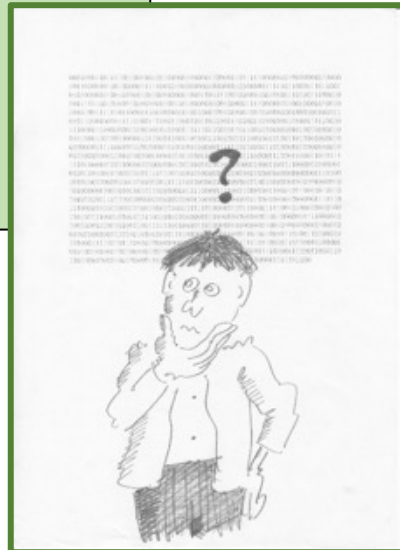


# Correlations and their statistical significance - cont

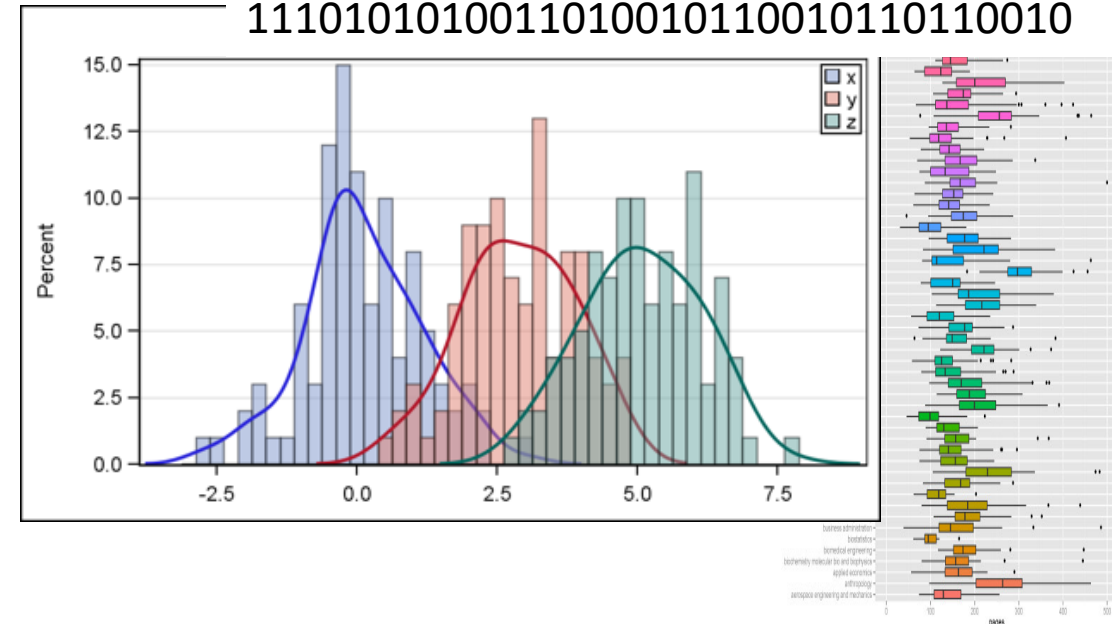
## Statistics and data analysis

Zohar Yakhini

IDC, Herzeliya



0010011101010100101010100100100010  
1010100010101111101011010011001001  
11101010100111010010110010110110010



# Outline

- Multinomial Correlations
- CI – an example
- Spurious correlations and intro to multiple testing

# Multinomial Distribution

- Let  $X \sim \text{MNom}(N, P)$  ,  $X = (X_1, X_2, \dots, X_d)$  .
- Example - roll a die  $N$  times, counts numbers on each face.
- What is  $d$ ?
- What is  $P$ ?
- What is  $E(X_i)$  ?
- What is  $V(X_i)$  ?
- Are these random variables collectively independent?
- Pairwise independent?

# Multinomial Distribution - covariances

Let  $X \sim \text{MNom}(N, P)$ ,  $X = (X_1, X_2, \dots, X_d)$ .

$$\text{Var}(X_i + X_j) = V(X_i) + 2\text{Cov}(X_i, X_j) + V(X_j)$$

Now observe that  $X_i + X_j \sim \text{BiNom}(N, p_i + p_j)$  and therefore, from the above identity we get:

$$\begin{aligned} 2\text{Cov}(X_i, X_j) &= \text{Var}(X_i + X_j) - V(X_i) - V(X_j) = \\ &= N[(p_i + p_j)(1 - p_i + p_j) - p_i(1 - p_i) - p_j(1 - p_j)] \\ &= -2Np_i p_j \end{aligned}$$

# Multinomial Distribution - correlations

Let  $X \sim \text{MNom}(N, P)$ ,  $X = (X_1, X_2, \dots, X_d)$ .

We saw that  $\text{Cov}(X_i, X_j) = -Np_i p_j$

Using the fact that  $\forall i \text{ } \text{Var}(X_i) = Np_i(1 - p_i)$  we can conclude that

$$\rho(X_i, X_j) = \frac{-\sqrt{p_i p_j}}{\sqrt{(1 - p_i)(1 - p_j)}}$$

For the case of a fair die we get, for example,  $\rho(X_1, X_2) = -\frac{1}{5}$ .

