# Lecture 07 – Correlation & Autocorrelation
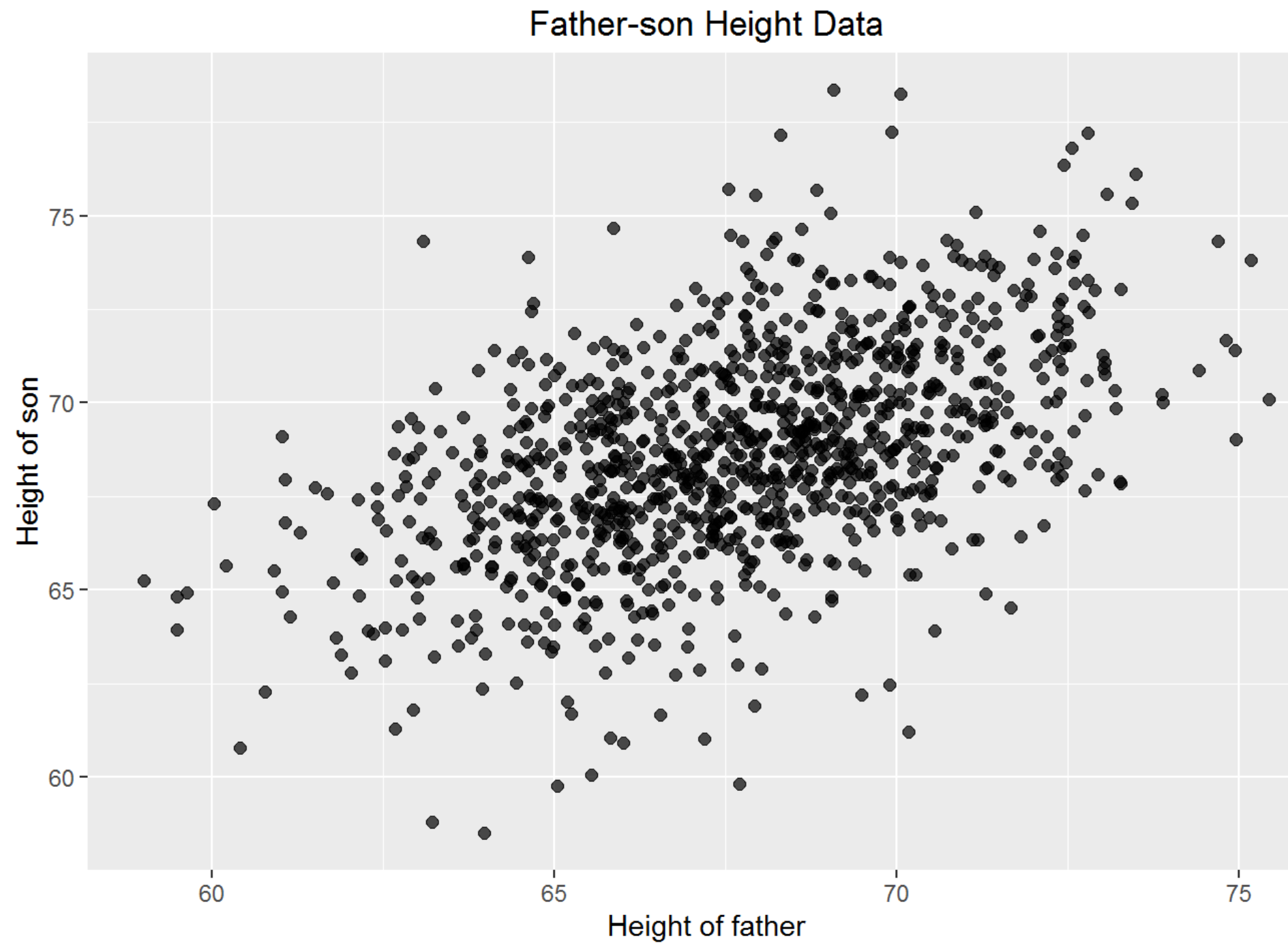
1

# Recap

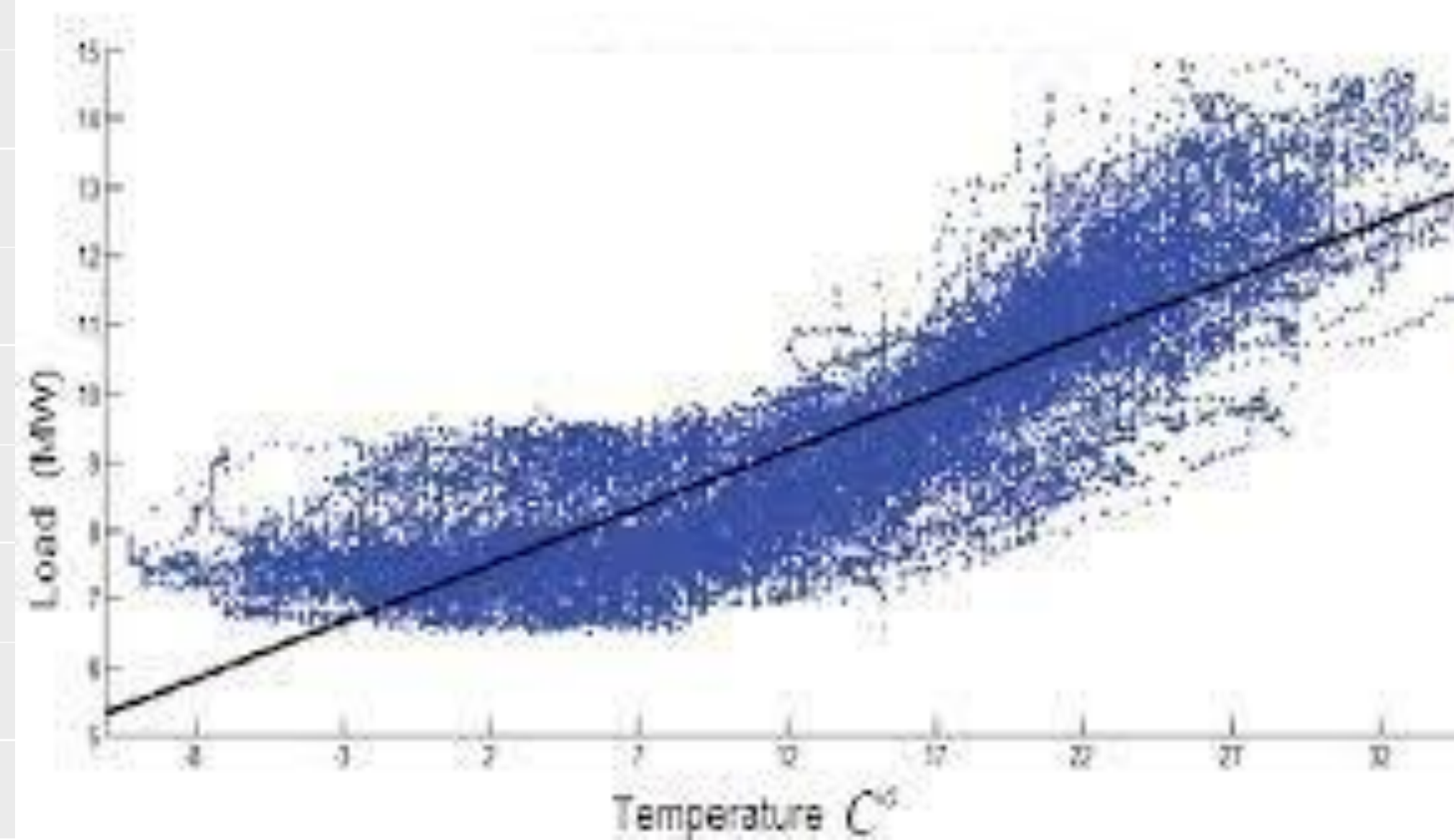- RMS, SD, Chebyshev inequality

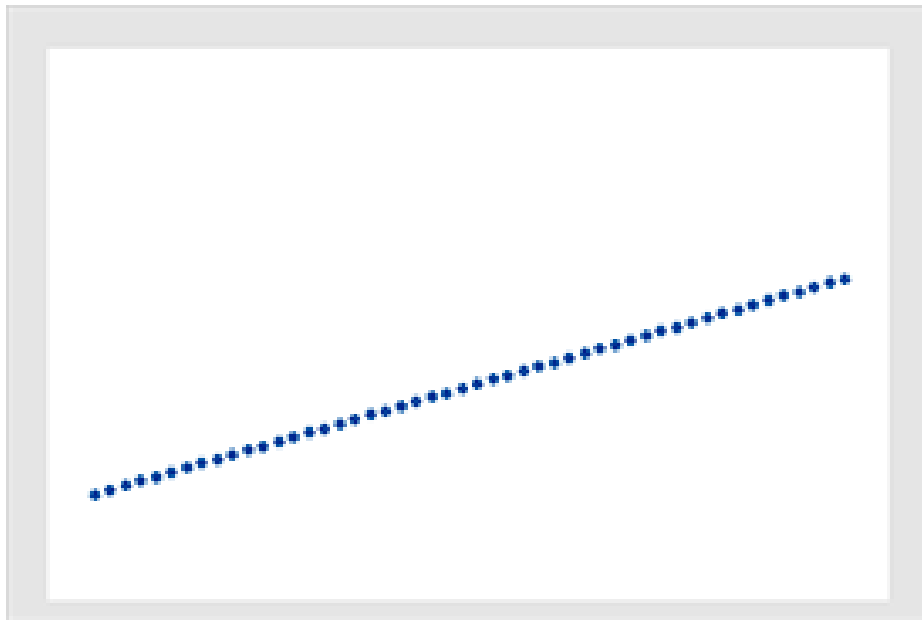# 1. Correlation

