

# Correlation, Model Combination

Jindřich Libovický (reusing materials by Milan Straka)

 November 20, 2025

After this lecture you should be able to

- Explain and implement different ways of measuring correlation: Pearson's correlation, Spearman's correlation, Kendall's  $\tau$ .
- Decide if correlation is a good metric for your model.
- Measure inter-annotator agreement and draw conclusions for data cleaning and for limits of your models.
- Use correlation with human judgment to validate evaluation metrics.

# Covariance

Given a collection of random variables  $\mathbf{x}_1, \dots, \mathbf{x}_N$ , we know that

$$\mathbb{E} \left[ \sum_i \mathbf{x}_i \right] = \sum_i \mathbb{E} [\mathbf{x}_i].$$

But how about  $\text{Var} \left( \sum_i \mathbf{x}_i \right)$ ?

$$\begin{aligned} \text{Var} \left( \sum_i \mathbf{x}_i \right) &= \mathbb{E} \left[ \left( \sum_i \mathbf{x}_i - \sum_i \mathbb{E}[\mathbf{x}_i] \right)^2 \right] \\ &= \mathbb{E} \left[ \left( \sum_i (\mathbf{x}_i - \mathbb{E}[\mathbf{x}_i]) \right)^2 \right] \\ &= \mathbb{E} \left[ \sum_i \sum_j (\mathbf{x}_i - \mathbb{E}[\mathbf{x}_i]) (\mathbf{x}_j - \mathbb{E}[\mathbf{x}_j]) \right] \\ &= \sum_i \sum_j \mathbb{E} \left[ (\mathbf{x}_i - \mathbb{E}[\mathbf{x}_i]) (\mathbf{x}_j - \mathbb{E}[\mathbf{x}_j]) \right]. \end{aligned}$$

We define **covariance** of two random variables  $\mathbf{x}, \mathbf{y}$  as

$$\text{Cov}(\mathbf{x}, \mathbf{y}) = \mathbb{E} \left[ (\mathbf{x} - \mathbb{E}[\mathbf{x}]) (\mathbf{y} - \mathbb{E}[\mathbf{y}]) \right].$$

Then,

$$\text{Var} \left( \sum_i \mathbf{x}_i \right) = \sum_i \sum_j \text{Cov}(\mathbf{x}_i, \mathbf{x}_j).$$

Note that  $\text{Cov}(\mathbf{x}, \mathbf{x}) = \text{Var}(\mathbf{x})$  and that we can write covariance as

$$\begin{aligned} \text{Cov}(\mathbf{x}, \mathbf{y}) &= \mathbb{E} \left[ (\mathbf{x} - \mathbb{E}[\mathbf{x}]) (\mathbf{y} - \mathbb{E}[\mathbf{y}]) \right] \\ &= \mathbb{E} [\mathbf{x}\mathbf{y} - \mathbf{x}\mathbb{E}[\mathbf{y}] - \mathbb{E}[\mathbf{x}]\mathbf{y} + \mathbb{E}[\mathbf{x}]\mathbb{E}[\mathbf{y}]] \\ &= \mathbb{E} [\mathbf{x}\mathbf{y}] - \mathbb{E} [\mathbf{x}] \mathbb{E} [\mathbf{y}]. \end{aligned}$$

# Correlation

Random variables  $x, y$  are **uncorrelated** if  $\text{Cov}(x, y) = 0$ ; otherwise, they are **correlated**.

Note that two *independent* random variables are uncorrelated, because

$$\begin{aligned}\text{Cov}(x, y) &= \mathbb{E} \left[ (x - \mathbb{E}[x]) (y - \mathbb{E}[y]) \right] \\ &= \sum_{x, y} P(x, y) (x - \mathbb{E}[x]) (y - \mathbb{E}[y]) \\ &= \sum_{x, y} P(x) (x - \mathbb{E}[x]) P(y) (y - \mathbb{E}[y]) \\ &= \left( \sum_x P(x) (x - \mathbb{E}[x]) \right) \left( \sum_y P(y) (y - \mathbb{E}[y]) \right) \\ &= \mathbb{E}_x [x - \mathbb{E}[x]] \mathbb{E}_y [y - \mathbb{E}[y]] = 0.\end{aligned}$$

However, dependent random variables can be uncorrelated – random uniform  $x$  on  $[-1, 1]$  and  $y = |x|$  are not independent ( $y$  is completely determined by  $x$ ), but they are uncorrelated.

