



# Applied Machine Learning

Dr. Harikrishnan N B

Computer Science and Information Systems

**BITS Pilani**  
Pilani Campus



# **SE ZG568 / SS ZG568, Applied Machine Learning Lecture No. 6 [23- Feb-2025]**

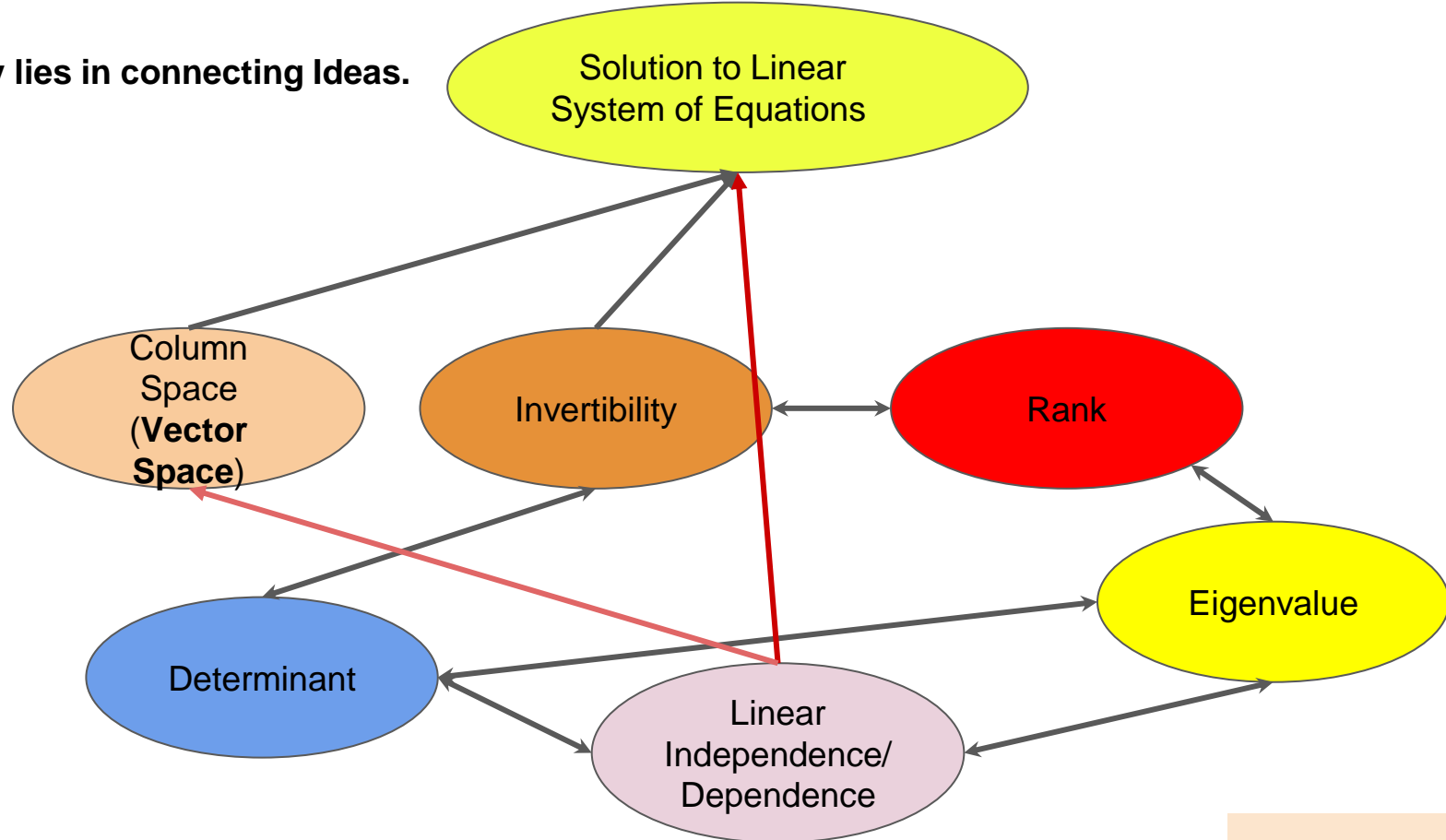
# Recap

---



Basics of Linear Algebra,  
Row Picture, Col Picture, Algebraic Way  
Solution to System of Linear Equations  
Inverse of a Matrix  
Linear Regression,  
**PCA**

Beauty lies in connecting Ideas.



## Practical Challenges and Important Points

**When can we apply  $A = X \Lambda X^{-1}$  ?**

- A should be a square matrix
- When A has 'n' linearly independent eigenvectors, then  $X^{-1}$  always exist.

**What happens when A is Symmetric (  $A^T = A$  )?**

- The eigenvectors of a symmetric matrix A can be chosen as **ORTHONORMAL**. So in this case **X is orthonormal**.
- For an **ORTHONORMAL** matrix X, the inverse is its transpose  **$X^{-1} = X^T$**

$$A = X \Lambda X^{-1}$$

$$A = X \Lambda X^T$$

- The eigenvectors of a symmetric matrix  $A$  can be chosen as **ORTHONORMAL**. So in this case  **$X$  is orthonormal. Why?**

## Practical Challenges and Important Points

**When can we apply  $A = X \Lambda X^{-1}$  ?**

- A should be a square matrix
- When A has 'n' linearly independent eigenvectors, then  $X^{-1}$  always exist.

**What happens when A is Symmetric (  $A^T = A$  )?**

- The eigenvectors of a symmetric matrix A can be chosen as **ORTHONORMAL**. So in this case **X is orthonormal**.
- For an **ORTHONORMAL** matrix X, the inverse is its transpose  $X^{-1} = X^T$

$$A = X \Lambda X^{-1}$$

$$A = X \Lambda X^T$$





## Practical Challenges

**What if  $A$  is not a square matrix?**

- We cannot apply Spectral Decomposition.

**Don't Worry!!!**



**Singular Value Decomposition works for any Matrix.**



## A Few more steps to PCA

What all minimum can we say about this data?

Day 1 in knee 1      Day 2 in knee 2

$f_1$  → Height  
 $f_2$  → Weight

|   |   |   |   |   |   |
|---|---|---|---|---|---|
| X | 1 | 2 | 3 | 4 | 5 |
| Y | 1 | 5 | 4 | 6 | 7 |

**X, Y** are the features

What all minimum can we say about this data?

