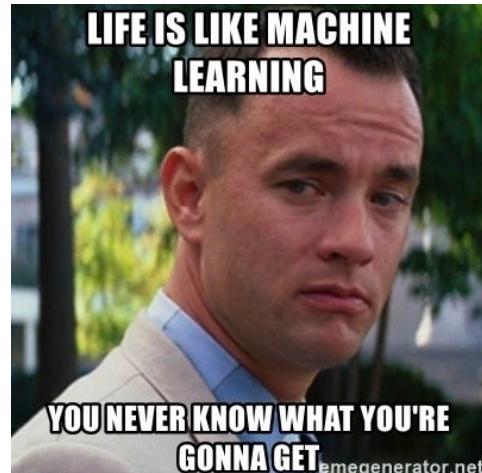
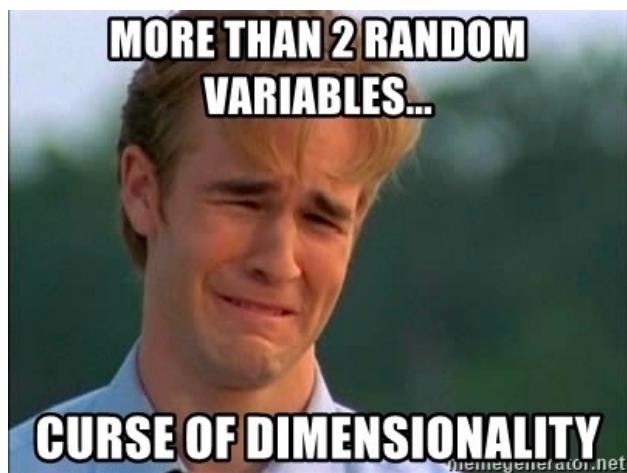
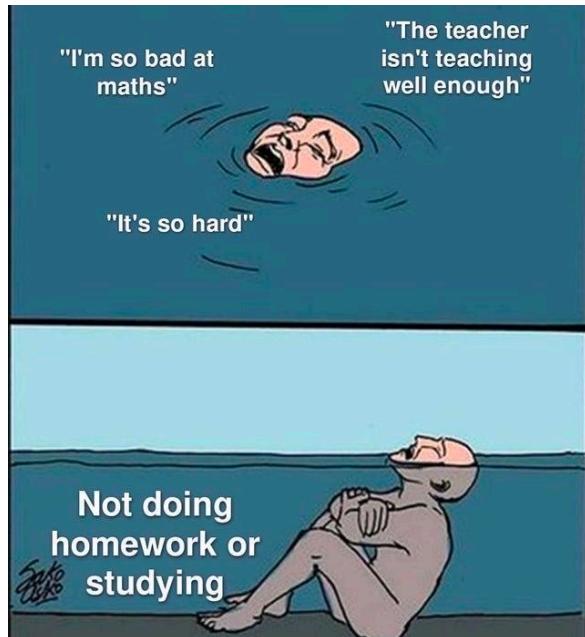


March 11, 2023

DSML : Math for ML.

Optimization 5: Principal Component Analysis.



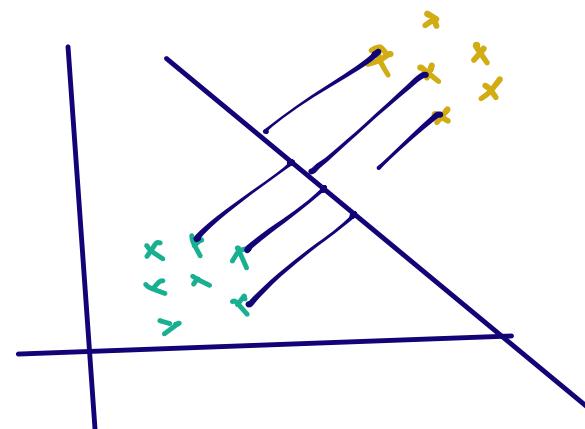
Recap:

- (a) Classification - choosing \bar{w} and w_0 .
- (b) Brute force: very inefficient.
- (c) Alternative: gradient descent
- (d) Functions, limits, derivatives.
- (e) Partial derivatives and gradients.
- (f) Gradient Descent in action.

Today:

- (a) Variants of gradient descent.
- (b) Dimensionality Reduction - PCA.

Recap:



$$\lambda(\omega; \bar{w}, w_0) = -\underbrace{\left(\sum_{i=1}^n (\bar{w}^\top \bar{x}_i + w_0) y_i \right)}_{\text{margin}} + \lambda \underbrace{(\|\bar{w}\| - 1)}_{\text{Regularizer}}$$

* Regularizer.

\downarrow
 λ controls
the importance

$$\bar{w}^{(t+1)} = \bar{w}^{(t)} - \eta \nabla_{\bar{w}} \lambda(\omega; \bar{w}, w_0)$$

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

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

Variants of Gradient Descent

Vanilla G.D.:

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

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

$n \rightarrow$ size of the dataset. $\mathcal{O}(n)$.

$D = \{(\bar{x}_i, y_i)\}_{i=1}^n$ \downarrow affects training time greatly.

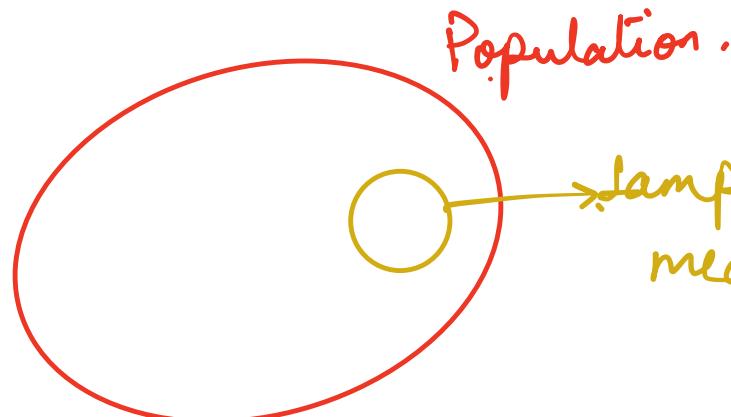
"10⁷ → 10 million images".
"Stochastic Gradient Descent"

"Stochastic Gradient Descent".

