

LIFE IS LIKE MACHINE  
LEARNING

YOU NEVER KNOW WHAT YOU'RE  
GONNA GET

e-mega-generator.net

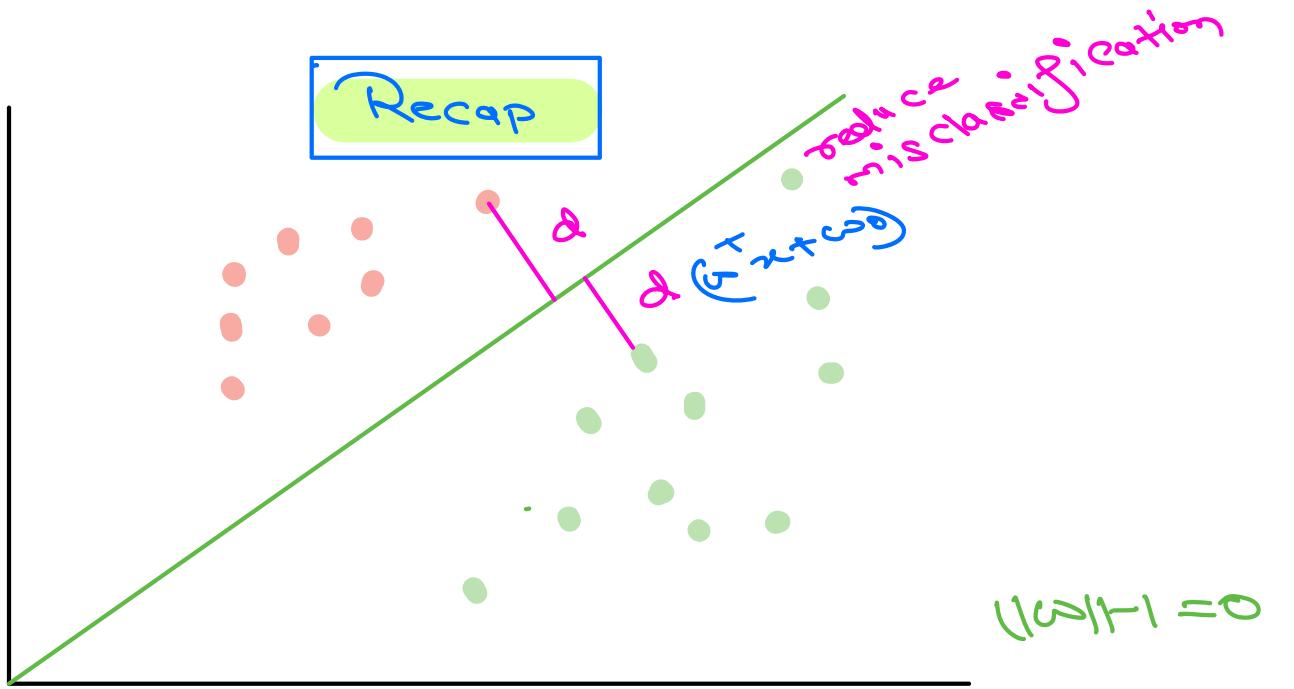
YOU SEE THIS? THIS  
IS TOO MANY DIMENSIONS

WATCH THIS HE'S GONNA  
SAY "USE PCA" THIS IS  
STANDARD MACHINE LEARNING STUFF

USE PCA

## Agenda

- ⑤ Recap
- ⑥ Variants of Gradient Descent
  - ⑦ Vanilla GD
  - ⑦ SGD
  - ⑦ Batch GD
- ⑥ Dimensionality Reduction
- ⑥ PCA



$$L(D, \bar{w}, w_0) = -\frac{\bar{w}^T \bar{x}_i + w_0}{||\bar{w}||} y_i \quad \text{S.t. } ||w|| = 1$$

+ \lambda (||w|| - 1)

$\lambda (\epsilon_q = 0)$

Gradient Descent

①  $\bar{w}^{(t+1)} = \bar{w}^{(t)} - \eta \nabla_{\bar{w}} L(D, \bar{w}, w_0)$  regularization

$\downarrow \bar{w}^T \bar{x}_1 + w_0 = 1.4$

$y_1$	$x_1 \dots x_d$	$\bar{w}_1$
$y_2$	$x_2$	$\bar{w}_2$
$\vdots$		
$y_n$	$x_n \quad x_n^d$	$\bar{w}_n$

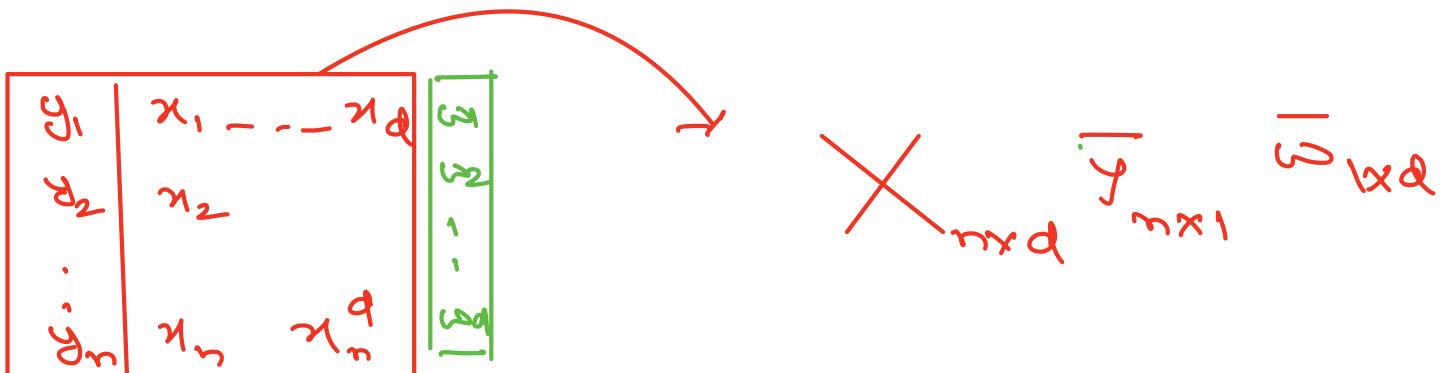
$\bar{w}^{(t+1)} = \bar{w}^{(t)} - \eta \left( -\sum_{i=1}^n y_i \bar{x}_i + \frac{\lambda \bar{w}}{||\bar{w}||} \right)$

