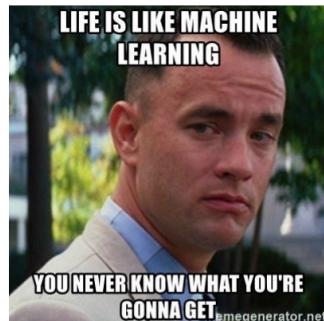
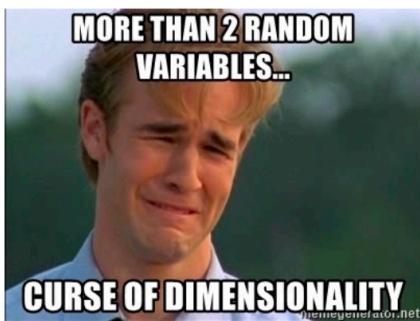
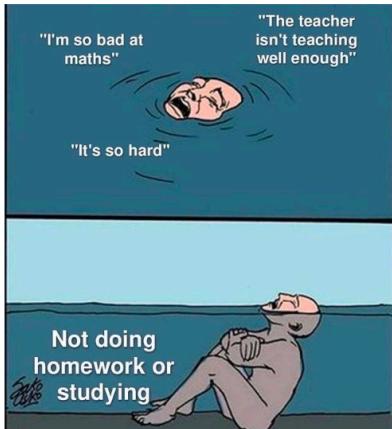


## Content

- Recap
- Variants of gradient descent
  - Vanilla GD
  - SGD
  - Batch GD
- Dimensionality reduction
- Input output relationship
- Example height vs age
- Visual understanding of PCA
- Implementing PCA
- Mathematical approach of optimizing PCA
- What is eigen vector
- PCA coding implementation

# Optimization 5: Principal Component Analysis.



Recap:

- (a) Classification - choosing  $\vec{w}$  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

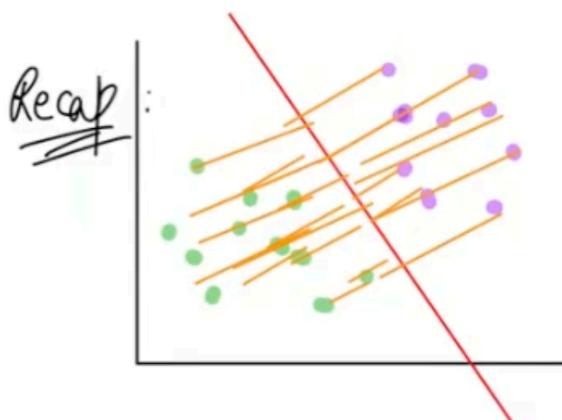
### Instructors note

- Use the last lecture's scribble note and recap the topics

Our classification problem was that we want to find the best lines that separates red points from green points such that the distances between the points and the lines is maximum

- to do that we learnt about the loss function and we also learnt how to update the parameters using the gradient descent rule

Double-click (or enter) to edit

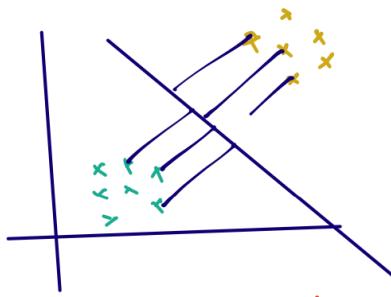


$$\text{L}(\vartheta, \bar{w}, w_0) = -\left(\sum_{i=1}^n \bar{w}^T \bar{x}_i + w_0\right) y_i + \lambda (\|\bar{w}\| - 1)$$

Q. How is the variable  $\lambda$  effecting the overall loss function?

- This is very important variable and is called as **regularizer**
- $\lambda$  controls the importance of the term  $\|\bar{w}\| - 1$ 
  - So if the value of  $\lambda$  is very large then  $\|\bar{w}\|$  becomes very very important when optimizing for the loss
    - because we want the loss function to be minimum and since  $\lambda$  is very large while optimizing we'll ensure that the  $\|\bar{w}\|$  becomes small so that when we multiply  $\lambda$  with  $\|\bar{w}\| - 1$  the values reduces thus reducing our loss function
  - And if the value of  $\lambda$  becomes very small then the importance of the term  $\|\bar{w}\| - 1$  reduces

Don't worry if you do not exactly understand this now we are just introducing the term later on we'll learn in detail how it affect the loss function

Recap:\* Regularizer.