$$\text{Mean}(X) = \frac{1}{N} \sum_{i=1}^N x_i$$

$$\text{Variance}(X) = \frac{1}{N-1} \sum_{i=1}^N (x_i - \mu_x)^2$$

$$\text{Cov}(X, Y) = \frac{1}{N-1} \sum_{i=1}^N (x_i - \mu_x)(y_i - \mu_y)$$

| X        | Y       | var(X) | var(Y) | cov(X,Y) | cov(Y,X) |
|----------|---------|--------|--------|----------|----------|
| 1        | 1       |        |        |          |          |
| 2        | 5       |        |        |          |          |
| 3        | 4       |        |        |          |          |
| 4        | 6       |        |        |          |          |
| 5        | 7       |        |        |          |          |
| Mean (X) | Mean(Y) | var(X) | var(Y) | cov(X,Y) | cov(Y,X) |
| ?        | ?       | ?      | ?      | ?        | ?        |

$$\text{Mean}(X) = \frac{1}{N} \sum_{i=1}^N x_i$$

$$\text{Variance}(X) = \frac{1}{N-1} \sum_{i=1}^N (x_i - \mu_x)^2$$

$$\text{Cov}(X, Y) = \frac{1}{N-1} \sum_{i=1}^N (x_i - \mu_x)(y_i - \mu_y) \quad \checkmark$$

$$\mu_x = \text{Mean}(X) = \frac{1+2+3+4+5}{5} = 3$$

$$\mu_y = \text{Mean}(Y) = \frac{1+5+4+6+7}{5} = \frac{23}{5} = 4.6$$

$$\text{Var}(X) = 2.5, \text{Var}(Y) = 5.3$$

| $\text{Var}(X)$ | $\text{Var}(Y)$     |
|-----------------|---------------------|
| $(1-3)^2 = 4$   | $(1-4.6)^2 = 12.96$ |
| $(2-3)^2 = 1$   | $(5-4.6)^2 = 0.16$  |
| $(3-3)^2 = 0$   | $(4-4.6)^2 = 0.36$  |
| $(4-3)^2 = 1$   | $(6-4.6)^2 = 1.96$  |
| $(5-3)^2 = 4$   | $(7-4.6)^2 = 5.76$  |
| $\frac{10}{4}$  |                     |

$$\text{Cov}(X, Y) = \frac{\cancel{4 \times 12.96} (1-3) \times (1-4.6) + (2-3) \times (5-4.6) + (3-3) \times (4-4.6) + (4-3) \times (6-4.6) + (5-3) \times (7-4.6)}{4} = 3.25$$

$$\text{Mean}(\mathbf{X}) = \frac{1}{N} \sum_{i=1}^N x_i$$

$$\text{Variance}(\mathbf{X}) = \frac{1}{N-1} \sum_{i=1}^N (x_i - \mu_x)^2$$

$$\text{Cov}(\mathbf{X}, \mathbf{Y}) = \frac{1}{N-1} \sum_{i=1}^N (x_i - \mu_x)(y_i - \mu_y)$$

| X | Y |
|---|---|
| 1 | 1 |
| 2 | 5 |
| 3 | 4 |
| 4 | 6 |
| 5 | 7 |

What all minimum can we say about this data?

$$\text{Mean}(X) = \frac{1}{N} \sum_{i=1}^N x_i$$

$$\text{Variance}(X) = \frac{1}{N-1} \sum_{i=1}^N (x_i - \mu_x)^2$$

$$\text{Cov}(X,Y) = \frac{1}{N-1} \sum_{i=1}^N (x_i - \mu_x)(y_i - \mu_y)$$

| X        | Y       | var(X) | var(Y) | cov(X,Y) | cov(Y,X) |
|----------|---------|--------|--------|----------|----------|
| 1        | 1       |        |        |          |          |
| 2        | 5       |        |        |          |          |
| 3        | 4       |        |        |          |          |
| 4        | 6       |        |        |          |          |
| 5        | 7       |        |        |          |          |
| Mean (X) | Mean(Y) | var(X) | var(Y) | cov(X,Y) | cov(Y,X) |
| 3.0      | 4.6     | 2.5    | 5.3    | 2.25     | 2.25     |

$$D = D_{kfa} = \begin{array}{|c|c|c|c|} \hline X & 2 & 4 & \dots \\ \hline y & 3 & 5 & \dots \\ \hline \end{array}$$


$$D_{mean} = \begin{array}{|c|c|c|c|} \hline X - \mu_x & 2 - \mu_x & 4 - \mu_x & \dots \\ \hline y - \mu_y & 3 - \mu_y & 5 - \mu_y & \dots \\ \hline \end{array}_{2 \times N}$$

$$\frac{1}{N-1} \left( D_{mean} D_{mean}^T \right)_{2 \times 2}$$

$$\frac{1}{N-1} \left( D_{mean} D_{mean}^T \right) \Rightarrow \begin{array}{|c|} \hline (x_1 - \mu_x) \quad x_2 - \mu_x \quad \dots \quad x_n - \mu_x \\ \hline (y_1 - \mu_y) \quad y_2 - \mu_y \quad \dots \quad (y_n - \mu_y) \\ \hline \end{array}_{2 \times N}$$

$$\frac{1}{N-1} \left[ \sum_{i=1}^n (x_i - \mu_x)^2 \right]_{2 \times 2}$$

$$\frac{1}{N-1} \left[ \begin{array}{|c|} \hline (x_1 - \mu_x) \quad y_1 - \mu_y \\ x_2 - \mu_x \quad y_2 - \mu_y \\ \vdots \quad \vdots \\ x_n - \mu_x \quad y_n - \mu_y \\ \hline \end{array} \right]_{N \times 2}$$

$$\frac{1}{N-1} \begin{bmatrix} \text{Var}(x) & \text{Cov}(x, y) \\ \text{Cov}(y, x) & \text{Var}(y) \end{bmatrix} \rightarrow \sum_{i=1}^N (x_i - \mu_x)(y_i - \mu_y)$$



## Variance- Covariance Matrix

$$A = U \Lambda U^T$$

$$\frac{1}{N-1} D_{\text{mean}} D_{\text{mean}}^T$$

$$\begin{matrix} X & & Y \\ X & \begin{bmatrix} \text{var}(X) & \text{cov}(X, Y) \\ \text{cov}(Y, X) & \text{var}(Y) \end{bmatrix} \\ Y & \end{matrix}$$

Recall the properties of a symmetric matrix!!!

$$\begin{matrix} & X & Y & Z \\ X & \text{var}(X) & \text{cov}(X, Y) & \text{cov}(X, Z) \\ Y & \text{cov}(Y, X) & \text{var}(Y) & \text{cov}(Y, Z) \\ Z & \text{cov}(Z, X) & \text{cov}(Z, Y) & \text{var}(Z) \end{matrix}$$

- Variance - Covariance Matrix is symmetric.  $\text{cov}(X, Y) = \text{cov}(Y, X)$
- The diagonal entries represents variance
- The off-diagonal entries represents the correlation of X and Y