②  $w_0^{(t+1)} = w_0^{(t)} - \eta \left( -\sum_{i=1}^n y_i \right) \frac{\partial L}{\partial w_0}$

$t+1 \rightarrow \text{Timestep } \tau / \text{iteration}$

## Variants of Batch Gradient Descent

**Epoch** : A full scan/pass through dataset



① Vanilla GD : Full Batch GD

10 million DP

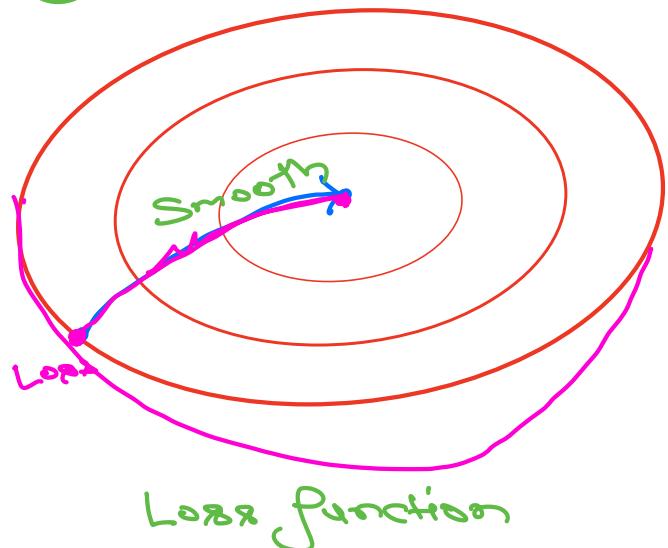
$$\bar{\omega}^{(t+1)} = \bar{\omega}^{(t)} - \eta \left( -\sum_{i=1}^n y_i \bar{x}_i + \eta \frac{\bar{\omega}}{\|\bar{\omega}\|} \right)$$

$$\bar{\omega}^{(t+1)} = \bar{\omega}^{(t)} - \eta \left( Y \cdot \bar{x} + \eta \frac{\bar{\omega}}{\|\bar{\omega}\|} \right)$$

$$\bar{x} = \bar{X} \cdot \bar{\omega}$$

	$x_1$	$x_2$	$x_3$	.....	$x_d$	$y$
1						
2						
3						

$d$ -dimension



\* Full Dataset is processed

Fastest but Requires High Computation

## ② SGD : Stochastic Gradient Descent

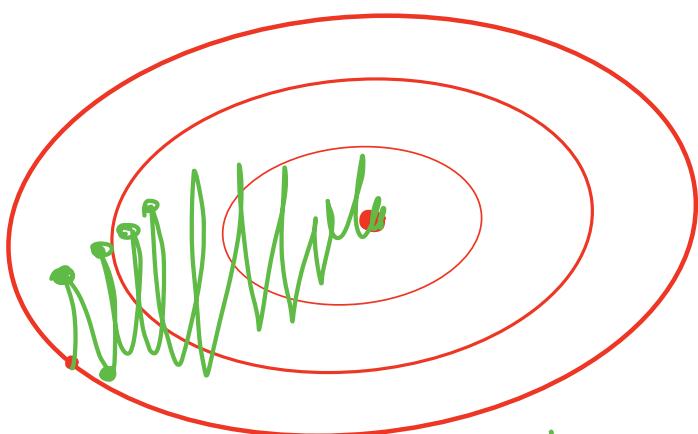
Pick  $i$  index Randomly

$$\bar{w}^{(t+1)} = \bar{w}^{(t)} - \eta \left( - \sum_{i=1}^n y_i \bar{x}_i + \gamma \frac{\bar{w}}{\|\bar{w}\|} \right)$$

Stochastic = Randomness

$$\bar{w}^{(t+1)} = \bar{w}^{(t)} - \eta \left( - y_i \bar{x}_i + \gamma \frac{\bar{w}}{\|\bar{w}\|} \right)$$

where  $i$  is picked randomly from  $(1, n)$



Loss Function



diss very  
high

## ③ BGD : Batch Gradient Descent

Mini Batch

Batch-size → Based on  
Hyperparameters Compute / Budget

Stochastic →  $1 < \kappa < N$  ← Vanilla

$$\bar{\omega}^{(t+1)} = \bar{\omega}^{(t)} - \eta \left( -\sum_{i=1}^{\kappa} y_i \bar{x}_i + \frac{\bar{\omega}}{\|\bar{\omega}\|} \right)$$

$10 <$

$\kappa = 10$  samples

How many iteration will be required for 1 Epoch?

$$10 // 5 \Rightarrow 2$$

Batch-Size 5

for i in range(10)  
rand\_idx = [---]