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

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

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

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

## ▼ Variants of Gradient descent

### ▼ Vanilla gradient descent

Q. What will be the time it takes to find the minimum value using the gradient descent update rule that we saw above? Let's say that we have 10 Million datapoints

Hint: What is n in the formula?

- It is the number of data points or the size of the datasets
- For one updation of  $\bar{w}$  and  $w_0$  we'll have to go through all the data points to find the summation value.
- Thus the update process of the parameters will be very very slow plus it'll require lot of computational resources.
- Thus the time complexity of this process would be  $O(n)$

Vanilla GI

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

$$\omega_0^{(t+1)} = \omega_0^{(t)} - \eta \left( -\sum_{i=1}^n y_i \right) \quad O(n)$$

$n \rightarrow$  size of dataset.

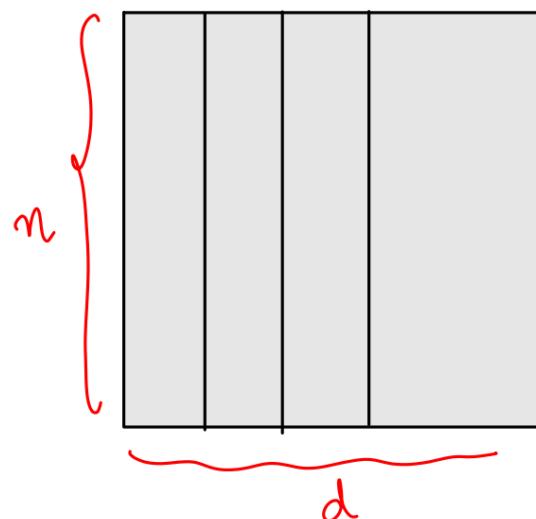
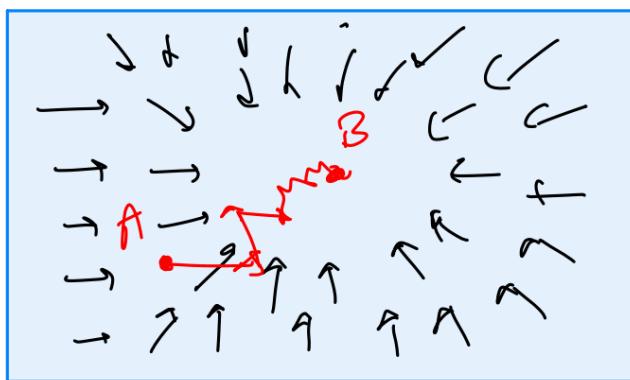
$10^7$  - 10 million datapoints.

## ✓ Stochastic Gradient descent

We saw in previous class a plot of gradients which were pointing towards the minimum value of the function and do you also recall that we start randomly with a point and then we take a step towards the minimum value

Q. Do you agree that for the vanilla gradient descent algo for each step we are looping over the entire dataset?

- Yes



So to solve this one idea was that instead of taking the entire dataset let's just take only one data point at random from the dataset and quickly update our parameters

- Since we are taking just one datapoint for updation the time complexity would be  $O(1)$

But there is one problem with this

- Let's say that the dataset that we have is our entire population and the point that we pick is a sample that we pick from the population

Q. If we have just one sample can we approximate the population mean?

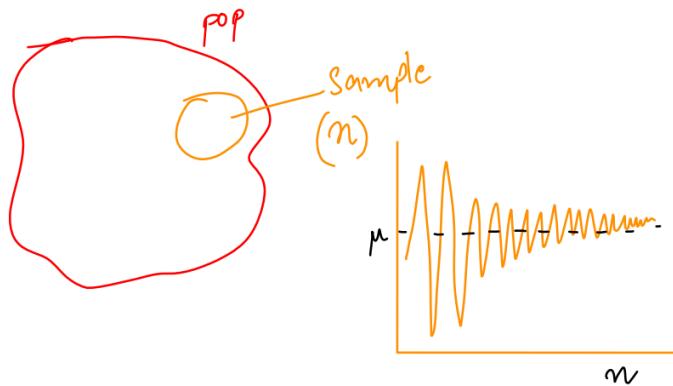
- we will be able to find some answer but we won't be much confident if it's the right value because the number of sample is just one
- Also do you remember that as we increase the sample size the variance was also decreasing ie the sample mean was approaching closer and closer to the population mean

There is also another problem which is for one point the direction would be towards the minima but for the next point it can take any direction and thus it would take us many iterations to reach out minimum value

*Stochastic Gradient Descent*

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

$$w_0^{(t+1)} = w_0^{(t)} - \eta \left( -y_i^* \right)$$



$w_0^{(t+1)} = w_0^{(t)} - \eta (-y_i^*)$

$O(1)$

We pick one point  $(\bar{x}_i, y_i^*)$  randomly from our dataset

## Batch gradient descent

So can we take leverage this concepts that we learnt in probability and statistics here?

