



# Principal Component Analysis

# Agenda

- Unsupervised Learning Approach
- Dimensionality Reduction
- Pre-requisite for fitting ML algorithms

# Principal Component Analysis (PCA)

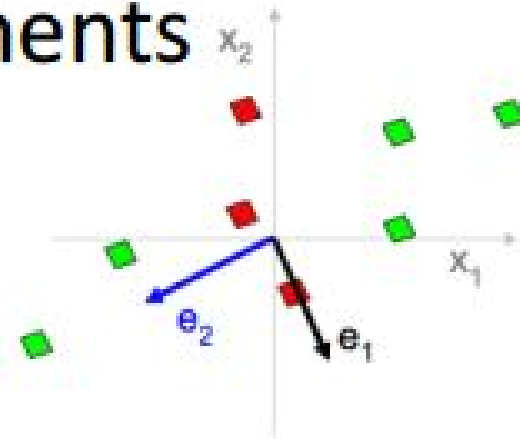
- Explain/summarize the underlying variance-covariance structure of a large set of variables through a few linear combinations of these variables

# Applications

- Uses:
  - Data Visualization
  - Data Reduction
  - Data Classification
  - Trend Analysis
  - Factor Analysis
  - Noise Reduction
- Examples:
  - How many unique “sub-sets” are in the sample?
  - How are they similar / different?
  - What are the underlying factors that influence the samples?
  - Which time / temporal trends are (anti)correlated?
  - Which measurements are needed to differentiate?
  - How to best present what is “interesting”?
  - Which “sub-set” does this new sample rightfully belong?

# Principal components

- Compute covariance matrix  $\Sigma$ 
  - covariance of dimensions  $x_1$  and  $x_2$ :
    - do  $x_1$  and  $x_2$  tend to increase together?
    - or does  $x_2$  decrease as  $x_1$  increases?
  - covariance: measure of variability



$$\begin{array}{c}
 \begin{matrix} x_1 & x_2 \\ x & \begin{pmatrix} 2.0 & 0.8 \\ 0.8 & 0.6 \end{pmatrix} \\ x_2^T & \end{matrix}
 \end{array}
 \rightarrow \text{var}(x_2) = \frac{1}{N} \sum_{j=1}^n (x_{2,j} - \mu_2)^2$$

$$\text{cov}(x_1, x_2) = \frac{1}{N} \sum_{j=1}^n (x_{1,j} - \mu_1)(x_{2,j} - \mu_2)$$

- Find the basis of  $\Sigma$ 
  - find vectors  $e_i$  which aren't turned by  $\Sigma$ 
    - $\Sigma e_i = \lambda_i e_i$ : eigenvalue / eigenvector
  - 1<sup>st</sup> PC: "longest"  $e_i$  (has largest  $\lambda_i$ ), 2<sup>nd</sup> PC: next longest, ...

$$\begin{array}{c}
 \lambda_1 \begin{pmatrix} 0.26 \\ 2.42 \end{pmatrix} \quad x \begin{pmatrix} e_1 & e_2 \\ 0.4 & -0.9 \\ -0.9 & -0.4 \end{pmatrix} \\
 \lambda_2
 \end{array}$$