# Correlation

- Father-son heights

- Temperature-Electric bill

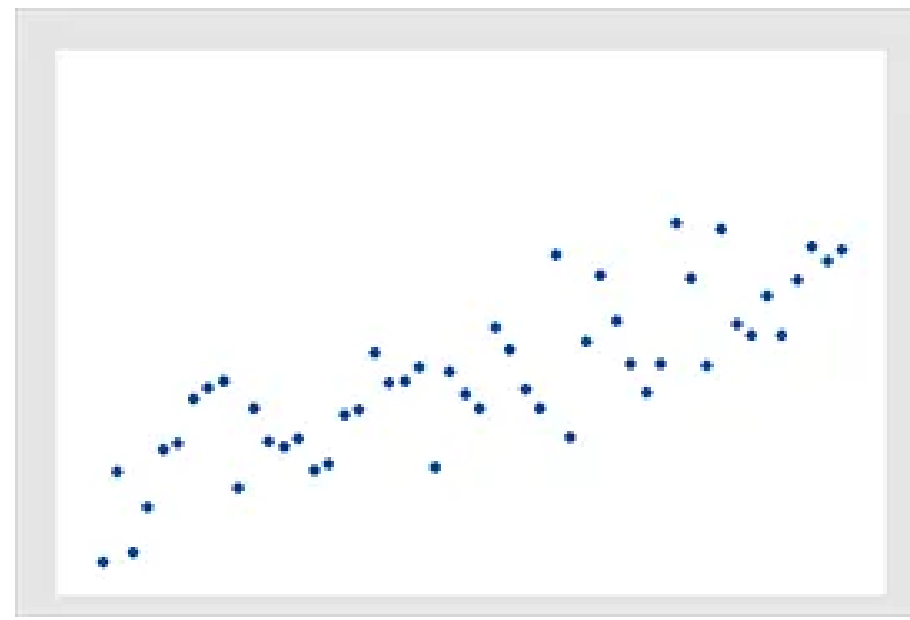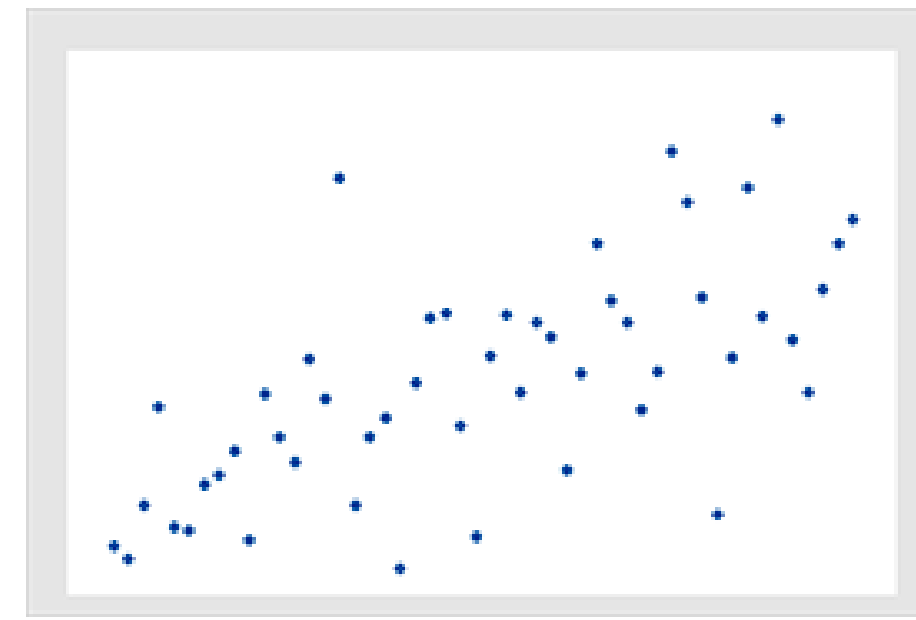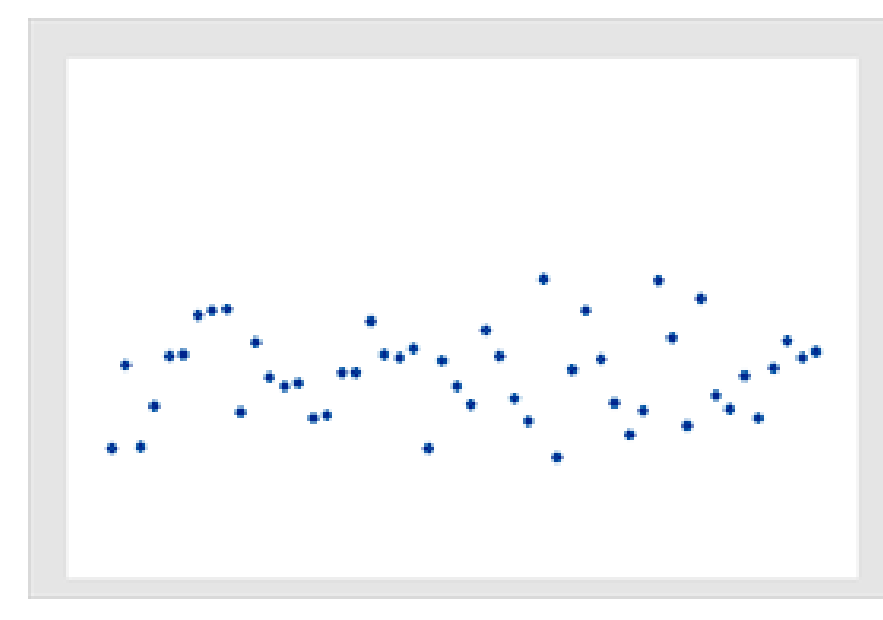# Correlation strength & coefficients