## What does Variance - Covariance Matrix signifies?

Case I

$$\begin{bmatrix} 2 & 3 \\ 3 & 5 \end{bmatrix}$$

Case II

$$\begin{bmatrix} 2 & -3 \\ -3 & 5 \end{bmatrix}$$

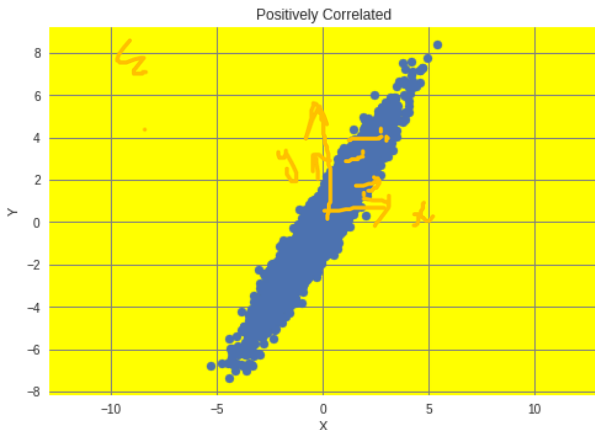
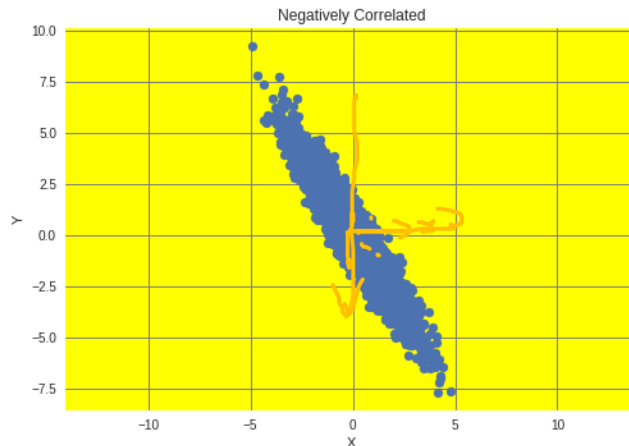
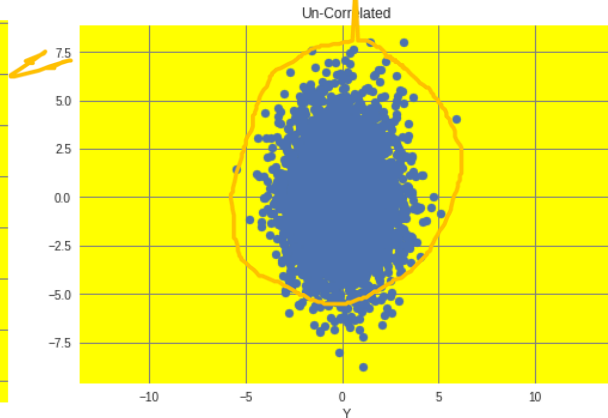
Case III

$$\begin{bmatrix} 2 & 0 \\ 0 & 5 \end{bmatrix}$$

Note: In all cases mean is (0,0)

$$p(x, y) = p(x) \cdot p(y)$$
 Case 2  

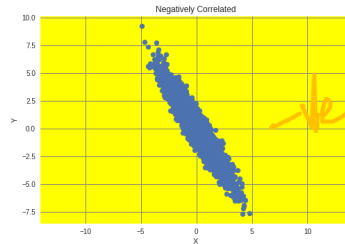
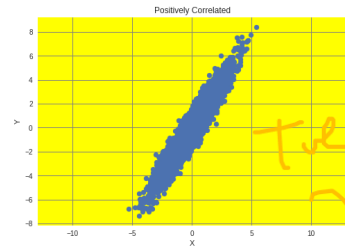
$$x^2 + y^2 = 1$$

$$\begin{bmatrix} 2 & 3 \\ 3 & 5 \end{bmatrix}$$

$$\begin{bmatrix} 2 & -3 \\ -3 & 5 \end{bmatrix}$$

$$\begin{bmatrix} 2 & 0 \\ 0 & 5 \end{bmatrix}$$


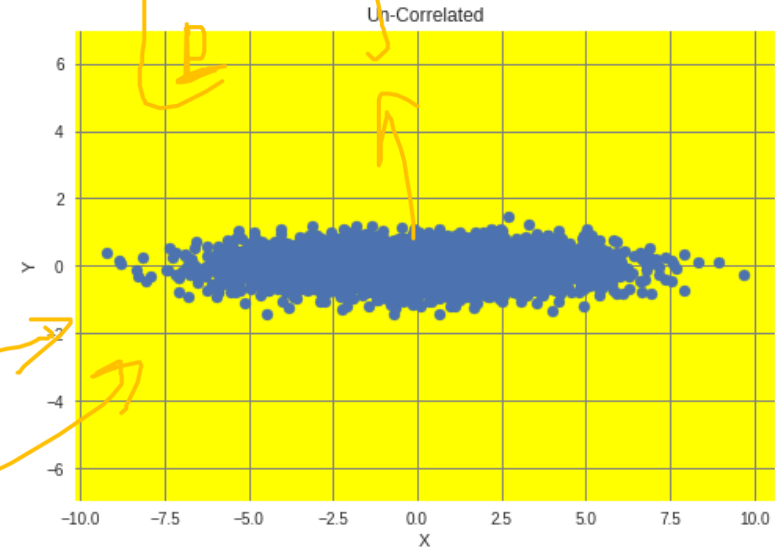
Harikrishnan N B

## So what ~~do~~ PCA do ?

- Principal Component Analysis (PCA) makes the data **UNCORRELATED**.



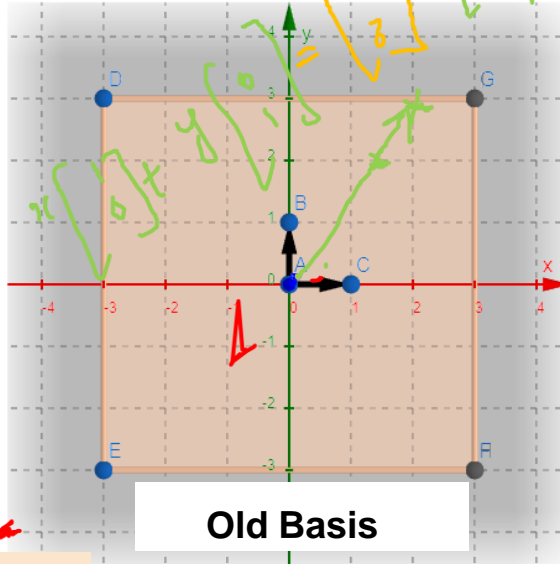
PCA achieves this by  
**Change of Basis**



# Change of Basis



My life is full of struggle.



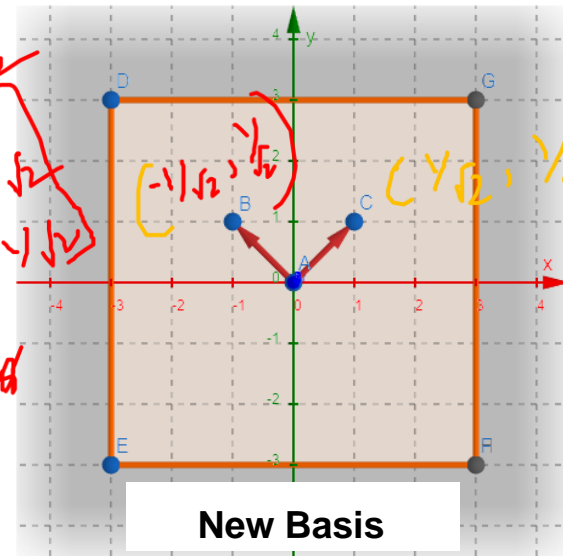
Old Basis

IT depends on the way you look at life. Try to see life from a different **POINT OF VIEW**



Handwritten red notes:

- $b_1 \checkmark$  and  $b_2 \checkmark$  with arrows pointing to the new basis vectors.
- A transformation matrix:  $\begin{bmatrix} 1/\sqrt{2} & -1/\sqrt{2} \\ 1/\sqrt{2} & 1/\sqrt{2} \end{bmatrix}$
- The text "Basis Vectors" with two arrows pointing towards the new basis graph.