- * Pick one (\bar{x}_i, \bar{y}_i) randomly from our dataset.
Use that to do the following update:

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

$$w_0^{(t+1)} = w_0^{(t)} - \eta \cdot (-\bar{y}_i) \quad O(1).$$



Sample, calculate "approximate" mean & confidence intervals.

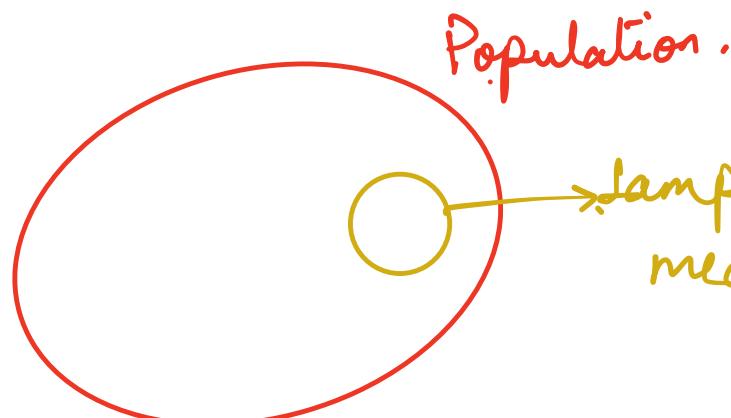
"Batch - G. D." ↵

- * Pick any k (\bar{x}_i, y_i) randomly from our dataset.
Use that to do the following update:

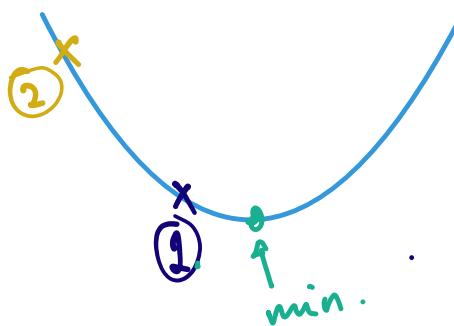
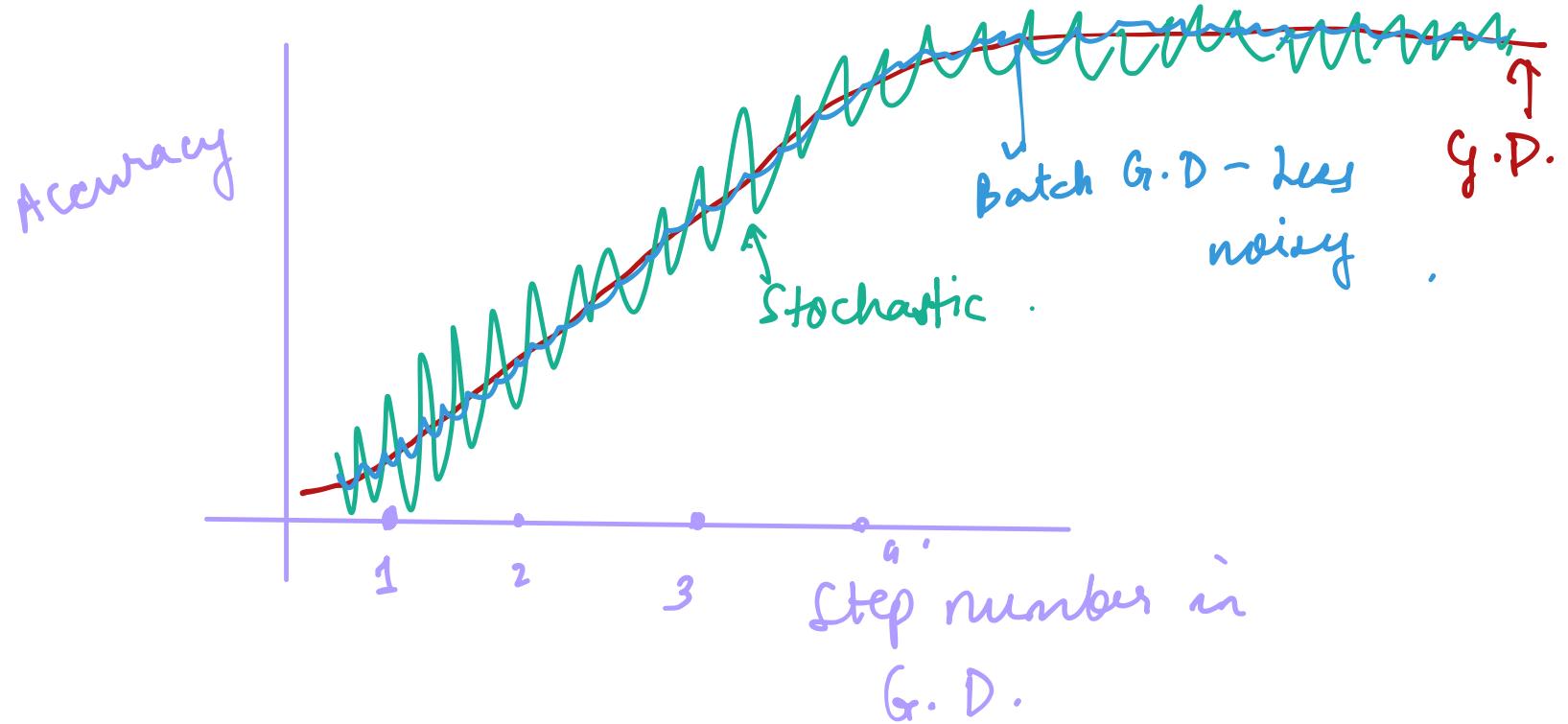
$$\bar{w}^{(t+1)} = \bar{w}^{(t)} - \eta \left(-\sum_{i=1}^k y_i \bar{x}_i + \lambda \cdot \frac{\bar{w}}{\|\bar{w}\|} \right)$$