# Pearson correlation coefficient

There are several ways to measure correlation of random variables  $x, y$ .

**Pearson correlation coefficient**, denoted as  $\rho$  or  $r$ , is defined as

$$\rho \stackrel{\text{def}}{=} \frac{\text{Cov}(x, y)}{\sqrt{\text{Var}(x)} \sqrt{\text{Var}(y)}}$$

$$r \stackrel{\text{def}}{=} \frac{\sum_i (x_i - \bar{x})(y_i - \bar{y})}{\sqrt{\sum_i (x_i - \bar{x})^2} \sqrt{\sum_i (y_i - \bar{y})^2}},$$

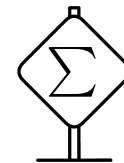
where:

- $\rho$  is used when the full expectation is computed (population Pearson correlation coefficient);
- $r$  is used when estimating the coefficient from data (sample Pearson correlation coefficient);
  - $\bar{x}$  and  $\bar{y}$  are sample estimates of the respective means.



# Pearson correlation coefficient

The value of Pearson correlation coefficient is in fact normalized covariance, because its value is always bounded by  $-1 \leq \rho \leq 1$  (and the same holds for  $r$ ).



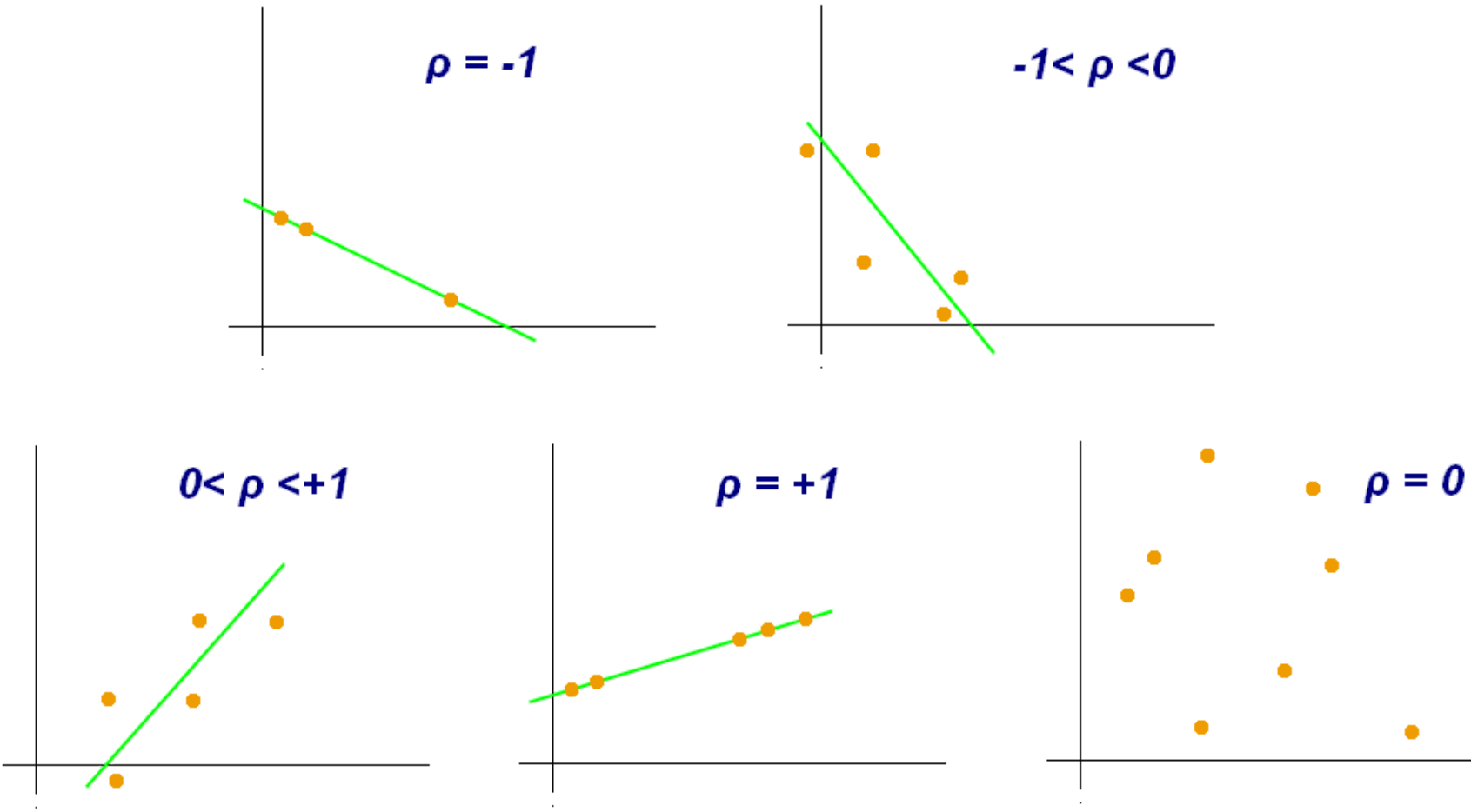
The bound can be derived from

$$\begin{aligned} 0 &\leq \mathbb{E} \left[ \left( \frac{(x - \mathbb{E}[x])}{\sqrt{\text{Var}(x)}} - \rho \frac{(y - \mathbb{E}[y])}{\sqrt{\text{Var}(y)}} \right)^2 \right] \\ &= \mathbb{E} \left[ \frac{(x - \mathbb{E}[x])^2}{\text{Var}(x)} \right] - 2\rho \mathbb{E} \left[ \frac{(x - \mathbb{E}[x])}{\sqrt{\text{Var}(x)}} \frac{(y - \mathbb{E}[y])}{\sqrt{\text{Var}(y)}} \right] + \rho^2 \mathbb{E} \left[ \frac{(y - \mathbb{E}[y])^2}{\text{Var}(y)} \right] \\ &= \frac{\text{Var}(x)}{\text{Var}(x)} - 2\rho \cdot \rho + \rho^2 \frac{\text{Var}(y)}{\text{Var}(y)} = 1 - \rho^2, \end{aligned}$$

which yields  $\rho^2 \leq 1$ .

# Pearson correlation coefficient

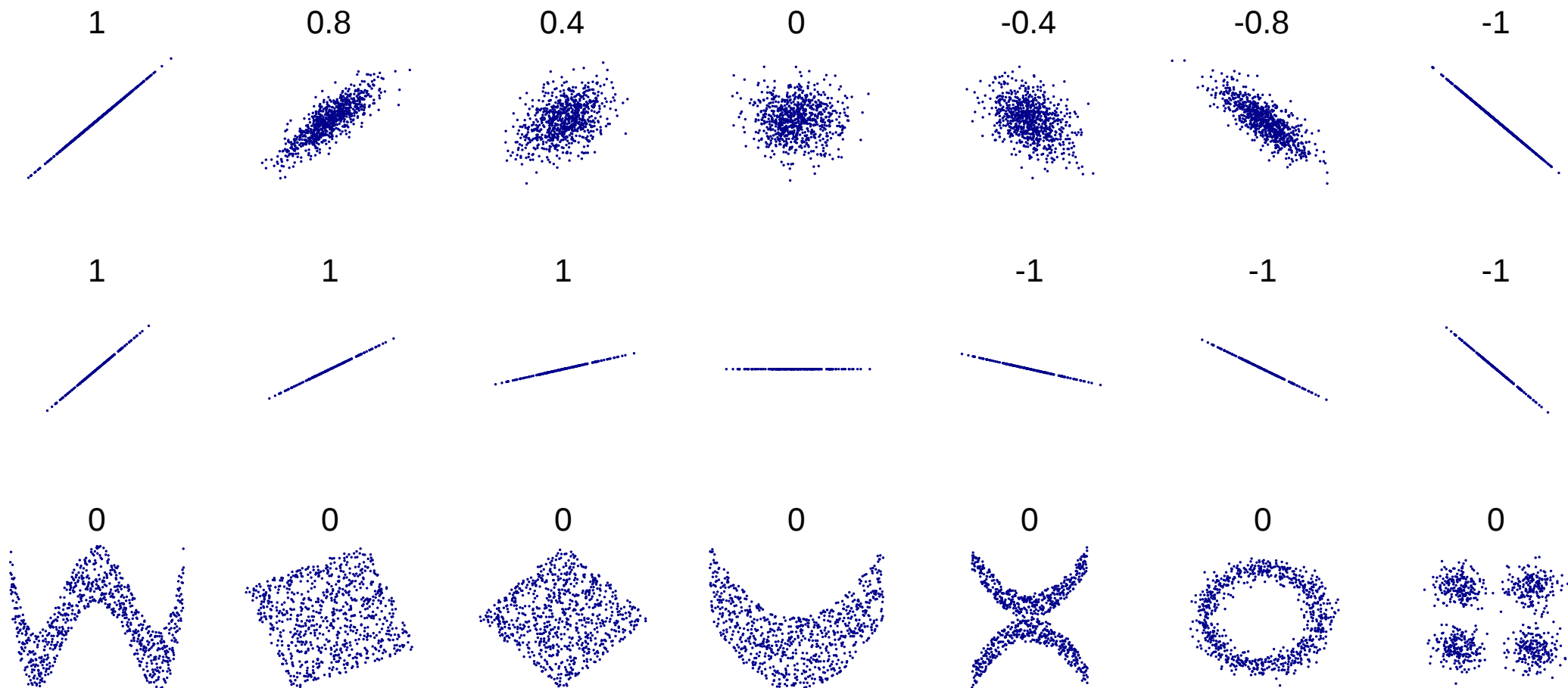
Pearson correlation coefficient quantifies **linear dependence** of the two random variables.