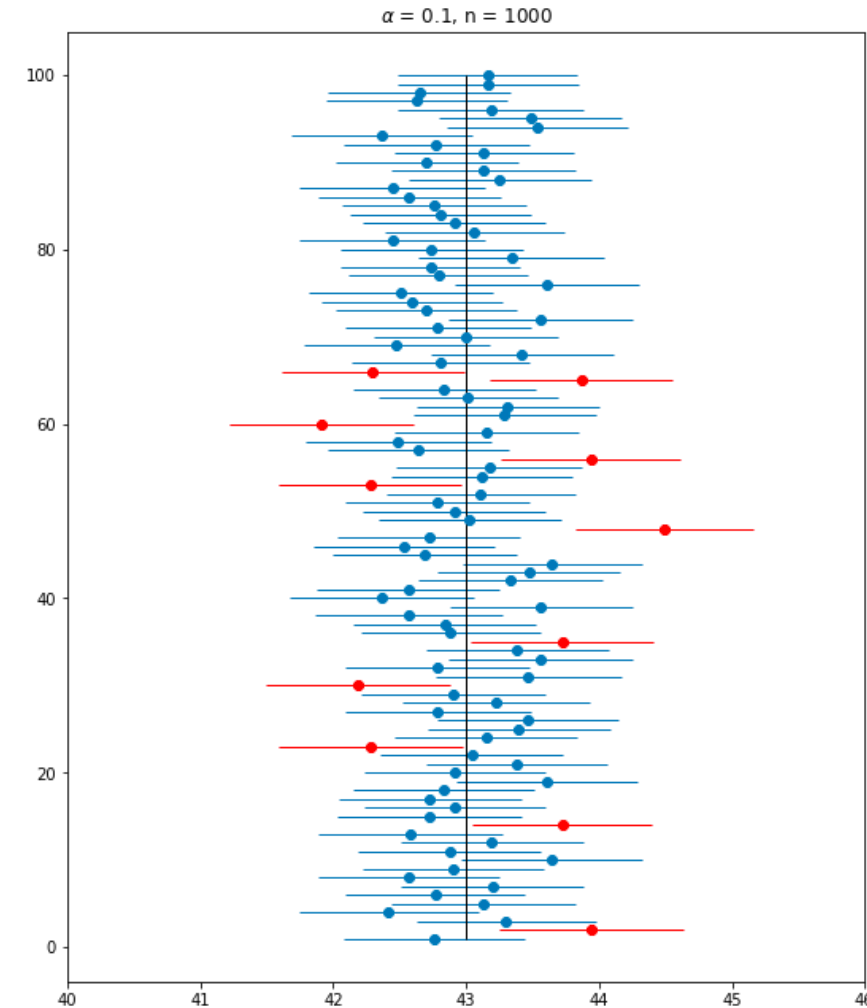
# Confidence intervals

$$P(\hat{p}(A) - \gamma \cdot \hat{\sigma}(A) \leq p \leq \hat{p}(A) + \gamma \cdot \hat{\sigma}(A)) \approx 2\Phi(\gamma) - 1$$

where  $\hat{p}(A)$  is the empirical proportion

$$\text{and } \hat{\sigma}(A) = \frac{\sqrt{\hat{p}(A)(1-\hat{p}(A))}}{\sqrt{n}}$$

$$P\left(p \in \hat{p}(A) \pm \Phi^{-1}\left(1 - \frac{\alpha}{2}\right) \cdot \hat{\sigma}(A)\right) \approx 1 - \alpha$$



# CIs for correlations? based on the Fisher Transform

When drawing  $n$  samples from a bivariate normal distribution with correlation coefficient  $\rho$  and computing the empirical Pearson correlation on the sample,  $\hat{\rho}_n$ , we have:

$$\Sigma = \begin{pmatrix} V(X) & Cov(X, Y) \\ Cov(X, Y) & V(Y) \end{pmatrix}$$

$$Cov(X, Y) = \rho(X, Y)\sigma(X)\sigma(Y)$$

$$P\left(F(\rho) \in F(\hat{\rho}_n) \pm \frac{1}{\sqrt{n-3}} \Phi^{-1}\left(1 - \frac{\alpha}{2}\right)\right) = 1 - \alpha$$

And also:

$$P\left(F(\rho) \geq F(\hat{\rho}_n) - \frac{1}{\sqrt{n-3}} \Phi^{-1}(1 - \alpha)\right) = 1 - \alpha$$

## How do we use this?

Suppose that we obtained  $\hat{\rho} = 0.8$  when calculated for grades in two exams in a class of 20 students.

What can you say about the correlation between the grades in these two exams, in the entire world population?

Assume that the grades in these two exams follow a bivariate normal joint distribution.

Also, assume, of course, that our 20 students were sampled from the relevant (entire world) population.



## How do we use this?

we obtained  $\hat{\rho} = 0.8$  when calculated for grades in two exams in a class of 20 students.

Therefore, working with  $\alpha = 0.05$ , say, we have

$$F(\rho) \in F(0.8) \pm \frac{1}{\sqrt{17}} \Phi^{-1}(0.975) = 1.099 \pm 0.243 \cdot 1.96 = [0.624, 1.574]$$

To be useful, we now want to convert this to a CI for  $\rho$ .