What do we essentially want?

1. We want to move towards the right direction ie towards the minima
2. We want to reach minimum quickly

So in vanilla gradient descent it takes too much time to update each parameter whereas in SGD we update our weight vectors very quickly but it would take many steps some in right direction and some in wrong directions before reaching the minimum value.

Also it was observed that it is better to update the weights more often in comparison to increasing the intake of sample size

- thus we want to update the weights quickly and it has more importance in the real world rather than using the entire dataset for the update step
- but now if we are using just one data point to update the weights we are consuming lot of time
- So the solution to this is to use  $k$  data points together when updating the weights that way we'll be getting better directions at the same time updating the weights frequently
  - Let say if we have 100 data points in our dataset and  $k=10$  then we'll update weights 10 times whereas if we use the SGD process then we are updating weights 100 times.
  - This is called as **Mini-Batch gradient descent** or **Batch gradient descent**

### *Mini Batch Gradient Descent*

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

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

$n = \text{total \# data points}$

$k=1$  SGD

$1 < k < n$  Mini-batch GD

$k=n$  GD

In the above update rules

- if  $k=1$  then we have our stochastic gradient descent
- if  $1 < k < n$  then we have batch gradient descent
- if  $k=n$  then we have our vanilla gradient descent

### What is Epoch?

- Using our entire dataset or iterating over our entire dataset to update the weights once is called as one epoch
  - So if we have 100 data points and if we have a batch size of 10 then we'll update our weights 10 times once this is referred to as 10 iterations
  - once the 10 iterations are done we have essentially gone through the entire dataset thus we call that as one epoch.
  - In the next epoch we'll again go through the entire dataset with the batch size of 10

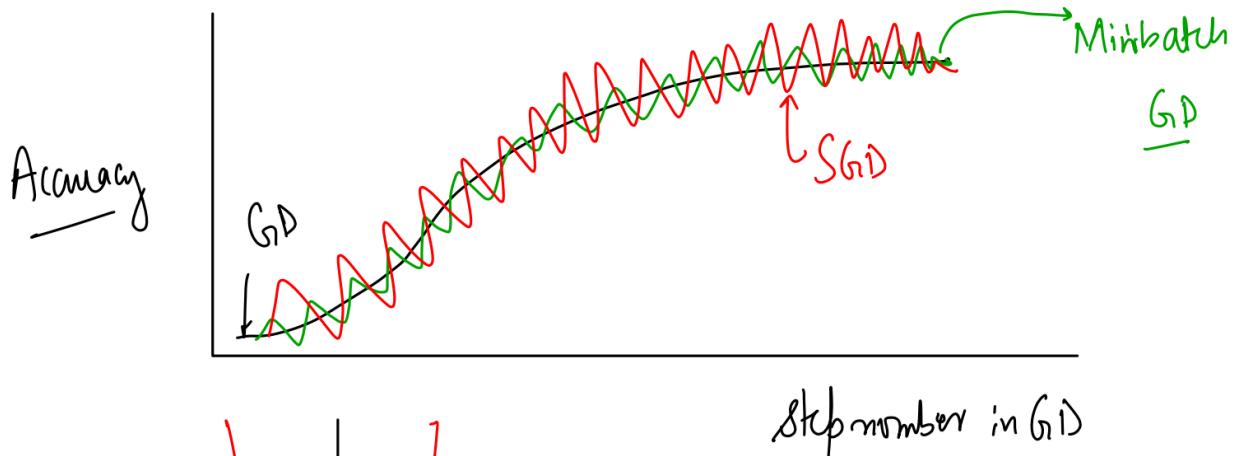
Q. If we want to perform 10 epoch where the batch size is 10 and number of data points is 100 then how many iterations we would have?