[https://upload.wikimedia.org/wikipedia/commons/3/34/Correlation\\_coefficient.png](https://upload.wikimedia.org/wikipedia/commons/3/34/Correlation_coefficient.png)

# Pearson correlation coefficient

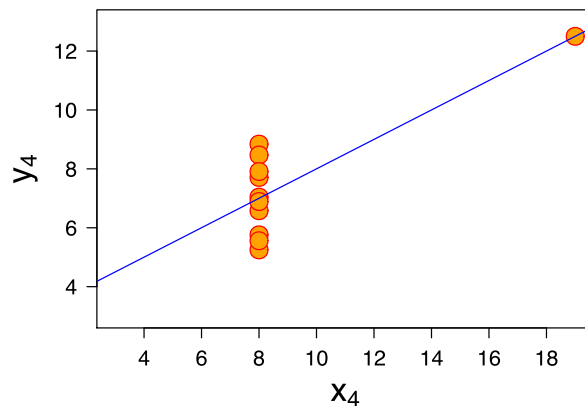
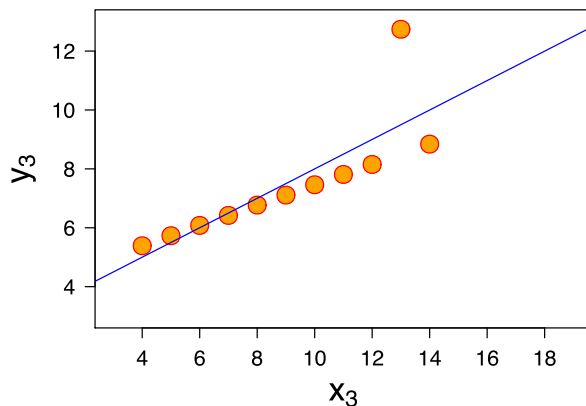
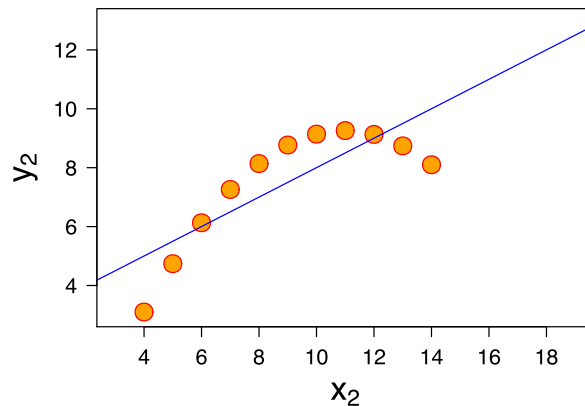
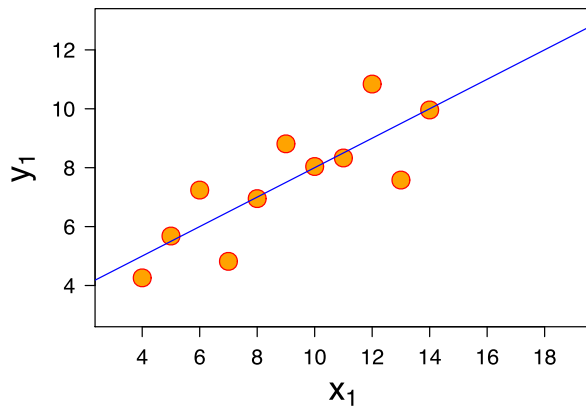
Pearson correlation coefficient quantifies **linear dependence** of the two random variables.