$$w_0^{(t+1)} = w_0^{(t)} - \eta \cdot \left(-\sum_{i=1}^k y_i \right)$$

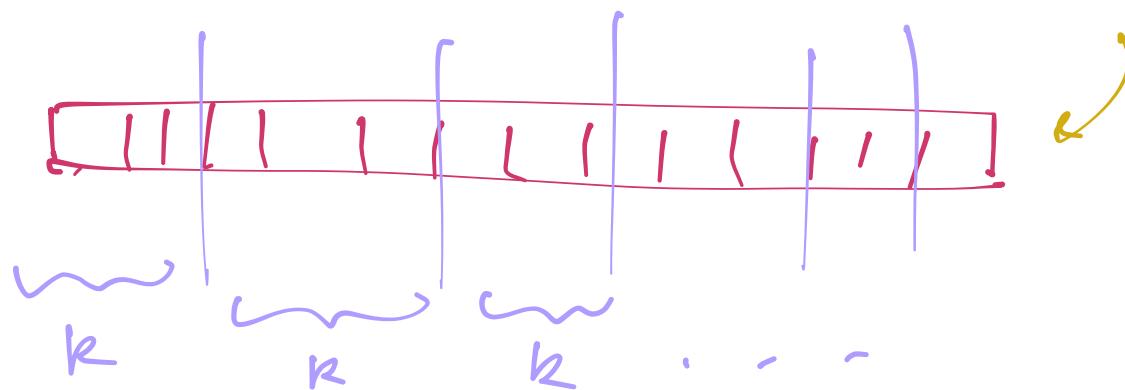
$R \approx 32 \rightarrow 256$.



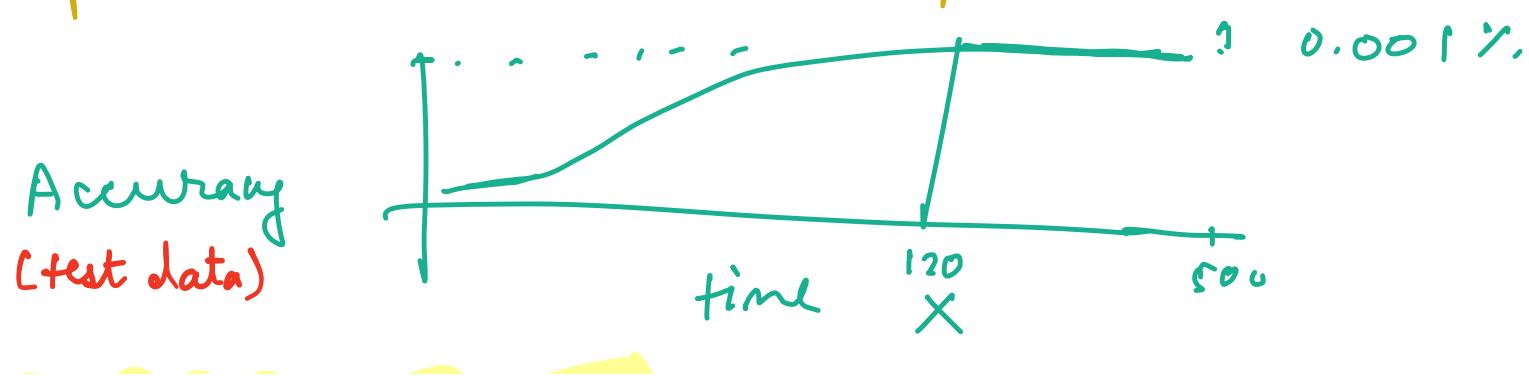
Sample, calculate "approximate" mean & confidence intervals.



- * Random permutation of the data is used to implement Batch G.D.



- * Using the whole dataset once to update our model \rightarrow 1 "epoch".



Dimensionality Reduction.

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

$d \gg 2 \rightarrow$ Very difficult to visualize !!

Examples:

- ① Flattening.
- ② Feature engineering.
- ③ t-SNE.
- ④ Autoencoders.

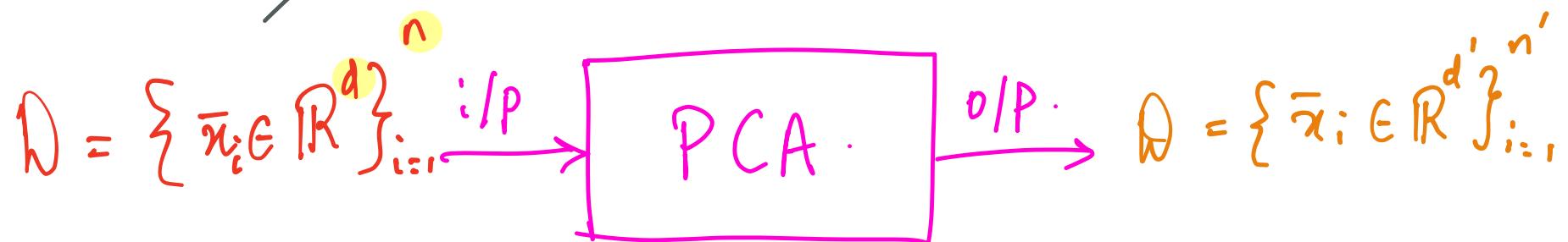
Problems of high dimensional data:

- ① Visualization is tough.
- ② Training time increases.
- ③ Computational costs (time/memory).
- ④ Difficult to work with (Math is harder).