Since the Fisher Transform,  $F$ , is monotone and invertible we get:

$$\rho \in F^{-1}([0.624, 1.574]) = [F^{-1}(0.624), F^{-1}(1.574)] = [0.554, 0.918] .$$

And so we can state that  $0.554 \leq \rho \leq 0.918$  with 95% confidence.

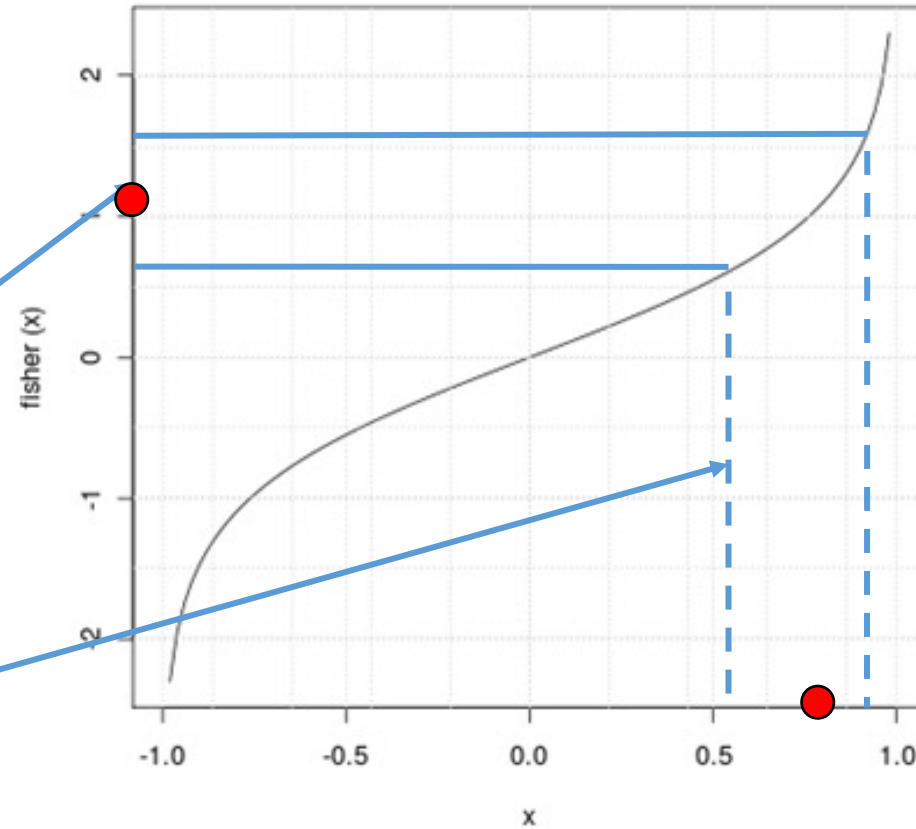
# Inverting the Fisher Transform

Note: the CI around  $\rho$  is NOT symmetric

$$r = F^{-1}(u) = \frac{e^{2u} - 1}{e^{2u} + 1}$$

CI for  $F(\rho)$ , obtained from a normal cdf

CI for  $\rho$ , obtained from the above by taking the inverse image, under F.