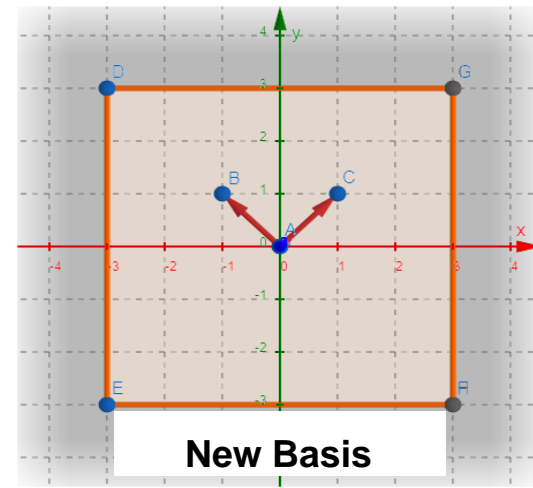
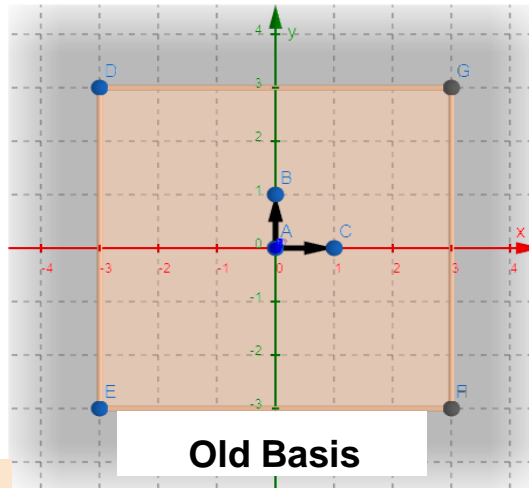
New Basis

## Recall

**Dimension of a Vector space** - Every vector space has a dimension. Dimension is the number of basis vectors required to span the vector space.

**Properties of Basis Vectors** -

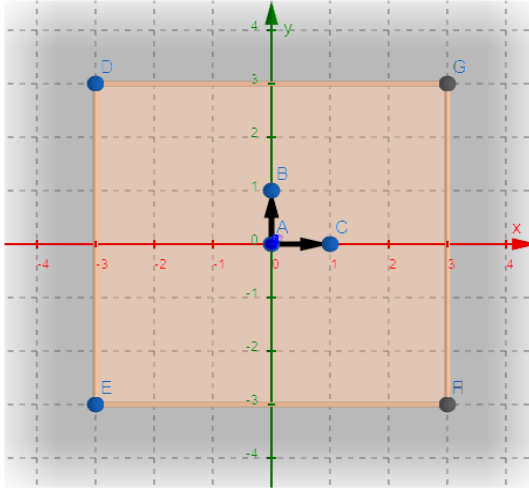
- Basis vectors has to be linearly independent.
- Basis vectors should span the vector space.



# Example of Change of Basis

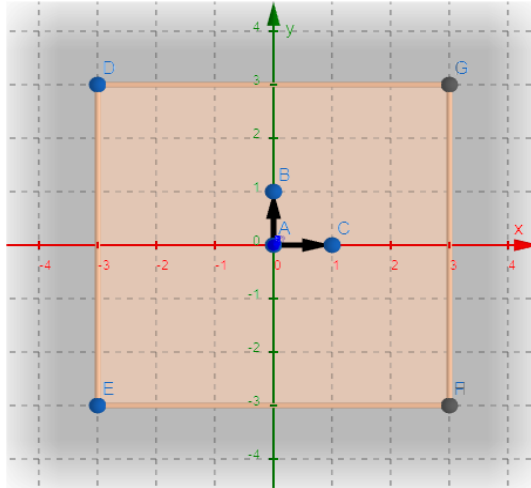
To represent a point (2,3) in old basis and new basis- How to understand this?

$$\text{Old basis} = \left\{ \begin{bmatrix} 1 \\ 0 \end{bmatrix}, \begin{bmatrix} 0 \\ 1 \end{bmatrix} \right\}$$



# Example of Change of Basis

To represent a point (2,3) in old basis and new basis- How to understand this?



$$\text{Old basis} = \left\{ \begin{bmatrix} 1 \\ 0 \end{bmatrix}, \begin{bmatrix} 0 \\ 1 \end{bmatrix} \right\}$$

$$2 \begin{bmatrix} 1 \\ 0 \end{bmatrix} + 3 \begin{bmatrix} 0 \\ 1 \end{bmatrix} = \begin{bmatrix} 2 \\ 3 \end{bmatrix}$$

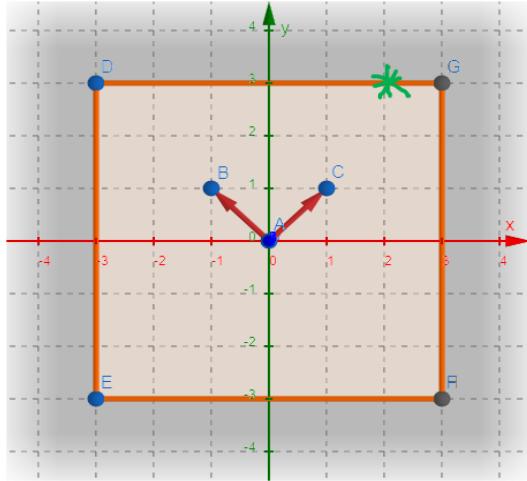
Handwritten notes in red ink show the vectors  $b_1 = \begin{bmatrix} 1 \\ 0 \end{bmatrix}$  and  $b_2 = \begin{bmatrix} 0 \\ 1 \end{bmatrix}$  being combined to form the vector  $\begin{bmatrix} 2 \\ 3 \end{bmatrix}$ . Below this, a diagram illustrates the linear combination:

$$\left[ \begin{array}{|c|} \hline b_1 \\ \hline \end{array} \right] \begin{array}{|c|} \hline b_2 \\ \hline \end{array} \begin{array}{|c|} \hline x \\ \hline y \end{array} = \begin{array}{|c|} \hline 2 \\ \hline 3 \end{array}$$

The vectors  $b_1$  and  $b_2$  are shown in orange boxes,  $x$  and  $y$  in a brown box, and the resulting vector  $\begin{bmatrix} 2 \\ 3 \end{bmatrix}$  in a red box. Red arrows indicate the mapping from the matrix equation above to this diagram.