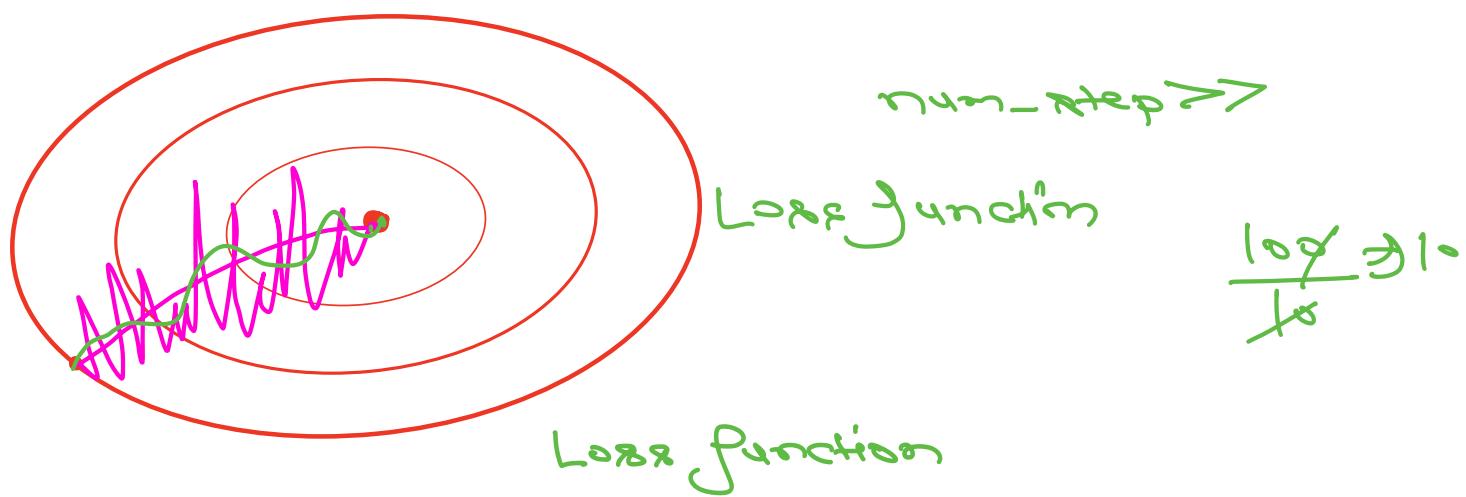
$$10 \times 5 \Rightarrow 50$$

$$5 \times 2 = 10$$

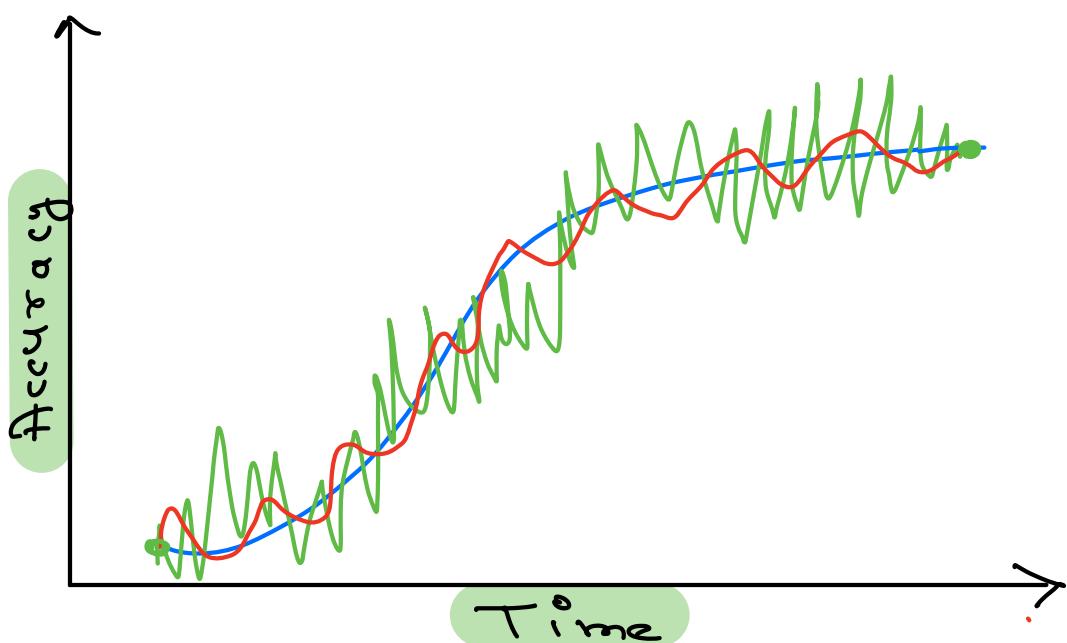
$$\omega = \omega - \eta \nabla L$$

$$* \text{Epoch} = \frac{\text{Total Datasize}}{\text{Batch-size}} \Leftrightarrow \frac{100}{10} = 10$$