[https://upload.wikimedia.org/wikipedia/commons/d/d4/Correlation\\_examples2.svg](https://upload.wikimedia.org/wikipedia/commons/d/d4/Correlation_examples2.svg)

# Pearson correlation coefficient

The four displayed variables have the same mean 7.5, variance 4.12, Pearson correlation coefficient 0.816 and regression line  $3 + \frac{1}{2}x$ .



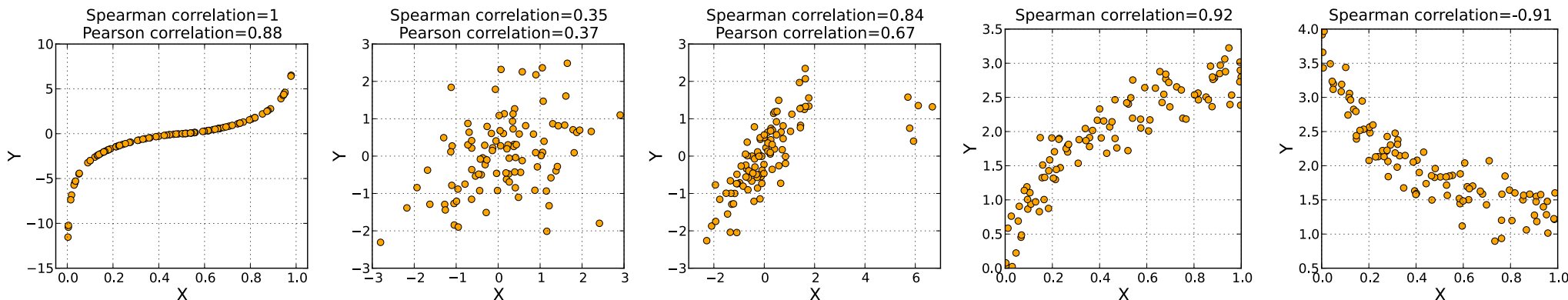
[https://upload.wikimedia.org/wikipedia/commons/e/ec/Anscombe%27s\\_quartet\\_3.svg](https://upload.wikimedia.org/wikipedia/commons/e/ec/Anscombe%27s_quartet_3.svg)

# Nonlinear Correlation – Spearman's $\rho$

To measure also nonlinear correlation, two coefficients are commonly used.

## Spearman's rank correlation coefficient $\rho$

Spearman's  $\rho$  is Pearson correlation coefficient measured on **ranks** of the original data, where a rank of an element is its index in sorted ascending order.



[https://upload.wikimedia.org/wikipedia/commons/4/4e/Spearman\\_fig{1,2,3,5,4}.svg](https://upload.wikimedia.org/wikipedia/commons/4/4e/Spearman_fig{1,2,3,5,4}.svg)