Dimensionality Reduction algorithms help us deal with these issues - Eg: Principal Component Analysis (PCA).

Input / Output Relationship -

unlabelled data.



$\geq, \leq, =$

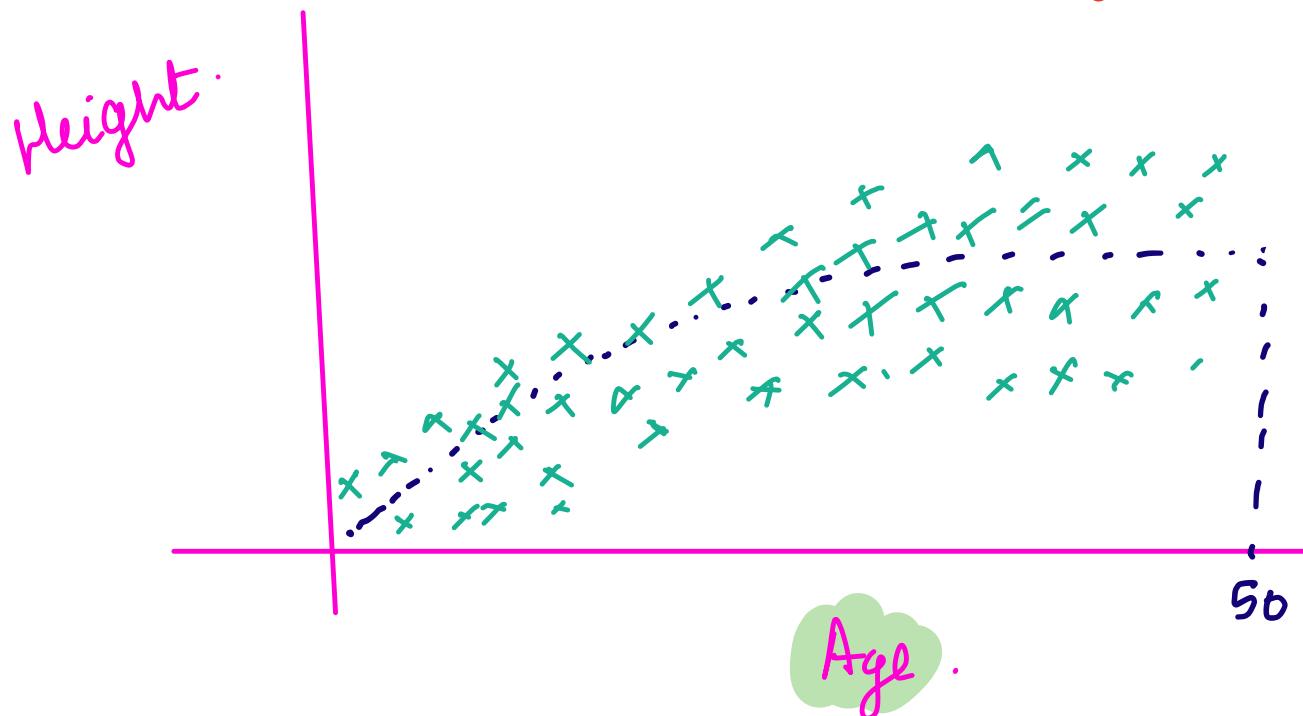
$$d \geq d'$$
$$n = n'$$

This is important
because we don't
want to lose essential
information.

Example 2:

Height vs. Age

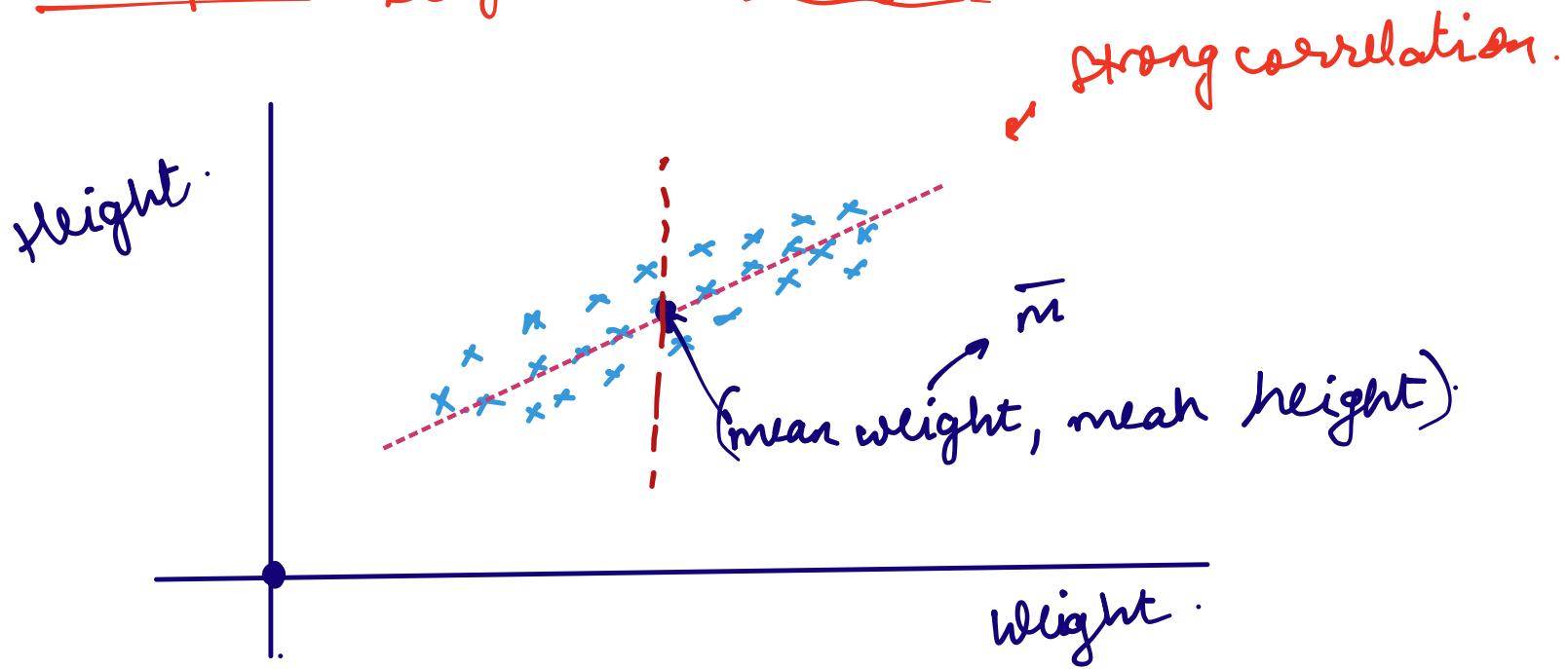
Strong correlation \Rightarrow magnitude of C.C. > 0.5 .



