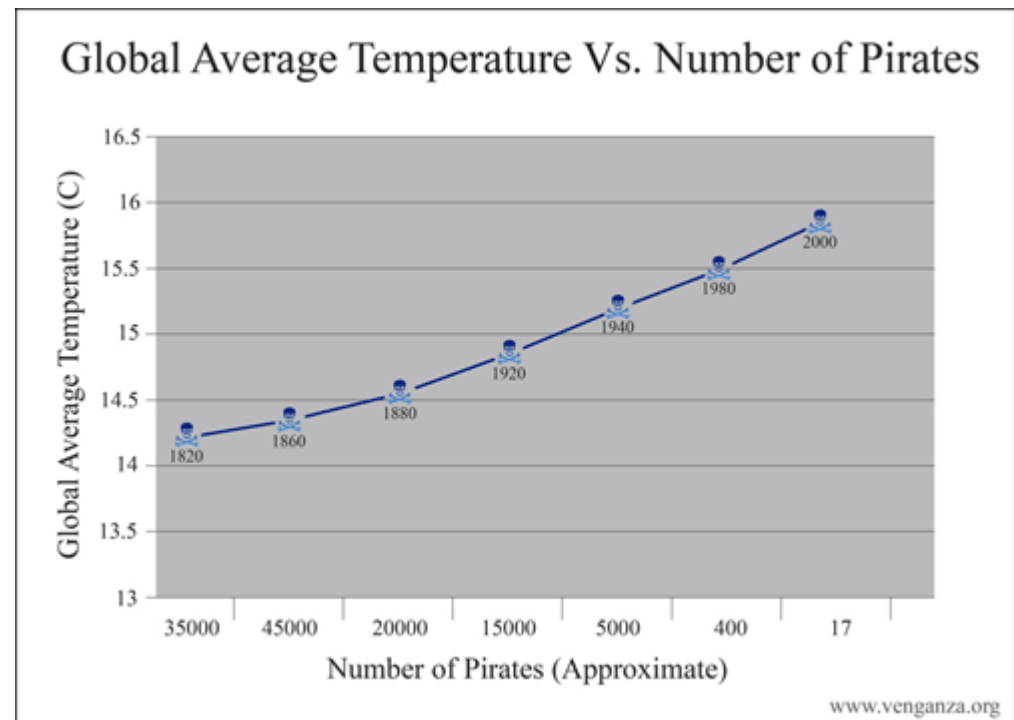


Physics 509: Non-Parametric Statistics and Correlation Testing

Scott Oser
Lecture #19



What is non-parametric statistics?

Non-parametric statistics is the application of statistical tests to cases that are only semi-quantitative. This includes:

1) When the data themselves are not quantitative. This can take various forms:

- A. simple classification or labelling of data
- B. rankings of data, without exact quantitative values
- C. “interval” data, in which relative differences between quantities have exact meaning, but the absolute values are meaningless

2) When the underlying distribution from which the data is drawn is not known. In this case the data may be quantitative, but we want to develop statistics that make no assumptions about or are at least insensitive to the details of the parent distribution.

Why use non-parametric statistics?

Obviously sometimes we have no choice---we may be handed non-quantitative data to start with.

While the physical sciences makes less use of non-parametric statistics than for example the social sciences do, you should be knowledgeable about and able to evaluate these other techniques.

Even if your data is perfectly quantitative and you *think* you know the underlying distribution from which it is drawn, a non-parametric statistic can cover your ass in case you're wrong about that!

General conclusion: non-parametric statistics make few model assumptions, and so simultaneously are more robust but less powerful than equivalent parametric tests.

Bayesian non-parametric statistics?

Bayesian methods are not easily made non-parametric. To use Bayes' theorem, you must have a quantitative model of how the data is distributed under various hypotheses and be prepared to evaluate each hypothesis quantitatively.

But often if you don't know the underlying distribution, you can get away with parametrizing the underlying distribution to cover all likely cases and try fitting for the shape of the distribution. Obviously this produces a lot of nuisance parameters and hence Occam factors (and so weakens the power of the test), but it can be done.

