

Statistics Refresher II

18th JRC Annual training on Composite Indicators and Scoreboards

William Becker

Joint Research Centre

Fundamentally, we are interested in how the random variable x_1 depends on the random variable x_2 .

Covariance

$$\operatorname{cov}(x_1, x_2) = E[(x_1 - \mu_1)(x_2 - \mu_2)]$$





Fundamentally, we are interested in how the random variable x_1 depends on the random variable x_2 .

Covariance

$$cov(x_1, x_2) = E[(x_1 - \mu_1)(x_2 - \mu_2)]$$





Fundamentally, we are interested in how the random variable x_1 depends on the random variable x_2 .

Covariance

$$cov(x_1, x_2) = E[(x_1 - \mu_1)(x_2 - \mu_2)]$$





Fundamentally, we are interested in how the random variable x_1 depends on the random variable x_2 .

Covariance

$$cov(x_1, x_2) = E[(x_1 - \mu_1)(x_2 - \mu_2)]$$





Fundamentally, we are interested in how the random variable x_1 depends on the random variable x_2 .

Covariance

$$cov(x_1, x_2) = E[(x_1 - \mu_1)(x_2 - \mu_2)]$$





Fundamentally, we are interested in how the random variable x_1 depends on the random variable x_2 .

Covariance

$$cov(x_1, x_2) = E[(x_1 - \mu_1)(x_2 - \mu_2)]$$





Fundamentally, we are interested in how the random variable x_1 depends on the random variable x_2 .

Covariance

$$cov(x_1, x_2) = E[(x_1 - \mu_1)(x_2 - \mu_2)]$$



Fundamentally, we are interested in how the random variable x_1 depends on the random variable x_2 .

Covariance

$$cov(x_1, x_2) = E[(x_1 - \mu_1)(x_2 - \mu_2)]$$

Correlation

$$R(x_1, x_2) = \operatorname{corr}(x_1, x_2) = \frac{\operatorname{cov}(x_1, x_2)}{\sigma_1 \sigma_2}$$

Standardises covariance so that $R \in [-1,1]$: 1 or -1 is perfect correlation, 0 is no correlation. Allows comparability.





Measures of dependence Correlation

$$R(x_{1}, x_{2}) = \operatorname{corr}(x_{1}, x_{2}) = \frac{\operatorname{cov}(x_{1}, x_{2})}{\sigma_{1}\sigma_{2}}$$

$$\operatorname{Coefficient of determination: } \mathbb{R}^{2}$$

$$R^{2}(x_{1}, x_{2}) = [\operatorname{corr}(x_{1}, x_{2})]^{2} = \left[\frac{\operatorname{cov}(x_{1}, x_{2})}{\sigma_{1}\sigma_{2}}\right]^{2}$$

$$\operatorname{Coefficient of determination: } \mathbb{R}^{2}$$

 R^2 is a measure of linear dependence. $R^2 \in [0,1]$: higher values indicate stronger dependence.





Very strong positive correlation R = 0.97

No positive correlation R = 0.02











12 JRC-COIN © | 2021 JRC Week on Composite Indicators and Scoreboards

Size of correlation	Interpretation
0.9 to 1 (-0.9 to -1)	Very high positive (negative) correlation
0.7 to 0.9 (-0.7 to -0.9)	High positive (negative) correlation
0.5 to 0.7 (-0.5 to -0.7)	Moderate positive (negative) correlation
0.3 to 0.5 (-0.3 to -0.5)	Low positive (negative) correlation
0 to 0.3 (0 to -0.3)	Negligible correlation

BUT: check for statistical significance!

"Given the sample size (number of points, countries, ...), is the correlation significantly different from zero?"

A quick test for significance is $|R| \ge \frac{2}{\sqrt{N}}$

E.g. Europe (N≈30) -> **0.35** Worldwide (N≈130) -> **0.18**

(p<0.05)



13 JRC-COIN © | 2021 JRC Week on Composite Indicators and Scoreboards



What happens to correlation when there are outliers?

correlation:

0.02 > 0.84

Be careful with outliers.

Always plot your data!



Finally

Correlation is *linear* dependence. Extensions or variations of correlation include:

- Spearman rank correlation (correlation coefficient between ranks)
- Kendall rank correlation (slightly more sophisticated version of Spearman)
- Nonlinear dependence measures: correlation ratio (first order sensitivity index) and friends

But

Linear correlation is usually a pretty good approximation. Using both is a safe bet.

Thank you



© European Union 2021

Unless otherwise noted the reuse of this presentation is authorised under the <u>CC BY 4.0</u> license. For any use or reproduction of elements that are not owned by the EU, permission may need to be sought directly from the respective right holders.



16 JRC-COIN © | 2021 JRC Week on Composite Indicators and Scoreboards