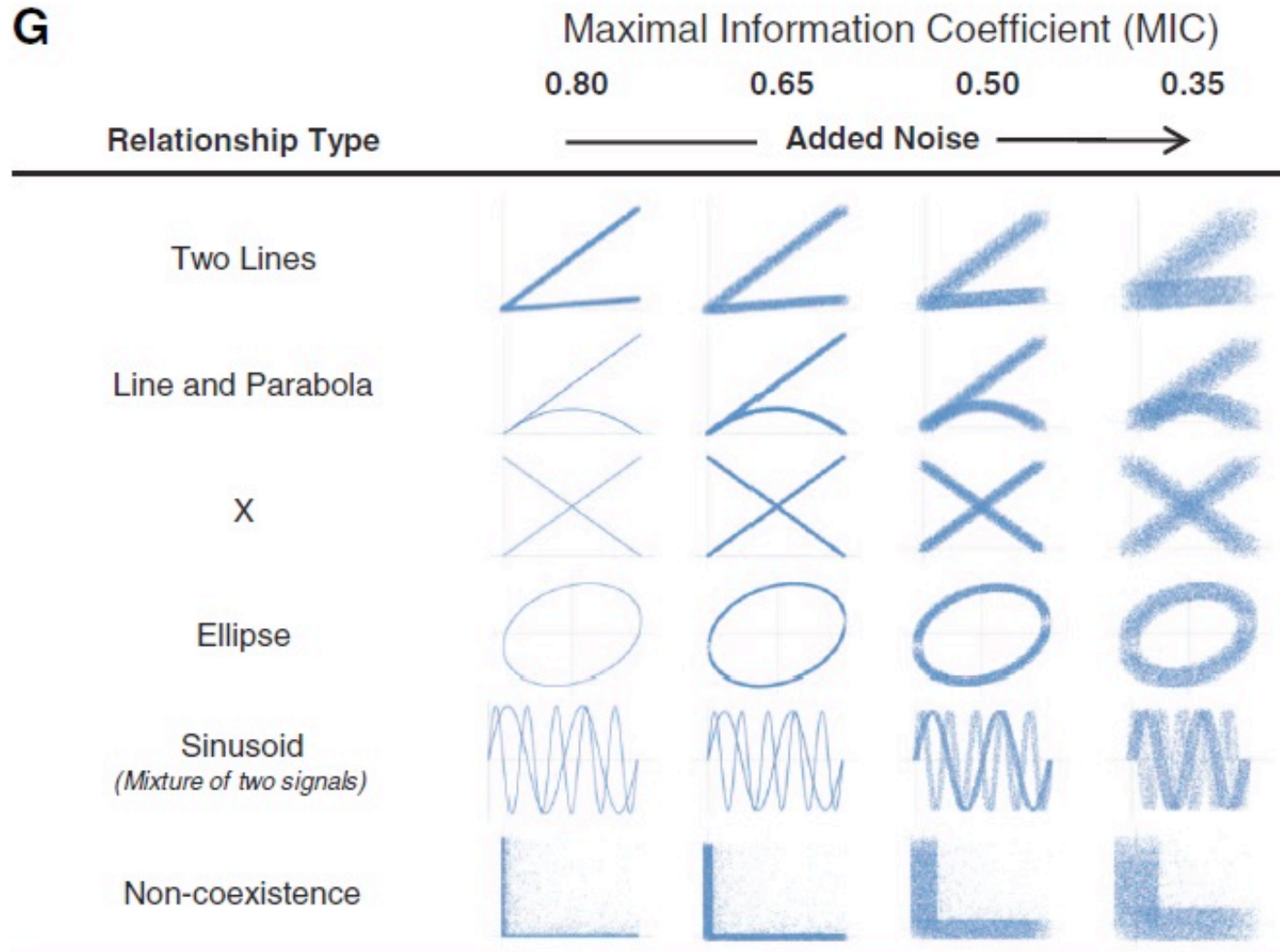
## p-Values for correlations

Built in functions typically compute p-values for Pearson's correlation using normal approximations (+Fisher), under the null model of  $X$  and  $Y$  following an independent standard bivariate normal distribution.