Goodness of fit: the run test

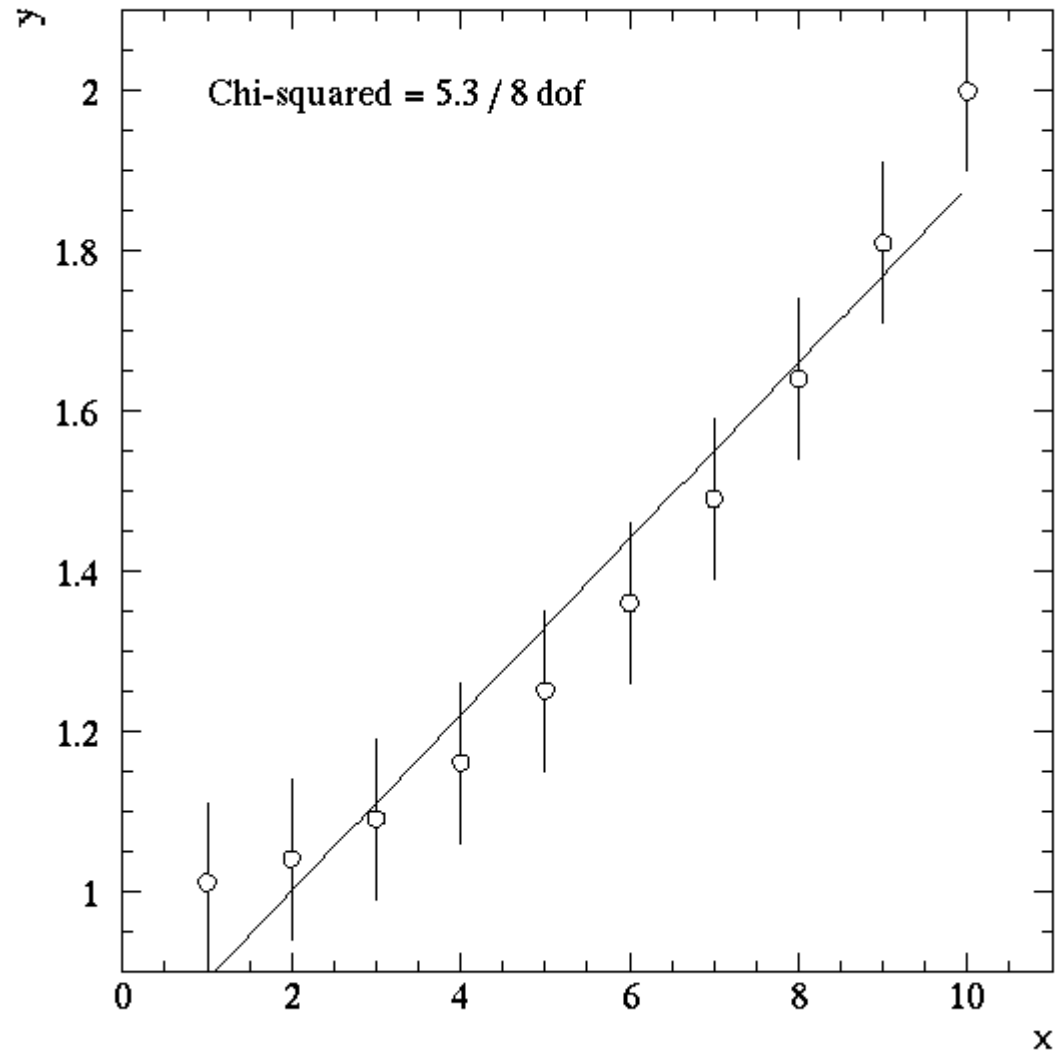
In spite of the excellent χ^2 , you should be suspicious of the linear fit to this data.

Notice the pattern of how some points are above the line, and others below:

HHLLLLLHH

There's an obvious pattern here, and in fact a parabolic fit would be much better.

Notice that a χ^2 test makes no use of the sign of the deviations since they are all squared.



Calculating the run test statistic

In the run test you simply count how many different “runs” there are in the data. In this data sample there are 3 runs: AABBBBBBAA

And here's one with 11 runs: ABAAABBAABBBABA

Given how many high data points and low data points there are, we can predict the probability of observing r runs:

$$P(r|r \text{ even}) = \frac{2 \binom{N_A - 1}{r/2 - 1} \binom{N_B - 1}{r/2 - 1}}{\binom{N}{N_A}}$$

$$P(r|r \text{ odd}) = \frac{\binom{N_A - 1}{(r-3)/2} \binom{N_B - 1}{(r-1)/2} + \binom{N_A - 1}{(r-1)/2} \binom{N_B - 1}{(r-3)/2}}{\binom{N}{N_A}}$$

The run test statistic: approximate formula

The mean and variance for the run test statistic are:

$$\langle r \rangle = 1 + \frac{2 N_A N_B}{N}$$
$$V(r) = \frac{2 N_A N_B (2 N_A N_B - N)}{N^2 (N - 1)}$$

If N_A and N_B are >10 , the distribution is approximately Gaussian.

For the case of AABBBBBBAA, there are just 3 runs. The probability of getting 3 or fewer runs is 0.87%.

The run test is complementary to a chi-squared goodness of fit, and provides additional information on the goodness of fit, although it is a pretty weak test.

The two-sample problem: the run test

The same kind of run test that provided a goodness of fit can be used to compare two data sets. Parametrically we could ask:

- 1) Do these distributions have the same mean? (t-test)
- 2) Do these distributions have the same variance? (F-test)

These tests only compare two specific features of the distributions, and rely on the assumption of Gaussianity to get their significance!