## Kendall rank correlation coefficient $\tau$

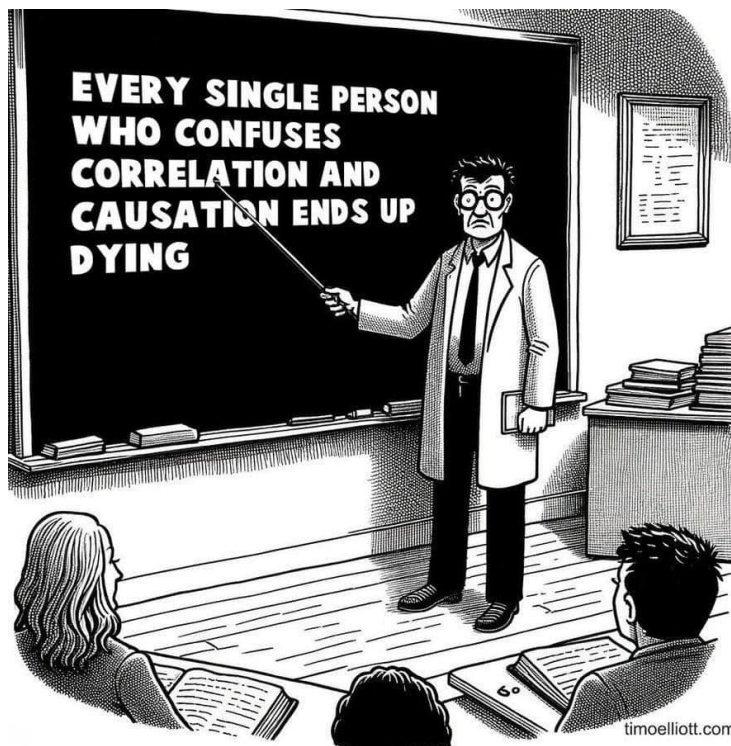
Kendall's  $\tau$  measures the amount of *concordant pairs* (pairs where  $y$  increases/decreases when  $x$  does), minus the *discordant pairs* (where  $y$  increases/decreases when  $x$  does the opposite):

$$\begin{aligned}\tau &\stackrel{\text{def}}{=} \frac{|\{\text{pairs } i \neq j : x_j > x_i, y_j > y_i\}| - |\{\text{pairs } i \neq j : x_j > x_i, y_j < y_i\}|}{\binom{n}{2}} \\ &= \frac{\sum_{i < j} \text{sign}(x_j - x_i) \text{sign}(y_j - y_i)}{\binom{n}{2}}.\end{aligned}$$

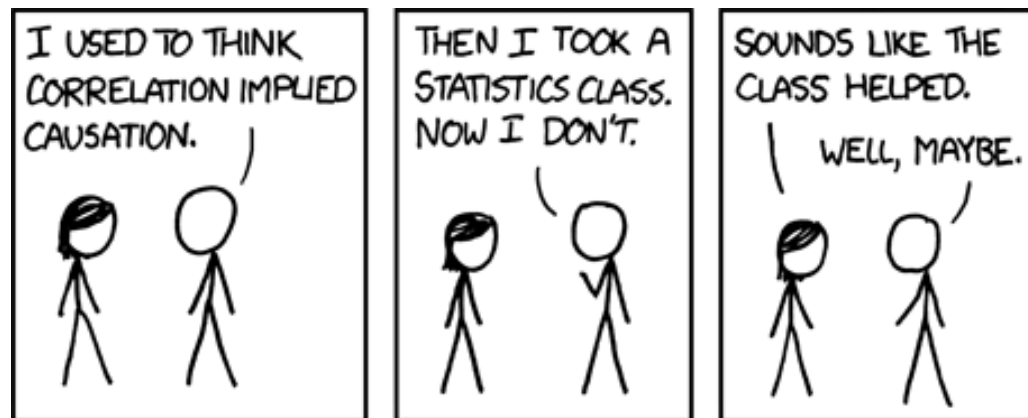
There is no clear consensus on whether to use Spearman's  $\rho$  or Kendall's  $\tau$ . When there are no/few ties in the data, Kendall's  $\tau$  offers two minor advantages –  $\frac{1+\tau}{2}$  can be interpreted as a probability of a concordant pair, and Kendall's  $\tau$  converges to a normal distribution faster.

As defined, the range of Kendall's  $\tau \in [-1, 1]$ . However, if there are ties, its range is smaller – therefore, several corrections (not discussed here) exist to adjust its value in case of ties.

# Correlation is not causation



<https://timoelliott.com/blog/cartoons/yet-more-analytics-cartoons>



<https://xkcd.com/552/>

# Correlation in Machine Learning



In ML, correlation is commonly used as

- Evaluation metric for some tasks;
- Measuring data annotation quality;
- Assessing the quality of automatic metrics by comparing them to human judgment.

