RAND3 + RAND4$$

$$SUM1234 = SUM12 + SUM34$$

The test run attempts to predict SUM1234 from RAND0 - RAND9, SUM12, and SUM34. The output is shown below. Brief explanatory comments are interspersed.

```
*****
*
* Computing relevance minus redundancy for optimal predictor subset *
*      12 predictor candidates                                     *
*      12 best predictors will be printed                         *
*      5 predictor bins                                           *
*      5 target bins                                              *
*      100 replications of complete Monte-Carlo Permutation Test *
*
*****
```

Initial candidates, in order of decreasing mutual information with SUM1234

| Variable | MI |
|----------|--------|
| SUM34 | 0.2877 |
| SUM12 | 0.2610 |
| RAND3 | 0.1307 |
| RAND4 | 0.1263 |
| RAND1 | 0.1129 |
| RAND2 | 0.1085 |
| RAND8 | 0.0015 |
| RAND5 | 0.0014 |
| RAND6 | 0.0012 |
| RAND7 | 0.0010 |
| RAND0 | 0.0008 |
| RAND9 | 0.0006 |

28 Predictors having Max Relevance, Min Redundancy

| Predictors so far | Relevance | Redundancy | Criterion |
|-------------------|-----------|------------|-----------|
| SUM34 | 0.2877 | 0.0000 | 0.2877 |

We see from the table above that the first candidate chosen is the one which has maximum mutual information with the target. Naturally this would be either SUM12 or SUM34, and it happens to be the latter. Then, in the table below we see that SUM12 has the largest relevance (its mutual information with the target) and essentially no redundancy with SUM34 (again, no surprise). This gives it the highest selection criterion and it is chosen.

Additional candidates, in order of decreasing relevance minus redundancy

| Variable | Relevance | Redundancy | Criterion |
|----------|-----------|------------|-----------|
| SUM12 | 0.2610 | 0.0014 | 0.2596 |
| RAND1 | 0.1129 | 0.0016 | 0.1112 |
| RAND2 | 0.1085 | 0.0009 | 0.1076 |
| RAND6 | 0.0012 | 0.0007 | 0.0005 |
| RAND0 | 0.0008 | 0.0009 | -0.0000 |
| RAND8 | 0.0015 | 0.0017 | -0.0002 |
| RAND5 | 0.0014 | 0.0016 | -0.0002 |
| RAND9 | 0.0006 | 0.0008 | -0.0002 |
| RAND7 | 0.0010 | 0.0012 | -0.0003 |
| RAND3 | 0.1307 | 0.3154 | -0.1847 |
| RAND4 | 0.1263 | 0.3158 | -0.1895 |

| Predictors so far | Relevance | Redundancy | Criterion |
|-------------------|-----------|------------|-----------|
| SUM34 | 0.2877 | 0.0000 | 0.2877 |
| SUM12 | 0.2610 | 0.0014 | 0.2596 |

Now we come to an important observation. One might think that the next candidate selected would be either RAND1, RAND2, RAND3, or RAND4, the four components of the SUM1234 target. However, the table on the next page shows that these four candidates actually fall at the bottom of the list! This is because they have so much redundancy with SUM12 and SUM34 (taken as a group) that they will not be chosen next. In fact, RAND6, which has no relationship whatsoever with any of the other variables, is chosen based only on its tiny random relevance and slightly smaller random redundancy.

Predictors having Max Relevance, Min Redundancy 29

Additional candidates, in order of decreasing relevance minus redundancy

| Variable | Relevance | Redundancy | Criterion |
|----------|-----------|------------|-----------|
| RAND6 | 0.0012 | 0.0009 | 0.0003 |
| RAND0 | 0.0008 | 0.0008 | 0.0000 |
| RAND8 | 0.0015 | 0.0015 | 0.0000 |
| RAND9 | 0.0006 | 0.0008 | -0.0002 |
| RAND5 | 0.0014 | 0.0017 | -0.0003 |
| RAND7 | 0.0010 | 0.0013 | -0.0004 |
| RAND3 | 0.1307 | 0.1581 | -0.0274 |
| RAND4 | 0.1263 | 0.1585 | -0.0322 |
| RAND1 | 0.1129 | 0.1527 | -0.0398 |
| RAND2 | 0.1085 | 0.1485 | -0.0399 |

| Predictors so far | Relevance | Redundancy | Criterion |
|-------------------|-----------|------------|-----------|
| SUM34 | 0.2877 | 0.0000 | 0.2877 |
| SUM12 | 0.2610 | 0.0014 | 0.2596 |
| RAND6 | 0.0012 | 0.0009 | 0.0003 |

But now that the selected set's redundancy with the remaining candidates has been 'diluted' by the inclusion of the unrelated RAND6, RAND1-RAND4 jump to the top of the list due to their relatively large relevance but lessened redundancy.

Additional candidates, in order of decreasing relevance minus redundancy

| Variable | Relevance | Redundancy | Criterion |
|----------|-----------|------------|-----------|
| RAND3 | 0.1307 | 0.1058 | 0.0249 |
| RAND4 | 0.1263 | 0.1061 | 0.0202 |
| RAND1 | 0.1129 | 0.1021 | 0.0107 |
| RAND2 | 0.1085 | 0.0995 | 0.0090 |
| RAND0 | 0.0008 | 0.0010 | -0.0002 |
| RAND9 | 0.0006 | 0.0009 | -0.0003 |
| RAND5 | 0.0014 | 0.0017 | -0.0003 |
| RAND8 | 0.0015 | 0.0018 | -0.0004 |
| RAND7 | 0.0010 | 0.0015 | -0.0006 |

| Predictors so far | Relevance | Redundancy | Criterion |
|-------------------|-----------|------------|-----------|
| SUM34 | 0.2877 | 0.0000 | 0.2877 |
| SUM12 | 0.2610 | 0.0014 | 0.2596 |
| RAND6 | 0.0012 | 0.0009 | 0.0003 |
| RAND3 | 0.1307 | 0.1058 | 0.0249 |

30 Predictors having Max Relevance, Min Redundancy

There is little point in continuing to show the inclusion steps. We now jump to the final table that lists all candidates in the order in which they were selected, along with associated p-values.

-----> Final results predicting SUM1234 <-----

| Final predictors | Relevance | Redundancy | Criterion | Solo pval | Group pval |
|------------------|-----------|------------|-----------|-----------|------------|
| SUM34 | 0.2877 | 0.0000 | 0.2877 | 0.010 | 0.010 |
| SUM12 | 0.2610 | 0.0014 | 0.2596 | 0.010 | 0.010 |
| RAND6 | 0.0012 | 0.0009 | 0.0003 | 0.570 | 0.010 |
| RAND3 | 0.1307 | 0.1058 | 0.0249 | 0.010 | 0.010 |
| RAND4 | 0.1263 | 0.0797 | 0.0465 | 0.010 | 0.010 |
| RAND1 | 0.1129 | 0.0617 | 0.0511 | 0.010 | 0.010 |
| RAND2 | 0.1085 | 0.0505 | 0.0581 | 0.010 | 0.010 |
| RAND8 | 0.0015 | 0.0014 | 0.0001 | 0.320 | 0.010 |
| RAND5 | 0.0014 | 0.0014 | -0.0001 | 0.340 | 0.010 |
| RAND7 | 0.0010 | 0.0014 | -0.0004 | 0.650 | 0.010 |
| RAND0 | 0.0008 | 0.0013 | -0.0004 | 0.850 | 0.010 |
| RAND9 | 0.0006 | 0.0012 | -0.0006 | 0.980 | 0.010 |

Two different p-values are printed for each predictor candidate. The *Solo pval* is the same quantity printed in the Univariate test. This is the probability that, if the predictor has no actual mutual information with the target, a mutual information (Relevance here) as large as that obtained could have occurred. Understand that this quantity considers each candidate in isolation, not involving any other candidates. Note how nicely this reveals the uselessness of the third candidate chosen, RAND6.

The *Group pval* considers the associated candidate along with every prior candidate. It tests the null hypothesis that the group of candidates selected so far, on average, has no mutual information with the target.

Regrettably, I am not aware of any way of computing what would be an especially useful p-value, that which tests the null hypothesis that selecting the candidate provides no additional (non-redundant) relevance. Such a p-value would be valuable for determining when to stop including additional candidates in the selected subset. The problem appears to be that the test statistic at any step is strongly dependent on the relevance of those predictors already selected. If anyone knows of a way around this problem, I would love to hear about it.

Hidden Markov Models with Target Correlation

When working with time series data, the developer need not assume a direct relationship between predictors and a target. Sometimes it is better to posit an underlying condition, the *state* of the process under study, which impacts both the predictors and the target. This process is assumed to exist at all times in exactly one of two or more possible states. The state at any given time impacts the distribution of associated variables. Some of these variables may be observable at the present time (predictors), while others may be unknown at the present time but be of great interest (targets). Our goal is to use measured values of the observable variables to determine (or make an educated guess at) the state of the process, and then use this knowledge to estimate the value of an unobservable variable which interests us.

It is vital to distinguish this application from ordinary classification methods which are not restricted to time series data. In simple classification, one measures some predictor variable(s) and makes a class decision, which in turn may imply likely values of other (probably unmeasurable) variables. But a hidden Markov model assumes a sequential process with an important property: the probability of being in a given state at an observed time depends on the process's state at the prior observed time. In other words, a hidden Markov model has memory, while ordinary classification does not.

This memory is immensely useful in some applications. For example, it may prevent whipsaws. Suppose a certain state tends to be persistent in real life. Ordinary classification will suffer if there is large random noise in the observed variables, which may snap the decision back and forth at the whim of chance. But the memory inherent in a hidden Markov model will tend to hold its decision in a persistent state even as noise in the measured variables tries to whip the decision back and forth. Of course, the downside of this memory is a tendency toward delayed decisions; the model may need several observed values to confirm a state change. But this is often a price well worth paying, especially in high-noise situations.

One application of a hidden Markov model is the prediction of a financial market. Perhaps the developer assumes that it is always in either a bull market (a long-term up-trend), a bear market (a long-term down-trend) or a flat market (no long-term trend). By definition, bull and bear markets cover an extended time period; one does not go from a bull to a bear market in one day, and then return to a bull market the next day. Such direction changes are just short-term fluctuations in a more extensive move. If one were to use frequent observations to make daily predictions of whether the market is in a bull

32 Hidden Markov Models with Target Correlation

or bear state, these decisions could reverse ridiculously often. One is better off taking advantage of the memory of a hidden Markov model to stabilize behavior.

Specifying the Test Parameters

When the user clicks *Tests / Hidden Markov model*, a dialog similar to that shown will appear. The various parameters are described below.

Hidden Markov Model correlation with a target

| Predictors | Target |
|-------------|-------------|
| DEP_RANDOM0 | DEP_RANDOM0 |
| DEP_RANDOM1 | DEP_RANDOM1 |
| DEP_RANDOM2 | DEP_RANDOM2 |
| DEP_RANDOM3 | DEP_RANDOM3 |
| DEP_RANDOM4 | DEP_RANDOM4 |
| DEP_RANDOM5 | DEP_RANDOM5 |
| DEP_RANDOM6 | DEP_RANDOM6 |
| DEP_RANDOM7 | DEP_RANDOM7 |
| DEP_RANDOM8 | DEP_RANDOM8 |
| DEP_RANDOM9 | DEP_RANDOM9 |
| RAND0 | RAND0 |
| RAND1 | RAND1 |
| RAND2 | RAND2 |
| RAND3 | RAND3 |
| RAND4 | RAND4 |
| RAND5 | RAND5 |
| RAND6 | RAND6 |
| RAND7 | RAND7 |
| RAND8 | RAND8 |
| RAND9 | RAND9 |

Cancel OK

Dimension (1-3) 1

Number of states 2

Monte-Carlo Permutation Test

☒ Complete

☐ Cyclic

Replications 0

Max printed 100

The leftmost column is used to specify the set of predictor candidates. Multiple candidates can be selected by dragging the mouse cursor across a block, or by clicking the first candidate in a block, holding the Shift key, and clicking the last candidate in the block. Individual candidates can be toggled on and off by holding the Ctrl key while clicking on the variable.

The *Target* column is used to specify the target variable. This variable is ignored when the models are computed; rather it plays a role in selecting the 'best' model.

The *Dimension* must be 1, 2, or 3. This is the number of predictor variables that will be used by the hidden Markov model.

The *Number of states* is exactly that, the number of states in which the process can exist. It must be at least two, and it typically is small, rarely more than four. Execution time blows up rapidly as the number of states increases.