The two-sample run test on the other hand is perfectly general. Consider two sets of data that are graded on the same scale:

A = 1.3, 4.8, 4.9, 6.2, 9.8, 12.1

B = 0.2, 2.2, 3.9, 5.0, 7.8, 10.8

The only thing we assume is that events of both types can all be ordered together.

The two-sample problem: the run test

A = 1.3, 4.8, 4.9, 6.2, 9.8, 12.1

B = 0.2, 2.2, 3.9, 5.0, 7.8, 10.8

Now merge the list, keeping track of the order of A's and B's:

BABBAABABABA

Just calculate the observed number of runs and compared to the expected number.

The test works best when $N_A \approx N_B$.

It is an extremely general test, and as expected isn't all that powerful.

An alternate to the two-sample run test is to do a KS test between the two sets of data.

The sign test for the median

An evil physics professor reports that the median grade on the midterm was 72% while the average was 66%, but doesn't show the whole distribution. As a TA you poll the students in your tutorial section, who report receiving the following grades:

80, 36, 70, 40, 52, 92, 79, 66, 62, 59

Is the performance of your tutorial section representative of the class as a whole? If your section did better than the class as a whole, ask for a raise. If they did worse, you can complain to the department for having put the worst students in your section!

For the students in your section, the average score was 63.6% while the median was 64%. Is this significant?

The sign test for the median

Your first inclination may be to look at the average. But how is it distributed? Without knowing the distribution you have no idea, and with $N=10$ you can't rely on the central limit theorem to bail you out. So you cannot evaluate a significance!

But 50% of all scores should be above the median if the distributions are identical. Look at the scores and compare each to the median:

80,	36,	70,	40,	52,	92,	79,	66,	62,	59
+	-	-	-	-	+	+	-	-	-

7 of 10 students performed worse than the median. Use a binomial distribution to get the probability of getting a result this bad, or worse:

$$P = \frac{1}{2^{10}} \left(1 + \binom{10}{1} + \binom{10}{2} + \binom{10}{3} \right) \approx 0.172$$

You're off the hook!

The (Wilcoxon-)Mann-Whitney Test, aka the U test, aka the rank sum test

Question: do two distributions have the same median? If you wanted a completely general way to ask whether two distributions are the same, you could use a run test or a KS test. But what if you just want to know if the medians are different? This is a more restricted question and so a more powerful test exists.

You start out as in the run test:

A = 1.3, 4.8, 4.9, 6.2, 9.8, 12.1

B = 0.2, 2.2, 3.9, 5.0, 7.8, 10.8

Now merge the list, keeping track of the order of A's and B's:

BABBAABABABA

Keep track of any ties.

The (Wilcoxon-)Mann-Whitney Test, aka the U test, aka the rank-sum test

Now add up the numerical ranks of all of the A's:

BABBAABABABA

$$U_A = 2 + 5 + 6 + 8 + 10 + 12 = 43$$

If there are any ties, assign each position the average of the tied ranks (ex. 1st, 2nd, 3.5th, 3.5th, 5th)

Equivalently you could add up all of the B ranks. Note that $U_A + U_B = N(N+1)/2$

If the medians are identical, then we expect $U_A \approx N_A(N+1)/2$ and $U_B \approx N_B(N+1)/2$.

What is the distribution of U_A ?

The distribution of U, the Wilcoxon-Mann-Whitney rank-sum statistic

What is the distribution of U_A ?

1. Simple solution: Monte Carlo it!
2. Use lookup tables
3. For large N_A and N_B (>12), there is an approximate Gaussian distribution:

$$\langle U_A \rangle = \frac{N_A(N+1)}{2}$$

$$\text{Var}(U_A) = \frac{N_A N_B (N+1)}{12}$$

Note that if you apply the Gaussian approximation, always remember that U_A is an integer. To get the probability that $U < U_A$, add +0.5 to U_A first. To get the probability that $U > U_A$, subtract -0.5 instead. Then use a lookup table from a normal distribution to get the probability.

The Matched Pairs Sign Test

Consider the following survey:

A. Professor Oser is an effective instructor.

(1) strongly agree (2) agree (3) neutral (4) disagree (5) strongly disagree

B. I would take a class from Professor Oser again.

(1) strongly agree (2) agree (3) neutral (4) disagree (5) strongly disagree

C. Professor Oser's lectures were clear and well-organized.

(1) strongly agree (2) agree (3) neutral (4) disagree (5) strongly disagree

The quiz is scored by adding up the numbers from all three questions.

This quiz is given to 12 students in the middle of the term, and then again right after the midterm is handed back.

The Matched Pairs Sign Test

After the midterm the mean score has actually gotten slightly worse!

But can you really trust such a non-quantitative system? It's all very subjective with significant student-by-student variation.

Notice that 10 of 12 scores actually got better. The probability of this sort of improvement happening by chance is only 1.9%.

Using matched pairs (“before” vs. “after”) removes student-by-student variability. A sign test gives a non-parametric indication of how the midterm improved the class' mood.