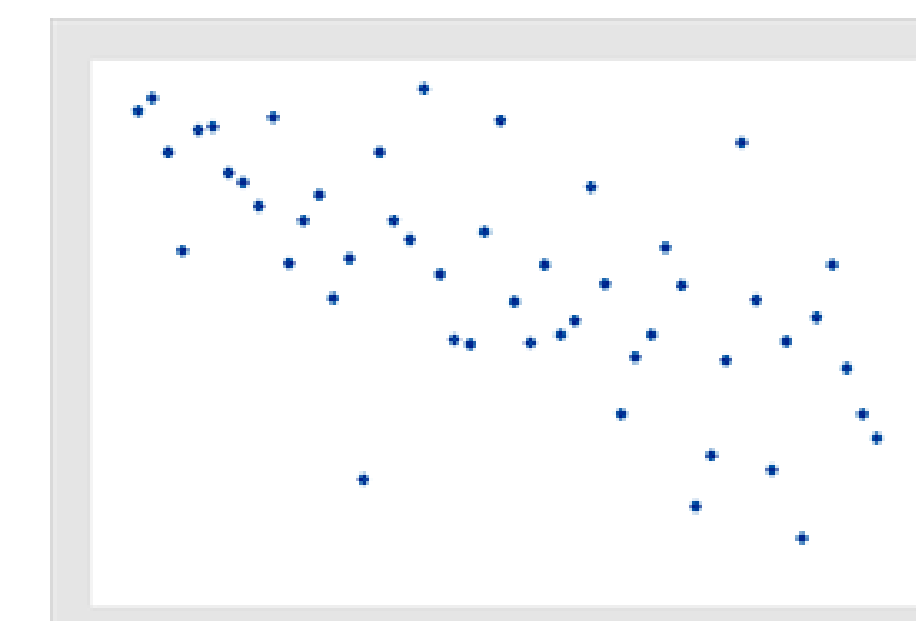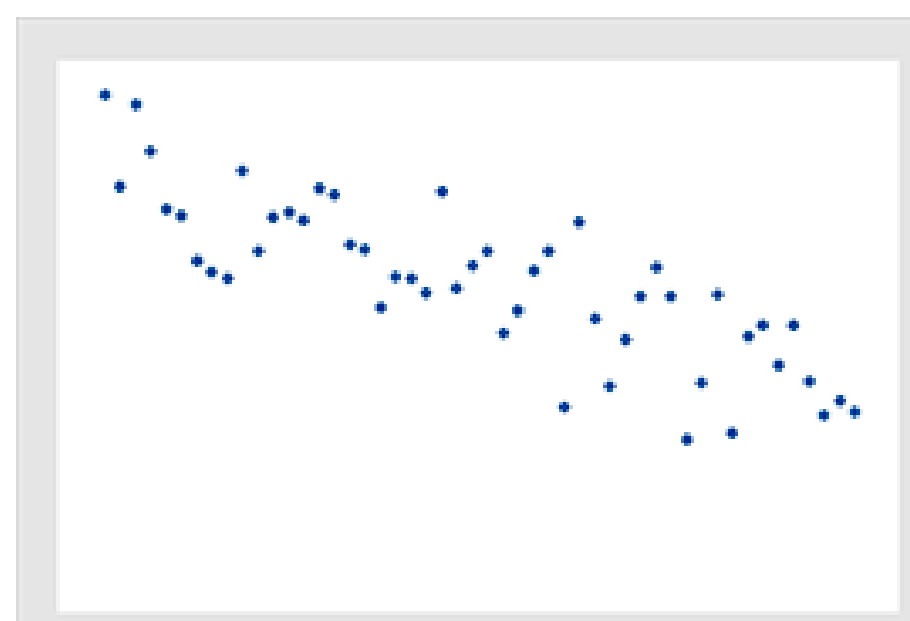- Very Strong      Strong      Moderate      None



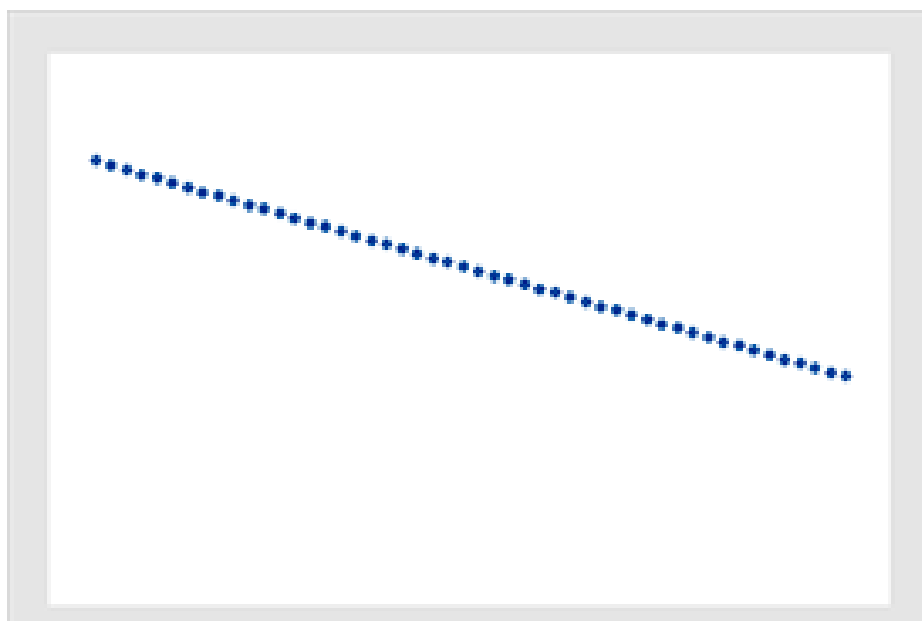-      1         0.8         0.6         0
-      -1        -0.8        -0.6