- 100

### Accuracy vs Step number of Gradient descent

when we plot accuracy vs step number

- after certain number of steps the curve plateaus ie the accuracy does not improve much no matter how hard we try.
- So this smooth curve would be for our vanilla gradient descent
- For Mini batch SGD it'll be less smoother there will be lot of variations along the way
- For SGD we'll have too much variance along the way compared to mini batch sgd



## ▼ Dimensionality reduction

We had written our dataset mathematically which is  $D = \{x_i, y_i : x_i \in \mathbb{R}^d\}_{i=1}^n$

- Here d is the dimension of the data point
- We can't visualize more than 3 dimensions

Let's see what are the problems with high dimensions ie as dimensions increases

- Visualization is tough
- Training time increases
  - We saw a calculation in previous classes that when we increased only one feature it would take 2 years for the computation.
- Computation resources requirement increases
- If we have too many dimensions the math required to solve becomes complex

So dimensionality reduction helps us solve this problem. One such algorithm that we'll study is **Principal component analysis**

## Dimensionality Reduction

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

$d > 3$  difficult to visualize

Problems with higher dimensional data -

- ① Visualisation is tough
- ② Training time increases.
- ③ Computation resources requirement increases.
- ④ Difficult to work math around such dataset.

Q. Don't you think when we say that we want to reduce the dimensions of the data essentially we'll be lossing our information of the dataset?

- Yes
- Thus when we'll learn about PCA we'll look at how smartly it reduces the dimensions without lossing the essence of the data
  - here the term essence of the data means that we are not lossing out much of the information present in the dataset

## Input Output relationship

Let's see on a high level what does PCA does exactly

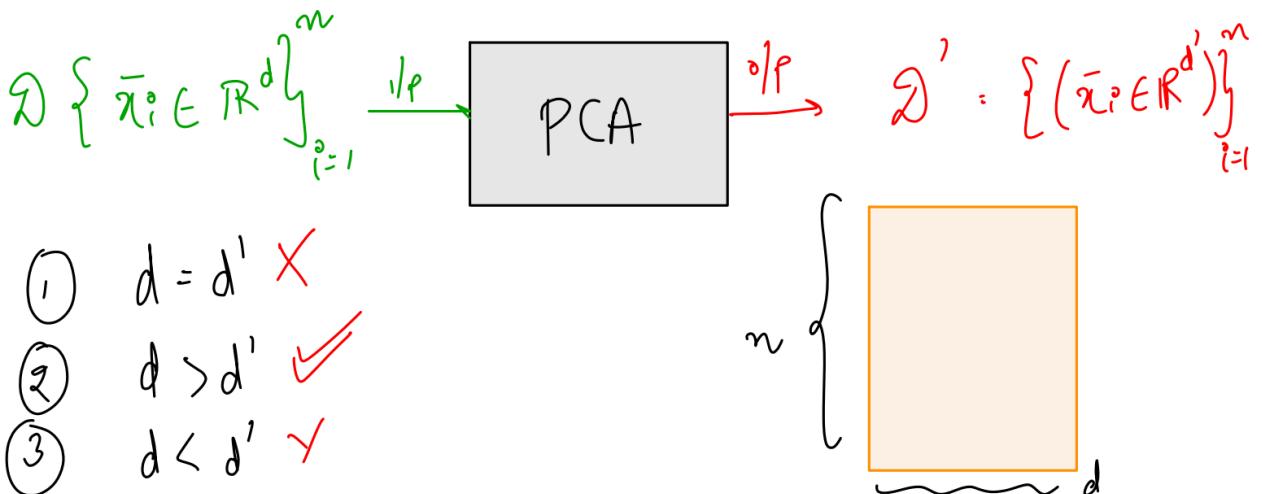
Let's consider PCA as a black box model for now ie we don't know what it does internally

- So the input to this PCA is a dataset with data points belonging to  $d$  dimensions
- And the output would be another dataset where the datapoint belongs to  $d'$  dimensions which is different from  $d$

Q. Choose the correct option for  $d'$

1.  $d' > d$
  2.  $d' < d$
  3.  $d' = d$
- Correct answer is 2

Dimensionality Reduct<sup>ion</sup> algorithms help us to deal with these problems Eg: Principal Component Analysis-



## Example Height vs Age

For PCA to work perfectly there is a condition that it needs to be satisfied let's understand what that condition is