The user must choose either *Complete* or *Cyclic* permutations and the number of replications to perform. Please refer to the discussions of this issue earlier in this document. However, because hidden Markov models virtually always are applied to serially correlated data, cyclic permutation is the default.

Max printed is the maximum number of models printed in the log file.

WARNING... This test can be extremely slow. While threads are being initialized for the first set of models, the ESCape key is ignored. After that, ESCape is polled only at widely spaced intervals. Then, when waiting for the final threads to complete, ESCape is again ignored. For a few thousand cases, 2 dimensions, and 2 states, the complete test should run in a few minutes or less on modern computers. But if there are many thousands of cases, 3 dimensions, and 4 or more states, the test could require several hours to complete. If you get in over your head, you may need to use Task Manger to force a shutdown of the program. Sorry about that, but as of yet I have not been able to figure out an efficient way to interrupt threads that are in the middle of extensive computation without inducing significant overhead, which just makes the situation worse.

Operation of This Test

The *Hidden Markov Model* test operates in two completely separate steps. In the first step, every possible combination of predictor candidates is used to fit a hidden Markov model. Let N be the number of candidates specified by the user (selected from the list in the left column of the dialog). If the dimension is specified to be 1, then each candidate is used alone, resulting in N models, one for each candidate. If the dimension is 2, then there are $N*(N-1)/2$ models, one for each possible pair of candidates. If the dimension is 3, then there are $N*(N-1)*(N-2)/6$ models, one for each possible trio. It must be emphasized that these models are optimized without regard to the target variable; the target plays no role whatsoever in the development of the models.

After this (potentially large!) set of hidden Markov models has been found, the relationship between each of them and the user-specified target variable is found. The relationship between a model and the target is defined as the multiple-R (the multivariate correlation coefficient) between the vector of state probabilities and the target. In other words, for a given model, each case will have associated with it a vector giving the probability that this observation is in each possible state. These state probability vectors are regressed on the target variable using ordinary multiple linear regression.

Details of the best (most highly correlated) model are printed. Then the models (up to *Max printed* of them) are listed in descending order of relationship with the target. The multiple-R is printed for each. If Monte-Carlo replications were specified, solo and unbiased p-values are printed for each model. The *solo p-value* is the probability that, if there were actually no relationship between the state (as defined by that model) and the target, we could have obtained a multiple-R at least as large as we did obtain. The *unbiased p-value* for the best model is the probability that if *none* of the models were related to the target, the best among them would have a multiple-R at least as large as that obtained. Subsequent unbiased p-values are upper bounds on similarly defined probabilities. This issue is discussed in detail earlier in this document.

Note that exact results will not in general be replicated if runs are repeated. This is because training a hidden Markov model relies on random number generation, and Windows' scheduling of training threads is rarely consistent. The competing models will receive their random numbers in different orders during different runs, resulting in slightly different solutions being obtained. In rare cases, a 'satisfactory' solution will not be obtained at all. But the probability of this happening depends on how well the data is explained by a hidden Markov model. Data which is almost entirely random noise will have the highest probability of leading to disappointing or unstable models.

An Example of Hidden Markov Models

This section demonstrates a revealing example of the algorithm using synthetic data to clarify the presentation. The variables in the dataset are as follows:

RAND0 - RAND9 are independent (within themselves and with each other) random time series. These are the predictor candidates.

SUM12 = RAND1 + RAND2. This is the target variable.

I chose to use two predictors and allow four states in the models. The program fits a hidden Markov model to each of the $(10-9)/2=45$ pairs of predictor candidates. Not surprisingly, the model based on *RAND1* and *RAND2* has the highest correlation with *SUM12*. Its means and standard deviations for each state are printed first:

Means (top number) and standard deviations (bottom number)

| State | RAND1 | RAND2 |
|-------|---------------------|---------------------|
| 1 | 0.06834 0.48729 | -0.66014 0.21358 |
| 2 | -0.73466 0.17187 | 0.07687 0.54038 |
| 3 | -0.02272 0.39033 | 0.35902 0.39555 |
| 4 | 0.73542 0.17546 | 0.08884 0.52133 |

RAND1 and *RAND2* are totally random (they exist in only one state), so attempting to fit a hidden Markov model to them should be extremely unstable. Indeed, in ten runs of this test, twice the program found solutions in which the means of the states were all nearly zero, indicating no differentiation between states. But most of the time it came up with a pattern essentially identical to the one shown above. This solution is remarkably similar to a sort of principal components decomposition: *RAND1* distinguishes between State 2 and State 4, while *RAND2* distinguishes between State 1 and State 3. Thus, knowledge of which of the four states the process is in provides great information about *SUM12*.

36 Hidden Markov Models with Target Correlation

Next we see the transition probabilities. The figure in Row i and Column j is the probability that the process will transition from State i to State j . Not surprisingly, they are almost all identical. The relatively small discrepancies are just due to random variation in the data.

Transition probabilities...

| | 1 | 2 | 3 | 4 |
|---|--------|--------|--------|--------|
| 1 | 0.2638 | 0.2037 | 0.3494 | 0.1830 |
| 2 | 0.2438 | 0.1945 | 0.3638 | 0.1979 |
| 3 | 0.2130 | 0.1682 | 0.4174 | 0.2014 |
| 4 | 0.2404 | 0.2148 | 0.3272 | 0.2176 |

Further properties of each state are then printed:

Percent is the percentage of cases in which this state has the highest probability. The sum of these quantities across all states may not reach 100 percent, because cases in which there is a tie for the highest probability are not counted. If the data is continuous, this should almost never happen.

Correlation is the ordinary correlation coefficient between the target and the membership probability for this state. On first consideration it might be thought that the beta weight in the linear equation predicting the target from the state probabilities would be the better quantity to print. But the beta weight is not printed at all due to the fact that such weights are notoriously unstable and hence uninformative. Suppose there is very high correlation between the membership probabilities of two states, a situation which is especially likely to happen if the user specifies more states than actually exist in the process. Then both of these probabilities could be highly correlated with the target, while they might actually have opposite signs for their beta weights!

Target mean is the mean of the target when this state has the highest membership probability. Cases in which there is a tie for maximum (almost impossible for continuous data) do not enter into this calculation.

Target StdDev is the standard deviation of the target when this state has the highest membership probability. Cases in which there is a tie for maximum (almost impossible for continuous data) do not enter into this calculation.

| State | Percent | Correlation | Target mean | Target StdDev |
|-------|---------|-------------|-------------|---------------|
| 1 | 23.76 | -0.53350 | -0.54538 | 0.45473 |
| 2 | 21.73 | -0.52368 | -0.71809 | 0.56342 |
| 3 | 34.03 | 0.38210 | 0.35674 | 0.47747 |
| 4 | 20.48 | 0.62840 | 0.92173 | 0.49069 |

The reader should look back at the table of RAND1 and RAND2 means for each of the four states and confirm that the correlations and target means shown in the table above make sense. We also see that the state membership probabilities conform with the transition matrix. As expected for random series, the target standard deviations are all about the same.

Last but not least is the list of models, sorted in descending order of their multiple-R with the target. As expected (or at least hoped), the models involving either RAND1 or RAND2 appear first, and they are all extremely significant. As soon as these two variables are exhausted, multiple-R plunges and significance is lost. The remainder of this table is not shown here, but this situation continues.

-----> Hidden Markov Models correlating with SUM12 <-----

| Predictor 1 | Predictor 2 | Multiple-R | Solo pval | Unbiased pval |
|-------------|-------------|------------|-----------|---------------|
| RAND1 | RAND2 | 0.8896 | 0.0010 | 0.0010 |
| RAND1 | RAND3 | 0.6937 | 0.0010 | 0.0010 |
| RAND1 | RAND5 | 0.6680 | 0.0010 | 0.0010 |
| RAND0 | RAND1 | 0.6619 | 0.0010 | 0.0010 |
| RAND1 | RAND9 | 0.6604 | 0.0010 | 0.0010 |
| RAND1 | RAND8 | 0.6590 | 0.0010 | 0.0010 |
| RAND2 | RAND5 | 0.6579 | 0.0010 | 0.0010 |
| RAND0 | RAND2 | 0.6554 | 0.0010 | 0.0010 |
| RAND2 | RAND9 | 0.6493 | 0.0010 | 0.0010 |
| RAND1 | RAND7 | 0.5870 | 0.0010 | 0.0010 |
| RAND1 | RAND4 | 0.5845 | 0.0010 | 0.0010 |
| RAND2 | RAND4 | 0.5756 | 0.0010 | 0.0010 |
| RAND2 | RAND3 | 0.5721 | 0.0010 | 0.0010 |
| RAND2 | RAND7 | 0.5667 | 0.0010 | 0.0010 |
| RAND2 | RAND6 | 0.5648 | 0.0010 | 0.0010 |
| RAND2 | RAND8 | 0.5623 | 0.0010 | 0.0010 |
| RAND1 | RAND6 | 0.3938 | 0.0010 | 0.0010 |
| RAND3 | RAND9 | 0.0307 | 0.1110 | 0.8760 |

A More Practical Example of Hidden Markov Models

This section demonstrates an example of hidden Markov models using actual data, in this case an application that predicts future movement of a financial market. There are five candidates for predictor variables and a single target:

CMMA_5 is the current closing price of the market, minus its 5-day moving average. This shows the degree to which the market just (as of the end of the current day) departed from its recent price level.

CMMA_10 is a similar quantity, but based on the 10-day moving average.

CMMA_20 is a similar quantity, but based on the 20-day moving average.

LIN_ATR_7 is the slope of the best-fit straight line connecting the prices over the most recent 7 days, normalized by average true range. This indicates the short-term price trend in the market.

LIN_ATR_15 is a similar quantity, but based on the 15-day trend.

DAY_RETURN_1 is the market change over the next day, normalized by average true range. This variable serves as the target, as it represents the future change of the market price.

This example specifies that two predictors will be used by the model, and three states are possible. The model that correlates most highly with the target uses *CMMA_5* and *CMMA_20* as predictors. The means and standard deviations of these variables are shown for each of the three states:

Means (top number) and standard deviations (bottom number)

| State | CMMA_20 | CMMA_5 |
|-------|----------------------|-----------------------|
| 1 | -20.81845 9.42582 | -15.87819 16.57821 |
| 2 | 24.57826 8.25328 | 17.83951 15.22672 |
| 3 | 3.57633 7.27092 | 2.36846 17.76842 |

The three states are highly distinct in terms of their predictor distributions. CMMA_20, in particular, has means that are widely separated relative to their standard deviations. We see that State 1 is characterized by today's price being much lower than recent prices, State 2 is characterized by today's price being much higher than recent prices, and State 3 is characterized by today's price being about the same as recent prices. This sounds almost too 'sensible' to be believed, but numerous reruns of the test consistently produced similar results.

The transition probability matrix, shown below, reveals several interesting properties. First, we see that states have considerable persistence; there is about a 90 percent probability that tomorrow will remain in the same state as today. What is also interesting is that it is nearly impossible for the market to transition between States 1 and 2 without going through State 3, and in fact probably staying in State 3 for some time. In fact, the probability of going from State 1 to State 2 is zero to at least four digits!

