

Confidence intervals for rank order statistics: Somers' D , Kendall's τ_a and their differences

Roger Newson (Guy's, King's and St Thomas' School of Medicine)

- Somers' D and Kendall's τ_a
- Why rank order statistics?
- Why confidence intervals?
- The program `somersd`
- An example (from the `auto` data)

Parameters behind “non-parametric” statistics

Suppose (X_1, Y_1) and (X_2, Y_2) are sampled independently from the same bivariate population. Kendall's τ_a (the signed difference covariance) is

$$\tau_{XY} = E[\text{sign}(X_1 - X_2) \text{sign}(Y_1 - Y_2)],$$

and is the difference between the two probabilities of concordance and discordance.

Somers' D (the corresponding regression coefficient) is

$$D_{YX} = \tau_{XY} / \tau_{XX},$$

and is the difference between the corresponding *conditional* probabilities, given unequal X -values.

These parameters are estimated by corresponding sample statistics, which are used to test null hypotheses. If X is binary, then ranksum tests $H_0 : D_{YX} = 0$.

Why rank order statistics?

No statistical method is simply “robust”. However, Kendall’s τ_a and Somers’ D are robust to:

- Extreme values. (These often “throw” classical regression and correlation coefficients.)
- Non-linearity. (They are not affected by monotonic transformations, and can be ± 1 for perfect non-linear relationships.)

They are useful as a preliminary to homing in on a particular regression model. (And losing a section of the skeptical public, who do not believe the model assumptions.)

Why confidence intervals?

- They *might* discourage people from arguing that a high P -value proves a null hypothesis.
- For continuous data, we often have Greiner's correspondence between Kendall's τ_a and Pearson's ρ ,

$$\rho = \sin\left(\frac{\pi}{2} \tau\right).$$

So, given a CI for τ , we can define an “outlier-resistant” CI for ρ .

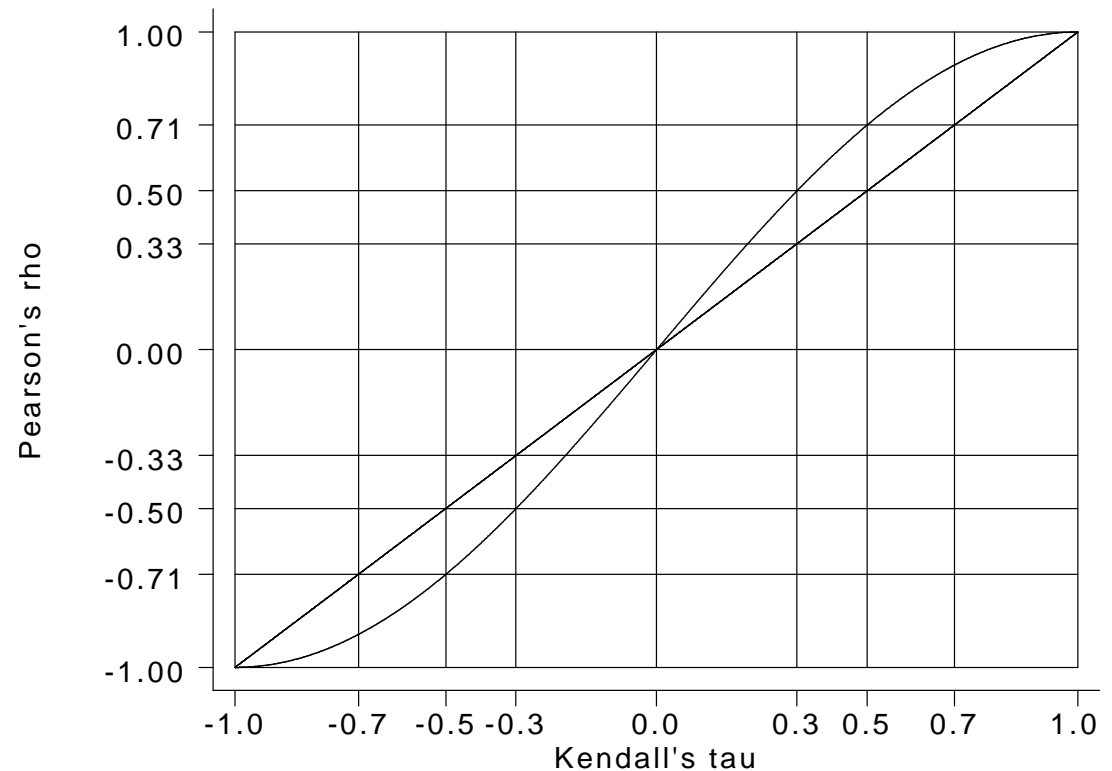
- We might want to know CIs for differences, such as

$$\tau_{XY} - \tau_{WY} \quad \text{or} \quad D_{YX} - D_{YW},$$

where Y is an outcome variable and W and X are competing predictor variables. This is because a larger τ_a cannot be secondary to a smaller τ_a .

Greiner's relation: $\rho = \sin\left(\frac{\pi}{2}\tau\right)$

- This holds under the bivariate normal distribution, and many other continuous distributions.
- It therefore saves the labour of hunting for a pair of transformations.
- Kendall's τ is less vulnerable than Pearson's ρ to “outliers”.
- *However*, τ is “less impressive”. ($\tau = \pm 1/3$ when $\rho = \pm 1/2$, and $\tau = \pm 1/2$ when $\rho = \pm 1/\sqrt{2}$.)



somersd - A program to calculate confidence intervals for rank order statistics

- Calculates jackknife confidence intervals for either Somers' D or Kendall's τ_a , between one variable X and a list of others $Y^{(1)} \dots Y^{(p)}$.
- Offers a choice of transformations: Fisher's z , Daniels' arcsine, Greiner's ρ , and the z -transform of Greiner's ρ .
- A `cluster()` option is available (mostly for measuring intra-class correlation, eg between twin sisters).
- Estimation results are saved as for a model fit, so differences can be estimated using `lincom`.