There isn't a strong correlation
between age & height.

\rightarrow In such cases, it is not possible to do D.R without losing information.

Example 2 : Height vs. Weight

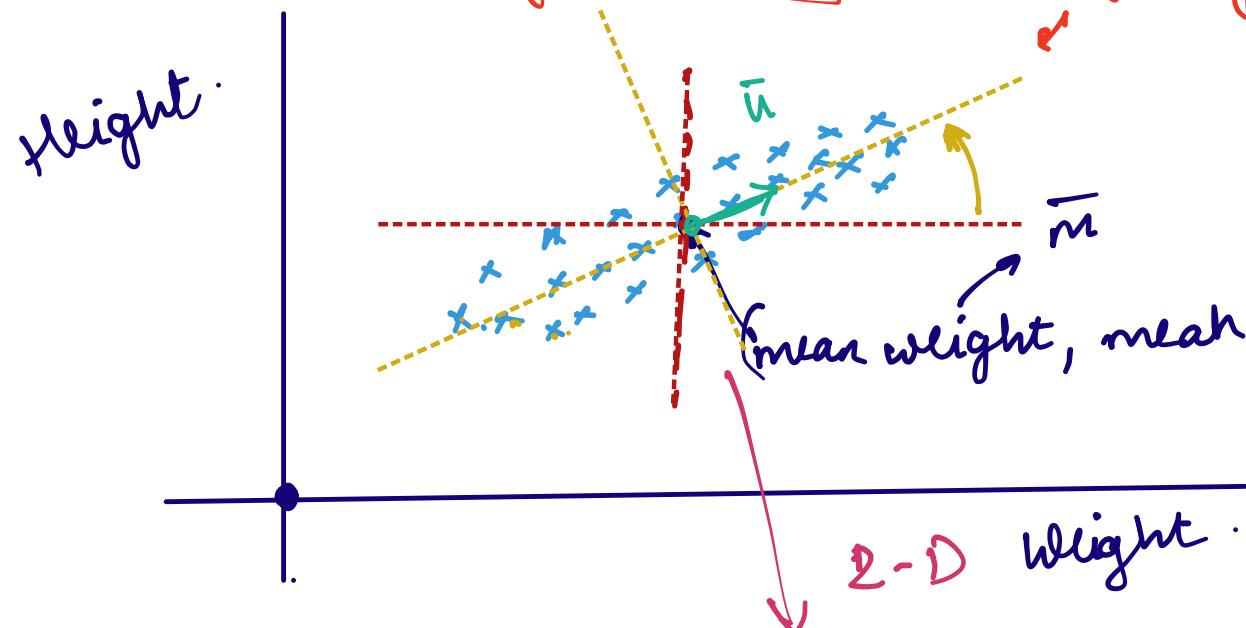


→ In such cases, Dimensionality reduction is possible.