## New Basis Representation



$$\text{New Basis} = \left\{ \begin{bmatrix} \frac{1}{\sqrt{2}} \\ \frac{1}{\sqrt{2}} \end{bmatrix}, \begin{bmatrix} \frac{-1}{\sqrt{2}} \\ \frac{1}{\sqrt{2}} \end{bmatrix} \right\}$$



$$x \begin{bmatrix} \frac{1}{\sqrt{2}} \\ \frac{1}{\sqrt{2}} \end{bmatrix} + y \begin{bmatrix} \frac{-1}{\sqrt{2}} \\ \frac{1}{\sqrt{2}} \end{bmatrix} = \begin{bmatrix} 2 \\ 3 \end{bmatrix}$$



$$b_1 = \begin{bmatrix} 1 \\ 1 \end{bmatrix}$$

$$b_2 = \begin{bmatrix} 1 \\ -1 \end{bmatrix}$$

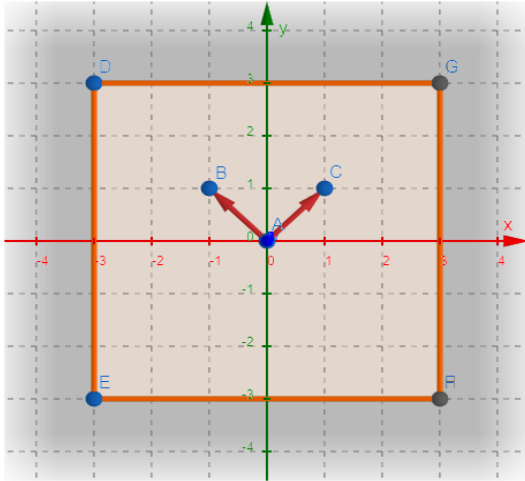
$$\begin{bmatrix} 1 & 1 \\ 1 & -1 \end{bmatrix} \begin{bmatrix} x \\ y \end{bmatrix} = \begin{bmatrix} a \\ b \end{bmatrix}$$

$$\begin{bmatrix} x \\ y \end{bmatrix} = \begin{bmatrix} 1 & 1 \\ 1 & -1 \end{bmatrix}^{-1} \begin{bmatrix} a \\ b \end{bmatrix}$$

$$\begin{bmatrix} 1 & 1 \\ 1 & -1 \end{bmatrix} \cdot \begin{bmatrix} 2 \\ 4 \end{bmatrix} = \begin{bmatrix} 2 \\ 0 \end{bmatrix}$$

$$x \begin{bmatrix} 1 \\ 1 \end{bmatrix} + y \begin{bmatrix} 1 \\ -1 \end{bmatrix} = \begin{pmatrix} 2 \\ 0 \end{pmatrix}$$

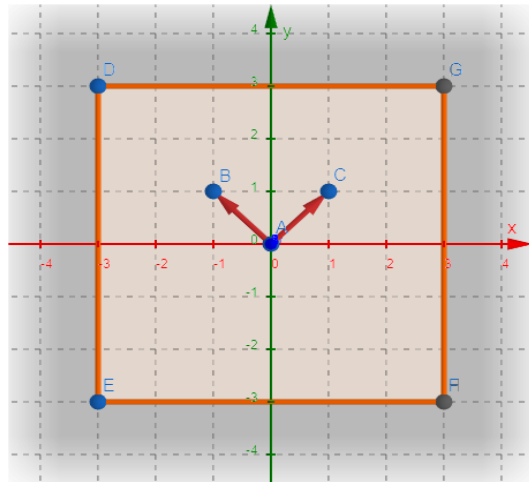
## Finding x and y for representing (2,3) using new basis



$$x \begin{bmatrix} \frac{1}{\sqrt{2}} \\ \frac{1}{\sqrt{2}} \end{bmatrix} + y \begin{bmatrix} \frac{-1}{\sqrt{2}} \\ \frac{1}{\sqrt{2}} \end{bmatrix} = \begin{bmatrix} 2 \\ 3 \end{bmatrix}$$

$$\begin{bmatrix} \boxed{\frac{1}{\sqrt{2}}}_{b_1} & \boxed{\frac{-1}{\sqrt{2}}}_{b_2} \end{bmatrix} \begin{bmatrix} x \\ y \end{bmatrix} = \begin{bmatrix} 2 \\ 3 \end{bmatrix}$$

$$\begin{bmatrix} \frac{1}{\sqrt{2}} & \frac{-1}{\sqrt{2}} \\ \frac{1}{\sqrt{2}} & \frac{1}{\sqrt{2}} \end{bmatrix} \begin{bmatrix} x \\ y \end{bmatrix} = \begin{bmatrix} 2 \\ 3 \end{bmatrix}$$



$$\begin{bmatrix} \frac{1}{\sqrt{2}} & \frac{-1}{\sqrt{2}} \\ \frac{1}{\sqrt{2}} & \frac{1}{\sqrt{2}} \end{bmatrix} \begin{bmatrix} x \\ y \end{bmatrix} = \begin{bmatrix} 2 \\ 3 \end{bmatrix}$$

$$\vec{x} = \begin{bmatrix} x \\ y \end{bmatrix}$$

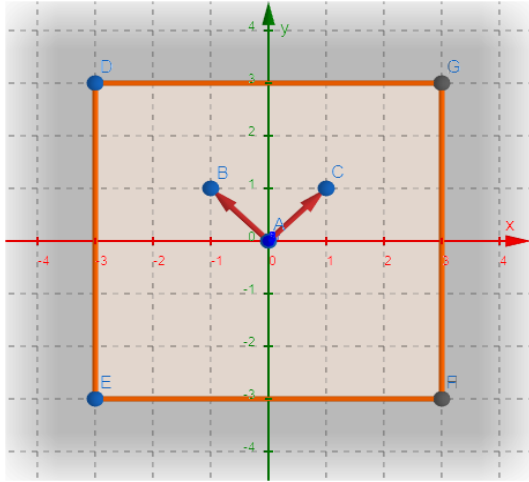
$$\begin{aligned} &\downarrow \\ P\vec{x} &= \vec{y} \\ P^{-1}P\vec{x} &= P^{-1}\vec{y} \\ \vec{x} &= P^{-1}\vec{y} \end{aligned}$$