# Correlation coefficient

## Simple form

$$\rho = \frac{(a - \bar{a})^T (b - \bar{b})}{\|a - \bar{a}\|\|b - \bar{b}\|}$$

**Denominator for normalizing between -1 and 1**

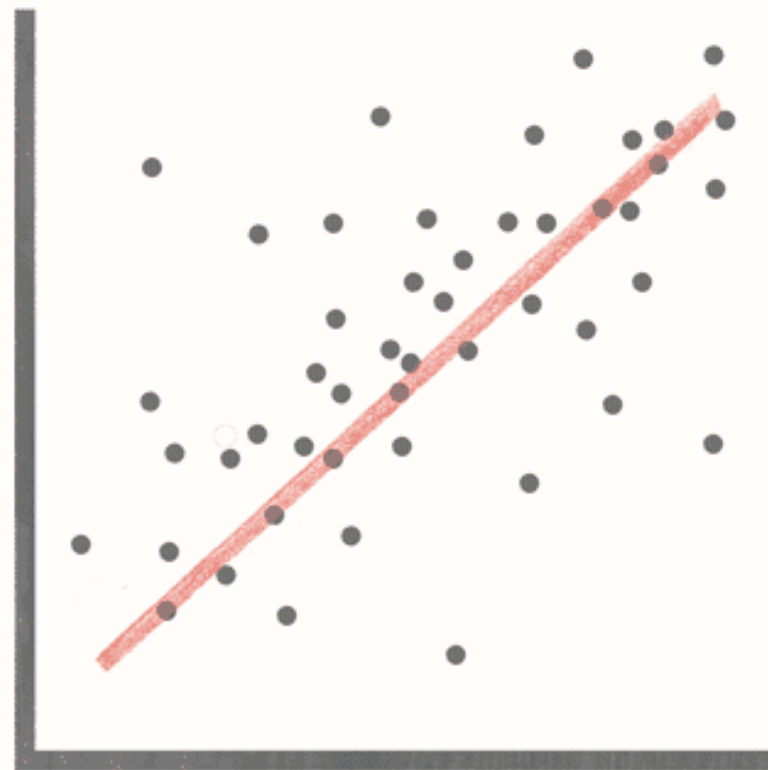- Why not $\quad \rho = \dfrac{a^T b}{\|a\|\|b\|}$

## Nasty form

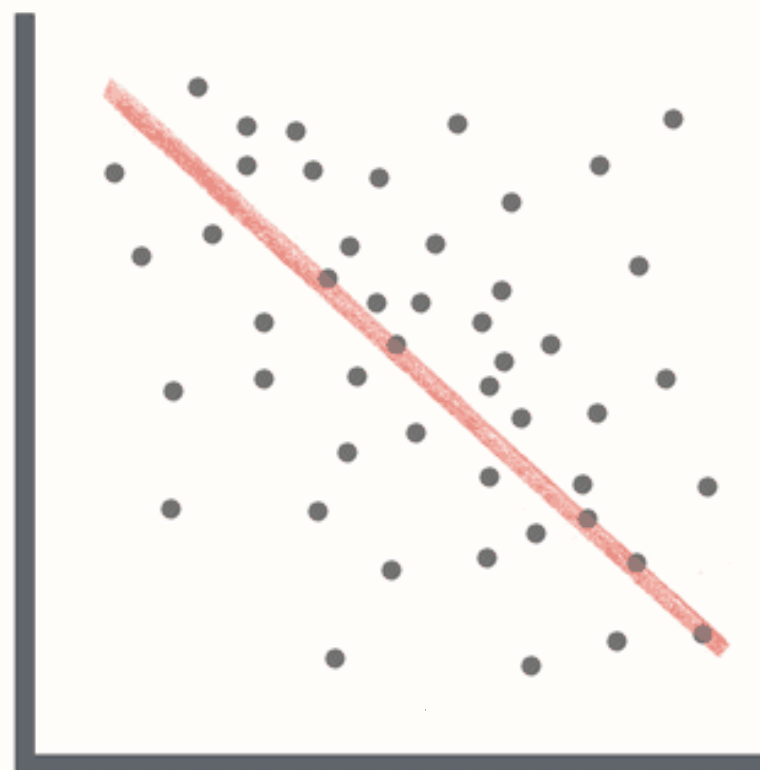$$\rho = \frac{(a - \frac{1}{n}\mathbf{1}^T a \mathbf{1})^T (b - \frac{1}{n}\mathbf{1}^T b \mathbf{1})}{\|a - \frac{1}{n}\mathbf{1}^T a \mathbf{1}\|\|b - \frac{1}{n}\mathbf{1}^T b \mathbf{1}\|}$$