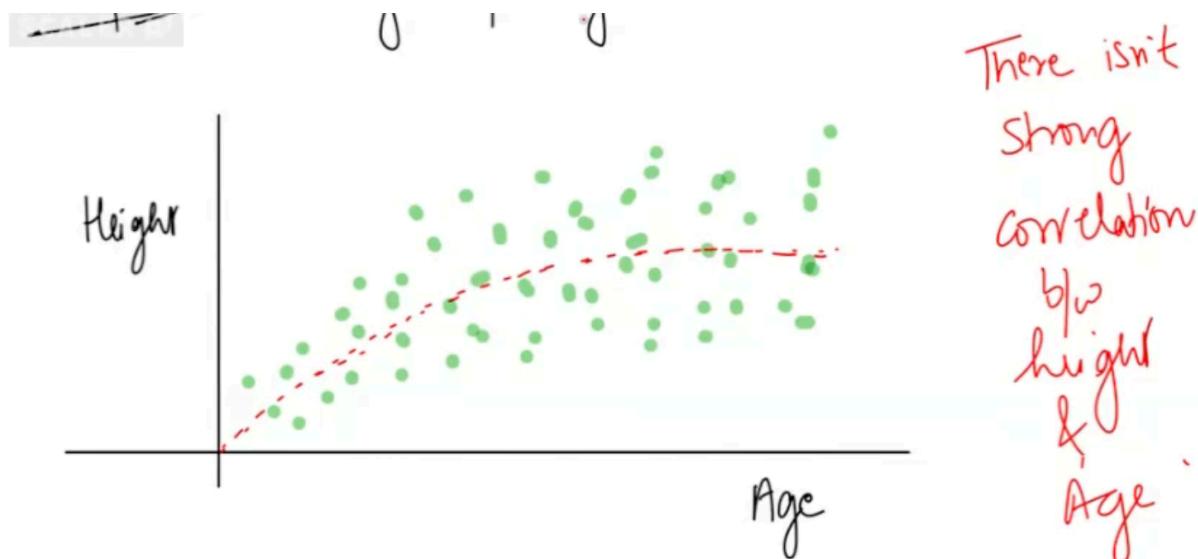
Let's see we have a plot of height vs age

- when we are babies we have a lesser height and as we grow our height increase upto a certain age after that our height stops increasing

Q. If were to find out pearson correlation value would it be a big value or a low value?

- There isn't a strong correlation between height and age because it is a non linear function
- So in such cases dimensionality reduction is not possible without losing the essence of data.

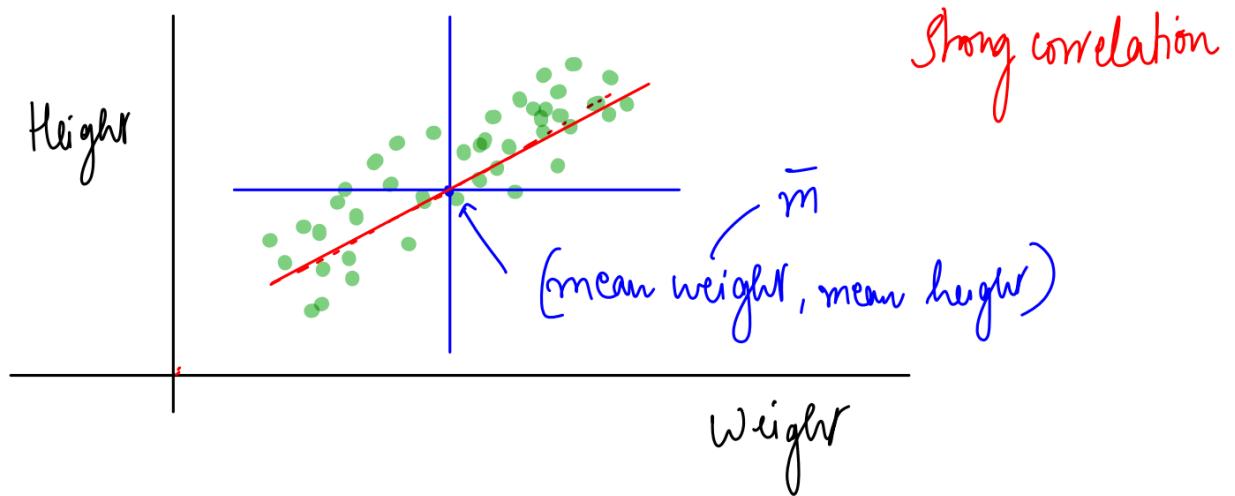
- There should be strong correlation between the 2 variables for dimensionality reduction to work



In such cases, DR is not possible without losing information of data.