Interesting to note that the two students who became more negative had bad scores to start!

Student	Before	After
1	10	15
2	8	7
3	5	4
4	4	3
5	8	15
6	9	7
7	10	9
8	6	5
9	5	4
10	7	6
11	11	10
12	6	5
Total	7.4± 2.2	7.5 ± 3.9

χ^2 test on contingency tables

A political scientist surveys the party affiliations of physicists and economists:

	Liberal	Conservative	NDP
Economists	45	63	12
Physicists	50	21	35

Is this data consistent with the hypothesis that party affiliation is identical between the two groups? $P(X,Y)=f(X)g(Y)$

Form a contingency table: figure out how many you expect in each box, using the average distributions. For example, out of 226 total entries, 95 are Liberals. So of the 120 economists we expect $120 \cdot (95/226) = 50.44$ to be Liberals.

	Liberal	Conservative	NDP
Economists	50.4	44.6	25.0
Physicists	44.6	39.4	22.0

χ^2 test on contingency tables

	Liberal	Conservative	NDP
Economists	45/50.4	63/44.6	12/25.0
Physicists	50/44.6	21/39.4	35/22.0

Form a χ^2 between the observed and expected values:

$$\chi^2 = \frac{(45 - 50.4)^2}{50.4} + \frac{(63 - 44.6)^2}{44.6} + \frac{(12 - 25.0)^2}{25.0} + \frac{(50 - 44.6)^2}{44.6} + \frac{(21 - 39.4)^2}{39.4} + \frac{(35 - 22.0)^2}{22.0} = 31.9$$

This should follow a χ^2 distribution with $(\text{rows}-1)(\text{columns}-1) = 2$ d.o.f.

Prob ($\chi^2 > 31.9$) = 1.2×10^{-7} . Highly significant!

By the way, this data is all made up.

Fisher exact test for 2x2 contingency table

A χ^2 test fails when the number of events per bin is small (any table entry less than 5 makes the whole test suspect). For the special case of a 2x2 table, there is however an exact formula that works for any sample size. Consider the hypothesis that A correlates with B, such as in this table:

		A		
		No	Yes	
B	Yes	2	7	9
	No	8	2	10
		10	9	19

0	9	9	1	8	9	2	7	9	3	6	9	4	5	9	...	8	1	9	9	0	9
10	0	10	9	1	10	8	2	10	7	3	10	6	4	10		2	8	10	1	9	10
10	9	19	10	9	19	10	9	19	10	9	19	10	9	19		10	9	19	10	9	19

We list all the possibilities from most correlated to most anticorrelated, as above.

Fisher exact test for 2x2 contingency table

We can now calculate the probability of any particular configuration from combinatorical considerations:

a	b	a+b
c	d	c+d
a+c	b+d	N

$$P = \frac{(a+b)!(c+d)!(a+c)!(b+d)!}{N!a!b!c!d!}$$

1	8	9
9	1	10
10	9	19

$$P = \frac{9!10!10!9!}{19!1!9!8!1!} = 9.7 \times 10^{-4}$$

2	7	9
8	2	10
10	9	19

$$P = \frac{9!10!10!9!}{19!2!8!7!2!} = 0.0175$$

You have the option of doing either a one-sided test (is the data skewed in a particular direction) or a two-sided test (is the data skewed from the null hypothesis expectation).

Testing linear correlation: Gaussian distributions

Suppose that X & Y are two random variables drawn from a 2D Gaussian distribution. We wish to test whether the correlation coefficient $\rho=0$.

A parametric test: calculate an estimator for ρ and see if it is consistent with zero.

$$r = \frac{\sum_{i=1}^N (X_i - \bar{X})(Y_i - \bar{Y})}{\sqrt{\sum_{i=1}^N (X_i - \bar{X})^2 \sum_{i=1}^N (Y_i - \bar{Y})^2}}$$

$$t = \frac{r \sqrt{N-2}}{\sqrt{1-r^2}}$$

is distributed as a t distribution with $N-2$ d.o.f.
if $\rho=0$

Spearman's Correlation Coefficient

Spearman's correlation coefficient is a powerful non-parametric method for testing for correlation. You don't have to assume a particular distribution for X or Y .

We start with N pairs of (x,y) .

First, rank all of the x measurements and all of the y measurements from 1 to N . Let X_i be the rank of the i^{th} value of x (so X_i is an integer between 1 and N), and similarly for Y_i .

Calculate Spearman's correlation coefficient:

$$r_s = 1 - 6 \frac{\sum (X_i - Y_i)^2}{N^3 - N}$$

Using Spearman's Correlation Coefficient

Spearman's correlation coefficient returns a value between -1 and 1. You either have to use a lookup table to tell you what value is significant given any given number of data pairs, or else if $N > 30$ you can use an asymptotic approximation to a t distribution with $N-2$ degrees of freedom.

$$t_s = \frac{r_s \sqrt{N-2}}{\sqrt{1-r_s^2}}$$

Spearman's correlation coefficient is a very efficient non-parametric test. This means that it's only a little less powerful than the classical method of calculating r , the estimator of ρ . For example, for Gaussian data the parametric test requires ~91 data points to reject the null hypothesis at the same significance that Spearman's test achieves with 100 data points.