**After all, dot product measures similarity**
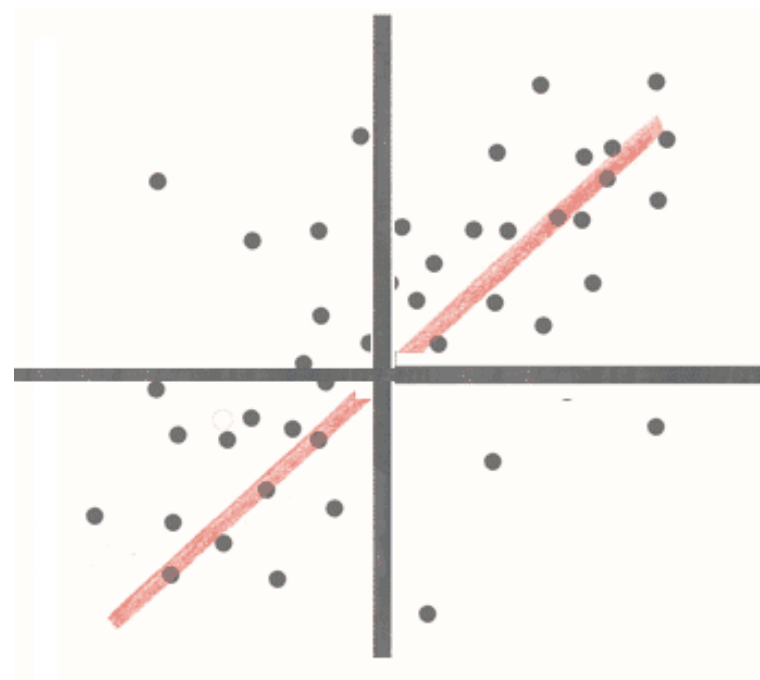
Positive Correlation