Dimensionality reduction works well when we have strong correlation as shown below for height vs weight plot

- here height and weight are linearly correlated with each other



In such cases, DR is possible

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

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

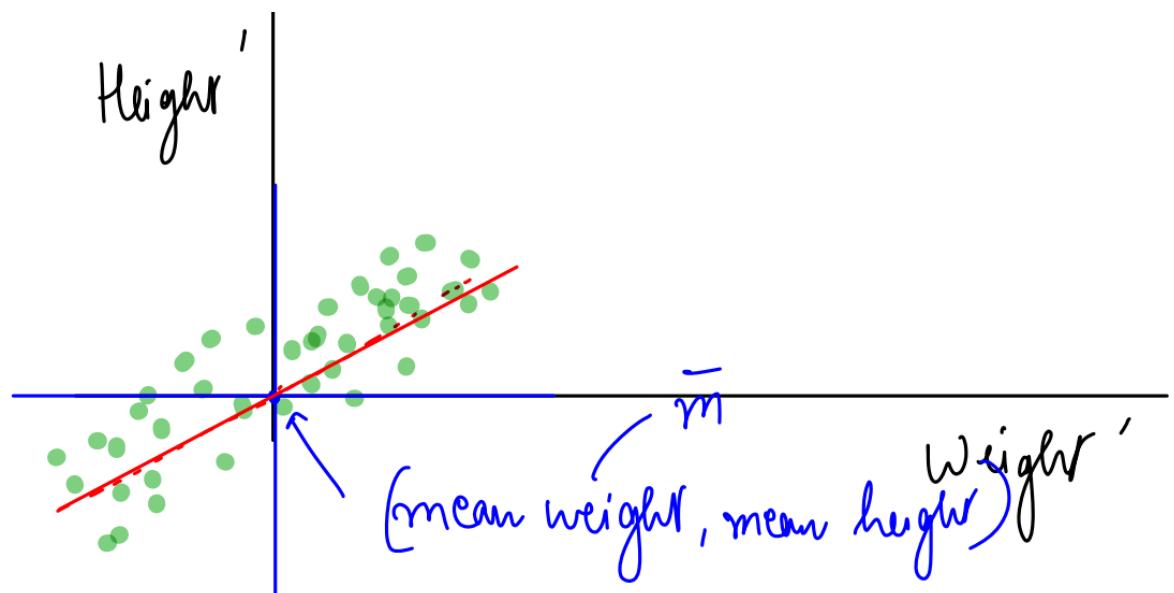
## ▼ Visual understanding of PCA

Now let's learn the interrals of what a PCA does

- First is we'll shift the origin to the mean values of the dataset ie we'll find out what the mean height and mean weight and we can represent this mean point as  $\bar{m}$

Q. How do we shift the origin to the mean point  $\bar{m}$

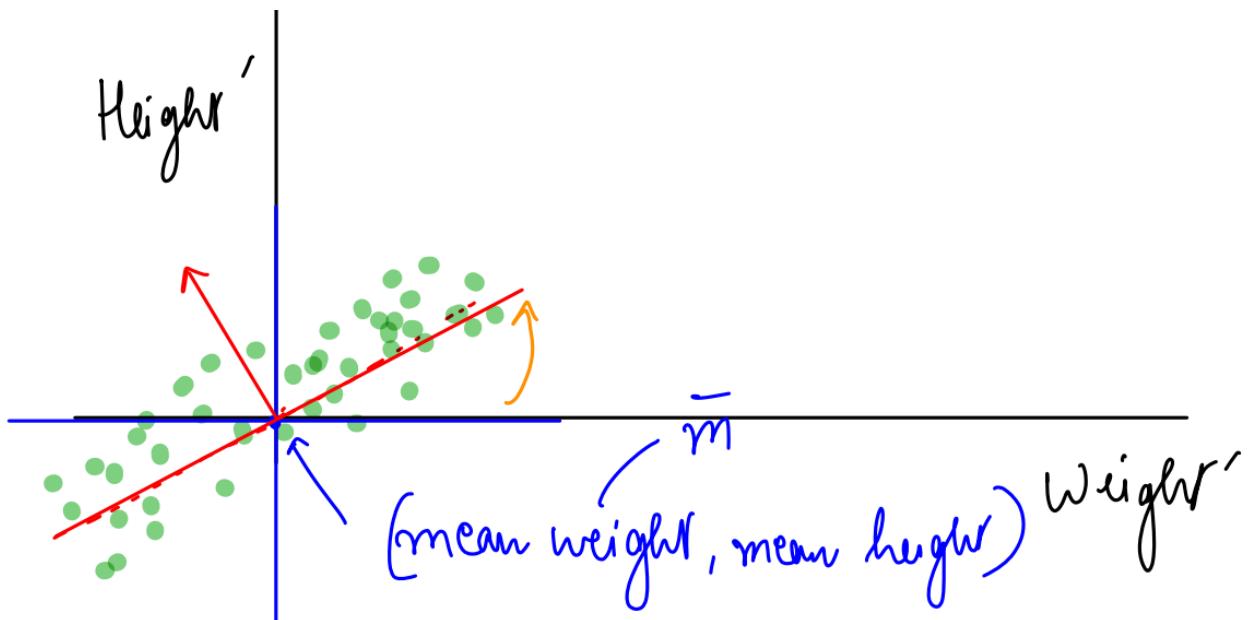
- We subtract each of our points in the dataset with the mean point that way the original mean points cancel out and we'll get a new origin which sits at the mean point  $\bar{m}$