Comparison of Spearman's Correlation Test to the Parametric Test for fake data sets

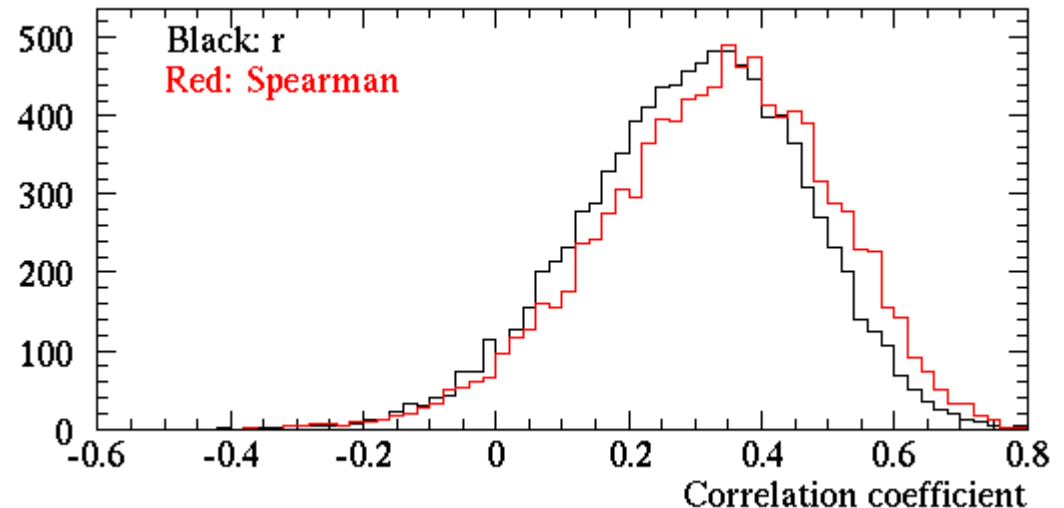
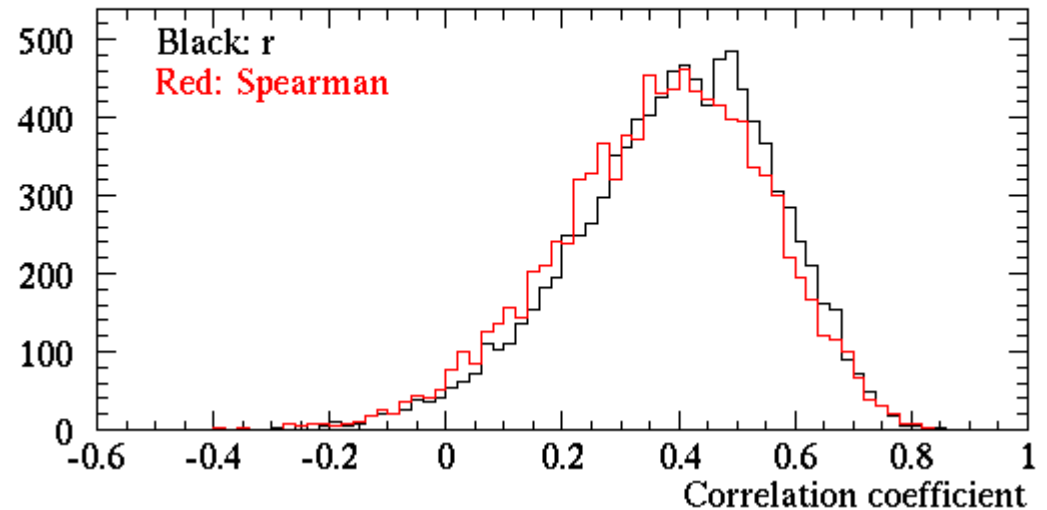
Top figure: comparison of Pearson's to Spearman's correlation coefficient for linear covariance:

Fraction detected at 95% CL by:
Pearson's (parametric): 65%
Spearman's test: 60%

Bottom figure: correlation is $y \propto x^{1/3}$.

Fraction detected at 95% CL by:
Pearson's (parametric): 42%
Spearman's test: 50%

Spearman's test actually was better for this non-linear correlation: it only checks whether y varies monotonically with x , not just linear dependence.



Do sumo wrestlers cheat?