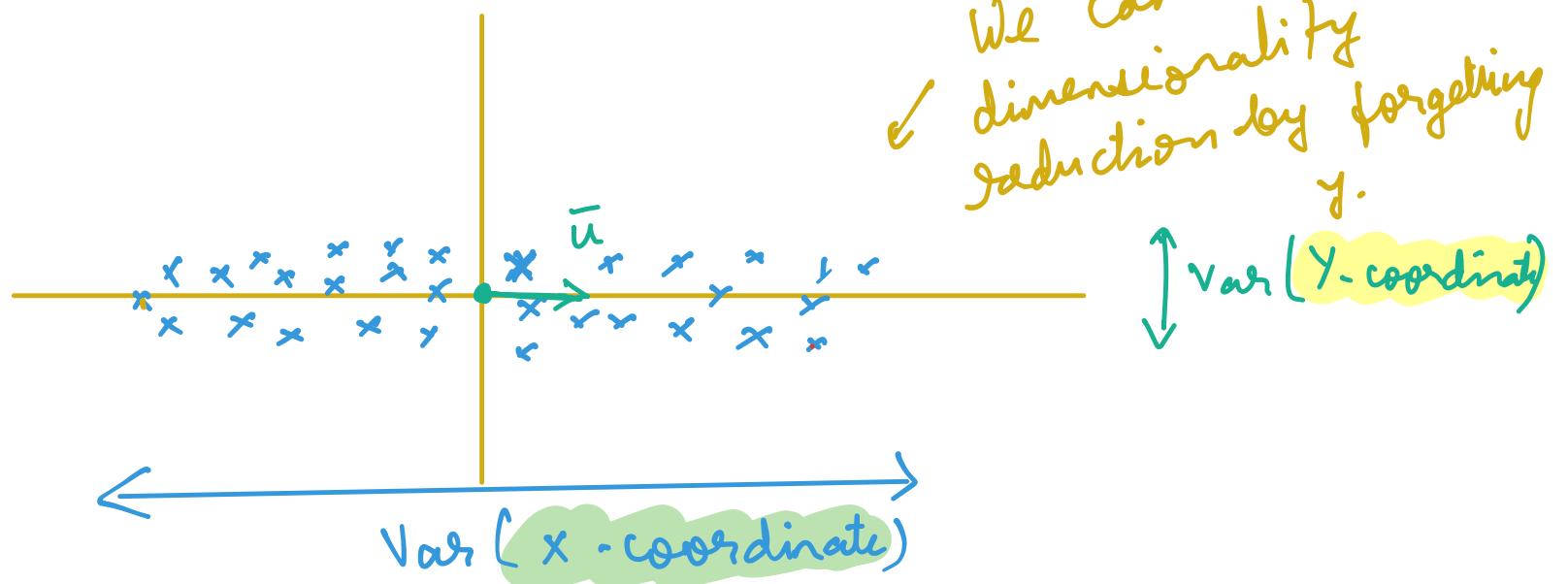
$$D = \left\{ \bar{x}_i \in \mathbb{R}^d \right\}_{i=1}^n$$

$$D' = \left\{ (\bar{x}_i - \bar{m}) \right\}$$

Visual understanding of PCA



Strong correlation.



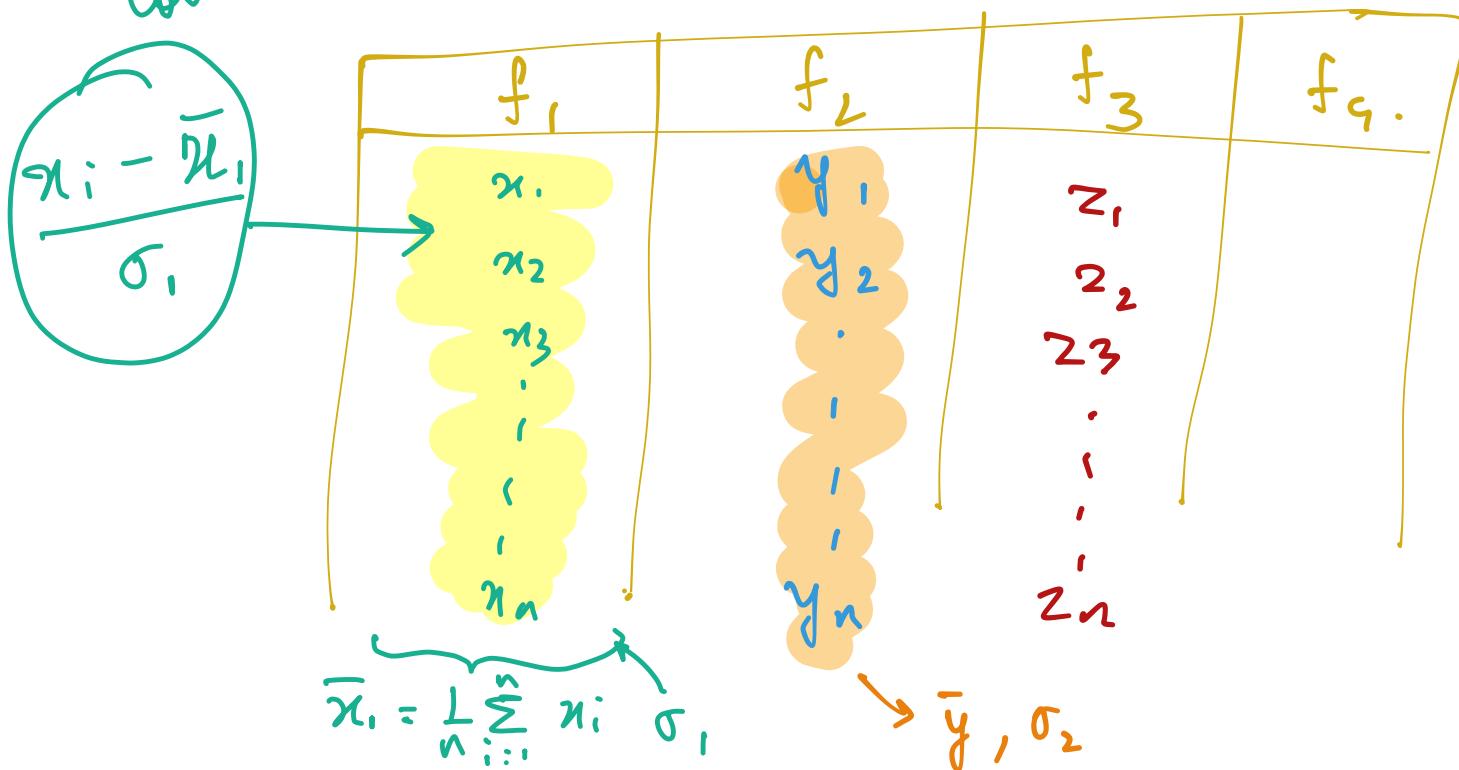
We can achieve dimensionality reduction by forgetting y.