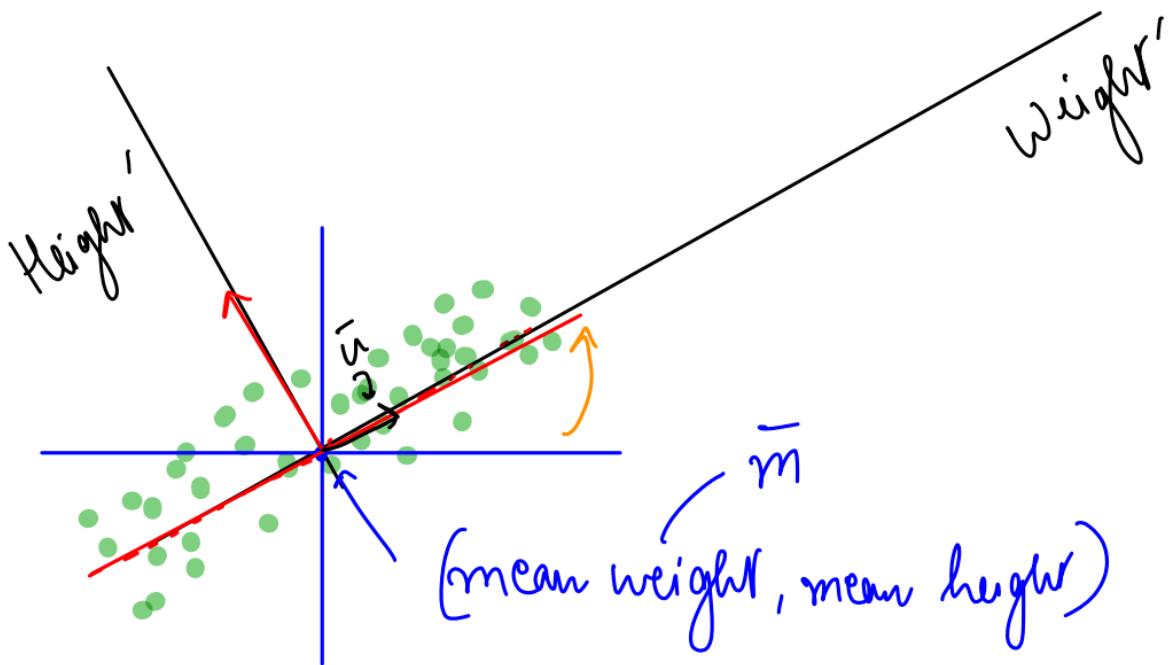
$$\mathcal{D}' : \{ (\bar{x}_i - \bar{m}) \}$$

We can now create a norm vector which is perpendicular to our Red line

Next what we want to do is we want to rotate our original axis such that they align with the red line.