Data taken from

“Winning Isn't Everything:
Corruption in Sumo
Wrestling”, by Mark Duggan
and Steven D. Levitt

<http://www.nber.org/papers/w7798>

Structure of sumo tournaments

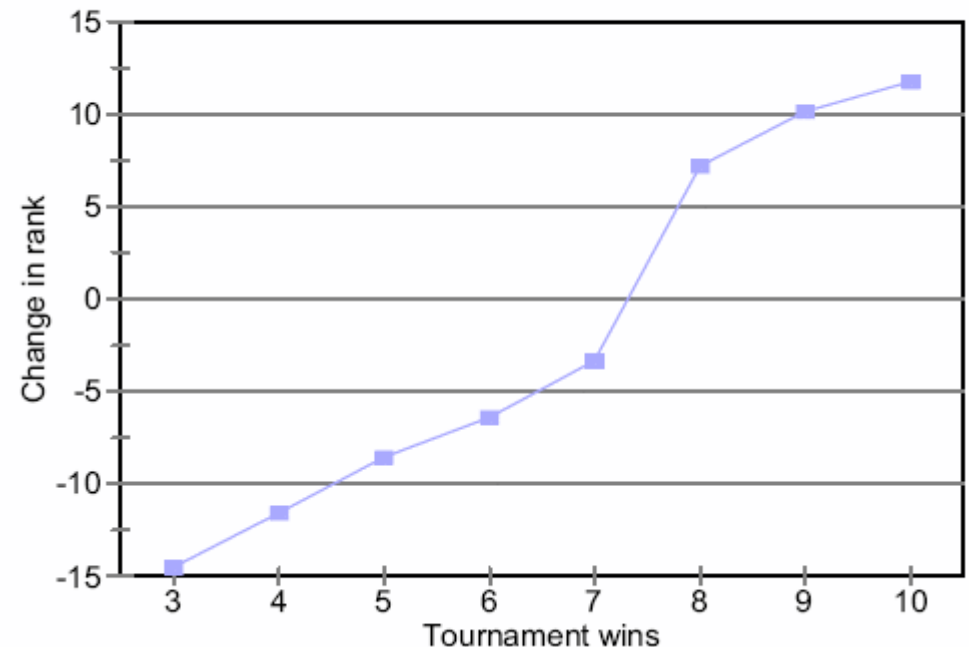
A sumo tournament features ~70 wrestlers, each of whom compete in 15 randomly assigned matches over the course of the tournament.

All wrestlers who finish with a winning record (8-7 or better) are guaranteed to rise in the rankings, while a losing record makes you drop in the ranking.

Ranking -> earnings

The plot at the right shows there is a non-linear “kink”---you gain more from winning your 8th match than from any other match.

Figure 1: Payoff to Tournament Wins



Bribe incentive: if your record is 7-7 on the last day, pay opponent to lose.

Sumo data

The following table shows the outcomes of 558 matches in which exactly one opponent had a 7-7 record on the final day.

We don't know the expected distribution of winning percentages---for example, how often do we expect a 7-7 wrestler to beat a 5-9 wrestler?

But we can safely assume that $P < 0.5$ for the 7-7 wrestler to beat a wrestler with a better record!

Opponent's Record	Matches	Wins	Winning Percentage
0 to 2	32	23	71.9
3 to 5	99	77	77.8
6	101	77	76.2
8	201	160	79.6
9	79	58	73.4
10 to 14	46	18	39.1
Total	558	413	74.0

The 7-7 wrestler won 236 of 326 matches against opponents with better records. At best we'd expect him to win half.

Probability of this: $>8\sigma$!

Sumo “binomial” distribution

The following graph shows the distribution of how many matches each wrestler wins, compared to the binomial expectation

Do 7-7 wrestlers just try harder?

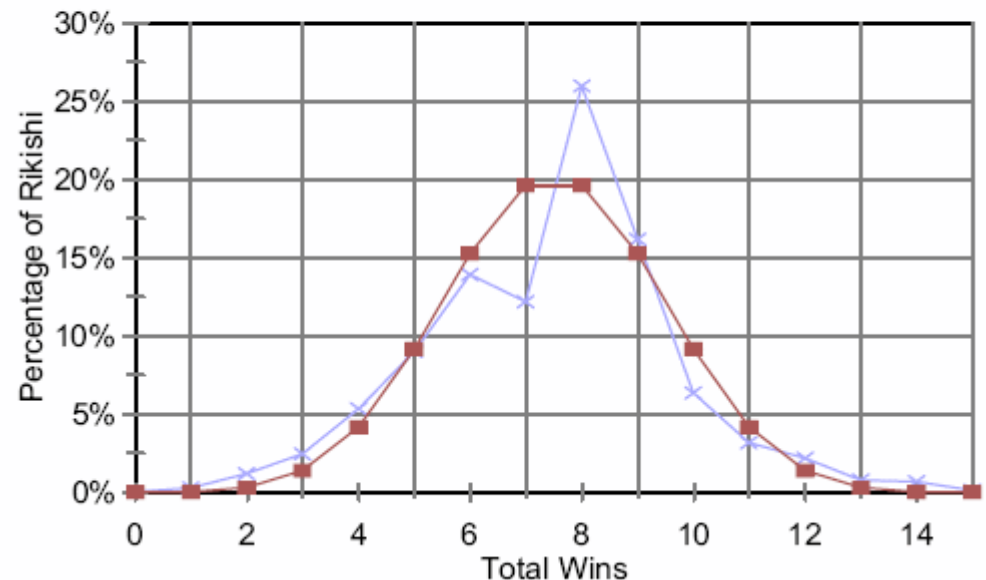
A. In rematches, the winning wrestler *loses* much more than expected. Payback?