Implementing PCA

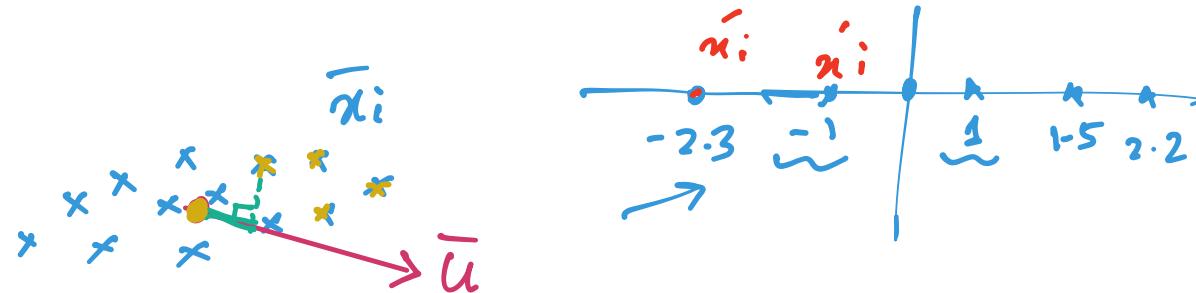
Step 1 : Column standardization .

$$\underline{\bar{x}_i - \bar{\mu}}$$

σ .
Column standardization - $\frac{y_i - \bar{y}}{\sigma_2}$



Step 2: Find a vector \bar{u} such that the variance of \bar{x} along \bar{u} is maximized.



Projection of \bar{x}_i' on \bar{u} :
$$\frac{\bar{x}_i^T \bar{u}}{\|\bar{u}\|}$$

$$\text{Var}(\bar{x}_i') = \frac{1}{n} \sum_{i=1}^n \left(\frac{\bar{x}_i^T \bar{u}}{\|\bar{u}\|} \right)^2$$

$$\left[\begin{array}{c} x'_1 \\ -2.2, \quad x'_2 \\ 1.3, \quad x'_3 \\ 1.7, \quad x'_4 \\ 0.86, \quad x'_5 \\ -3.4 \end{array} \right]$$

↓
add. → not going to help.

var? $\frac{1}{n} \sum_{i=1}^n (x_i' - \bar{\mu})^2$

$$\Rightarrow \bar{u}^* = \underset{\bar{u}}{\operatorname{Argmax}} \frac{1}{n} \sum_{i=1}^n (\bar{x}_i^\top \bar{u})^2$$

$$\text{s.t. } \|\bar{u}\|^2 = 1$$

↓ Lagrange multipliers

$$\bar{u}^* = \underset{\bar{u}}{\operatorname{Argmax}} \frac{1}{n} \sum_{i=1}^n (\bar{x}_i^\top \bar{u})^2 + \lambda (\|\bar{u}\|^2 - 1)$$

↳ We can stop here and apply Gradient Ascent to find \bar{u} .

But, finding a math solution is more efficient.

$$\bar{u}^* = \underset{\bar{u}}{\operatorname{argmax}} \sum_{i=1}^n (\bar{x}_i^\top \bar{u})^2 + \lambda (\|\bar{u}\|^2 - 1)$$

-charge - !
 how we are
 seeing this -

$$X = [\bar{x}_1^\top \quad \bar{x}_2^\top \quad \dots \quad \bar{x}_n^\top]^{n \times d}$$

$$\mathcal{W} = \{ \bar{x}_i \in \mathbb{R}^d \}_{i=1}^n$$

$$X \bar{u} = \bar{v}$$

$$\begin{bmatrix} \bar{x}_1^\top \bar{u} \\ \bar{x}_2^\top \bar{u} \\ \bar{x}_3^\top \bar{u} \\ \vdots \\ \bar{x}_n^\top \bar{u} \end{bmatrix} \rightarrow \|\bar{v}\|^2$$

This matrix is called the data matrix.
or X .

$$\bar{u}^* = \underset{\bar{u}}{\operatorname{argmax}} \underbrace{\frac{1}{n} \cdot \bar{u}^\top X^\top X \bar{u}}_{f(\bar{u})} + \lambda (\|\bar{u}\|^2 - 1).$$

$$\bar{v} = X \bar{u}$$

$$\begin{aligned} \bar{v}^\top \bar{v} &= (X \bar{u})^\top (X \bar{u}) \\ &= \underbrace{\bar{u}^\top X^\top X \bar{u}} \end{aligned}$$

Why did we do this? to make taking the gradient easier !!

→ Fact 1: $f(\bar{u}) = \bar{u}^\top S \bar{u}$ $S \rightarrow \text{any matrix}$

then $\nabla_{\bar{u}} f(\bar{u}) = (S + S^\top) \bar{u}$

$$f(\bar{u}) = \frac{1}{n} \bar{u}^\top \underbrace{x^\top x \bar{u}}_{\text{Fact 1: } \downarrow} + \lambda \underbrace{\bar{u}^\top \bar{u} - 1}_{-1}$$

$$\nabla_{\bar{u}} f(\bar{u}) = \frac{1}{n} \left((x^\top x) + (x^\top x)^\top \right) \bar{u} + \lambda (2 \bar{u})$$

$$\rightarrow \text{Since } (x^\top x)^\top = x^\top x. + \lambda (2 \bar{u})$$

$$= \frac{2}{n} x^\top x \bar{u} + \lambda (2 \bar{u})$$

$$\nabla_{\bar{u}} \bar{u}^\top \bar{u} = 2 \bar{u}$$

$$\frac{\partial x^2}{\partial x} = 2x$$

$$\nabla_{\bar{u}} f(\bar{u}) = 2 \left(\frac{1}{n} X^T X \bar{u} + \lambda \bar{u} \right) = 0..$$

$$\Rightarrow \underbrace{\frac{1}{n} X^T X}_{\text{I}} \bar{u} = -\lambda \bar{u}$$

$$\Rightarrow X^T X \bar{u} = \underbrace{-(n\lambda)}_{\text{II}} \bar{u}$$

$$\boxed{(X^T X) \bar{u} = \lambda' \bar{u}} \quad \textcircled{1}$$

Eigenvalue, eigenvector expression!!