Transition probabilities...

| | 1 | 2 | 3 |
|---|--------|--------|--------|
| 1 | 0.8978 | 0.0000 | 0.1022 |
| 2 | 0.0014 | 0.9095 | 0.0890 |
| 3 | 0.0711 | 0.0747 | 0.8542 |

The table of additional properties shows how these states relate to the target, the price change of the market the next day. We see that State 3, that corresponding to prices remaining fairly constant, is the most common, occurring almost 40 percent of the time. We also see at least one-day persistence of price movements into the future, as State 1, which corresponds to a pattern of today's closing price being far below recent prices, is associated with a negative price movement tomorrow. Similarly, State 2, which corresponds to a pattern of today's closing price being far above recent prices, is associated with an upward price movement tomorrow. Finally, it is noteworthy that the standard deviation of the target when in State 1 is almost fifty percent higher than when in the other two states. Thus, we can expect unusually large market turbulence when we have been in a pattern of prices closing far below their recent values. This agrees well with intuition, but it is nice to see it corroborated numerically.

| State | Percent | Correlation | Target mean | Target StdDev |
|-------|---------|-------------|-------------|---------------|
| 1 | 27.75 | -0.07034 | -0.05099 | 0.86047 |
| 2 | 32.41 | 0.06831 | 0.08906 | 0.60901 |
| 3 | 39.84 | -0.00049 | 0.02438 | 0.64007 |

40 Hidden Markov Models with Target Correlation

Finally, we have the list of models sorted according to their relationship to the target. The major take-away from this list is that the CMMA variables are much more important to predicting tomorrow's price movement than the linear trend variables. Also, the degree of significance of these relationships is impressive, usually the minimum obtainable from the 1000 Monte-Carlo replications performed.

| Predictor 1 | Predictor 2 | Multiple-R | Solo pval | Unbiased pval |
|-------------|-------------|------------|-----------|---------------|
| CMMA_20 | CMMA_5 | 0.0807 | 0.0010 | 0.0010 |
| CMMA_5 | LIN_ATR_7 | 0.0762 | 0.0010 | 0.0010 |
| CMMA_10 | CMMA_5 | 0.0689 | 0.0010 | 0.0010 |
| CMMA_10 | CMMA_20 | 0.0686 | 0.0010 | 0.0010 |
| CMMA_20 | LIN_ATR_7 | 0.0650 | 0.0010 | 0.0010 |
| CMMA_20 | LIN_ATR_15 | 0.0442 | 0.0010 | 0.0010 |
| CMMA_10 | LIN_ATR_7 | 0.0408 | 0.0010 | 0.0010 |
| CMMA_10 | LIN_ATR_15 | 0.0330 | 0.0020 | 0.0080 |
| CMMA_5 | LIN_ATR_15 | 0.0227 | 0.0480 | 0.1500 |
| LIN_ATR_15 | LIN_ATR_7 | 0.0168 | 0.1790 | 0.4750 |

Stationarity Test for Break in Mean

Stationarity in the mean is vital to most prediction schemes. If a predictor or target significantly changes its mean in the midst of a data stream, it would be foolish to assume that a prediction model will perform well on both sides of this break. Thus, we should always check for this sort of nonstationarity in all predictors and targets.

Even for applications in which series being evaluated are not being used as predictors or targets, this test is also useful. We may have a process whose performance is indicated by a numerical value. It may be the error rate of a prediction system, or cost savings achieved by a new manufacturing process. A classic example is following the performance of a market trading system. Suppose a previously profitable system suddenly deteriorates. We naturally wish to determine whether this falloff in performance is within historical norms or perhaps signifies something more serious.

This test is performed by clicking *Test / Stationarity break in mean*. The dialog box shown below will appear:

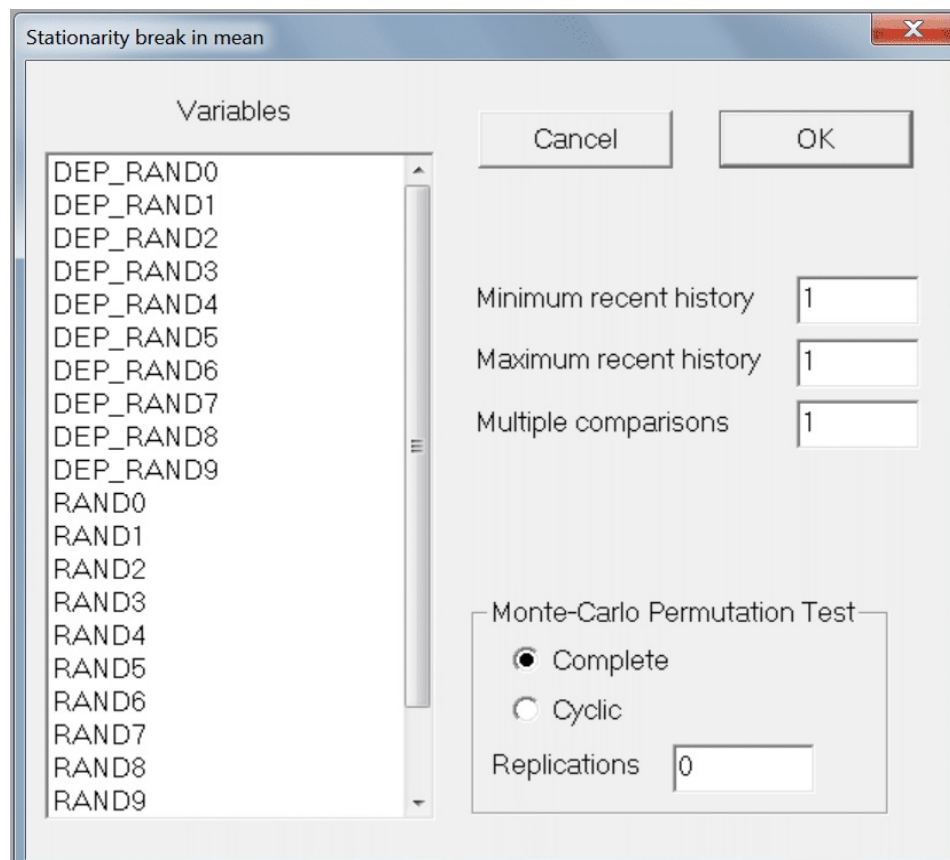


Figure 5: Dialog for stationarity test for break in mean

The user must select one or more variables. The user also specifies the range of recent history which will be searched for a break in the mean. The default of doing no search at all, but rather looking at only the most recent observation, allows the fastest detection of a change. However, it is also the least sensitive test, being based on a single observation relative to the rest of history. Employing a wider search range greatly increases the sensitivity of the test, at the price of delayed confirmation of a change in the mean.

The *multiple comparisons* field has a subtle but important function. Suppose you are performing a one-time test. You have one or more series which you plan to employ as predictors and/or targets in a modeling operation. You simply want to test whether any of them have a significant break in their mean. Then you can leave the *multiple comparisons* field at its default of one.

But now suppose you are monitoring incoming data from a series. For example, you may be assessing quarterly returns of a market trading system. Every time a new quarter rolls around you repeat the test. The statistical term for this repetition of the same test with different data is *multiple comparisons*. Its effect is to increase the chance that you will observe a statistically significant result, even though the effect you are looking for is not present. Sooner or later, random chance is going to present a significant result due to nothing more than luck.

The user can compensate for this effect by having the program adjust its p-values under the assumption that a specified number of tests will be performed. Of course, in real life it would be difficult to make an honest assessment in advance of exactly how much testing will be done. Still, this capability is better than blithely ignoring this vital issue! At a minimum, the user can see the effect of multiple tests on the computed p-values, and make a good-faith assessment of the number of tests that will be performed.

Because there will be huge correlation between successive test statistics due to overlap of the testing regions, ordinary multiple-comparison tests are invalid. For this special application, an ad hoc but reasonable methodology is followed. Look at Figure 6 below:

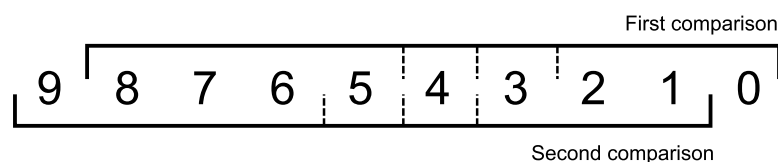


Figure 6: Testing for break in mean

Figure 6 illustrates the simple situation of testing with a range of 3 (*minimum recent history*) to 5 (*maximum recent history*) cases on the ‘recent’ side of the hypothetical break. It also shows 2 multiple comparisons. The dotted lines show the breakpoints tested.

For the original, unpermuted data only the ‘First comparison’ would be performed. Whichever of the three trial breakpoints produces the largest break will be the score as of the most current observation.

For all permutations, both comparisons will be performed. The null-hypothesis score will be the greatest of the six scores (three for each of the two comparisons). We then count how many of these null hypothesis scores equal or exceed the obtained score for each test. As per the usual Monte-Carlo permutation test, let there be m permutations, with k of them having a score equaling or exceeding the greatest score among the tests (which, strictly speaking, is not known until all tests are complete!). Then the p-value is $(k+1)/(m+1)$. This is the approximate probability that, if there were no break in the mean, we would have obtained a maximum break score across all tests that is at least as large as that actually observed.

There are several theoretical problems with this multiple-comparison test. Foremost, it is not strictly correct to keep re-evaluating the p-value on each test. By rights we should wait until all tests are complete and examine the maximum break across all tests. The computed p-value relates to this maximum break. Of course, in real life this would defeat the whole purpose of the test! We want to test on an ongoing basis. But I strongly suspect that, compared to other sources of random error, this is of minor consequence.

Also, the shifting of test windows probably does a good job of accounting for serial correlation in the test statistics, but I have no rigorous proof. Because each sequential test involves a massive overlap in the data that goes into the test, the test statistics will have similarly massive serial correlation. The algorithm illustrated in Figure 6 simulates what would happen in real life, but rigorous justification would be nice.

In short, the mathematical foundations of this test are shaky. Nonetheless, in a multiple-comparisons situation, this test is almost certainly far superior to failing to compensate in any way, and I have reasonable confidence that it is actually quite good. But be warned.

If the user sets the Monte-Carlo permutation test replications to zero or one, no MCPT will be performed, and only one column of results will be printed. This column is labeled $Z(U)$, and it is the absolute z-score corresponding to the Mann-Whitney U-test statistic for the difference in means between the data before and after the break point.

In the more usual situation of the user specifying a large number of replications (100-1000 or so), two additional columns are printed. The *Solo pval* for a variable is the approximate p-value for that variable considered in isolation; it is the probability that if the variable had no break in its mean we would have obtained a test statistic at least as large as was actually obtained.

If this quantity is not small, the developer should be inclined to believe that the variable does not have a significant mean break. Of course, this logic is, in a sense, accepting a null hypothesis, which is well known to be a dangerous practice. However, if a reasonable number of cases are present and a reasonable number of Monte-Carlo replications have been done, this test is powerful enough that failure to achieve a small p-value can be interpreted as the variable being decently stationary in its mean.

If more than one variable is specified, then the *Unbiased pval* column has a useful interpretation. When several variables are tested, chances are that one or more of them will, by sheer chance, have an usually large apparent mean break, even if in truth no such break exists. The *Unbiased pval* compensates for this effect.

The *Unbiased pval* is printed for all variables. For the first variable, the one having the greatest observed mean break, this is the approximate probability that, if none of the variables had a mean break, we could get a greatest mean break among them at least as large as that observed. For those other, lesser candidates, the *Unbiased pval* is an upper bound for the true unbiased p-value of the variable. Thus, a very small *Unbiased pval* for any candidate is a strong indication that the candidate has a significant mean break. Unfortunately, unlike the *Solo pval*, large values of the *Unbiased pval* are not necessarily evidence that the candidate is break-free. Large values, especially near the bottom of the sorted list, may be due to over-estimation of the true p-value. The author is not aware of any algorithm for computing correct unbiased p-values for any candidate other than that having the largest break. However, because this measure is conservative, it does have great utility in discovering nonstationary variables.

On a final note, be aware that having a *statistically* significant mean break does not equate to having a *practically* significant mean break. If the dataset is large, even a trivial mean break, something of no practical consequence, may show statistical significance. This test should be treated as a tool, a supplementary source of information, as opposed to the final arbiter of stationarity.

Serial Correlation and Cyclic Permutation

Like many other tests in *VarScreen*, the user may select either *complete* or *cyclic* permutation. In other tests, the cyclic method is useful for variables having significant serial correlation. However, with a test for a break in the mean, one must be cautious, as positive serial correlation can, in and of itself cause the mean to wander. Thus, any situation in which cyclic permutation is warranted is likely to also be a situation in which the mean will be inherently nonstationary, or at least appear so! The user will nearly always employ complete permutation. Still, in some special situations, cyclic permutation is appropriate. This should become more clear as we work through several examples.

We begin with the most basic situation in which the variables have negligible serial correlation and the user wishes to search nearly the entire extent of the data stream for a break in the mean. We'll use the same RAND0 through RAND9 that have appeared in prior demonstrations. These are random series, independent within themselves and with one another. There are 6300 cases. We decide to keep at least 50 cases on each side of the sought-after break in order to provide decent sensitivity, realizing that if a mean break happens in the outer 50 cases we will miss it. Thus, we specify a minimum of 50 and a maximum of 6250 recent cases. Since this is a one-shot test we leave the multiple comparisons parameter at its default of one. The following results are obtained with 100 iterations:

| Variable | Z (U) | Solo pval | Unbiased pval |
|----------|--------|-----------|---------------|
| RAND6 | 3.3257 | 0.0600 | 0.3900 |
| RAND8 | 2.8718 | 0.1200 | 0.8100 |
| RAND5 | 2.5216 | 0.2800 | 0.9900 |
| RAND9 | 2.5128 | 0.3600 | 0.9900 |
| RAND0 | 2.4845 | 0.3900 | 1.0000 |
| RAND3 | 2.3591 | 0.3900 | 1.0000 |
| RAND2 | 2.0359 | 0.6600 | 1.0000 |
| RAND7 | 2.0212 | 0.7500 | 1.0000 |
| RAND1 | 1.9353 | 0.7400 | 1.0000 |
| RAND4 | 1.3772 | 0.9800 | 1.0000 |

One of these variables, RAND6, manages through luck to get a solo p-value of 0.06. But its unbiased p-value of 0.39 tells us that this was almost certainly just a fluke from having tested ten variables.