- Learning to rank (e.g., document retrieval): we do not care about the actual values
  - Kendall's  $\tau$ , Spearman's correlation
  - When we want the correct items to rank before incorrect ones: precision (assuming fixed top- $k$ , typically at 5, 10), recall (often ill-defined), mean reciprocal rank

$$\text{MRR} = \frac{1}{N} \sum_{i=1}^N \frac{1}{\text{rank of the first relevant item}}$$

- Evaluating pair similarity: word embeddings, sentence embeddings
  - Similarity estimates from psycholinguistic experiments: scores for word/sentence pairs
  - Measure Pearson/Spearman correlation between embedding distances and similarity scores

# Inter-annotator agreement (1)

- Inter-annotator agreement can tell us
  - How well defined the task is
  - How reliable annotators/user ratings are
  - What data items are suspicious / difficult
- For continuous target values: Pearson's/Spearman's correlation
- For classification tasks: Cohen's  $\kappa$   
 $p_O$  is observed agreement,  $p_E$  expected agreement by chance

$$\kappa = \frac{p_O - p_E}{1 - p_E}$$

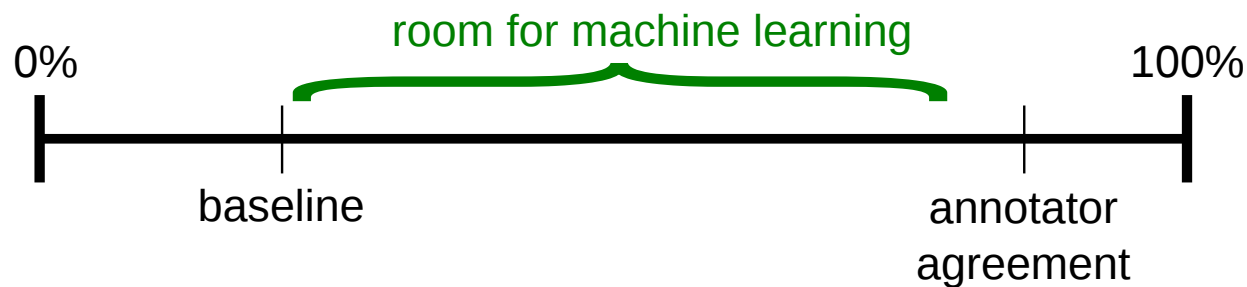


<https://www.surgehq.ai/blog/the-pitfalls-of-inter-rater-reliability-in-data-labeling-and-machine-learning>

# Inter-annotator agreement (2)

- Can be used to filter out confusing data points and unreliable annotators
- Not all outliers are noise! Low IAA can reveal cultural differences.

IAA sets natural upper boundary for ML performance. Performance over IAA is suspicious!

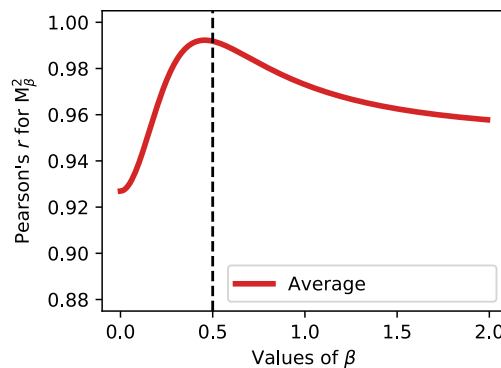


- Trivial baseline for classification: majority class, for regression average, or something based on simple rules
- Performance over IAA is more likely overfitting for the way the data is curated than super-human performance.

# Correlation with human judgment

For some tasks, it might not be clear how to measure the model performance:

**Grammar checking:** the  $\beta$  parameter



*J. Náplava, M. Straka, J. Straková, and A. Rosen. 2022. Czech Grammar Error Correction with a Large and Diverse Corpus. In TAACL, 10:452–467.*

**Machine translation:** evaluation is subjective by definition, we design metrics to correlate with human judgment.

- SoTA machine translation metrics are typically machine-learned.
- Different metrics might be suitable for different tiers of translation quality.
- There is an annual competition in MT quality and MT metric quality.

After this lecture you should be able to

- Explain and implement different ways of measuring correlation: Pearson's correlation, Spearman's correlation, Kendall's  $\tau$
- Decide if correlation is a good metric for your model
- Measure inter-annotator agreement and draw conclusions for data cleaning and for limits of your models
- Use correlation with human judgment to validate evaluation metrics