10 epochs ⚡  $10 * 10 = 100 \text{ steps}$



## Questions

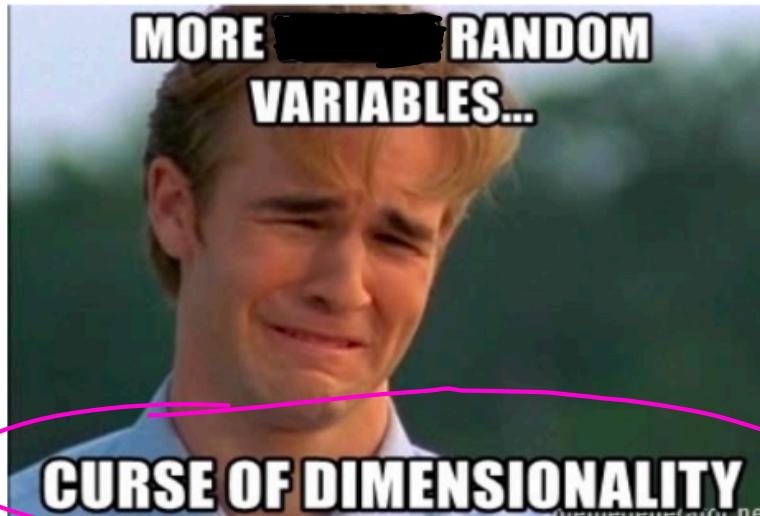


Blue ⚡ Vanilla

Green ⚡ Stochastic

Red ⚡ Mini-Batch

# Dimensionality Reduction



n datapoints

	$x_1$	$x_2$	$x_3$	-----	$x_d$	$y$
1						
2						
3						
$n$						

d-dimension

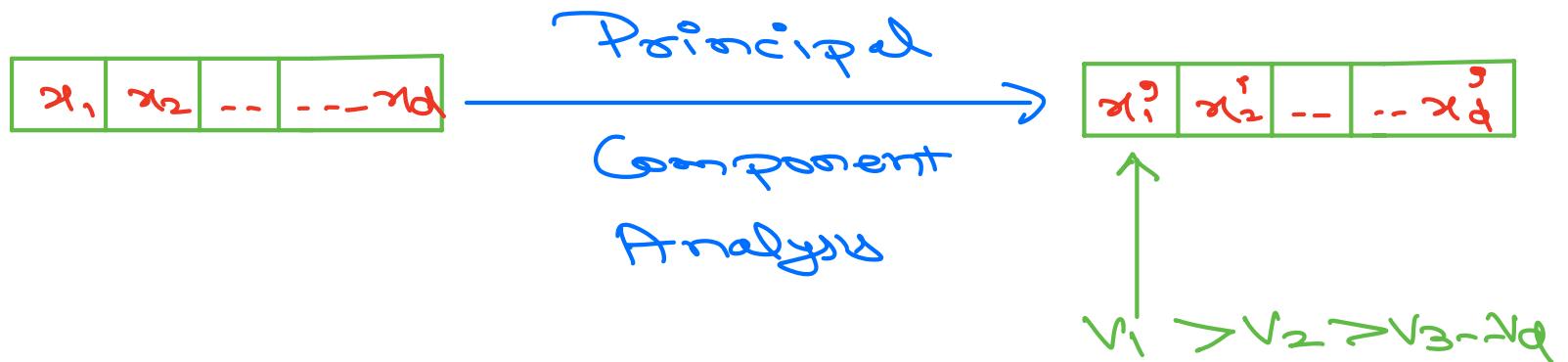
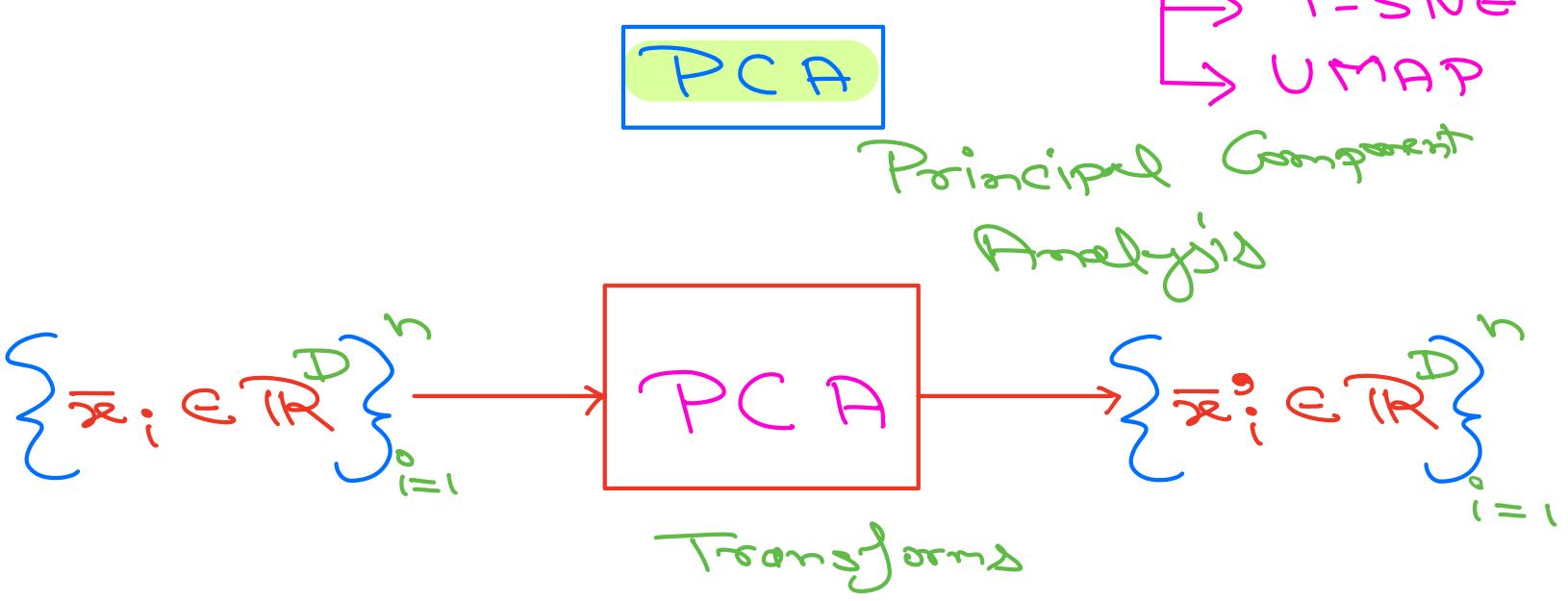
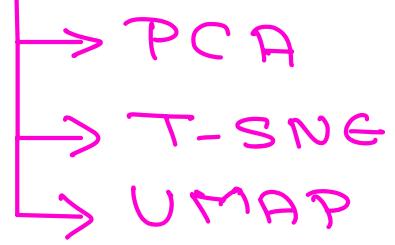
$d$  is high  
↓  
Curse of Dimensionality