But what if our variables have substantial positive serial correlation? It is vital that we do not attempt to perform an ‘across the extent’ test for a mean break, such as the 50-6250 test just shown. If we were to use complete permutation, the serial correlation in the null hypothesis runs would be destroyed, while the serial correlation in the unpermuted data would cause the mean to wander, making it virtually certain that a significant (probably *highly* significant!) break would be found, whether one truly exists or not. But cyclic permutation would not work here either. The only effect of the data rotation effected by cyclic permutation would be to shift the position of the break; the search for a break would still find it in nearly every permutation. So the null hypothesis distribution would be too large, resulting in overly large p-values.

The only sort of test we can do when the data has substantial serial correlation is to limit the searched range to a very small fraction of the number of cases, so that when the null hypothesis distribution is computed via cyclic permutation, only a tiny fraction of that distribution will find the original mean break in the search. The most common such situation is when we suspect that a series we are measuring has recently suffered a shift in mean beyond that which can be expected from any positive serial correlation.

So let’s suppose that we want to examine only the most recent ten cases out of 6300. We’ll use the DEP_RANDOM through DEP_RANDOM9 variables seen in other tests. These variables, while independent of one another, have large positive serial correlation. The incorrect approach is to use complete permutation, as this destroys the serial correlation in the null hypothesis. If we were to make this foolish mistake, we would get the result shown below. Remember that these variables have no break in the mean other than the wandering that is to be expected from positive serial correlation.

| Variable | Z (U) | Solo pval | Unbiased pval |
|-------------|--------|-----------|---------------|
| DEP_RANDOM3 | 4.2651 | 0.0100 | 0.0100 |
| DEP_RANDOM5 | 3.8843 | 0.0100 | 0.0100 |
| DEP_RANDOM9 | 3.6829 | 0.0100 | 0.0100 |
| DEP_RANDOM0 | 3.4434 | 0.0100 | 0.0100 |
| DEP_RANDOM6 | 3.3758 | 0.0100 | 0.0100 |
| DEP_RANDOM8 | 3.1717 | 0.0100 | 0.0200 |
| DEP_RANDOM7 | 3.1334 | 0.0100 | 0.0500 |
| DEP_RANDOM4 | 2.5973 | 0.0300 | 0.2300 |
| DEP_RANDOM1 | 2.1878 | 0.0900 | 0.6300 |
| DEP_RANDOM2 | 0.7855 | 0.9500 | 1.0000 |

It’s obvious that this is a crazy test. Even most of the unbiased p-values are tiny.

But if we switch to cyclic permutation, we will be testing whether the mean of the most recent few cases differs from the mean of the earlier cases more than is usual in a series with this level of serial correlation. The results are as follows:

| Variable | Z (U) | Solo pval | Unbiased pval |
|-----------|--------|-----------|---------------|
| DEP RAND3 | 4.2651 | 0.1000 | 0.6600 |
| DEP RAND5 | 3.8843 | 0.2000 | 0.8900 |
| DEP RAND9 | 3.6829 | 0.1600 | 0.9600 |
| DEP RAND0 | 3.4434 | 0.2400 | 0.9800 |
| DEP RAND6 | 3.3758 | 0.2000 | 0.9800 |
| DEP RAND8 | 3.1717 | 0.3400 | 0.9900 |
| DEP RAND7 | 3.1334 | 0.3000 | 0.9900 |
| DEP RAND4 | 2.5973 | 0.4500 | 1.0000 |
| DEP RAND1 | 2.1878 | 0.6600 | 1.0000 |
| DEP RAND2 | 0.7855 | 0.8800 | 1.0000 |

This is more reasonable. Even the solo p-value of the ‘worst’ variable is not terribly significant, and its unbiased p-value correctly confirms that nothing unusual is going on here.

Finally, suppose we want to perform this exact same test but with the understanding that a few more observations will be coming in and we will want to repeat the test. In particular, we agree that we will be performing this same test a total of five times, each time a new case arrives. So we specify 5 multiple comparisons and observe the following results:

| Variable | Z (U) | Solo pval | Unbiased pval |
|-----------|--------|-----------|---------------|
| DEP RAND3 | 4.2643 | 0.1300 | 0.8900 |
| DEP RAND5 | 3.8833 | 0.2600 | 0.9800 |
| DEP RAND9 | 3.6887 | 0.3600 | 0.9800 |
| DEP RAND0 | 3.4422 | 0.3300 | 0.9900 |
| DEP RAND6 | 3.3811 | 0.4800 | 0.9900 |
| DEP RAND8 | 3.1702 | 0.4600 | 1.0000 |
| DEP RAND7 | 3.1319 | 0.4300 | 1.0000 |
| DEP RAND4 | 2.5955 | 0.6200 | 1.0000 |
| DEP RAND1 | 2.1892 | 0.7600 | 1.0000 |
| DEP RAND2 | 0.7836 | 1.0000 | 1.0000 |

As a result of allowing multiple comparisons, the p-values have all increased somewhat. Also note that the Z(U) values have changed slightly. This is because the number of cases tested is slightly reduced for multiple comparisons, as illustrated in Figure 6 on Page 42.

FREL: Feature Weighting as Regularized Energy-Based Learning

The FREL algorithm (Yun Li et al, 'FREL: A Stable Feature Selection Algorithm', *IEEE Transactions on Neural Networks and Learning Systems*, July 2015) is a useful method for ranking, and even weighting, predictor variables in a classification application which is relatively low noise but is plagued by high dimensionality (numerous predictors) and small sample size. The implementation in *VarScreen* is strongly based on their innovative algorithm, but with significant modifications that I believe improve on the original version by providing more accurate and stable weights (at the cost of slower execution). My implementation also includes an approximate Monte-Carlo permutation test (MCPT) of the null hypothesis that all predictors have equal value, as well as an MCPT of the null hypothesis that the predictors, taken as a group, are worthless. Sadly, I am unable to devise a FREL-based MCPT of any null hypothesis concerning individual predictors taken in isolation.

The 'model' which inspires FREL is weighted nearest-neighbor classification. The distance between a test case having predictors $\mathbf{x} = \{x_1, \dots, x_K\}$ and a training-set case $\mathbf{t} = \{t_1, \dots, t_K\}$ is defined as the city-block distance between these cases, with each dimension having its own weight. This is defined as:

$$D(\mathbf{x}, \mathbf{t}) = \sum_k w_k |x_k - t_k|$$

Then, if one wishes to classify an unknown test case \mathbf{x} based on a training set, one would compute the distance between the test case and each member of the training set. The chosen class for the test case would be the class of the training case having minimum distance from the test case.

Of course, performing this classification presupposes that we know appropriate weights. The procedure can be inverted and used to find optimal weights, and we could then interpret the weights as measures of importance of the predictors (assuming that the predictors have commensurate scaling!). All we would do is define a measure of classification quality and then find weights that maximize this quality measure.

An approach to machine learning that is becoming more and more popular is *energy-based modeling*. One has a set of random variables, which in the current context would be predictors, and a prediction target or class membership. The model defines a scalar *energy* as a function of the values of these variables, sometimes called their *configuration*. This energy is a measure of the compatibility of the configuration, with small values of

energy corresponding to compatible configurations. If we have a known energy-based model and we wish to make an inference (a prediction or classification) based on specified values of the predictors, we fix the predictors and vary the target or class variable to identify the configuration that minimizes the energy.

In order to find a good energy-based model, we tune the parameters of the model in such a way that ‘correct’ configurations (as indicated by the training set) have small energy and ‘incorrect’ configurations have large energy.

Once the structure of the model is specified, in order to find optimal parameters we define a loss functional (a function of a function). The model is a function which maps configurations of variables to energy values, and the loss functional maps models to scalar loss values. In order to train the model, we find the version (parameters for the model family) which minimizes the loss functional.

The most common version of this latter operation, which we will do here, is to define a per-sample loss functional as a function of the model and a single case, and then average this per-sample measure across the entire training set.

This is a good time for a brief digression to make sure that two crucial issues are clear. First, many models, such as nearest-neighbor classification and some types of kernel regression, implicitly include the entire training set (or some other dataset) as a key component of the model. Do not confuse this with discussions of the training set related to training. It’s still just the model, and we need not explicitly mention the presence of the training set as part of the model. Second, do not confuse energy with loss. Energy is a measure of the compatibility of a given variable configuration with a model, and it is used to make a prediction. Loss is a measure of the quality of a model in a way that generally includes a training set, and it is used to find an optimal model.

The energy that a model M assigns to a hypothetical variable configuration $\{x, y\}$ can be conveniently written as $E(M, x, y)$. An extremely common and useful way to express the per-sample loss for a single training case $\{x^i, y^i\}$ is $L(y^i, E(M, x^i, Y))$, in which the term $E(M, x^i, Y)$ actually stands for multiple energy values, one for each possible value of y . In other words, the per-sample loss for a single case is a function of the true value of y for that case, and the energies given by the model for x associated with every possible y .

Note, by the way, that the distinction between *function* and *functional* become a bit murky here, depending on whether we think in terms of E being a hypothetical function or an observed number. In any case, the idea should be clear.

50 Feature Weighting as Regularized Energy-Based Learning

We are almost done presenting a general form of an effective loss function(al) for training an optimal (in the sense of the loss) model. We have seen the form of a per-sample loss, and stated that averaging this quantity over every sample in the training set is reasonable. The only remaining issue is that of *regularization*. This enables us to embed prior knowledge about the model in the final solution. Typically, this involves limiting the size of weights involved in the expression of the model, although other approaches are possible. With these things in mind, we can express the loss of a given model M for a given training set T and regularization function R as shown below. This is a scalar quantity which we will minimize in order to develop a good model.

$$L(M, T) = \frac{1}{K} \sum_k L[y^k, E(M, x^k, \Upsilon)] + R(M)$$

To review, a good model will fulfill two requirements: it will have low energy for correct configurations and high energy for incorrect configurations. Looked at another way, when a good model is presented with a set of predictors x , its energy will be low when it is simultaneously presented with the correct y for that x , and its energy will be high when it is simultaneously presented with any incorrect y .

It is tempting, and often appropriate, to consider only the first half of this two-part requirement: the model will have low energy for correct configurations. This is especially true for models in which fulfilling the first half automatically fulfills the second half. For example, suppose we have a regression equation as the model, and we define the energy associated with the model and a training case as the squared difference between the correct answer and the answer provided by the regression function. If the loss is just this energy, then averaged across the entire training set, the loss is the mean squared error (MSE). The optimal model is produced by minimizing the MSE, a venerable approach.

But for many model architectures, this halfway method is not a good approach. It is much better, if not mandatory, to explicitly take into account the second half of the requirement: the energy of incorrect answers should be large. And intuitively, we don't much care about easy situations, those incorrect answers that have huge energy. Even a weak model will do well with them. What we must worry about is those situations in which an incorrect answer has dangerously low energy. We want our model to be able to raise the energy of these problematic cases as much as possible above the energy of the correct answer.

This intuition leads to the following definition:

The *most offending incorrect answer* for a case, which we will call \tilde{y} , is the incorrect answer that has the lowest energy. This is the answer most likely to cause an error, because it is the incorrect answer that is most difficult for the model to distinguish from the correct answer. The second half of the training procedure discussed earlier, that incorrect answers should have large energy, is more general than is necessary. All we really care about is that the most offending incorrect answer has energy as large as possible, compared to the energy of the correct answer. The other incorrect answers are of relatively minor importance because they are easier for the model to avoid.

In particular, what we often want to maximize is the difference between the energy of the most offending incorrect answer and the energy of the correct answer. This will give us a model that is optimal in the sense of effectively handling the most difficult cases, while letting the easy cases slide.