The vector \bar{u} which maximizes $f(\bar{u})$ is the one which satisfies ①.

By observing ①, we notice that
 \bar{u} is nothing but the eigenvector
of the matrix $\underline{\underline{X^T X}}$.

What is an eigenvector?

"eigen" \rightarrow German word
 \approx identity.

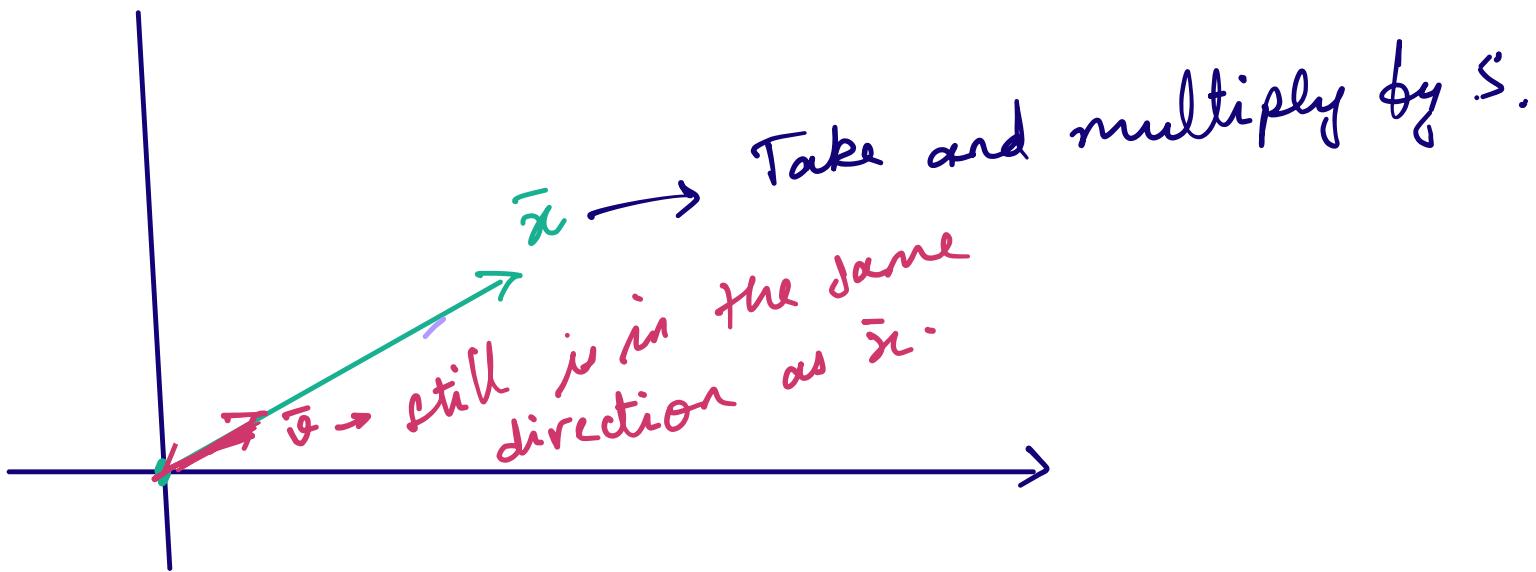
For any square matrix S ,
 \downarrow no. of rows = no. of cols.

same direction

if we find an \bar{x} such that

$$S \bar{x} = \underline{\lambda} \bar{x}$$

then \bar{x} is said to be an "eigenvector".



$$\underbrace{S_{d \times d}}_{\text{a } d \times d \text{ matrix}} \quad \underbrace{\bar{x}_{d \times 1}}_{\text{a } d \times 1 \text{ column vector}} = \underbrace{\bar{v}_{d \times 1}}_{\text{a } d \times 1 \text{ column vector}}$$

Result (of step 2) : The vector \bar{u} which captures max. variance is the eigenvector of the matrix $X^T X$.

How to find this? \rightarrow numpy

Each matrix has many eigenvectors,
each eigenvector has a corresponding
eigenvalue.

So, we pick the eigenvector corresponding
to the maximum eigenvalue for our
 \bar{u} .

Step 3: Other directions?

If we want to reduce

$$d \rightarrow d'$$

we simply pick the first d' eigenvectors with the largest corresponding eigenvalues.

$$3 \times 3 \rightarrow (\bar{x}_1, \overset{\checkmark}{\lambda}_1) \downarrow 0.5, (\bar{x}_2, \overset{\checkmark}{\lambda}_2) \downarrow 0.3, (\bar{x}_3, \overset{\checkmark}{\lambda}_3) \downarrow 0.2.$$

²

Pick (\bar{x}_1, \bar{x}_2) , $(0.5, 0.3)$,
these are the new ones.

Step 4: How to convert $\bar{x}_i \in \mathbb{R}^d$ to $\bar{x}_i' \in \mathbb{R}^{d'}$:

$$\frac{\bar{x}_i^\top \bar{u}_1}{\|\bar{u}_1\|} + \frac{\bar{x}_i^\top \bar{u}_2}{\|\bar{u}_2\|} \dots + \frac{\bar{x}_i^\top \bar{u}_{d'}}{\|\bar{u}_{d'}\|}.$$

↓
coordinate
along \bar{u}_1

↓
coordinate
along \bar{u}_2

⋮

Step 5: How to calculate the information loss?

$$\frac{(\lambda_1 + \lambda_2 + \dots + \lambda_{d'})}{\sum_{i=1}^d \lambda_i} \rightarrow \text{tells us the fraction of info retained.}$$