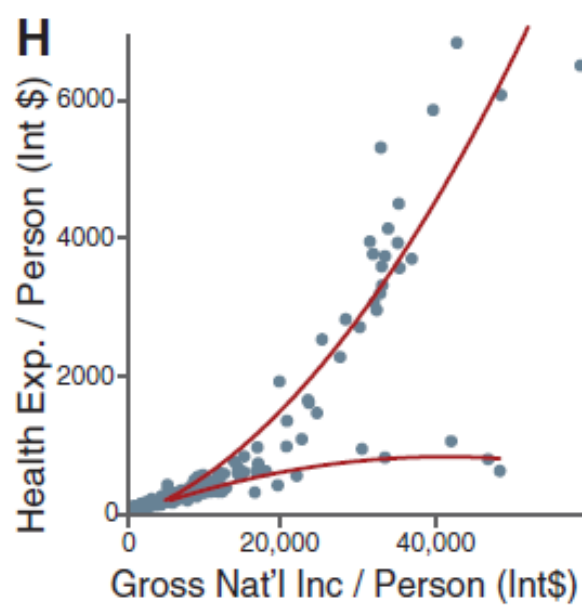
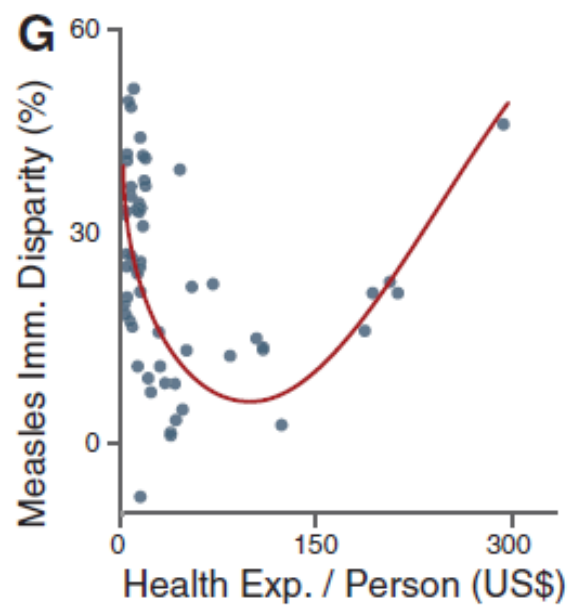
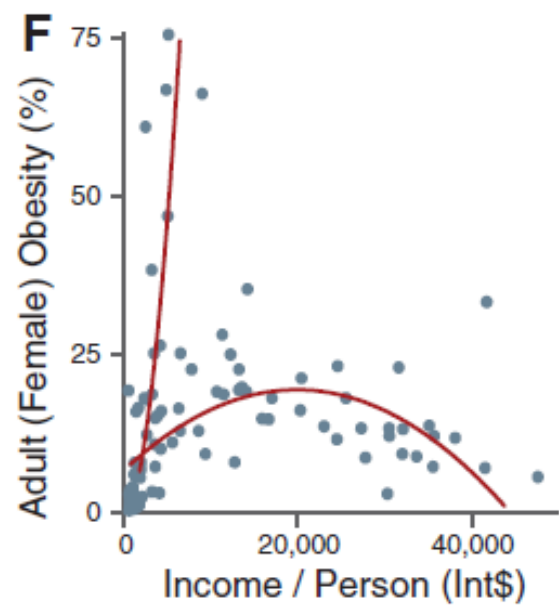
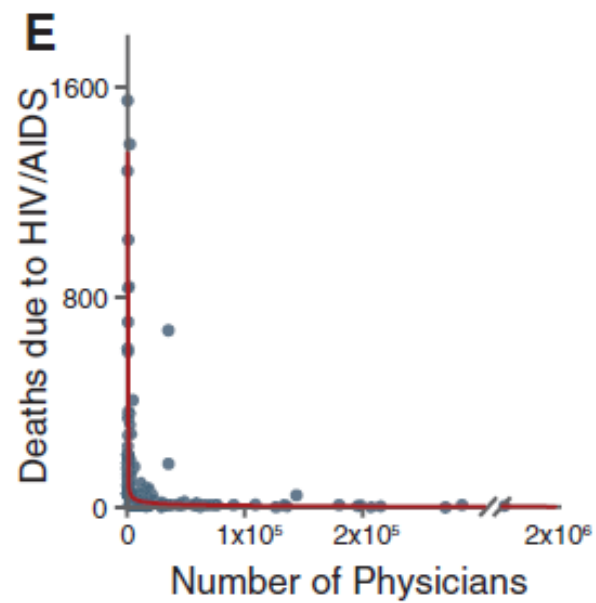
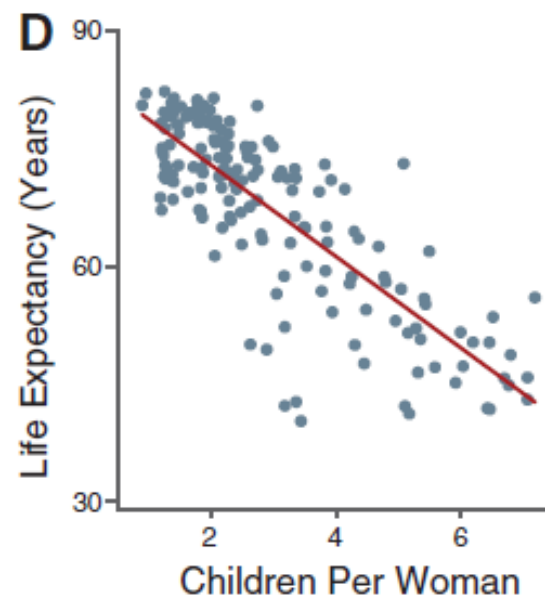
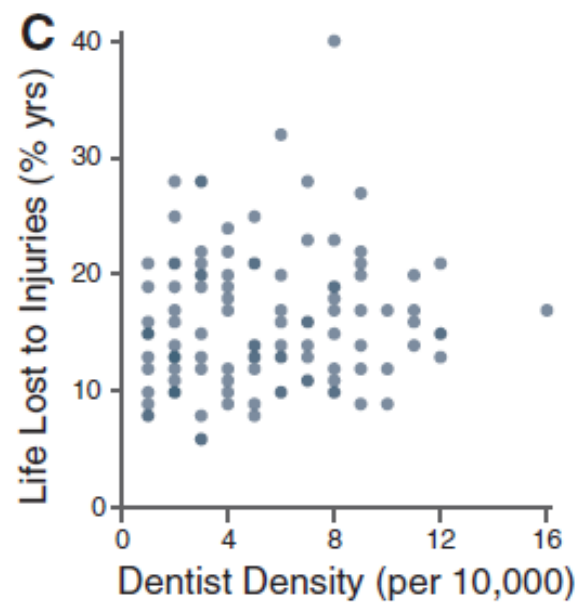
Built in functions compute p-values for Kendall's tau and Spearman's rho using either the exact permutation distributions (for small sample sizes), or large-sample normal approximations.

(what is the null model?)

# Correlations don't cover all association types



Reshef et al,  
Detecting Novel Associations  
in Large Data Sets,  
Science 2011