A popular per-sample loss criterion, and which is used in *VarScreen*, is the log loss shown below. Note how it is a monotonic function of the difference between the two energies, so optimizing either is equivalent to optimizing the other (for a single case, not averaged across the training set!).

$$\text{Loss}(\mathbf{M}, \mathbf{x}^i, \mathbf{y}^i) = \log(1 + \exp[E(\mathbf{M}, \mathbf{x}^i, \mathbf{y}^i) - E(\mathbf{M}, \mathbf{x}^i, \tilde{\mathbf{y}}^i)])$$

Now that a theoretical foundation is laid, we can apply these ideas to the specific model used in the FREL paper and *VarScreen*. Recall from the beginning of this section that we use weighted nearest-neighbor classification. Thus, in order to compute $E(\mathbf{M}, \mathbf{x}^i, \mathbf{y}^i)$ for training case i , we check all other training cases in the correct class, \mathbf{y}^i . The smallest distance is the energy for the correct class. Similarly, to compute $E(\mathbf{M}, \mathbf{x}^i, \tilde{\mathbf{y}}^i)$ we search all other training cases in an incorrect class and find the distance to the nearest. Of course, although this is simple to describe and implement, it can be horrendously slow to compute. The quantity being minimized is the average across the training set of the per-sample losses shown in the equation above. If there are n training cases and K predictors, a single evaluation of the grand loss function requires on the order of Kn^2 operations. Yikes! Luckily, FREL is most useful for situations in which the training set is small relative to the number of predictor candidates, so that squared term will hopefully not be a serious problem.

52 Feature Weighting as Regularized Energy-Based Learning

All that remains to be settled is the regularization. In any reasonable application, the energy of the incorrect answers will, on average, exceed that of the correct answers; otherwise the model would be worthless! For the loss function just shown applied to weighted nearest-neighbor classification, increasing the weights together will decrease the loss, because the term being exponentiated will become increasingly negative. Thus, naive minimization of the loss will result in the weights blowing up without bound. Thus, we are inspired to penalize large weights. This is common practice, even in situations in which this blowup is not natural. The reason is that in many models, large weights are associated with overfitting and poor out-of-sample performance. In *VarScreen* we use the common method of penalizing by the sum of the squares of the weights. The sum of their absolute values is also common and may be implemented in a future version of the program.

The optimal weights determined by minimizing regularized loss can be interpreted as measures of importance of the individual predictors. However, two issues must be considered. First, the scaling of the predictors obviously impacts the weights, so their scaling should be commensurate. *VarScreen* takes care of this by internally scaling per their standard deviation. Second, interpretation by the user is aided by normalizing the weights in some way for display. In *VarScreen* they are linearly normalized so as to sum to 100.

A frequently useful variation is to take many bootstrap samples from the dataset and compute the final weight estimate by averaging the estimates produced from each bootstrap sample. The sampling must be done without replacement, as nearest-neighbor algorithms are irreparably damaged when the dataset contains exact replications of cases. Bootstrapping FREL has at least two major advantages over doing one FREL analysis of the entire dataset:

- 1) Stability is usually improved. A critical aspect of any weighting scheme is that the computed optimal weights should be affected as minimally as possible by small changes in the dataset. Such changes might be inclusion or exclusion of a few training cases, or change might be effected by the addition of noise to the data. An average of bootstraps is much more robust against data changes compared to a single complete FREL processing.
- 2) Because run time of the FREL algorithm is proportional to the square of the number of cases, we can greatly decrease the run time by performing many iterations of a small sample.

For these reasons, bootstrapping is generally recommended.

FREL Operation in VarScreen

We've already discussed the mathematics behind the FREL implementation in *VarScreen*. This section covers the user interface. When the user clicks *Test / Regularized energy-based*, a dialog box appears. The following information must be supplied by the user:

The leftmost (*Predictors*) column is used to specify the set of predictor candidates. Multiple candidates can be selected by dragging the mouse cursor across a block, or by clicking the first candidate in a block, holding the Shift key, and clicking the last candidate in the block. Individual candidates can be toggled on and off by holding the Ctrl key while clicking on the variable.

The *Target* column is used to specify the target variable. This variable will be partitioned into two or more classes based on its values. FREL does not permit continuous targets.

Target bins specifies the number of bins into which the target will be categorized. The number of cases in each bin will be made as equal as possible.

Regularization factor traditionally prevents model weights from running away to problematic large values. However, in *VarScreen* this is a fairly non-critical parameter; even zero is acceptable. This is because the optimization algorithm in *VarScreen* inherently prevents weight runaway as part of its stability enforcement. In practical terms, the effect of the regularization factor is to control the relative spread of weights. Suppose that predictability is concentrated in just one or a few candidates. If the user specifies a small or zero value for this parameter, the computed weights will strongly reflect this focus. However, if a very large regularization factor is specified, the focus will be less intense; some of the weight will be redistributed away from the dominant predictors and given to predictors of lesser value. Intense focus on one or a few dominant predictors can, in some cases, be seen as a form of overfitting.

Bootstrap operation usually increases robustness of the weight estimates and also decreases runtime, a happy confluence of outcomes. By default, no bootstrapping is done. But the user can specify that a given number of *iterations* are performed, each having a specified *sample size*. The sample size must be large enough that each sample is virtually guaranteed to have a significant number of representatives from each target class. For the number of iterations, my own rough rule of thumb is that the product of the number of iterations times the sample size should be about twice the number of training cases.

54 Feature Weighting as Regularized Energy-Based Learning

A *Monte-Carlo permutation test* is a useful, though time consuming, way to test certain null hypotheses about the predictor candidates. It is vital to understand that these tests are radically different from the other permutation tests in *VarScreen*. For one thing, I am not aware of any way of performing a perfect individual-candidate MCPT with FREL; the best I can do is come up with a rough approximation that appears to work well in practice. More importantly, in other tests, the candidate predictors are handled individually, so the p-values (at least the solo tests) are independent. But FREL considers all candidates simultaneously. This dependence changes the nature of MCPT. One effect is for dominant candidates to ‘suck’ weight out of lesser candidates, thus reducing their apparent significance. But the important effect is to radically change the nature of the null and alternative hypotheses of the test.

In other *VarScreen* tests, the null hypothesis for each solo p-value is that the *individual* candidate is *worthless*, and that for the unbiased p-values is that *all* candidates are *worthless*, and the power of the test is in identifying *individual* candidates which have predictive power. But for FREL, the individual MCPT tests have no useful power in situations in which all candidates have equal predictive power, regardless of whether that power is tiny or large. The null hypothesis is still generated by making all candidates worthless, exactly as in other tests. But because of the joint estimation of weights, it is more intuitive (though not strictly correct!) to think of the null hypothesis as being that *all candidates have equal predictive power*, with the unbiased p-values compensating for the fact that we are testing numerous candidates, and any of them may be outstanding by random luck. In other words, these individual tests are related to the predictive power of each candidate *relative to their competitors*. Their *individual* predictive powers play no easily identifiable role in determining p-values.

With this in mind, we can look at the p-values of candidates at the top of the list, those ranked highest in terms of predictive power and having the largest weights, and consider the p-values as being the probability that if all candidates were truly *equal* in predictive power, the top-ranked candidates would have *outperformed the others* to the degree shown. Suppose we see a highly significant result for the single best candidate. It may be that this best candidate is *almost* worthless, and its competitors are *completely* worthless. Or it may be that this single candidate is *excellent*, while its competitors are merely *very, very good*. In either case we may see the best candidate having a highly significant p-value. Again, I emphasize that this interpretation is not strictly correct, but I believe that it is close enough, especially the unbiased p-values, to be effective indicators of the validity of the obtained results.

The sucking of weight from relatively poor predictors to good predictors has a peculiar and potentially confusing effect on the solo p-values. As we drop down the sorted list to the low-ranked candidates, we can see the solo p-values cover a wide range, jumping up and down between high and low significance randomly. This is illustrating in an exaggerated manner the fact that the p-values for worthless candidates in any statistical test have a uniform distribution, with all values being equally likely. This is yet another reason why we should focus on the unbiased p-values, ignoring the solo p-values except perhaps (and with great caution) for the few top-ranked candidates.

VarScreen does print one additional p-value, called the *Loss p-value*. This is a ‘grand’ measure of the ability of all predictors taken together to be effective at correct classification. The null hypothesis is that none of the candidates have any predictive power, and the Loss p-value is the probability that if this were so, we would have achieved a loss at least as low as that obtained. This p-value being small is a necessary condition for any of the individual p-values to be meaningful. If we cannot be reasonable certain that at least one of the candidates has predictive power, then there is no point in considering their relative power!

The user may specify several parameters for the MCPT:

Replications defaults to zero, in which case no Monte-Carlo permutation test is performed. However, if computer time permits, it is usually best to set this to at least 100, and perhaps as much as 1000, so that solo and unbiased p-values will be computed. Note that the minimum possible p-value is the reciprocal of the number of permutations. So, for example, if the user specifies 100 permutations, the minimum p-value that can appear is 0.01. Run time of this test is linearly related to the number of permutations.

The user must choose either *Complete* or *Cyclic* permutations. If the user is confident that there is no dependency as described earlier in this document, then *Complete* should be used; it is the traditional approach which does a complete random shuffle for each permutation. However, if there is dependency, this type of shuffling will produce underestimation of *p-values*, a very dangerous situation. If the dependency is serial (the data is a time series and the dependency is among samples close in time) then a slight improvement in the situation can be obtained by using *Cyclic* permutation. In this type of shuffle, the time order of the target is kept intact except at the ends by rotating the targets with end-point wraparound. Shuffling this way preserves most of the serial dependency in the permuted targets, which makes the algorithm more accurate. The *p-values* computed this way will generally be larger than those computed with complete shuffling, and hence less likely to lead to false rejection of the null hypothesis of no

56 Feature Weighting as Regularized Energy-Based Learning

predictive power. But be warned that the cure is far from complete; computed *p-values* will still underestimate the true values, just not as badly.

Note that in most cases it is legitimate to use *Cyclic* permutation instead of *Complete* when there is no dependency. However, if the dataset is small, *Cyclic* permutation will limit the number of unique permutations and hence increase the random error inherent in the process. As long as the dataset is large, some users may prefer to use *Cyclic* permutation even if it is assumed that there is no serial dependency; in case there really is hidden serial dependency, this is a cheap insurance policy. Still, the best practice is to make sure that the data does not contain dependency and then use *Complete* permutation. Relying on *Cyclic* permutation to take care of dependency problems is living dangerously. And if the dataset contains fewer than 1000 or so cases, use of *Cyclic* permutation is not recommended unless it is necessary to handle dependency.

CUDA Considerations

First, be aware that the default CUDA parameters (*Kernels* and *Granularity*) should be fine for nearly all applications and hardware. However, for users who wish to tweak operation (or those who must do so because of timeouts) the FREL dialog allows the user to specify two parameters.

Computation of the loss function entails two nested loops. The outer loop performs cross validation, letting each training case play the role of a test case, with these individual losses averaged across the entire training set. The inner loop passes through all cases other than the test case and finds the energy of the correct answer and that of the most offending incorrect answer. Since this latter operation also involves finding the weighted distance between cases, this results in a *lot* of mathematical operations.

Microsoft Windows has the infamous ‘feature’ of limiting the time during which CUDA computation can monopolize the video display in a contiguous stretch, typically two seconds. Therefore, the *CUDA Kernels* parameter lets the outer loop be broken up into multiple kernel launches. By default all computation is performed in a single launch, which is good, because launches have considerable overhead. But if the screen goes black and a message pops up that the display adapter has been reset, you will have to increase (as little as possible!) the *CUDA Kernels* parameter.

The *Granularity* parameter is more subtle and require an understanding of CUDA hardware to be fully appreciated. If the granularity is set to 1, each outer-loop case is assigned to a thread, and this single thread handles the entire inner loop. But CUDA devices prefer much finer granularity so that they can run thousands or even millions of threads simultaneously. Otherwise, vast amounts of hardware resources sit idle, a grievous waste. To avoid this, the inner loop for each outer-loop case is broken up into *Granularity* sub-tasks, where this parameter cannot exceed the number of cases. The bottom line is that a total of *Number of cases* times *Granularity* separate threads are executed. Users with a late-model extremely powerful CUDA processor may benefit from increasing the granularity beyond the default, perhaps even to its limit of the number of cases.

FSCA: Forward Selection Component Analysis

The algorithms provided here are greatly inspired by the paper “Forward Selection Component Analysis: Algorithms and Applications” by Luca Puggini and Sean McLoone, published in *IEEE Transactions on Pattern Analysis and Machine Intelligence*, December 2017, and widely available for free download on various Internet sites. However, I have made several small modifications that I believe make it somewhat more practical for real-life applications.