$$D = \left\{ (\bar{x}_i, y_i) : \bar{x}_i \in \mathbb{R}^d \right\}_{i=1}^n$$

## \* Issues with Higher Dimensions:

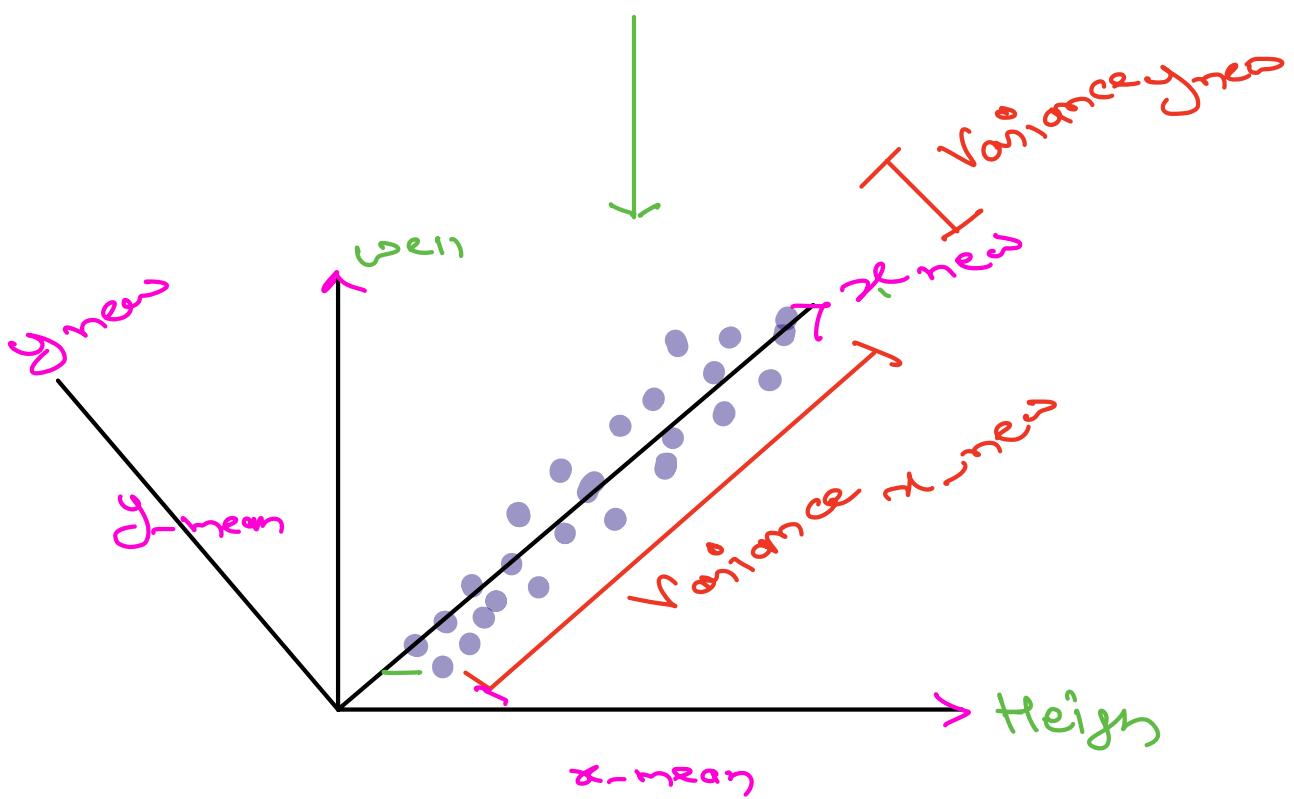
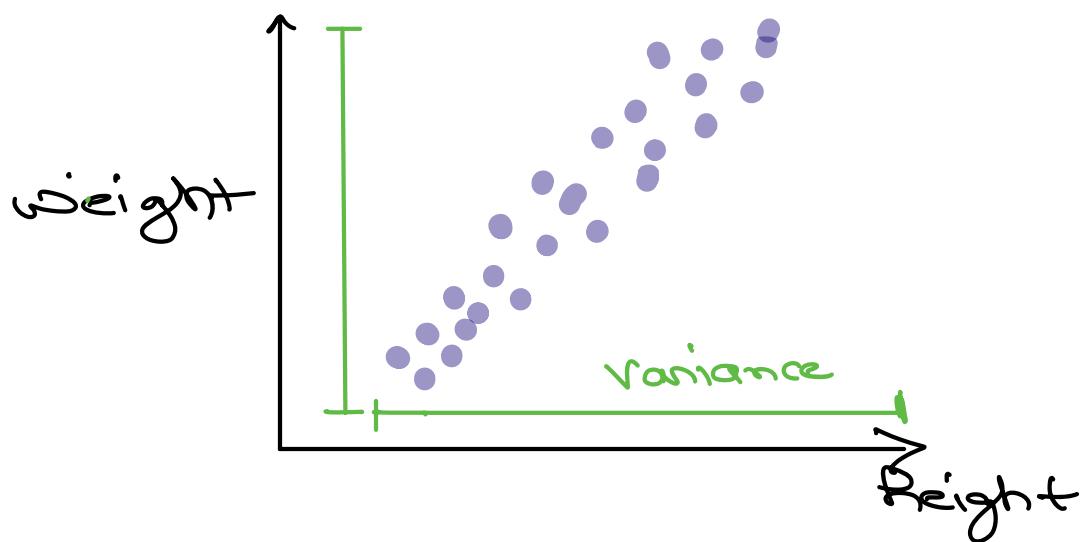
- ① Visualization is difficult (3d-vision EDA)
- ② Training and processing Time
- ③ Computational Cost
- ④ Difficult Maths

# Dimensionality Reduction Algo



\* Variance  $\Rightarrow$  The new dimensions capture variance in reverse sorted order