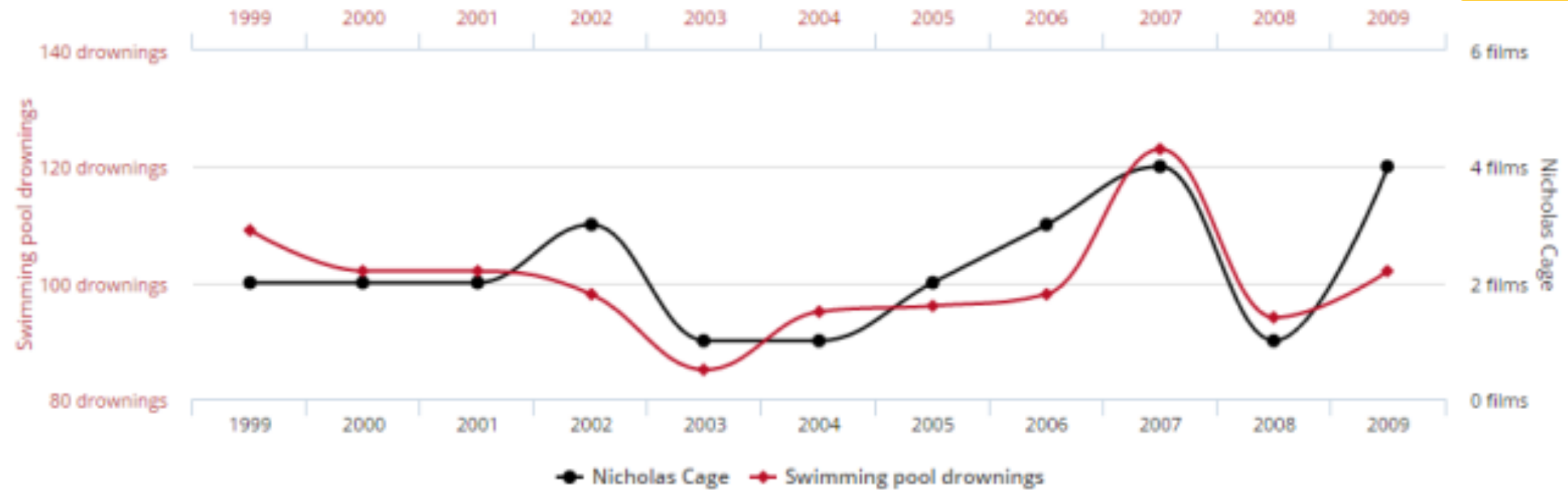


# Spurious Correlations



Number of people who drowned by falling into a pool  
correlates with  
Films Nicolas Cage appeared in

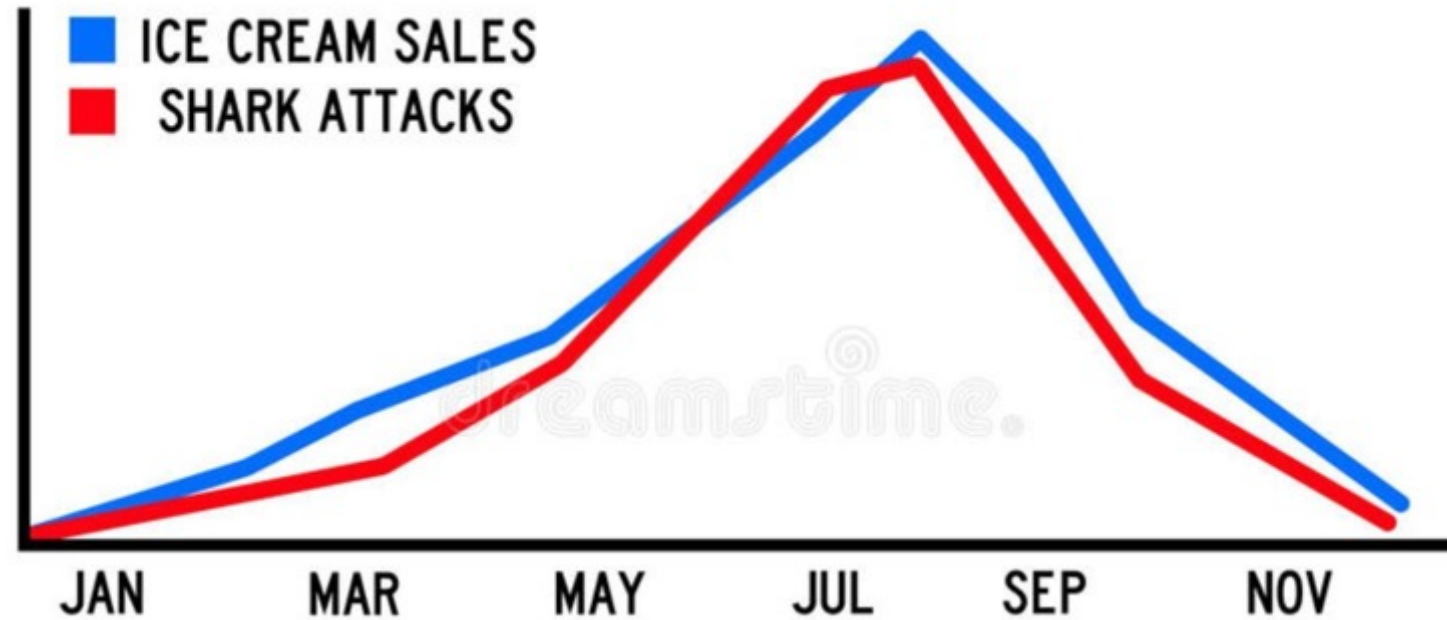
Correlation: 66.6% ( $r=0.666004$ )