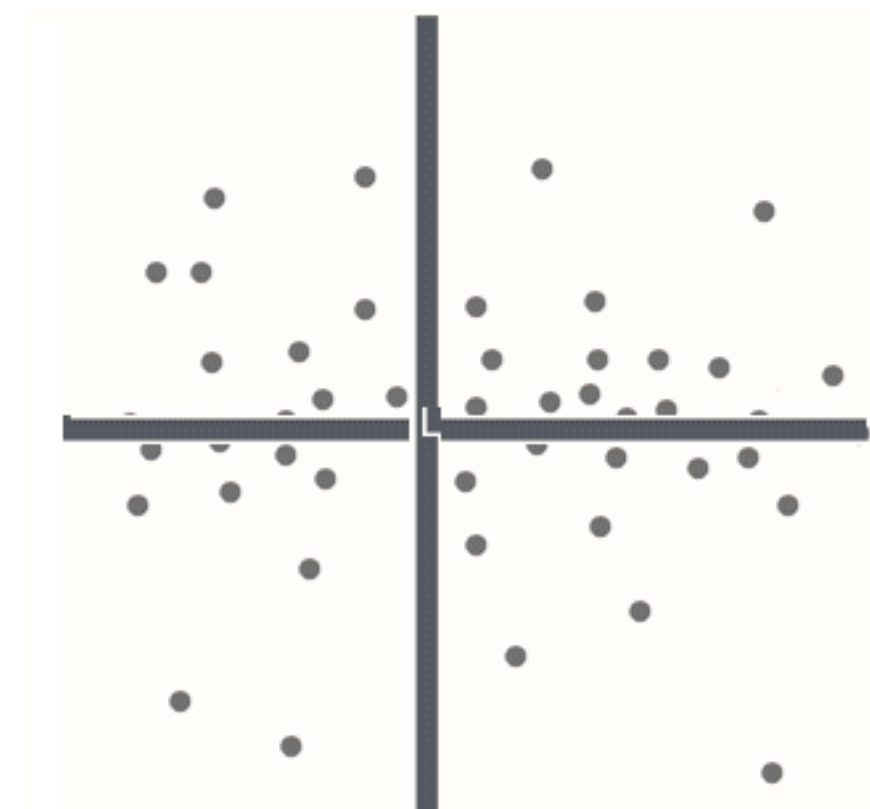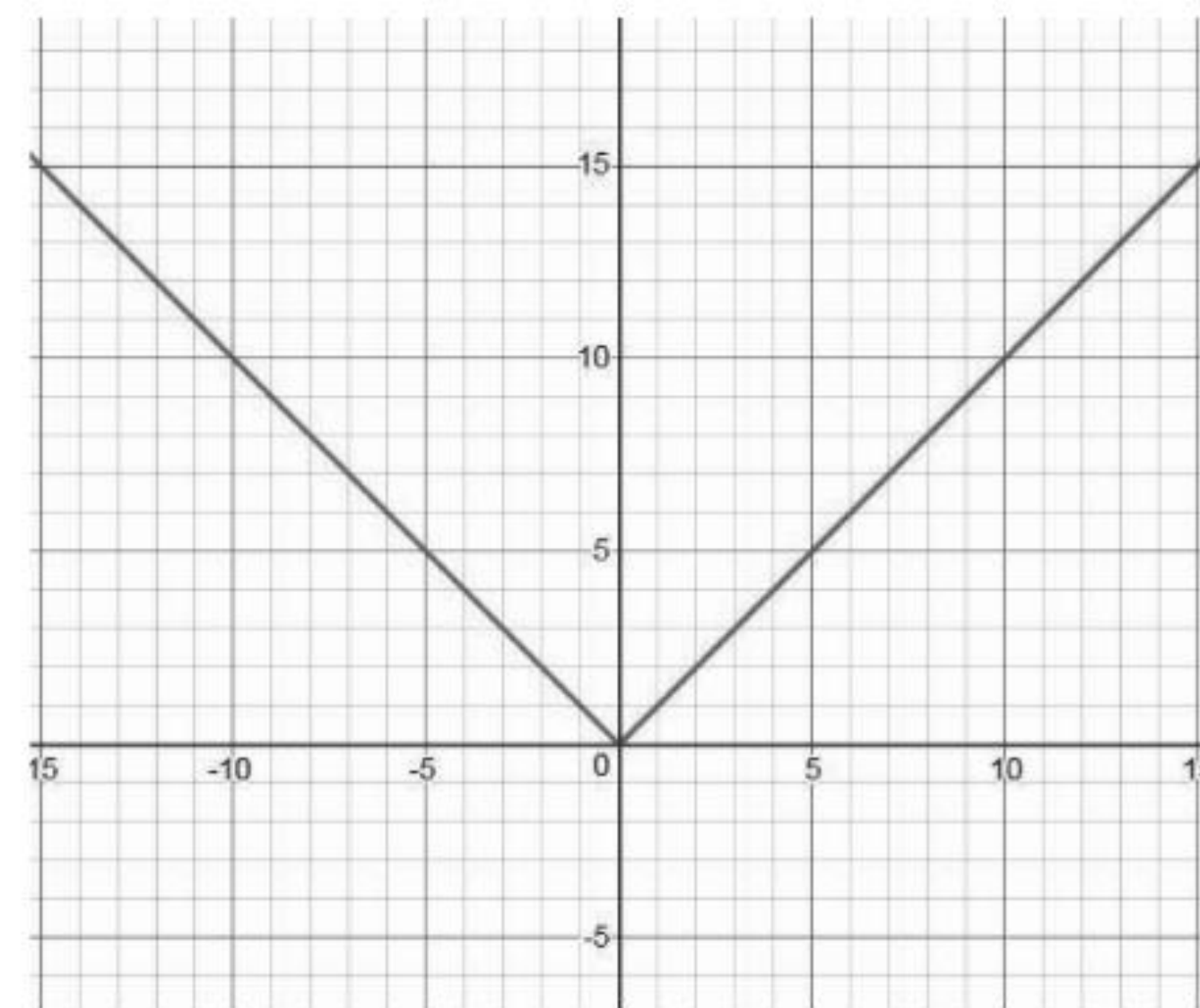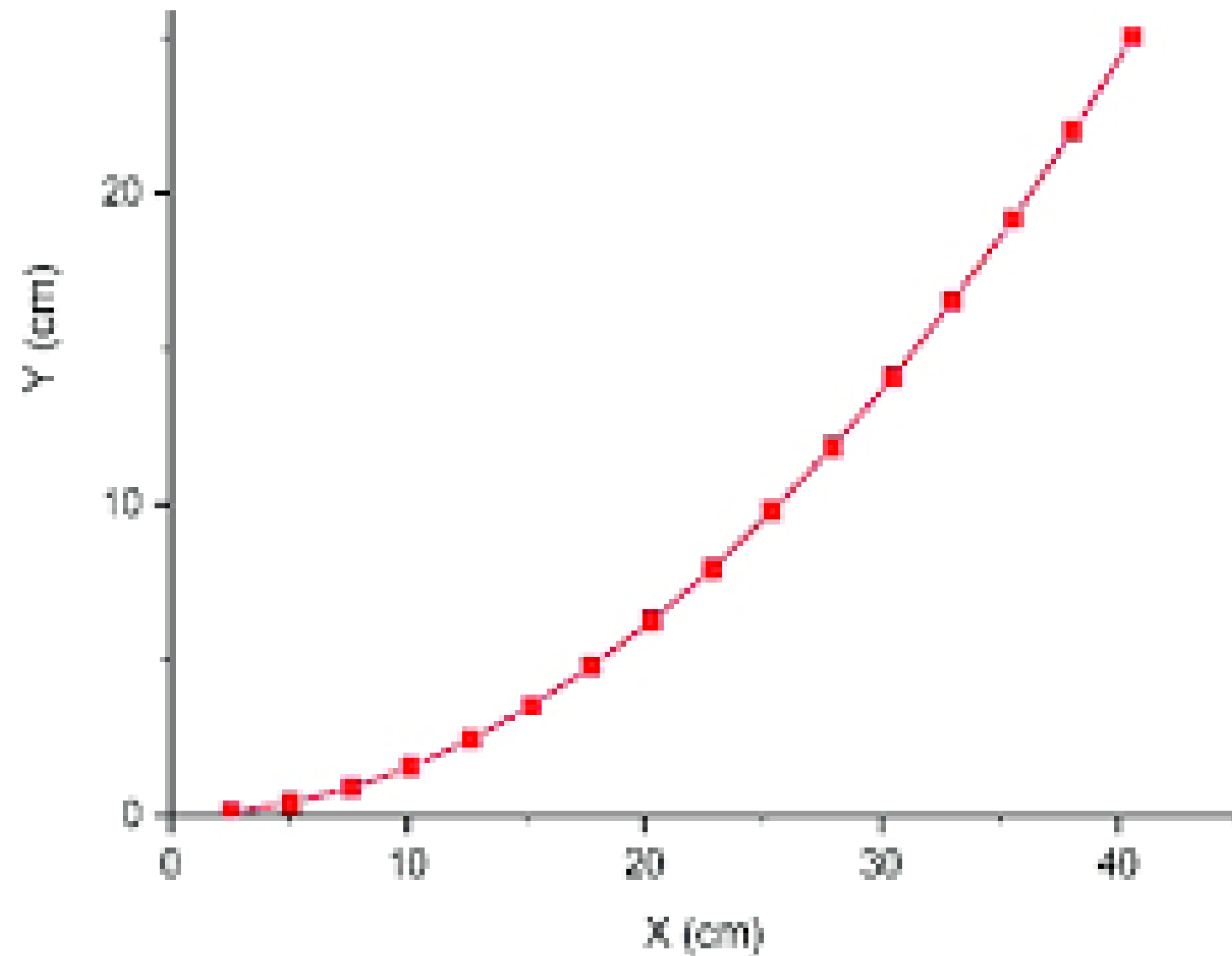
Negative Correlation