Example: Weight and fuel consumption in US and non-US cars (1)

In the auto data set, we use `somersd` (instead of `ranksum`) to assess US origin as a predictor of fuel consumption (gallons/mile) and weight (lbs.). We use Fisher's z -transform to compute asymmetric confidence limits for the two Somers' D s:

```
. somersd us gpm weight,tran(z)
```

```
Somers' D
```

```
Transformation: Fisher's z
```

```
Valid observations: 74
```

us	Coef.	Jackknife Std. Err.	z	P> z	[95% Conf. Interval]
gpm	.4937249	.1708551	2.890	0.004	.1588551 .8285947
weight	.9749561	.1908547	5.108	0.000	.6008878 1.349024

```
95% CI for untransformed Somers' D
```

	Somers_D	Minimum	Maximum
gpm	.45716783	.15753219	.67972072
weight	.75087413	.53768098	.87382282

We note that US cars (usually) are heavier than the rest, and consume more fuel per mile.

Example: Weight and fuel consumption in US and non-US cars (2)

Using `lincom`, we can show that US origin predicts weight better than it predicts fuel consumption. We estimate the difference (in z -units) between the two positive z -transformed Somers' D values:

```
. lincom weight-gpm
```

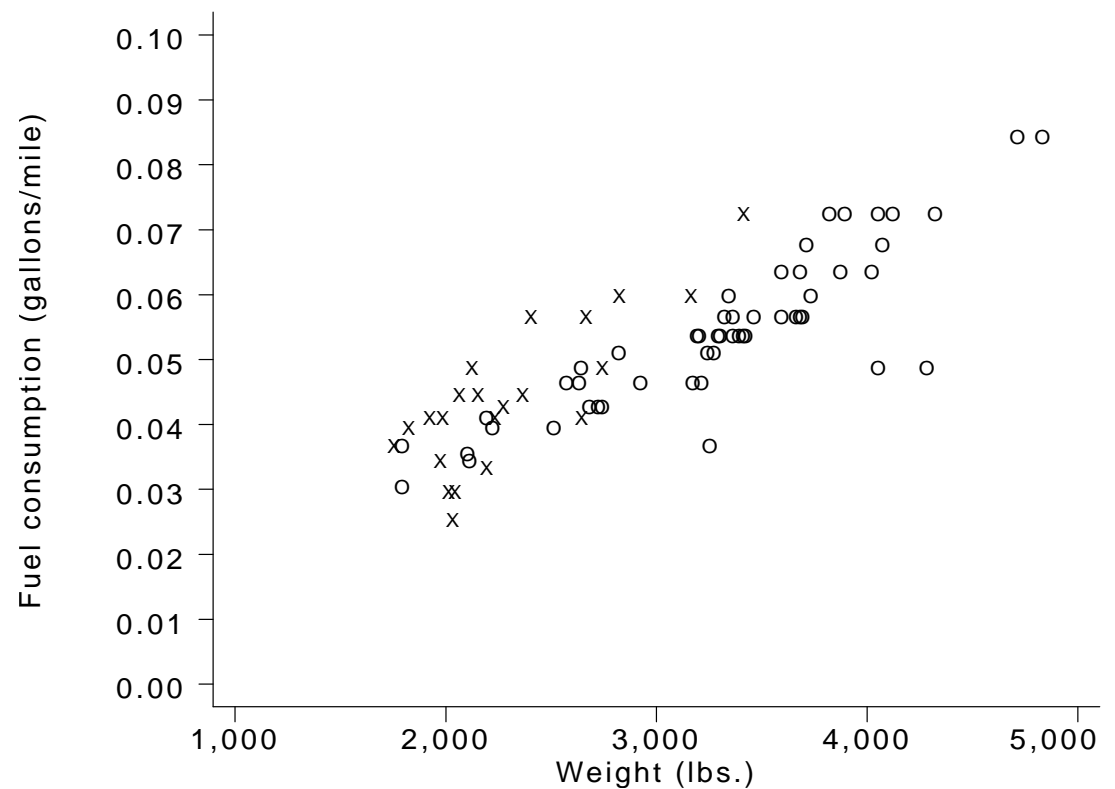
```
( 1) - gpm + weight = 0.0
```

us	Coef.	Std. Err.	z	P> z	[95% Conf. Interval]
(1)	.4812312	.1235452	3.895	0.000	.2390871 .7233753

The difference is positive. So, given two cars, one consuming fewer gallons to move more weight, the other consuming more gallons to move less weight, the *first* (the *more efficient* one) is more likely to be US-made.

Example: Weight and fuel consumption in US and non-US cars (3)

- Data points are US models (“O”) and non-US models (“X”).
- Given two models, one moving more mass with less gas, the other moving less mass with more gas, the *first* is more likely to be a US model.
- Using `somersd`, we can show this without contentious assumptions such as linearity, because the difference between the two Somers’ D s depends entirely on such discordant pairs.



Further reading

Newson, R. 2000. `somersd` - Confidence intervals for “non-parametric” statistics and their differences. *Stata Technical Bulletin* 55, in press. (To appear in May 2000.)

Edwardes, M. D. deB. 1995. A Confidence Interval for $\Pr(X < Y) - \Pr(X > Y)$ Estimated From Simple Cluster Samples. *Biometrics* 51: 571-578.

Kendall, M. G. and J. D. Gibbons. 1990. *Rank correlation methods*. 5th ed. New York: Oxford University Press.