B. Wrestlers who face each other more frequently show a larger effect. Collusion?

C. Certain sumo stables show much more effect than others.

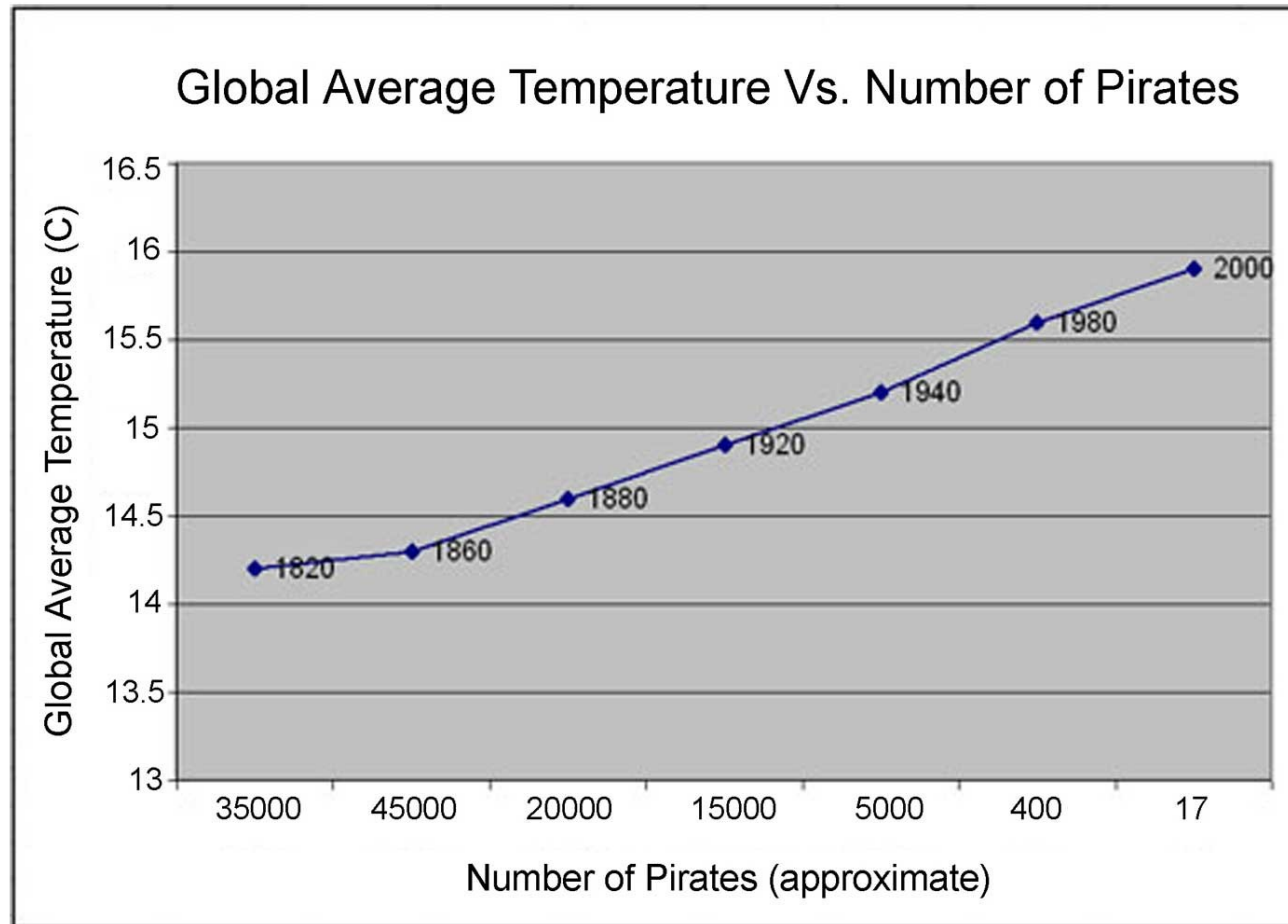
D. Sumo wrestlers about to retire do no better than average---no incentive to cheat then.

Figure 2: Wins in a Sumo Tournament
Actual vs. Binomial



Partial correlations & lurking third parameters

STOP GLOBAL WARMING: BECOME A PIRATE



WWW.VENGANZA.ORG

Partial correlations & lurking third parameters

It is observed that there is a strong correlation between the number of firefighters at a fire and the amount of damage done by the fire.

Should we send fewer firefighters to fires?