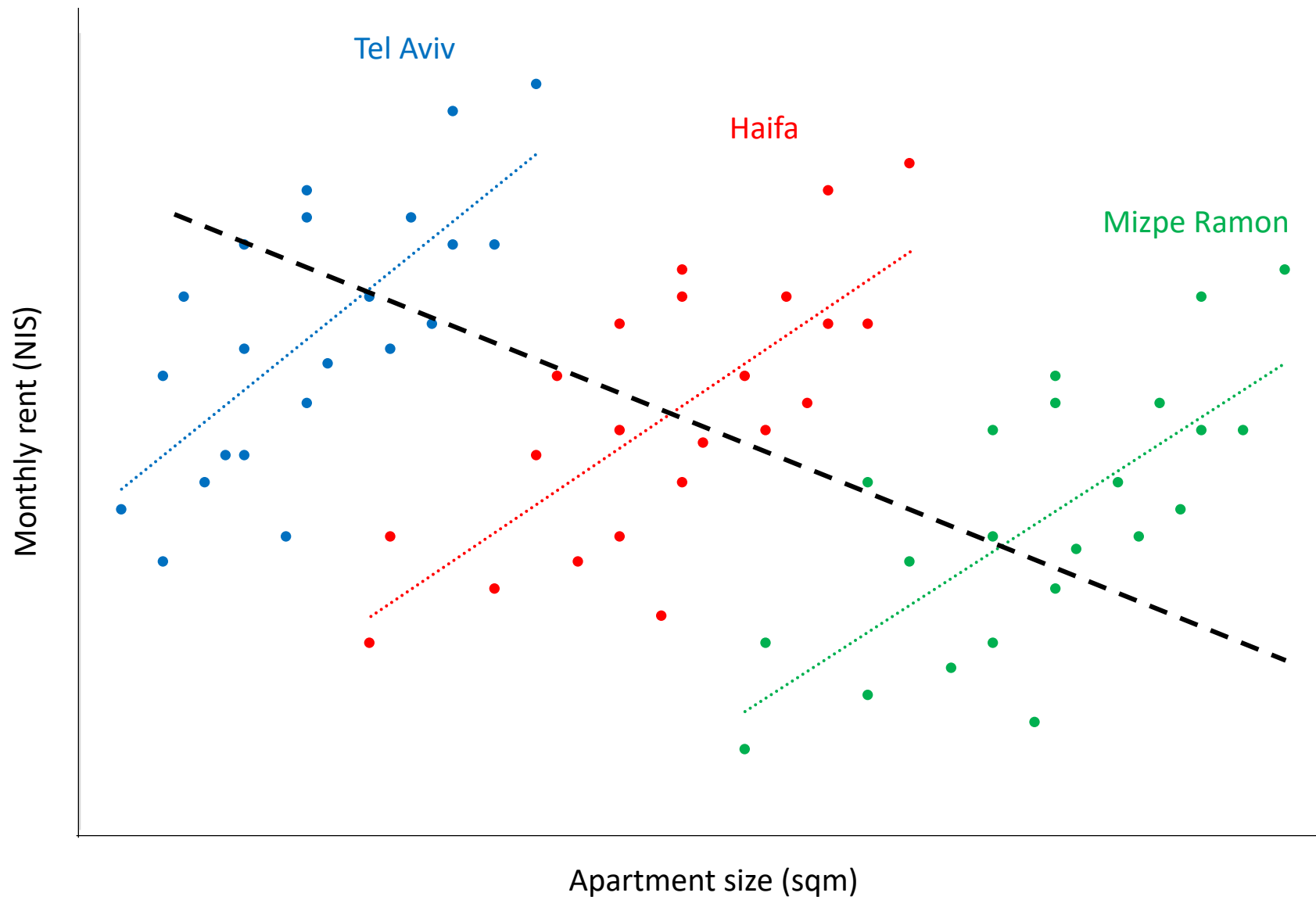


Data sources: Centers for Disease Control & Prevention and Internet Movie Database