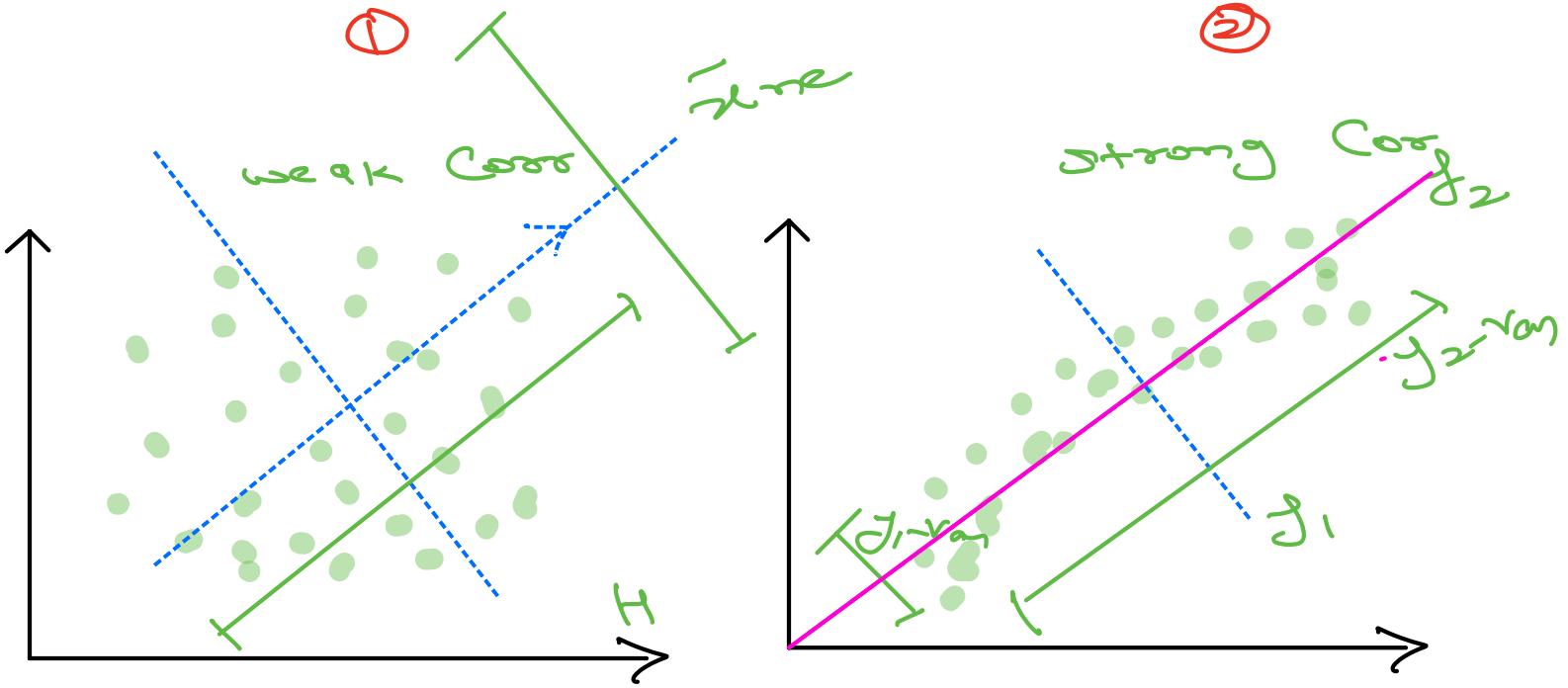
W	H	BPTI



$x$        $y$

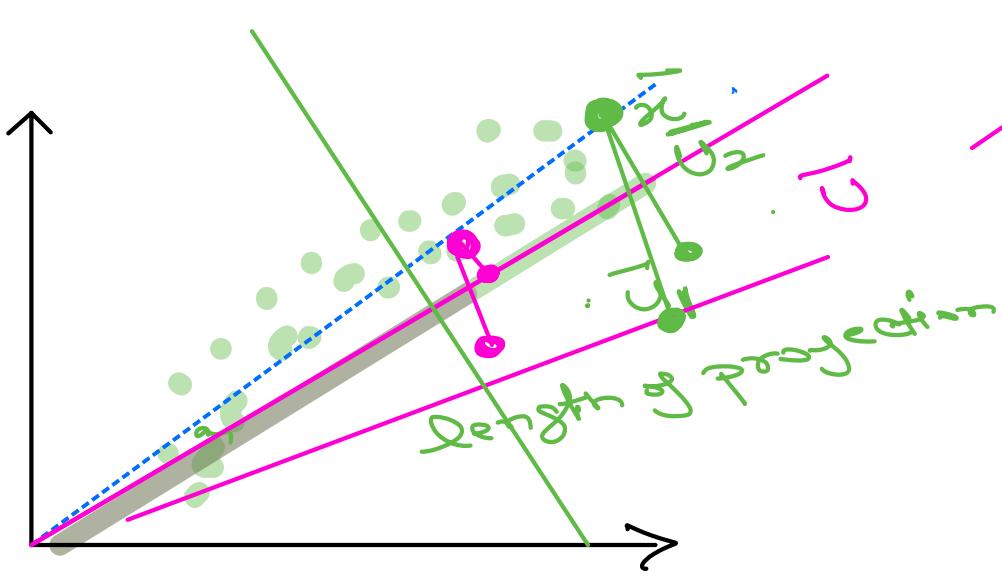
$x_{\text{new}}$        $y_{\text{new}}$   
 $\text{Var}_x > \text{Var}_y$