Suppose that we have the following random variables:

$$P(X) = \frac{1}{\sqrt{2\pi}\sigma} \exp\left[-\frac{1}{2} \frac{X^2}{\sigma^2}\right]$$

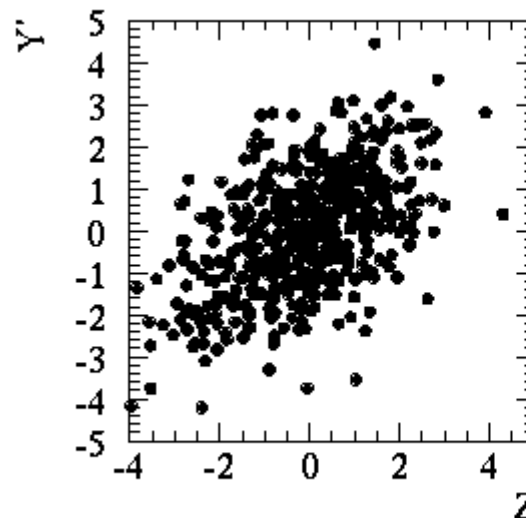
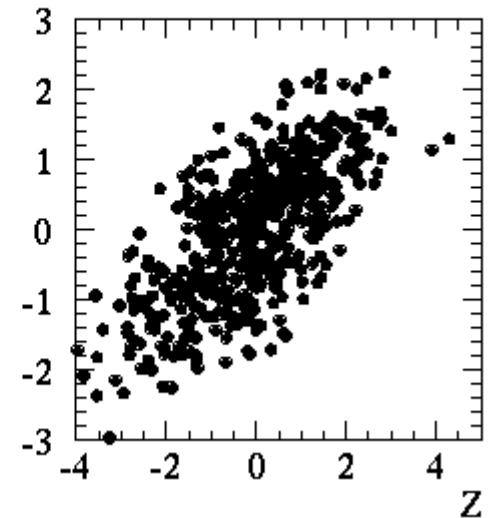
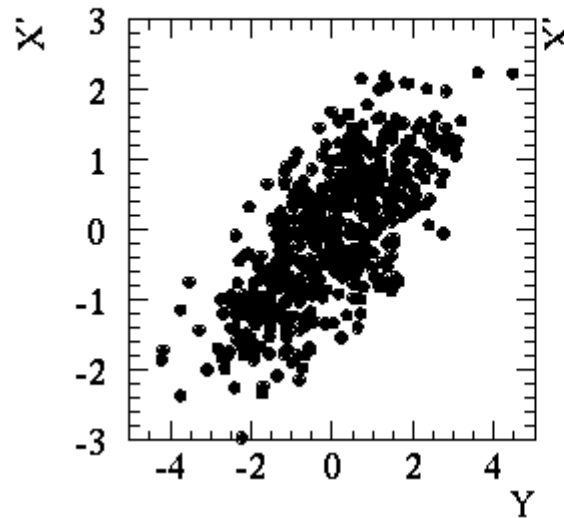
$$P(Y|X) = \frac{1}{\sqrt{2\pi}\sigma} \exp\left[-\frac{1}{2} \frac{(Y-X)^2}{\sigma^2}\right]$$

$$P(Z|X) = \frac{1}{\sqrt{2\pi}\sigma} \exp\left[-\frac{1}{2} \frac{(Z-X)^2}{\sigma^2}\right]$$

Partial correlations & lurking third parameters

All three variables are correlated with each other.

From looking at the scatterplot of Y vs. Z, we would conclude there was a significant correlation. But is there any causation---does larger Z cause Y to be larger as well?



First-order partial correlation coefficient

Calculate the regular correlation coefficients between all pairs of variables. Then calculate:

$$r_{xy.z} = \frac{r_{xy} - r_{xz}r_{yz}}{\sqrt{(1-r_{xz}^2)(1-r_{yz}^2)}}$$

This attempts to provide the “direct” correlation between x and y, not including any indirect correlation through z.

The standard error on the partial correlation coefficient is:

$$\sigma_{r_{xy.z}} = \frac{1 - r_{xy.z}^2}{\sqrt{N - m}}$$

where m is the number of variables involved.

Results with the first-order partial correlation coefficient

From Monte Carlo simulation with $N=25$

$$r_{xy} = 0.70 \pm 0.11$$

$$r_{xz} = 0.70 \pm 0.11$$

$$r_{yz} = 0.49 \pm 0.16$$

$$r_{xy.z} = 0.57 \pm 0.14$$

$$r_{xz.y} = 0.57 \pm 0.14$$

$$r_{yz.x} = 0.00 \pm 0.21$$

The partial correlation coefficients identify the fact that Y depends directly on X, and so does Z, but that Y does not directly depend on Z even though they have a significant correlation.

In other words, all of the correlation between Y and Z is explained by their mutual correlation with X.

Note that these results at some level assume linear correlations, and so this may not be a truly non-parametric test.