No Correlation

Positive Correlation

- Imagine two datasets with points **perfectly** along the curve
- What is the correl coeff?

# Which has more correlation coefficient ?

- Imagine two datasets where points are **perfectly** along these lines
- Which has more correlation coefficient?

# Correlation coefficient other forms

## In Linear Algebra

## In Statistics

$$\rho = \frac{(a - \bar{a})^T (b - \bar{b})}{\|a - \bar{a}\|\|b - \bar{b}\|}$$

$$Cov(a, b) = \frac{\Sigma_{i=1}^n (a - \bar{a})(b - \bar{b})}{n}$$

**Cov(a,a) = Var(a)**

$$\rho = Correl(a, b) = \frac{Cov(a, b)}{\sigma_a \sigma_b}$$

$$\rho = \frac{u^T v}{n} \quad \Longleftarrow \quad \rho = \frac{1}{n}\Sigma_{i=1}^n \left(\frac{a_i - \bar{a}}{\sigma_a}\right)\left(\frac{b_i - \bar{b}}{\sigma_b}\right)$$

$$where \quad u = \frac{a - \bar{a}}{\sigma_a} \quad v = \frac{b - \bar{b}}{\sigma_b}$$

10

# Covariance & Correlation in numpy

- Covariance: np.cov()
  - Returns covariance matrix
- Correlation Coefficient: np.corrcoef(a, b)
  - Two 1-D vectors passed
  - Supports only Pearson correlation coefficient
- There are two more – Spearman and Kendall
  - When are they used? - Reading assignment
- What is the relation to np.correlate()?
- Pandas corr()

# Visualizing correlation

- Seaborn pair plots and heatmaps

# Correlation is not causation

# Correlation is not causation

- Ice cream sale increases as summer heat increases
- Shark attack increases as summer heat increases
- Ice cream sale is highly correlated to shark attacks
- Wrong to conclude ice cream sale caused shark attack

# 2. Correlation between time series signals

# Auto correlation

- Correlation of data with itself
- Auto-correlation is always 1
- What if it is time lagged version of itself?