If we eliminate  $y_{\text{new}}$  we won't  
lose a lot of information



What is the Goal of PCA

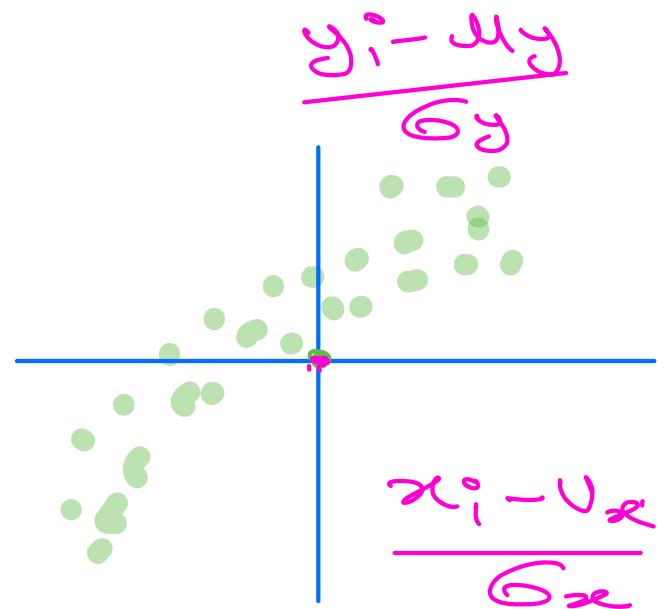
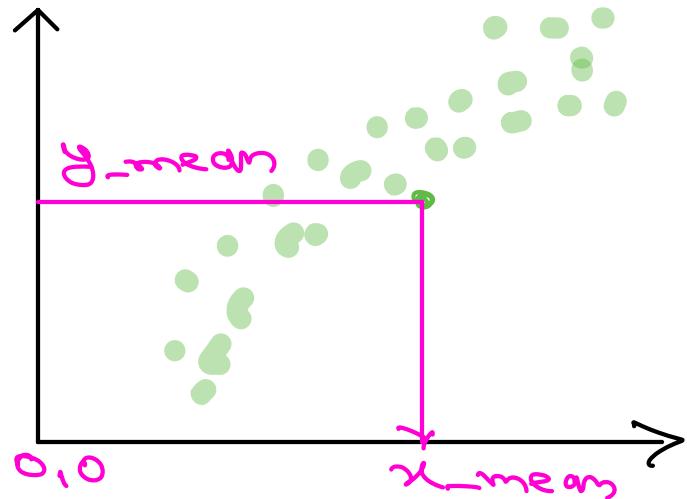
- ⑤ Find new axis such that we can eliminate some features



$$\frac{\|x \cdot j_1\|}{\|j_1\|}$$

# Visualizing PCA

Strong Cor



Step 1: Shift the Origin to the mean of features

- \* Normalization  $\rightarrow -1 \quad -1 \quad (0,0)$
- \* Standardization  $\rightarrow (\bar{x}, \bar{y})$

How can we do this Mathematically?

Objective  
maximize

$$\sum_{i=1}^n \left( \frac{\bar{x} \cdot \bar{v}}{\|\bar{v}\|} - u_i \right)^2$$

Step 2:

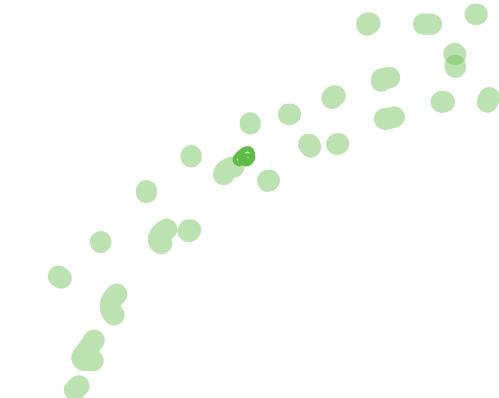
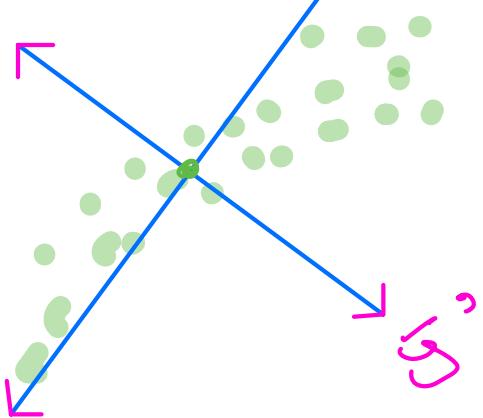
Rotate axis  
with objective