The technique of principal components has been used for centuries (or so it seems!) to distill the information (variance) contained in a large number of variables down into a smaller, more manageable set of new variables. Sometimes the researcher is interested only in the *nature* of the linear combinations of the original variables that provide new ‘component’ variables having the property of capturing the maximum possible amount of the total variance inherent in the original set of variables. In other words, principal components analysis can be viewed as an application of descriptive statistics. Other times the researcher wants to go one step further, computing and employing the principal components as predictors in a modeling application.

However, with the advent of extremely large datasets, several shortcomings of traditional principal components analysis have become problematic. The root cause of these problems is that traditional principal component analysis computes the new variables as linear combinations of *all* of the original variables. If you have been presented with thousands of variables, there can be issues with using all of them.

One possible issue is the cost of obtaining all of these variables going forward. Maybe the research budget allowed for collecting a huge dataset for initial study, but the division manager would look askance at such a massive endeavor on an ongoing basis. It would be a lot better if, after an initial analysis, you could request updated samples from only a much smaller subset of the original variable set.

Another issue is interpretation. Being able to come up with descriptive names for the new variables (even if the ‘name’ is a paragraph long!) is always good. It’s hard enough putting a name to a linear combination of a dozen or two variables; try understanding and explaining the nature of a linear combination of two thousand variables! So if you could identify a much smaller subset of the original set, such that this subset encapsulates the majority of the independent variation inherent in the original set, and then compute the new component variables from this smaller set, you are in a far better position to understand, name, and explain what these new variables represent.

Yet another issue with traditional principal components when applied to an enormous dataset arises is the all too common situation of groups of variables having large mutual correlation. For example, in the analysis of financial markets for automated trading systems, we may measure many families of market behavior: trends, departures from trends, volatility, and so forth. We may have many hundreds of such indicators, and among them we may have several dozen different measures of volatility, all of which are highly correlated. When we apply traditional principal components analysis to such correlated groups, an unfortunate effect of the correlation is to cause the weights within each correlated set to be evenly dispersed among the correlated variables in the set. So, for example, suppose we have a set of 30 measures of volatility that are highly correlated. Even if volatility is an important source of variation (potentially useful information) across the dataset (market history), the computed weights for each of these variables will be small, each measure garnering a small amount of the total 'importance' indication. As a result, we may examine the weights, see nothing but tiny weights for the volatility measures, and erroneously conclude that volatility does not carry much importance. When there are many such groups, and especially if they do not fall into obvious families, the possibility of intelligent interpretation becomes hopeless.

The algorithms presented here go a long way toward solving all of these problems. They work by first finding the single variable that does the best job of 'explaining' the total variability (all original variables) observed in the dataset. Roughly speaking, we say that a variable does a good job of explaining the total variability if knowledge of the value of that variable tells us a lot about the values of all of the other variables in the original dataset. So the best variable is the one that lets us predict the values of all other variables with maximum accuracy.

Once we have the best single variable, we consider the remaining variables and find the one that, in conjunction with the one we already have, does the best job of predicting all other variables. Then we find a third, and a fourth, et cetera. Application of this simple algorithm gives us an ordered set of variables selected from the huge original set, beginning with the most important, and henceforth with decreasing but always optimal importance (conditional on prior selections).

It is well known that a greedy algorithm such as the strictly forward selection just described can produce a sub-optimal set of variables. It is always optimal in a certain sense, but only in the sense of being conditional on prior selections. It can (and often does) happen that when some new variable is selected, a previously selected variable suddenly loses a good deal of its importance. Thus, the algorithms here optionally allow for continual refinement of the set of selected variables by regularly testing previously selected variables to see if they should be removed and replaced with some other

candidate. Unfortunately, we then lose the strict ordering-of-importance property that we have with strict forward selection, but we gain a more optimal final subset of variables. Of course, even with backward refinement we can still end up with a set of variables that is inferior to what could be obtained by testing every possible subset. However, the combinatoric explosion that results from anything but a very small universe of variables makes exhaustive testing impossible. So in practice, backward refinement is pretty much the best we can do.

When *FSCA* is selected from the *Create* menu, a dialog box will appear, from which the user makes the following specifications:

The leftmost (*Variables*) column is used to specify the universe of variables from which a subset will be selected. Multiple candidates can be selected by dragging the mouse cursor across a block, or by clicking the first candidate in a block, holding the Shift key, and clicking the last candidate in the block. Individual candidates can be toggled on and off by holding the Ctrl key while clicking on the variable.

The *Number of Components* specifies how many variables will be selected, although if the dataset contains extreme colinearity this number will be reduced as needed to prevent colinearity in the computed components. Setting this value to zero causes all variables to be selected. This, of course, runs counter to the primary purpose of this algorithm. On the other hand, it does let us see the universe of variables rank-ordered according to ability to reconstruct the complete dataset. This information is often interesting and useful. The number of components computed will always equal the number of variables selected.

Three algorithms for variable selection and corresponding component generation are available:

Principal Components of the traditional variety can be computed. This is a rather uninteresting option, but it is included for comparison purposes.

Forward selection, ordered uses strict forward selection; no backward refinement is done. As a result, the order in which variables are printed when the program is finished represents their order of importance in reproducing the entire dataset. In other words, the first variable in the list is the single most important. The second variable in the list is the one that, *given the value of the first variable selected*, is the most important among the remaining variables. The third is the one that, *given the values of the first two variables selected*, is best at reproducing the dataset. Et cetera.

Forward selection, refined combines forward selection with backward refinement. This generally improves the quality of the final subset of variables compared to the prior option, but backward refinement destroys the ordering of the variables. It can happen that the first variable selected, the single best, doesn't even make it to the final subset! At this time, this option (the slowest of the three) is the only one of the three that is multi-threaded for full use of multi-core CPUs.

All three of these options create a new set of variables in the database which can then be used in subsequent studies. If the user specified principal components, the variable names will be in the form *PrinCo_n_m*, while the other two options will produce variables named *FSCA_n_m*. In both cases, *n* refers to the sequence number in which they were computed as separate operations. The first time you run the algorithm, *n*=1. The second time, *n*=2, and so forth. In both cases, *m* is the component number, ranging from 1 through the number of variables in the selected subset.

For all three options, the newly computed variables will have zero mean, unit standard deviation, and they will be uncorrelated. The *VarScreen.log* file will provide information to allow the user to recreate the components with other data and programs, if desired.

For the *ordered* (no refinement) option, the log file will list the actual coefficients needed to convert *standardized* (zero mean, unit standard deviation) values of the original variables to the newly created component variables, also standardized. For the other two options, the log file will list the correlations between each component and the original variable, with the first column being the component that captures the most variance from the subset, the second column capturing the second-most variance, and so forth. If you require coefficients for computing standardized values of the components, just divide each correlation by the eigenvalue shown at the top of the table. Or you can use the correlations directly, without dividing by the eigenvalues, in which case you will get the same components, but they will not have unit standard deviations.

For all three options, the eigenvalues and eigenvectors of the correlation matrix of the universe will be printed first, with as many columns as variables/components specified by the user. This is followed by a list of the mean squared correlation of each variable in the universe with all other variables. Finally, the table of coefficients or component/variable correlations as described above is printed.

Here is an example of each of the two FSCA algorithms. For this example, the following variables are employed:

RAND0 - RAND6 are independent (within themselves and with each other) random time series.

SUM12 = RAND1 + RAND2

SUM34 = RAND3 + RAND4

SUM1234 = SUM12 + SUM34

When we run the FSCA algorithm using the option for strict ordering (no refinement), we first see the following results printed:

There are 6 unique (non-redundant) sources of variation
The number of components computed is therefore being reduced to this value.

Eigenvalues, cumulative percent, and principal component factor structure

| | | | | | | |
|------------|---------|---------|---------|---------|---------|---------|
| Eigenvalue | 2.988 | 1.986 | 1.052 | 1.015 | 0.987 | 0.972 |
| Cumulative | 33.195 | 55.263 | 66.957 | 78.240 | 89.203 | 100.000 |
| | | | | | | |
| RAND1 | 0.4835 | 0.4964 | -0.6476 | -0.1497 | -0.1080 | -0.2576 |
| RAND2 | 0.4597 | 0.5206 | 0.6390 | 0.1478 | 0.1037 | 0.2770 |
| RAND3 | 0.5246 | -0.4808 | -0.0470 | -0.2077 | 0.6690 | -0.0271 |
| RAND4 | 0.5175 | -0.4859 | 0.0620 | 0.2194 | -0.6661 | 0.0240 |
| RAND5 | -0.0198 | -0.0198 | -0.4669 | 0.4999 | 0.1474 | 0.7139 |
| RAND6 | 0.0020 | 0.0260 | 0.0233 | 0.7937 | 0.2265 | -0.5635 |
| SUM12 | 0.6800 | 0.7331 | -0.0090 | -0.0021 | -0.0036 | 0.0128 |
| SUM1234 | 0.9997 | 0.0239 | 0.0012 | 0.0040 | 0.0003 | 0.0073 |
| SUM34 | 0.7331 | -0.6800 | 0.0104 | 0.0076 | 0.0039 | -0.0023 |

We have 9 variables in the universe, but the program notes that there are only 6 unique sources of variation. This is not surprising, because the 3 sum variables are just combinations of the other variables. Since by definition the computed components must be independent, the program limits us to just 6 components.

The first eigenvector accounts for one-third of the total variation in the dataset, and it correlates almost perfectly with SUM1234, very highly with SUM12 and SUM34, and moderately highly with RAND1-RAND4. None of this should be surprising.

The second component is just the contrast between RAND1 and RAND2 versus RAND3 and RAND4. In conjunction with the first component, it gives us over 55 percent of the total variation. The remaining components are other contrasts as well as RAND5 and 6.

Next, we get a list of the mean squared correlation of each variable in the universe with all other variables:

Mean squared correlation of each variable with all others

| | |
|---------|-------|
| RAND1 | 0.091 |
| RAND2 | 0.088 |
| RAND3 | 0.096 |
| RAND4 | 0.095 |
| RAND5 | 0.000 |
| RAND6 | 0.000 |
| SUM12 | 0.181 |
| SUM1234 | 0.248 |
| SUM34 | 0.191 |

It is not surprising that RAND1-RAND4, along with their various sums, have positive mean squared correlations, while RAND5 and RAND6 have zero correlations.

Last of all we get the table of coefficients needed to compute the 6 components from the chosen 6 variables in the subset. Note that each component depends on only the corresponding ordered variable and all previously selected variables.

| Variable | 1 | 2 | 3 | 4 | 5 | 6 |
|----------|---------|---------|---------|---------|---------|---------|
| SUM1234 | 1.0000 | -0.9730 | 0.0181 | 0.0106 | 0.0047 | -1.4045 |
| SUM12 | -0.0000 | 1.3953 | -0.9696 | -0.0091 | -0.0131 | 0.9888 |
| RAND2 | 0.0000 | -0.0000 | 1.3842 | -0.0129 | 0.0380 | -0.0081 |
| RAND6 | 0.0000 | 0.0000 | -0.0000 | 1.0001 | -0.0169 | -0.0071 |
| RAND5 | 0.0000 | -0.0000 | 0.0000 | 0.0000 | 1.0007 | -0.0017 |
| RAND4 | -0.0000 | 0.0000 | -0.0000 | -0.0000 | -0.0000 | 1.4188 |

Observe that the best single variable for reproducing the entire universe of values is SUM1234, the sum of four other variables in the universe, and the first component is just this one variable (its coefficient is 1.0 and all other coefficients are 0.0).

The second variable selected is another sum variable, and the corresponding component's value is computed as that sum variable times 1.3953, minus the prior selected variable times 0.9730.

The third variable selected is a similar weighted sum, primarily based on RAND2. The next two components are essentially equal to the two completely independent variables, RAND6 and RAND5. Note that their coefficients are almost exactly 1, and all other coefficients are almost exactly 0. And the last component is a complex mix of other variables.

We use this same universe of variables to demonstrate the other FSCA option, forward selection combined with backward refinement. The initial information (eigenstructure and mean squared correlations) are the same as in the prior example, so we will skip straight to the interesting part, the log of variables being added and replaced:

```
Commencing stepwise construction with SUM1234
Added SUM12 for criterion=4.973085
  Replaced SUM1234 with SUM34 to get criterion = 4.973123
Added RAND2 for criterion=6.011605
  Replaced SUM12 with RAND1 to get criterion = 6.011623
Added RAND6 for criterion=7.011701
Added RAND5 for criterion=8.010402
Added RAND4 for criterion=8.999919
  Replaced SUM34 with RAND3 to get criterion = 8.999940
```

As in the prior example, the first variable selected is SUM1234. We then add SUM12, as in the prior example. (Both options will always select the same first two variables.) But then something interesting happens: SUM1234 is replaced by SUM34, giving us a two-variable set of SUM12 and SUM34. To me, this is prettier than SUM1234 and SUM12.