# Auto correlation

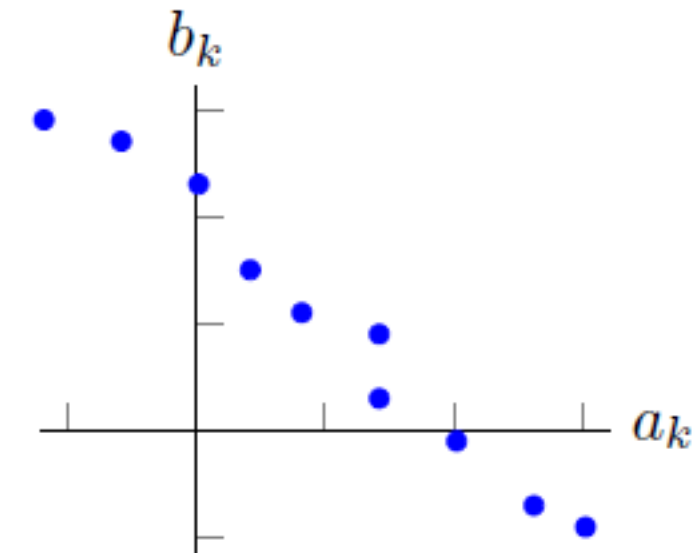- Correlation of time series signal with a lagged version of itself

$$r_k = \frac{\sum_{t=k+1}^{N}(y_t - \bar{y})(y_{t-k} - \bar{y})}{\sum_{t=1}^{N}(y_t - \bar{y})^2}$$

- Correlation of time series signal with a lagged version of itself is often useful

- Air pollution timeseries auto correlation

**Why is autocorrelation decreasing even for 12 month cycle with the passage of time?**



Autocorrelation Plot

- Notice the cyclical pattern in multiple of 12

# Correlation between time series
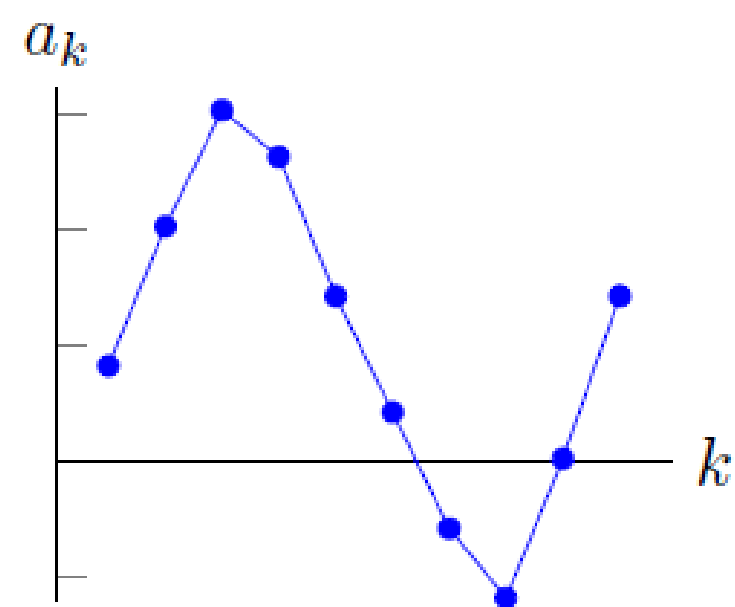
- Stock market correlation between stocks of same sector
- If they are not showing strong correlation, then something is wrong or some new opportunity
- Always analyze correlation with volatility (coefficient of standard deviation over a period)

# 2. Standard deviation of sum of vectors

# Standard deviation of the sum

- a and b are two equal sized vectors
- In statistical terms a and b are realizations of A and B

$$std(a + b) = \sqrt{std(a)^2 + std(b)^2 + 2\rho \, std(a) std(b)}$$

- Special cases:
  - Correlation coefficient = 1, -1 and 0

# Standard deviation of the sum

$$std(a + b) = \sqrt{std(a)^2 + std(b)^2 + 2\rho std(a)std(b)}$$

- Correlation coefficient = 1

$$std(a + b) = std(a) + std(b)$$

- Correlation coefficient = -1

$$std(a + b) = std(a) - std(b)$$

- Correlation coefficient = 0 (Uncorrelated features in ML)

$$std(a + b) = \sqrt{std(a)^2 + std(b)^2}$$

# Hedging investments

- Invest in two assets with same return ($\mu$) & risk($\sigma$)
- Asset returns over 5 year period in a and b vectors
- Hedged investment

$$c = \frac{a+b}{2} \qquad avg(c) = avg\left(\frac{a+b}{2}\right) = \mu \qquad \sigma_{\frac{a}{2}} = \sigma_{\frac{b}{2}} = \frac{\sigma}{2}$$

$$std(c) = \sqrt{std\left(\frac{a}{2}\right)^2 + std\left(\frac{b}{2}\right)^2 + 2\rho \, std(a) std(b)}$$

$$= \frac{\sqrt{2\sigma^2 + 2\rho\sigma^2}}{2} \qquad = \frac{\sigma}{\sqrt{2}}\sqrt{1+\rho}$$

**Two special cases: Rho = 0 and 1**

# Brief plan of what is next

- Remaining topics between chapter 1 & 3:
  - Linear Combinations