

MACHINE LEARNING WITH PYTHON

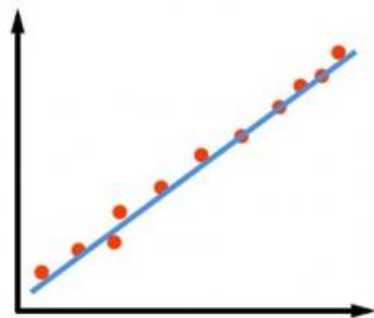
# FEATURE SELECTION

---

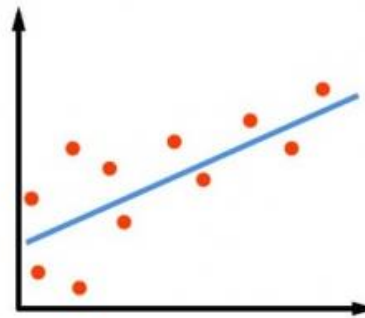
Themistoklis Diamantopoulos

# Correlation

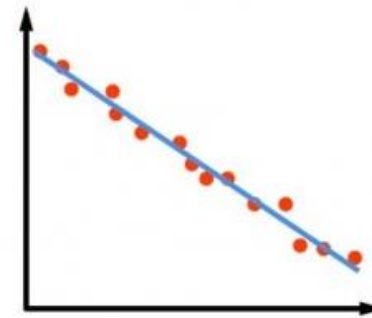
- Similar information



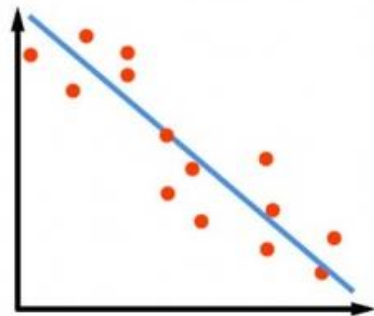
**STRONG POSITIVE CORRELATION**



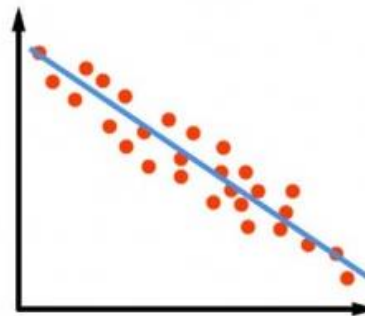
**WEAK POSITIVE CORRELATION**



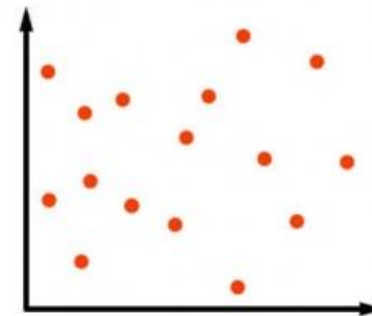
**STRONG NEGATIVE CORRELATION**



**WEAK NEGATIVE CORRELATION**



**MODERATE NEGATIVE CORRELATION**



**NO CORRELATION**

# Correlation Statistics

- Variables  $X, Y$

- Variance

$$\text{var}(X) = E \left[ (X - E[X])^2 \right]$$

- Covariance

$$\text{cov}(X, Y) = E \left[ (X - E[X])(Y - E[Y]) \right]$$

- Correlation

$$\text{corr}(X, Y) = \frac{\text{cov}(X, Y)}{\sqrt{\text{var}(X) \cdot \text{var}(Y)}}$$

Values from [-1, +1]  
+: Positive correlation  
-: Negative correlation  
0: No correlation

# Mutual Information

- Variables  $X, Y$
- Measures the dependence between random variables

$$I(X; Y) = \sum_{y \in Y} \sum_{x \in X} p(x, y) \log \left( \frac{p(x, y)}{p(x) p(y)} \right)$$

- Higher dependence  $\rightarrow$  Higher mutual information
- If the variables are independent then mutual information is 0

# Chi-Squared Statistic

$$\chi^2 = \sum_{i=1}^n \frac{(O_i - E_i)^2}{E_i}$$

$O_i$ : actual count

$E_i$ : expected count

	Play Chess	Don't Play Chess	Sum
Like Science Fiction	250 90	200 360	450
Don't Like Science Fiction	50 210	1000 840	1050
Sum	300	1200	1500

- Example:

$$\chi^2 = \frac{(250 - 90)^2}{90} + \frac{(50 - 210)^2}{210} + \frac{(200 - 360)^2}{360} + \frac{(1000 - 840)^2}{840} = 507.93$$

- Higher dependence  $\rightarrow$  Higher chi-squared