$$A = \begin{bmatrix} a & b \\ c & d \end{bmatrix}$$

$$A^{-1} = \frac{1}{|A|} \begin{bmatrix} d & -b \\ -c & a \end{bmatrix}$$

For **ORTHONORMAL MATRIX**,  $P^{-1} = P^T$

In our case the matrix P is ORTHONORMAL

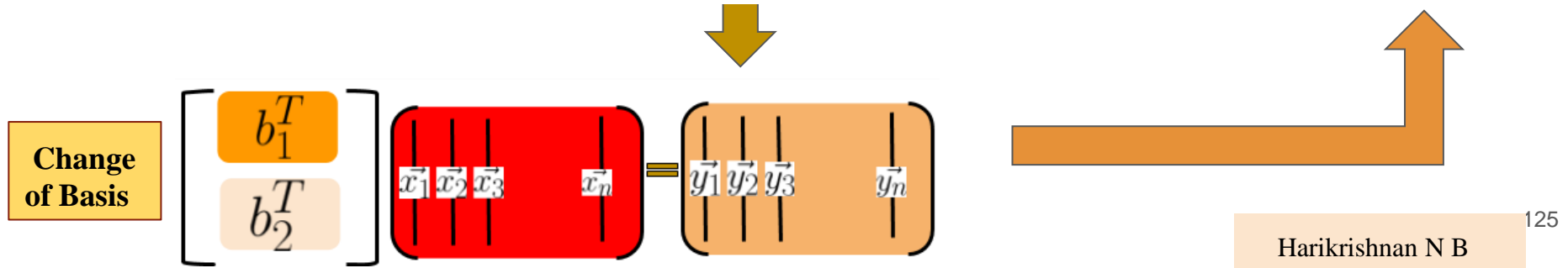
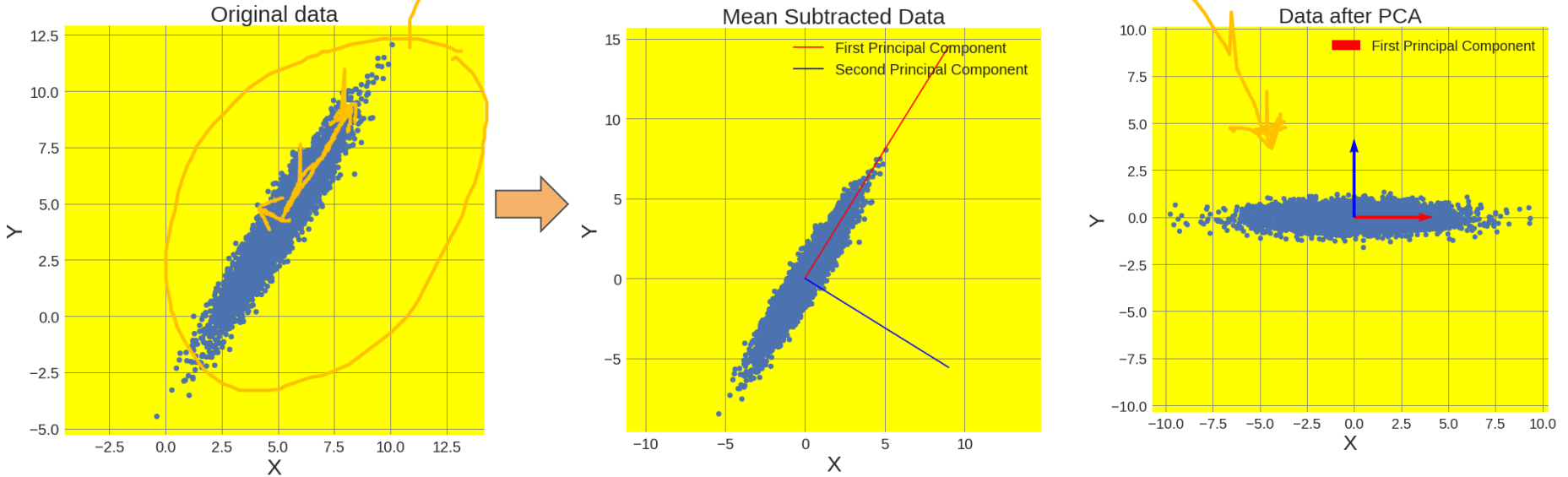


$$\vec{x} = P^{-1}\vec{y} = P^T\vec{y}$$

$$\begin{bmatrix} b_1^T \\ b_2^T \end{bmatrix} \begin{bmatrix} 2 \\ 3 \end{bmatrix} = \begin{bmatrix} x \\ y \end{bmatrix} \begin{bmatrix} \frac{1}{\sqrt{2}} & \frac{1}{\sqrt{2}} \\ \frac{-1}{\sqrt{2}} & \frac{1}{\sqrt{2}} \end{bmatrix} \begin{bmatrix} 2 \\ 3 \end{bmatrix} = \begin{bmatrix} \frac{5}{\sqrt{2}} \\ \frac{1}{\sqrt{2}} \end{bmatrix}$$

$$\frac{5}{\sqrt{2}} \begin{bmatrix} \frac{1}{\sqrt{2}} \\ \frac{1}{\sqrt{2}} \end{bmatrix} + \frac{1}{\sqrt{2}} \begin{bmatrix} \frac{-1}{\sqrt{2}} \\ \frac{1}{\sqrt{2}} \end{bmatrix} = \begin{bmatrix} 2 \\ 3 \end{bmatrix}$$

# Steps in PCA



$$D_{\text{mean}} = \begin{bmatrix} x_1 - \mu_x & x_2 - \mu_x & \dots & x_n - \mu_x \\ y_1 - \mu_y & y_2 - \mu_y & \dots & y_n - \mu_y \end{bmatrix}$$

$$Pb_{\text{mean}} = Z \rightarrow \text{linearly uncorrelated}$$

$$Pb_{\text{mean}} = Z \rightarrow \text{Weakly Uncorrelated}$$

$$\text{cov}(Z) = \frac{1}{N-1}$$

var-cov matrix of  $Z$

$$(AB)^T = B^T A^T$$

$$\frac{1}{N-1} Z Z^T$$

$$\frac{1}{N-1} PD_{\text{mean}} (Pb_{\text{mean}})^T$$

$$= \frac{1}{N-1} PD_{\text{mean}} D_{\text{mean}}^T P^T$$

$$= \frac{1}{N-1} P D_m D_m^T P^T \quad | \text{ call}$$

$$V^T V = I$$
  
~~$$V V^T = I$$~~  
$$V V^T = I$$

$$P \left[ \begin{array}{cc} 1 & D_m D_m^T \\ N-1 & \end{array} \right] P^T \rightarrow \text{se}$$

Symmetrie

$$A = \begin{bmatrix} \text{Var}(x) & \text{Cov}(x, y) \\ \text{Cov}(y, x) & \text{Var}(y) \end{bmatrix}$$

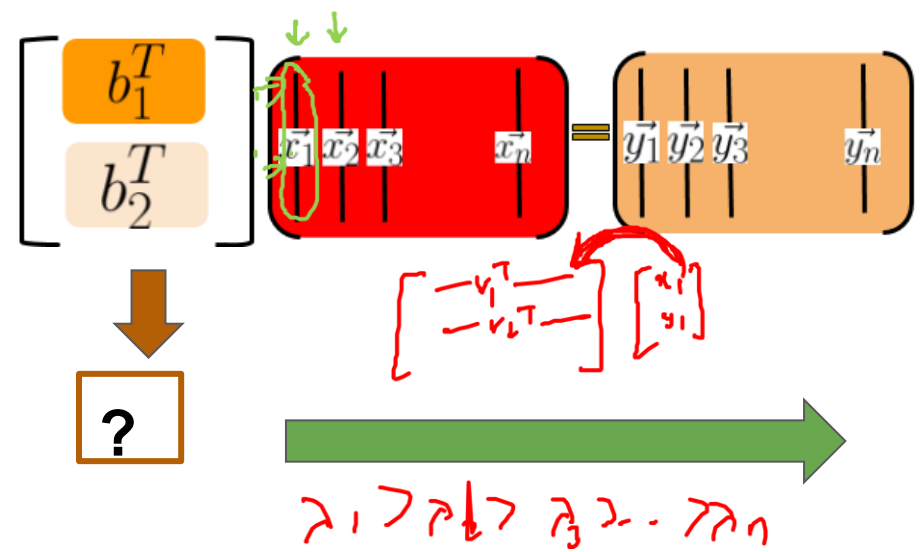
$$\text{Var-cov}(Z) = P \underbrace{V \Lambda V^T}_{W} P^T$$

$V$  is the eigen  $\rightarrow$  variance covariance matrix of  $Z$   
 $V$  is orthogonal.  
 that the anti-diagonal var-cov( $Z$ ) is 0.