tylervigen.com

# Correlation and causality

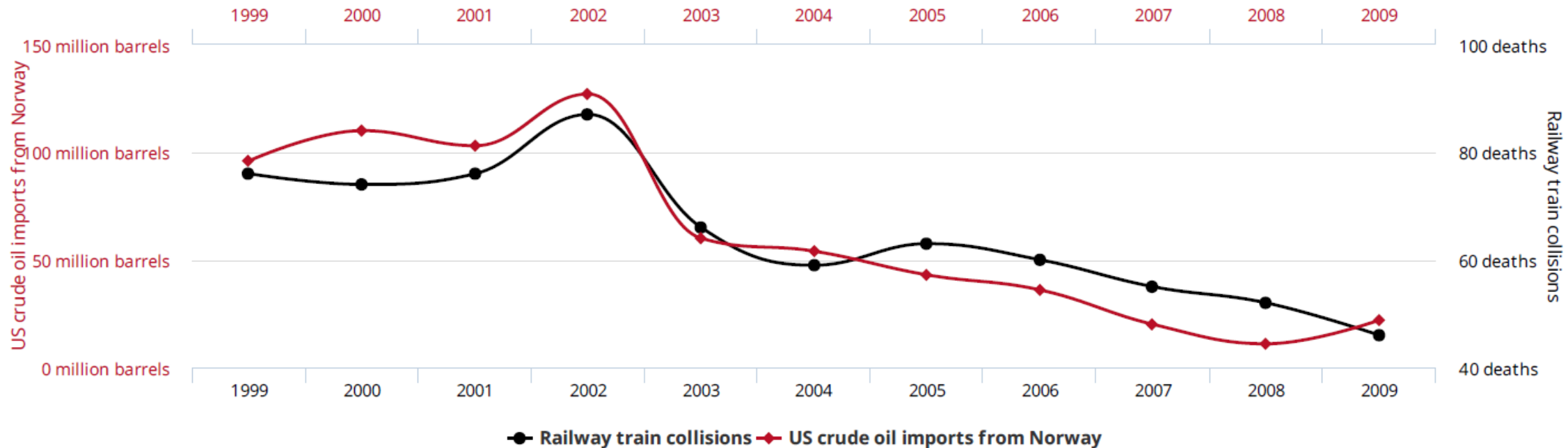


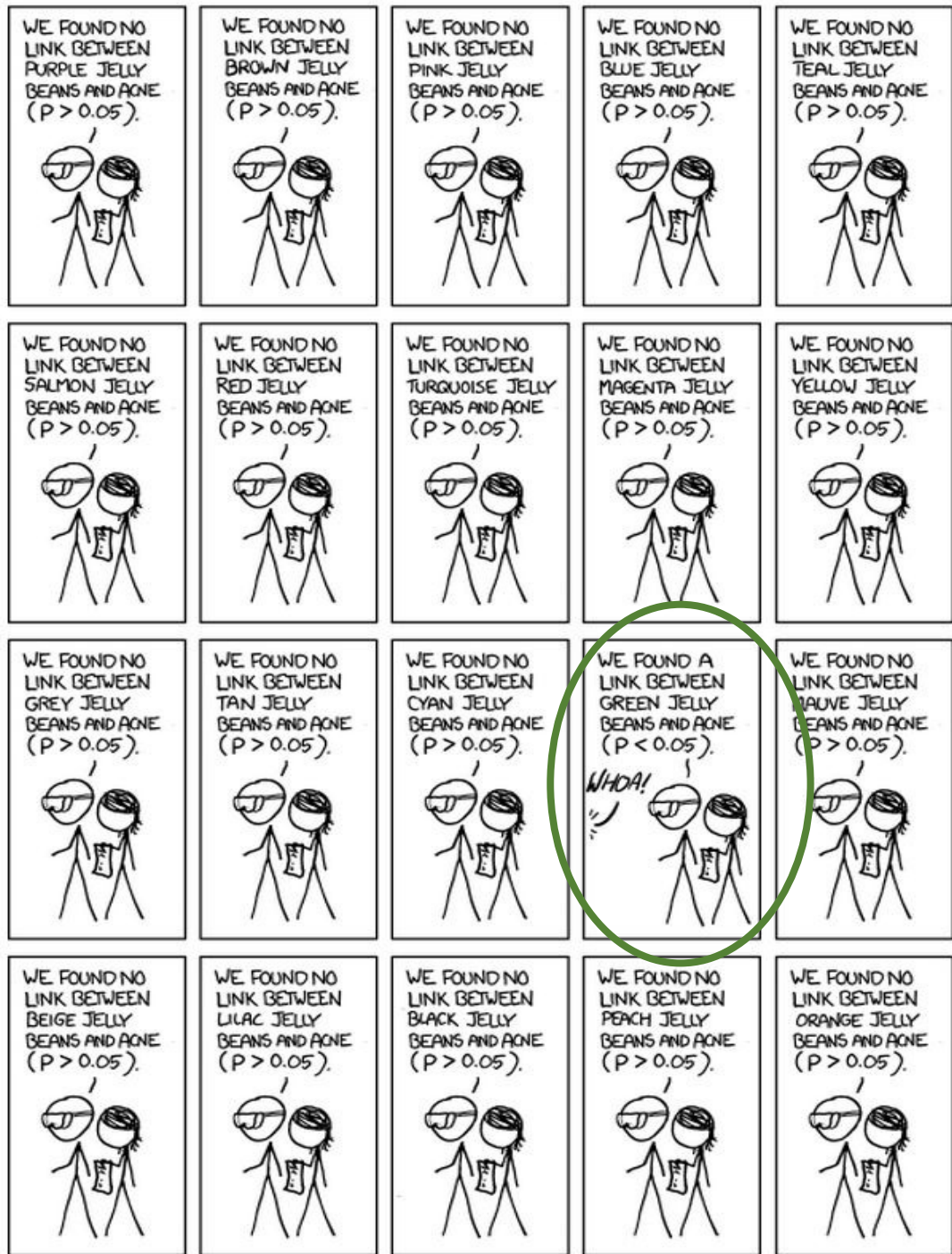




# Spurious correlations

## US crude oil imports from Norway correlates with Drivers killed in collision with railway train





# Summary

- Correlation measures serve to quantify relationships between different aspects of observed data
- Statistical assessment under a null model
- Pearson correlation is the classical and most popular correlation coefficients. It's a sample version of the population covariance.
- Pearson correlation affords confidence interval calculations.
- How can we capture more complex relationships between two variables/features/quantities?
- Between more than two variables/features/quantities?
- Finally – again, always assess presented correlations with a good measure of skepticism and attention to multiple testing