maximizing projecting

$$\max \left( \frac{\mathbf{x}^T \mathbf{u}}{\|\mathbf{u}\|} \right)^2$$

$\mathbf{u}^\perp$

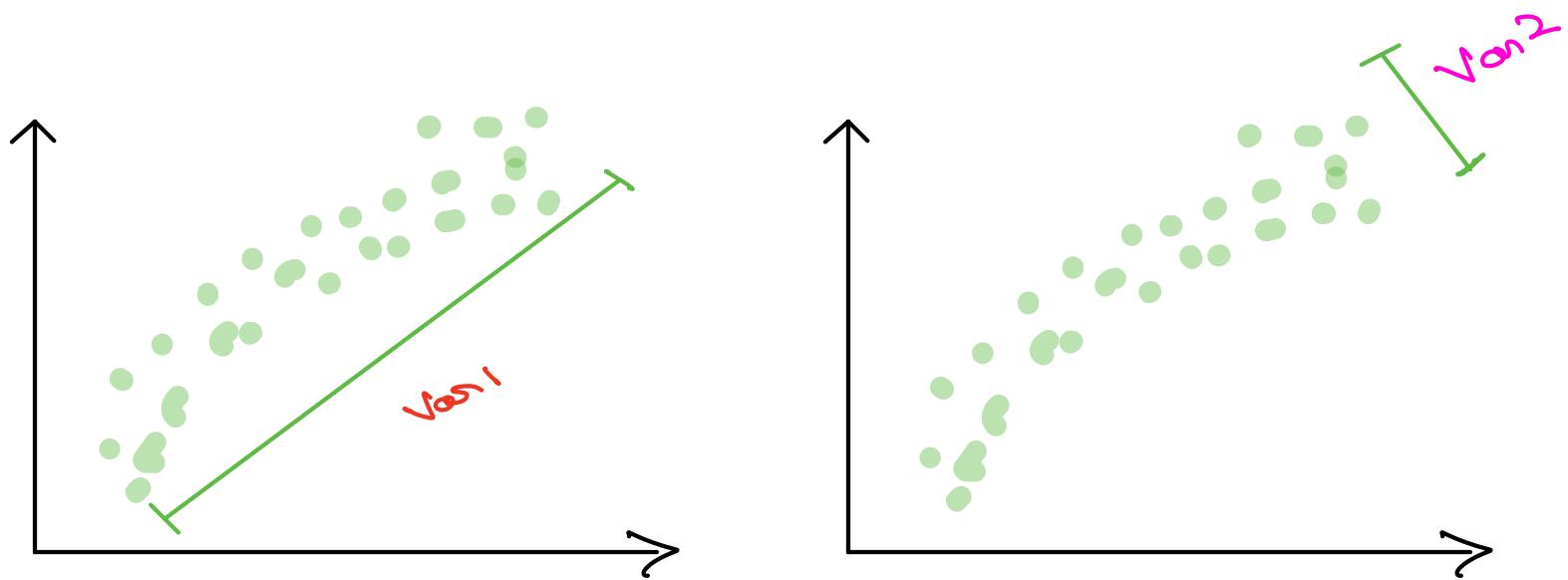
$$(\mathbf{x}^T \mathbf{u})^2 + \lambda (\|\mathbf{u}\|^2 - 1)$$



\* Note: PCA Returns new Orthogonal Features

feature at 90°  
with  
Each Other

## Implementing PCA



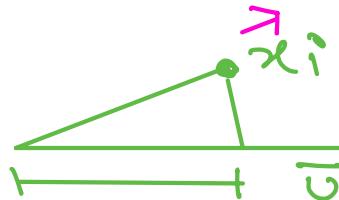
- \* Higher the Variance Captured on new axis more important the feature

Step 1 : Standardization

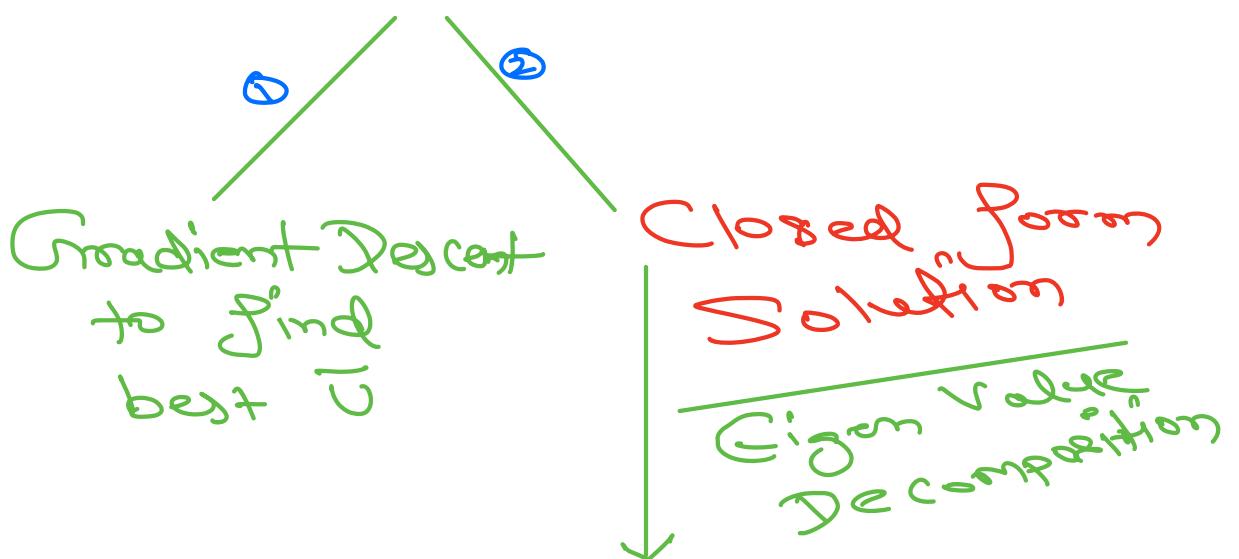
$$\frac{x_i - \bar{x}_i}{\sigma_x}$$

Step 2 : Maximize the Projections

$$U^* = \underset{U}{\operatorname{argmax}} \sum_{i=1}^n \frac{\bar{x}_i^T U}{\|U\|}$$



$$\bar{U}^* = \underset{\bar{U}}{\operatorname{argmax}} \frac{1}{n} \sum_{i=1}^n \bar{x}_i^T \bar{U} + \gamma (\|U\|_F - 1)$$



$$(X^T X) \bar{U} = ?$$

Covariance Matrix

$X \rightarrow$  Dataset

$\beta_1$	$\beta_2$	$\dots$	$\beta_d$
$x_1$			
$x_2$			
$x_n$			

$\lambda \rightarrow$  Eigen Value

$X \cdot X$

$\bar{U} \rightarrow$  Eigen Vector

- ④ 3 features Input  $\bar{x}_1, \bar{x}_2, \bar{x}_3$
- ⑤ 3 Eigen value  $\gamma_1, \gamma_2, \gamma_3$
- ⑥ 3 Eigen vector  $U_1, U_2, U_3$

$\lambda_{10} X$

$$\sum_{i=1}^n M_i^j = g$$

$$\sum_{i=1}^n M_i^{10} = g$$

④ % of

Information

DDY.

①

How much variation  $C_V$  store

$$\frac{\sigma}{\lambda_1 + \lambda_2 + \lambda_3}$$

$\lambda_1$