We then add RAND2, which immediately triggers the replacement of SUM12 with RAND1. After that we add the two totally independent variables, RAND6 and RAND5. Finally, we add RAND4, which triggers the replacement of SUM34 with RAND3. The final results are shown below:

Eigenvalues, cumulative percent, and selected principal component factor structure

| | | | | | | |
|------------|---------|---------|---------|---------|---------|---------|
| Eigenvalue | 1.056 | 1.023 | 1.002 | 0.988 | 0.983 | 0.948 |
| Cumulative | 17.600 | 34.646 | 51.343 | 67.811 | 84.194 | 100.000 |
| | | | | | | |
| RAND3 | 0.2269 | 0.5167 | 0.2772 | 0.5747 | -0.5221 | -0.0431 |
| RAND1 | 0.5751 | 0.0508 | -0.1047 | 0.3593 | 0.5829 | 0.4322 |
| RAND2 | -0.6877 | 0.0094 | 0.1597 | 0.2264 | 0.0044 | 0.6709 |
| RAND6 | -0.0862 | -0.6254 | 0.4725 | 0.4844 | 0.1757 | -0.3357 |
| RAND5 | 0.4228 | -0.5366 | 0.1187 | -0.2014 | -0.5355 | 0.4381 |
| RAND4 | 0.1210 | 0.2723 | 0.8070 | -0.4496 | 0.2300 | 0.0702 |

The final set of selected variables is intuitively more appealing than what we got with the strict ordering option, because it's just the individual random variables, without their various sums. Because replacement has destroyed any ordering of the subset, it makes the most sense to me to just compute the components as the principal components of the final subset. Note that the eigenvalues are all nearly equal, meaning that the components have no strong ordering either. Also note that the values in the table are the correlations between the components and the variables, and they can be converted to weights by dividing each column by the eigenvalue at the top of the column.

LFS: Local Feature Selection

Most common feature selection algorithms are primarily oriented toward favoring features that are at least somewhat predictive over the *entire domain* of the feature set. This predictivity may be nonlinear, and it may interact with other features, but such a predictor will be at a significant advantage over *more powerful but only locally predictive candidates* if the nature of its relationship to a target variable is at least somewhat consistent across the domain of all possible values of all candidate features.

This can be a major problem, because modern nonlinear models can obtain a lot of useful predictive information from variables whose power is limited to small areas of the domain, or whose predictive relationship changes significantly over the domain. But if our predictor selection algorithm fails to find such variables, focusing instead on more global candidates, we lose out on what may be valuable information.

For example, consider the XOR problem. Suppose we have two variables symmetric around zero, and we define two classes. A case is a member of Class 1 if both of our variables are positive or both negative, and it is in Class 2 if one is positive and the other negative. This classification problem can be solved with 100 percent accuracy by a simple rule, and modern nonlinear models should have no trouble achieving nearly perfect performance. Yet if we were to augment these two variables with numerous worthless predictor candidates and then try to identify the two true predictors, an amazing number of otherwise sophisticated predictor selection algorithms would fail to find them. Not only are the marginal distributions of both variables identical in both classes, but the relationship of each variable to the class depends completely on the value of the other variable, with the relationship reversing across the domain. This is a tough problem.

This same issue arises in applications that are closer to reality. For example, a common phenomenon in equity market prediction is that certain families of indicators have considerable predictive power in times of low market volatility, but become useless in times of high volatility. The presence of a large amount of high-volatility data in the dataset dilutes the predictive power of such variables and may put otherwise excellent indicators at a competitive disadvantage. And this problem arises in many other applications. The effectiveness of medical treatments can vary according to age, weight, and a potentially large number of other unknown conditions. Identification of vehicles and pedestrians by a self-driving car's control system can depend on features that are vital in some contexts and distracting clutter in others. We need a feature selection algorithm that is sensitive to predictive power that comes, goes, and even reverses, according to location in the feature domain.

In terms of modeling we can deal with inconsistent behavior by using sophisticated nonlinear models (which are prone to overfitting!), or by using different models in differing regimes (assuming that we know how to define these regimes!). But consider the pre-modeling stage, when we are searching for predictor candidates. We would like to have a predictor selection algorithm that can *automatically* find such regime-dependent behavior and identify powerful predictors, even if this power is localized.

The feature selection algorithm described in “Local Feature Selection for Data Classification” by Narges Armanfard, James P. Reilly, and Majid Komeili (*IEEE Transactions on Pattern Analysis and Machine Intelligence*, June 2016) fits the bill nicely. We’ll now present a condensed and intuitive overview of how it operates.

There are a large number of possible approaches to feature selection. We’ve seen some based on mutual information and uncertainty reduction, techniques that are effective at detecting highly nonlinear relationships. Some techniques actually train predictive models, and perform their feature selection by intelligently choosing inputs for these models. Early discriminant analysis methods involve the use of Mahalanobis distances to find dimensions of maximum separation when the predictors are highly correlated, optimally taking correlation into account. The *LFS* algorithm presented here is based on yet another approach, a concept akin to nearest neighbor classification, but taken to a much higher level of sophistication.

We begin with a simple example: we want to predict success in college, with students divided into two classes: those who graduate, and those who drop out. We measure four candidate predictors for each student in our study dataset, and standardize the values of these predictors (mean zero and standard deviation one) to put their variation on a level playing field. These candidate predictors are:

- 1) SAT score
- 2) High school grade point average
- 3) Circumference of thumb divided by circumference of index finger
- 4) Day of month student was born

Suppose we randomly choose two students, both in the *Success* class. For each of these four features, think about the average difference in predictor value we would see for these two students. But now suppose we randomly choose two students, one in the *Success* class, and one in the *Dropout* class. The expected difference between these two students would be about the same as it was for the ‘same class’ students for the third and fourth candidate predictors, but much larger for the first and second candidate predictors: a person who graduated would probably have a higher GPA and SAT score

than a dropout, leading to a large difference, while these two students would probably have similar finger sizes and birthdays, at least relatively speaking.

If we effectively estimated these expected differences throughout the dataset, looking at every pair of students, we would conclude that the first two candidate predictors are the ones we want, because the expected difference in these two features for students in different classes greatly exceeds the expected difference for students in the same class, while for the third and fourth candidates we observe about the same difference, regardless of whether the two students are in the same class or different classes.

Now, instead of looking at candidate predictors individually, let's look at them in pairs: 1 and 2, 1 and 3, et cetera. A good measure of the difference between two cases is the Euclidean distance between them. Let $x_k^{(i)}$ represent the value of variable k as measured for case i , and let $\mathbf{x}^{(i)}$ represent the vector of all variables for this case. Then the distance between case i and case j is given by the following Equation:

$$d_{ij} = \|\mathbf{x}^{(i)} - \mathbf{x}^{(j)}\| = \sqrt{\sum_k \left(x_k^{(i)} - x_k^{(j)}\right)^2}$$

It should be clear that the pair of variables consisting of the first two competitors will have the greatest expected inter-class distance between cases, the pair consisting of the last two competitors will have the least expected inter-class distance, and mixed pairs will have intermediate values.

Intuition can now guide us toward a good way to choose an effective set of candidate predictors. We look for a set that has a high contrast between expected intra-class distance (which we want to be small) and inter-class distance (which we want to be large). Neither quality alone is good. For example, if we find a set of candidates that produces large average inter-class separation between cases, but the expected separation between cases in the same class is also large, we have gained nothing; we cannot look at either quality in isolation. We must find a balance, a way to trade off the desirable quality of low intra-class separation with the also desirable quality of high inter-class separation. The LFS algorithm has an automated way to find the optimal tradeoff, a topic which we will return to later.

All that we've seen so far is good, and the algorithm just outlined would work fairly well in practice. However, it is missing the 'Local' component of the 'Local Feature Selection' algorithm. We still need a way to handle the problem of predictive power that varies across the domain of all features. For example, the distribution shown in Figure 7 would foil the algorithm just described.

In this example, we have two classes, one of which is split into two distinct subsets. Think about how the variable selection algorithm just described would perform when presented with this problem. Half of the cases in Class 2 would have excellent inter-class separation from Class 1 via X_1 , though no separation at all via X_2 . The other half would experience the opposite behavior, gaining great separation via X_2 but none via X_1 . If inter-class separation were the only consideration, the algorithm would pick up X_1 and X_2 easily.

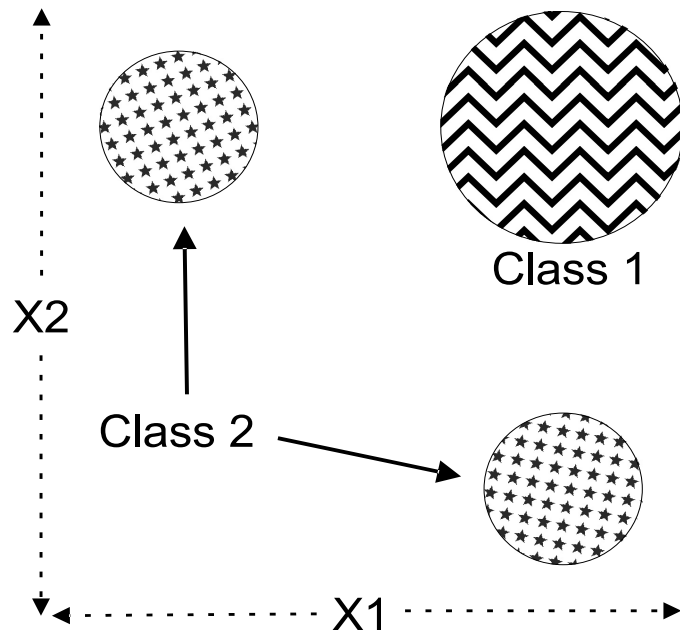


Figure 7: A job for Local Feature Selection

The problem lies with the intra-class separation. Cases that lie within either of the subsets of Class 2 would have nicely small separation. But if one case in Class 2 lies in one subset, and the other case lies in the other subset, the distance between them would be enormous, larger even than the inter-class separation! So the average intra-class separation for Class 2 would be so large that it would be nearly commensurate with the inter-class separation. It's unlikely that (X_1, X_2) would stand out as a set of effective predictors, even though this figure shows that they are fabulous.

The key element of the paper cited at the beginning of this section is that the problem shown in Figure 7 can be alleviated by weighting the distances with intelligently computed weights. The primary focus in the weighting scheme is that pairs of cases which are close are given higher weights than pairs which are distant, with the weighting dropping off exponentially with distance. It's somewhat more complicated than that, because the class memberships of the cases are taken into account, as well as global behavior of the distance metrics. The details are far too complex to get into here; see the cited paper if you are interested.

In order to get an idea of what's happening in regard to weights, the four histograms in Figure 8 show the weights generated from a test with data having the pattern shown in Figure 7.