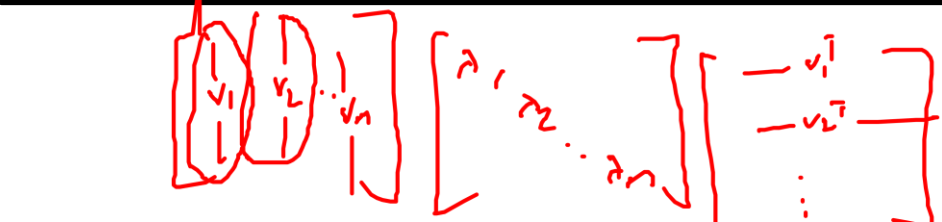
$[0]$   $\downarrow$   $p$  is unknown  
What should be my IP



What should be the NEW BASIS so that DATA is UNCORRELATED?



Rows of matrix P are the **eigenvectors** of the **variance-covariance matrix** of the **mean subtracted data**.



$$\begin{aligned}
 P\bar{X} &= \bar{Y} \\
 \text{var} \text{cov}(Y) &= \text{cov}(PX) \\
 \text{var} \text{cov}(PX) &= \frac{1}{N-1} (PX)(PX)^T \\
 \text{var} \text{cov}(PX) &= \frac{1}{N-1} PXX^T P^T \\
 \text{var} \text{cov}(PX) &= P \left( \frac{1}{N-1} XX^T \right) P^T \\
 \text{var} \text{cov}(PX) &= P \text{cov}(X) P^T \\
 \text{var} \text{cov}(PX) &= P(V\Lambda V^T) P^T \\
 P &= V^T \\
 \text{var} \text{cov}(PX) &= \Lambda
 \end{aligned}$$

## Some words about PCA

- PCA is “an orthogonal linear transformation that transfers the data to a new coordinate system such that the greatest variance by any projection of the data comes to lie on the first coordinate (*first principal component*), the second greatest variance lies on the second coordinate (*second principal component*), and so on.”

# Applications of PCA

number of features = 10

- Dimensionality Reduction
- Denoising
- Feature Extraction
- Image Compression
- EEG Analysis

$$P_{10 \times 10} = \begin{bmatrix} -v_1^T \\ -v_2^T \\ \vdots \\ -v_{10}^T \end{bmatrix}$$

$$P_{10 \times 10} D_{mean 10 \times 100} = Z_{10 \times 100}$$

↑  
Data is here

$Z[0:2, :]$   
 ↳ first two rows = 2  
 first row = 1  
 second row = 1  
 ⇒ first 2 PCs  
 ⇒ read PC

## Assumptions in PCA

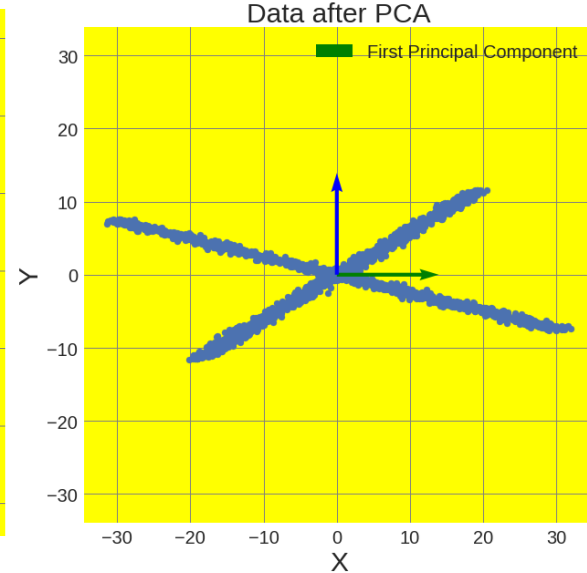
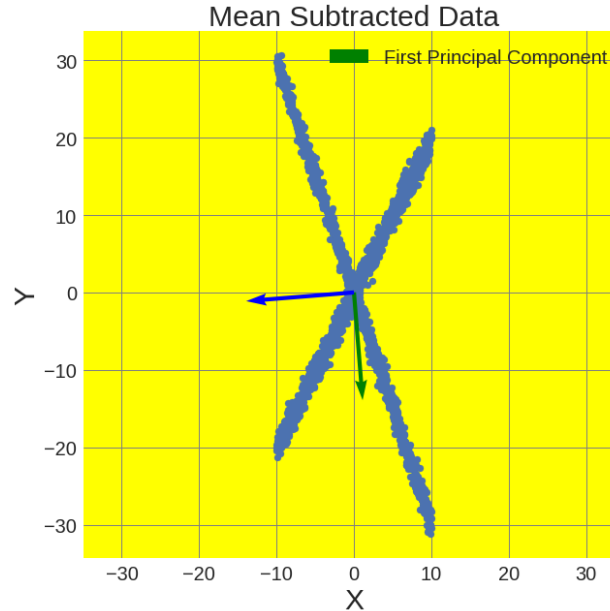
- Linearity
- Large variance have important structure
- Principal components are orthogonal

$$\begin{array}{c} \text{compact} \quad \text{size of cell} \\ \begin{bmatrix} f_1 & f_2 & \vdots \\ \vdots & \vdots & \ddots \end{bmatrix} \end{array}$$

$$\begin{array}{l} \vec{f}_1 \in [0, 1] \\ \vec{f}_2 \in [0, 1] \end{array}$$

## When does PCA fail?

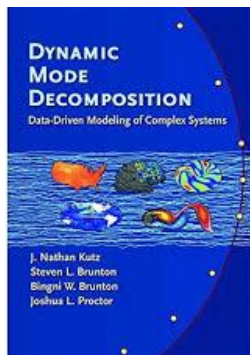
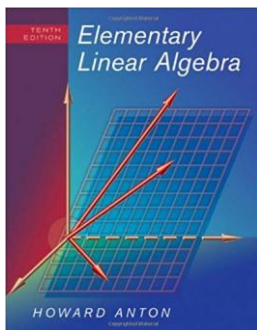
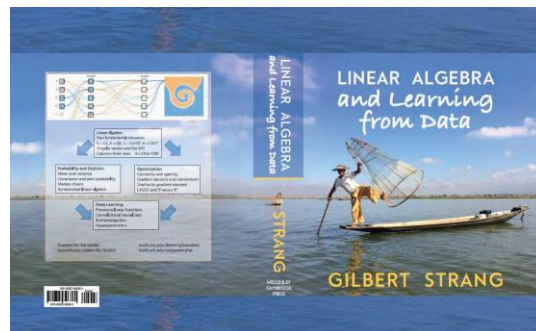
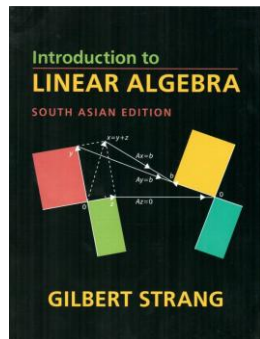
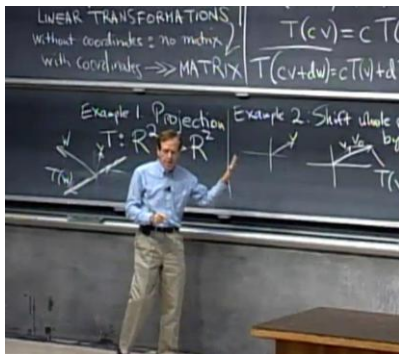
- Non-linearity
- Non-Gaussian
- Non-orthogonality