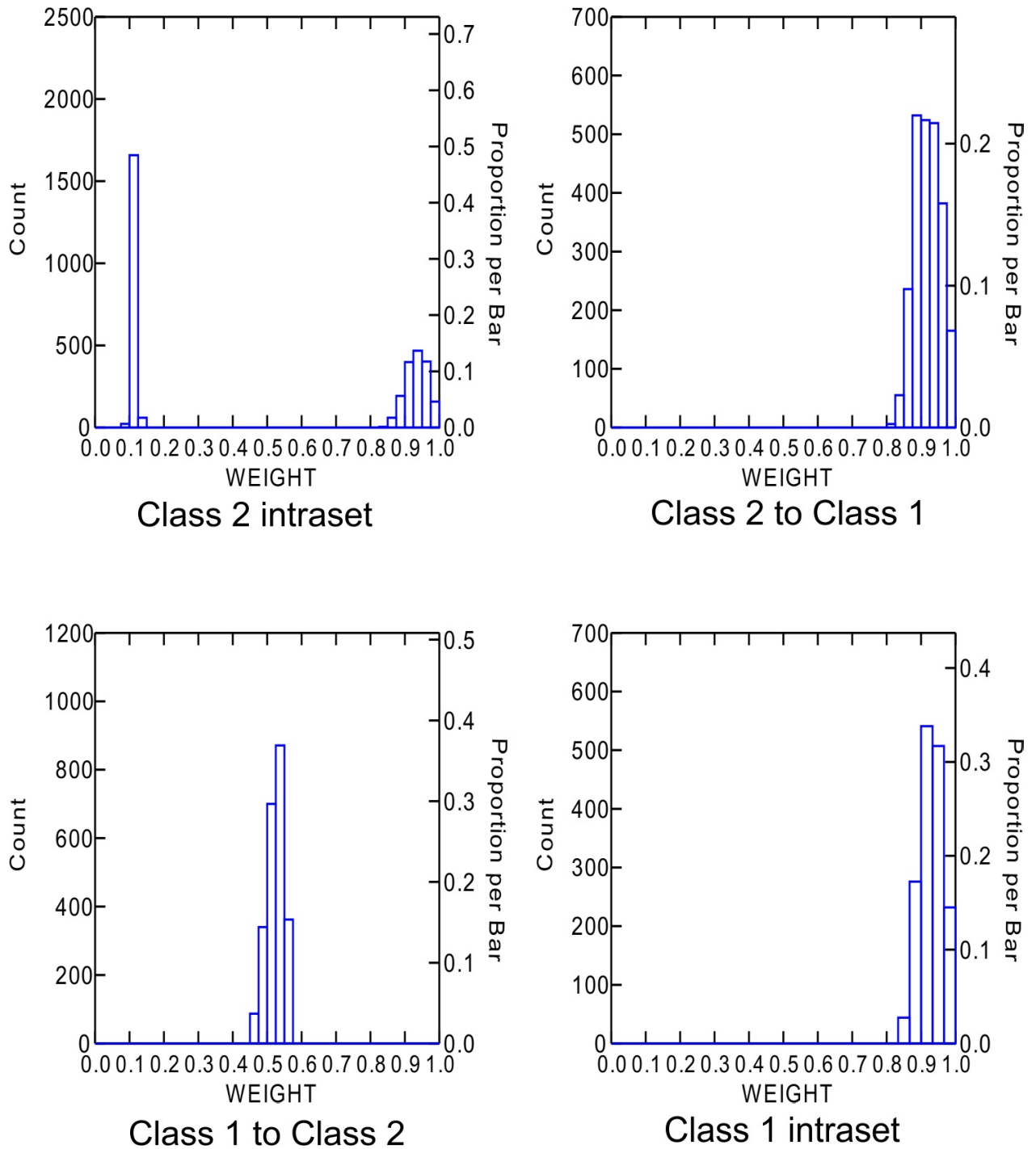


Figure 8: LFS weights for split-class example

The most interesting of these four histograms is the upper-left, which shows the weights for pairs of cases that are both in Class 2. We see that half of the weights are clustered near the maximum possible weight, one. These are the pairs of cases that are both in the same subset of Class 2. The other half of the weights are clustered near zero, the minimum possible. These are the pairs of cases that, while both in Class 2, are in different subsets of this class. So we see that when the intra-class separation (mean distance separating cases both in Class 2) is computed with weighted distances, pairs that span the two subsets are downplayed, thus providing a more realistic estimate of the intra-class separation.

The Class 1 intra-class weights are all close to one because this class is not split into subsets. Also, when we are considering cases in Class 2 and looking at their distances from cases in Class 1, we have full weighting. The weights are about 0.5 when we consider cases in Class 1 and look at the distances to cases in Class 2 (the weighting algorithm is asymmetric). Very roughly speaking, this is because there are two possible ways the difference can go. You can study the weight equations in the cited paper to see exactly how this comes to be.

What This Algorithm Reports

Because the algorithm performs optimal-candidate selection separately for each case, there is no practical way to report a single optimal candidate set, let alone a sorted list of subsets like we were able to achieve with some prior algorithms. Instead, it counts the number of times each candidate predictor appears in an optimal subset. For example, we might see that X2, X7, and X35 form an optimal subset for some region; X3, X7, and X21 form another optimal subset, X7 and X94 form another, and so forth. X7 appeared in an optimal subset 3 times, while each of the other subset members appeared just once. So it looks like X7 is on its way to becoming popular and heading up the popularity list.

This does not mean that X7 alone is valuable. In fact, it may be (and often is) that X7 alone is worthless; it's value is only in conjunction with other candidates. This is why LFS is superior to many other feature selection algorithms, which often rely on some form of stepwise selection and hence ignore individually worthless candidates. But *this property of reliance is not a problem*. The reason is that most modern prediction models, if given a list of the most popular predictors, can sort out the complex relationships between them and perform well. All they need is preprocessing to weed out the worthless candidates, so they are not overwhelmed.

Specifying the Test Parameters

When LFS is selected from the *Tests* menu, the following items must be specified:

The *leftmost column* specifies the set of predictor candidates. Multiple candidates can be selected by dragging the mouse cursor across a block, or by clicking the first candidate in a block, holding the Shift key, and clicking the last candidate in the block. Single candidates can be toggled on/off by holding the Ctrl key while clicking on the variable.

The *Target* column is used to select the target variable. The target is partitioned into bins that are as equal in size as possible. The user must specify the number of bins to employ for each, and unless the dataset is huge the default of three bins is frequently appropriate.

Max kept is the maximum number of variables ever employed in a metric space (subset of candidates). In general it is best to make this as small as possible, consistent with having enough variables simultaneously present to provide predictive power. In my experience, setting this to more than 5 is rarely, if ever needed. The default is 3.

Iterations is the number of LFS algorithm iterations to obtain good weight estimates. Run time is heavily impacted by this number. The point of diminishing returns is reached quickly; in many cases 2 is sufficient, and 3 is almost certainly more than enough for all but the most critical applications. The default is 3.

Binary random is the number of random trials employed to convert the floating-point usage flags to binary flags. More is better, but the default of 500 should be plenty for most applications, although if there are a great many variables this should be increased. It has a modest but not severe impact on run time for most applications.

Beta trials specifies the number of search points for optimizing the relative importance of intra-class versus inter-class separation discussed earlier in this section. The default of 20 should be sufficient for the vast majority of applications. It has a modest but not severe impact on run time for most applications.

Replications defaults to zero, in which case no Monte-Carlo Permutation Test is performed. However, it is usually best to set this to at least 100, and perhaps as much as 1000, so that solo and unbiased group p-values will be computed. Note that the minimum possible p-value is the reciprocal of the number of permutations. So, for example, if the user specifies 100 permutations, the minimum p-value that can appear is 0.01. Run time of this test is linearly related to the number of permutations.

The user must choose either *Complete* or *Cyclic* permutations if a Monte-Carlo Permutation Test is to be performed. If the user is confident that there is no dependency as described earlier in this document, then *Complete* should be used; it is the traditional approach which does a complete random shuffle for each permutation. However, if there is dependency, this type of shuffling will produce underestimation of *p-values*, a very dangerous situation. If the dependency is serial (the data is a time series and the dependency is among samples close in time) then a considerable improvement in the situation can be obtained by using *Cyclic* permutation. In this type of shuffle, the time order of the target is kept intact except at the ends by rotating the target with end-point wraparound. Shuffling this way preserves most of the serial dependency in the permuted target, which makes the algorithm more accurate. The *p-values* computed this way will generally be larger than those computed with complete shuffling, and hence less likely to lead to false rejection of the null hypothesis of no predictive power. But be warned that the cure is far from complete; computed *p-values* will still underestimate the true values, just not as badly.

Note that in most cases it is legitimate to use *Cyclic* permutation instead of *Complete* when there is no dependency. However, if the dataset is small, *Cyclic* permutation will limit the number of unique permutations and hence increase the random error inherent in the process. As long as the dataset is large, some users may prefer to use *Cyclic* permutation even if it is assumed that there is no serial dependency; in case there really is hidden serial dependency, this is a cheap insurance policy. Still, the best practice is to make sure that the data does not contain dependency and then use *Complete* permutation. Relying on *Cyclic* permutation to take care of dependency problems is living dangerously. And if the dataset contains fewer than 1000 or so cases, use of *Cyclic* permutation is not recommended unless it is necessary to handle dependency.

Important note: If you perform a Monte-Carlo permutation test, please see the discussion of solo and unbiased p-values that begins on Page 5 and continues onto the next page. That discussion covers vital issues related to what these figures mean, as well as when they are and are not valid.

CUDA note: As of Version 1.81, LFS will by default use CUDA-capable video hardware if present. This results in a speed increase of 1 or even 2 orders of magnitude if there are several thousand cases and not more than a few hundred variables. In other situations, CUDA may slow processing due to its overhead, and might better be disabled by clicking File/Use CUDA to make the check mark disappear.

An Example of Local Feature Selection

I created a dataset consisting of about 4000 cases and 10 variables, X0 through X9. Each random variable is uniformly distributed on $[-1, 1]$. Variables X3 and X4 determine the class. A case is in one class if X3 and X4 are both positive or both non-positive. The case is in the other class if one of these variables is positive and the other is not. This is a very difficult problem for many feature selection algorithms because the marginal distributions of these variables are identical for both classes, and the nature of the relationship between one of the variables with the class is determined by the value of the other variable. Here is the output of the LFS algorithm:

```
*****
*
* Computing Local Feature Selection for optimal predictor subset
*   10 predictor candidates
*   5 predictors at most will define a metric space
*   2 target bins
*   3 iterations of LFS algorithm
*  500 random trials for real-to-binary f conversion
*   20 trial values for beta optimization
*  100 replications of complete Monte-Carlo Permutation Test
*
*****
```

-----> Percent of times selected <-----

| Variable | Pct | Solo pval | Unbiased pval |
|----------|-------|-----------|---------------|
| X3 | 96.26 | 0.0100 | 0.0100 |
| X4 | 69.62 | 0.0100 | 0.0100 |
| X0 | 4.66 | 1.0000 | 1.0000 |
| X1 | 2.94 | 1.0000 | 1.0000 |
| X6 | 2.29 | 1.0000 | 1.0000 |
| X7 | 1.76 | 1.0000 | 1.0000 |
| X9 | 1.13 | 1.0000 | 1.0000 |
| X8 | 0.58 | 1.0000 | 1.0000 |
| X2 | 0.53 | 1.0000 | 1.0000 |
| X5 | 0.39 | 1.0000 | 1.0000 |

It's a little curious that X3 was selected somewhat more often than X4, when they have identical roles in predicting the class, but I've seen this happen often. It's undoubtedly a random occurrence that would change with a different random set of cases. What is certainly clear is that these two variables are selected vastly more often than their worthless competitors. Also, the computed solo and unbiased p-values are impressive, leaving no doubt about the conclusion reached by the algorithm.

Appendix: Version Updates

- 1.0 Univariate mutual information between predictor candidates and a single target

 Bivariate mutual information between a pair of predictor candidates and one or more target candidates
- 1.1 Added the option of uncertainty reduction instead of mutual information for bivariate mutual information, in order to accommodate targets with widely differing entropies
- 1.2 Peng, Long and Ding (2005) “Feature Selection Based on Mutual Information: Criteria of Max-Dependency, Max-Relevance, and Min Redundancy” algorithm implemented to select an optimal subset of predictors based on maximum relevance at predicting the target while simultaneously minimizing redundancy within the predictor set.
- 1.3 Hidden Markov models are defined using up to three predictors, without regard to a target. Then these models are sorted according to the multivariate correlation of their state probability vectors with a user-supplied target variable.
- 1.4 The univariate mutual information test now prints a new column: the probability that a selected candidate will have XVAL out-of-sample mutual information less than or equal to the median out-of-sample mutual information for all candidates.
- 1.5 One or more time series are examined for a break in their mean using the Mann-Whitney U-test. The user specifies how far in recent history to look back for a break. The test includes compensation for examining more than one series simultaneously, as well as compensation for repeating the test as time passes and new values for the series become available.
- 1.6 Feature Weighting as Regularized Energy-Based Learning (FREL): A recent development for feature ranking and weighting that is excellent for low-noise, high-dimension, small-sample-size applications.
- 1.7 Forward selection, as well as optional backward refinement, is used to find a relatively small subset of a very large set of variables such that the principal components of this subset capture the most variance possible from a subset of that size. This is valuable when faced with an extremely large set of predictors.

- 1.8 LFS (Local Feature Selection) for identifying predictors that are optimal in localized areas of the feature space but may not be globally optimal. Such predictors can be effectively used by nonlinear models but are neglected by many other feature selection algorithms.
- 1.81 CUDA computation added to the LFS algorithm, resulting in huge speed increase for problems with a large number of cases.
- 1.82 Fixed a serious non-thread-safe bug in a random number generator. Under certain unusual but possible conditions this could compromise Monte-Carlo permutation test results, especially for cyclic permutation.