Ref: <https://arxiv.org/abs/1404.1100>

# Interesting Materials

## Prof. Gilbert Strang



**Tutorial on PCA - ([Click here](#))**



# Learning - Function Approximation

---

## Problem Setting

- Set of possible instances  $X$ .
- Unknown target function  $f: X \rightarrow Y$
- Set of function hypotheses  $H = \{h \mid h: X \rightarrow Y\}$

## Input

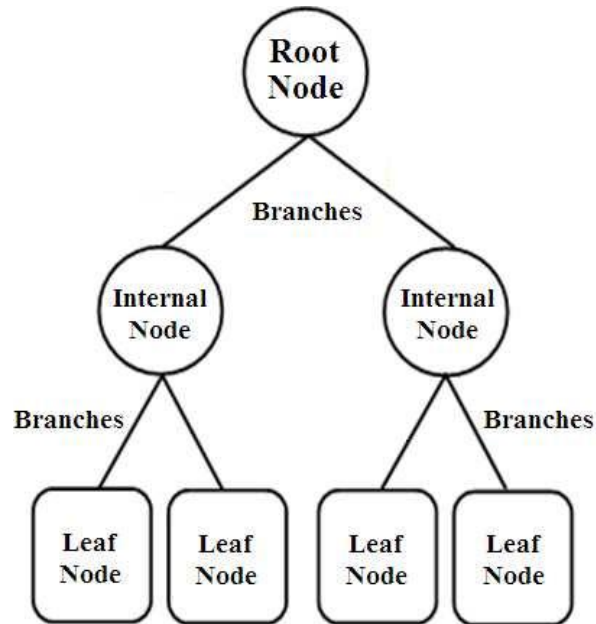
Training Examples of unknown target function  $f$ .

## Output

Hypothesis  $h \in H$  that best approximates target function  $f$ .

# Decision Tree

**Decision Tree Learning** is a method for approximating the target function ( $Y$ ), in which the **learned functions ( $h$ )** is represented by a **decision tree**.

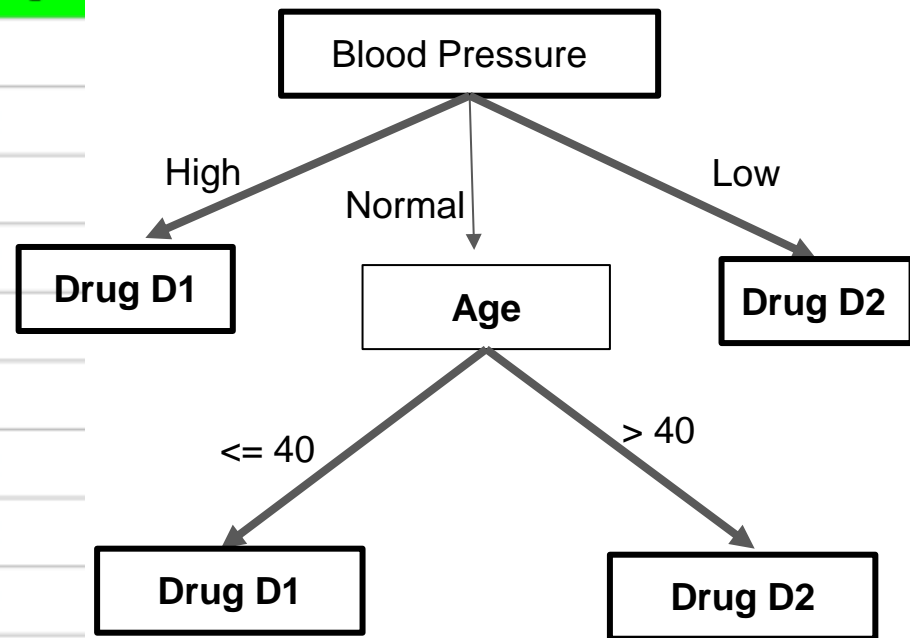


$$h: X \longrightarrow Y$$



| ID No | SEX | AGE | Blood Pressure | Drug |
|-------|-----|-----|----------------|------|
| 1     | M   | 20  | Normal         | D1   |
| 2     | F   | 73  | Normal         | D2   |
| 3     | M   | 37  | High           | D1   |
| 4     | M   | 33  | Low            | D2   |
| 5     | F   | 48  | High           | D1   |
| 6     | M   | 29  | Normal         | D1   |
| 7     | F   | 52  | Normal         | D2   |
| 8     | M   | 42  | Low            | D2   |
| 9     | M   | 61  | Normal         | D2   |
| 10    | F   | 30  | Normal         | D1   |
| 11    | F   | 26  | Low            | D2   |
| 12    | M   | 54  | High           | D1   |

| ID No | SEX | AGE | Blood Pressure | Drug |
|-------|-----|-----|----------------|------|
| 1     | M   | 20  | Normal         | D1   |
| 2     | F   | 73  | Normal         | D2   |
| 3     | M   | 37  | High           | D1   |
| 4     | M   | 33  | Low            | D2   |
| 5     | F   | 48  | High           | D1   |
| 6     | M   | 29  | Normal         | D1   |
| 7     | F   | 52  | Normal         | D2   |
| 8     | M   | 42  | Low            | D2   |
| 9     | M   | 61  | Normal         | D2   |
| 10    | F   | 30  | Normal         | D1   |
| 11    | F   | 26  | Low            | D2   |
| 12    | M   | 54  | High           | D1   |



# Decision Tree

- **ID3** (Iterative Dichotomiser 3)- **1986 Ross Quinlan**
- **C4.5** - is the successor to ID3 and removed the restriction that features must be categorical by dynamically defining a discrete attribute (based on numerical variables) that partitions the continuous attribute value into a discrete set of intervals.
- **C5.0** - is Quinlan's latest version release under a proprietary license. It uses less memory and builds smaller rulesets than C4.5 while being more accurate.
- **CART** - Classification and Regression Trees - s very similar to C4.5, but it differs in that it supports numerical target variables (regression) and does not compute rule sets. CART constructs binary trees using the feature and threshold that yield the largest information gain at each node.

scikit-learn uses an optimized version of the CART algorithm; however, the scikit-learn implementation does not support categorical variables for now.

<https://scikit-learn.org/stable/modules/tree.html#classification>

## ID3- Decision Tree algorithm for classification

# How do we construct such a tree?



| Day | outlook  | temperature | humidity | wind   | Decision |
|-----|----------|-------------|----------|--------|----------|
| 1   | sunny    | hot         | high     | weak   | No       |
| 2   | sunny    | hot         | high     | strong | No       |
| 3   | overcast | hot         | high     | weak   | Yes      |
| 4   | rainfall | mild        | high     | weak   | Yes      |
| 5   | rainfall | cool        | normal   | weak   | Yes      |
| 6   | rainfall | cool        | normal   | strong | No       |
| 7   | overcast | cool        | normal   | wtrong | Yes      |
| 8   | sunny    | mild        | high     | weak   | No       |
| 9   | sunny    | cool        | normal   | weak   | Yes      |
| 10  | rainfall | mild        | normal   | weak   | Yes      |
| 11  | sunny    | mild        | normal   | strong | Yes      |
| 12  | overcast | mild        | high     | strong | Yes      |
| 13  | overcast | hot         | normal   | weak   | Yes      |
| 14  | rainfall | mild        | high     | strong | No       |