We'll get something like as shown below

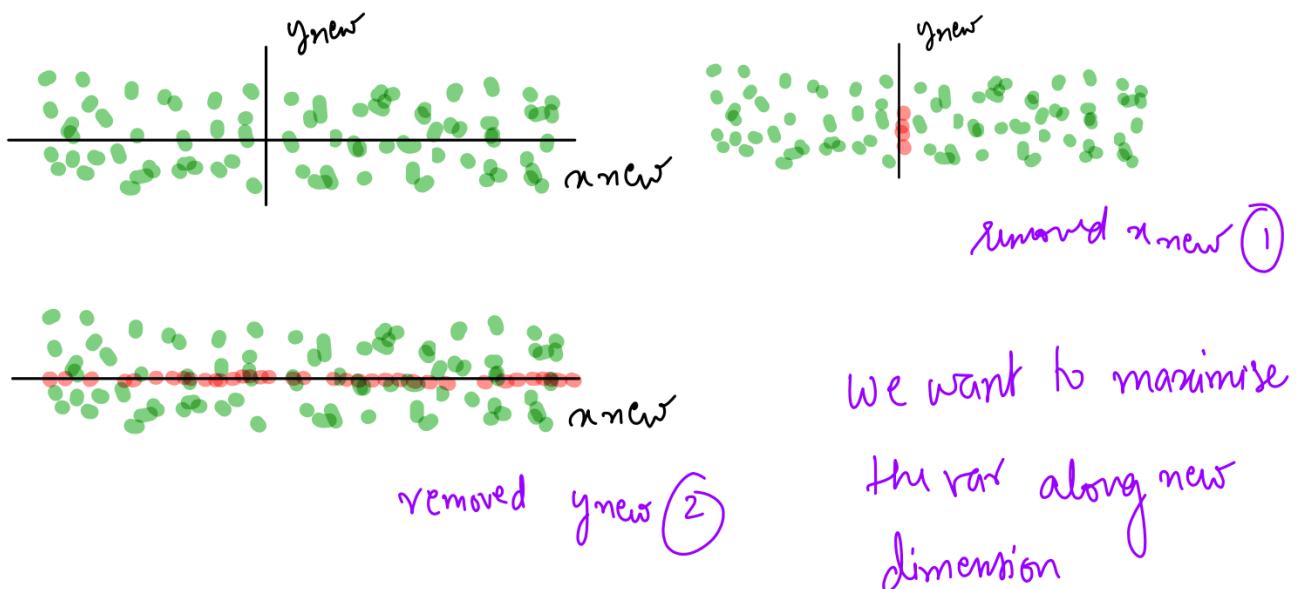
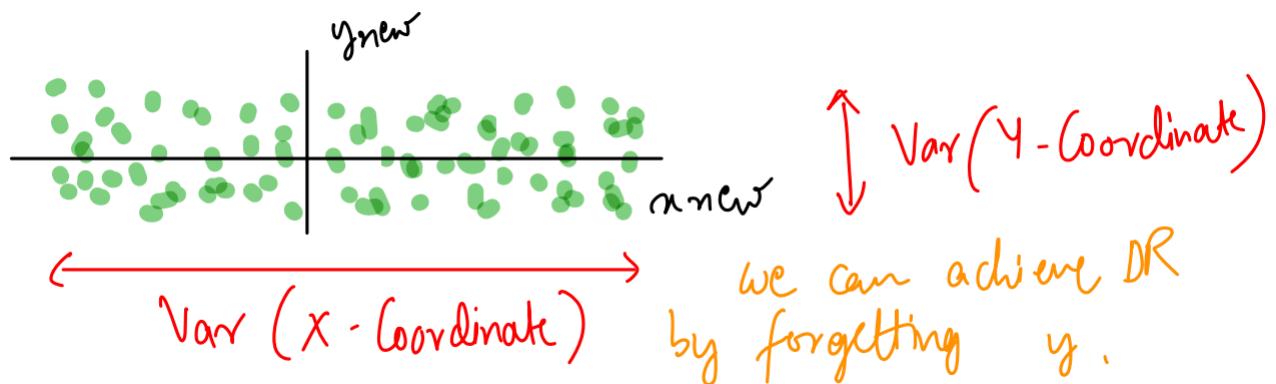


Let's call this new axis as  $x_{new}$  and  $y_{new}$

Next let's just focus on the green points and the new axis where do you think most of the information lies ?

- x axis, that is because let's say that the points to the left of the vertical black lines are peoples who had less height and less weight
- and points to the right of the vertical black lines are peoples who had more height and more weight

- so if we project them onto the x we'll still be able to differentiate between this 2 groups
- whereas if we were to project them onto the new y axis it would be very hard to differentiate between these 2 groups
- Mathematically we want to check what is the variance along the x axis and variance along the y axis
  - And we want to select that axis which has the maximum variance
- Thus we can achieve dimensionality reduction by forgetting the y axis

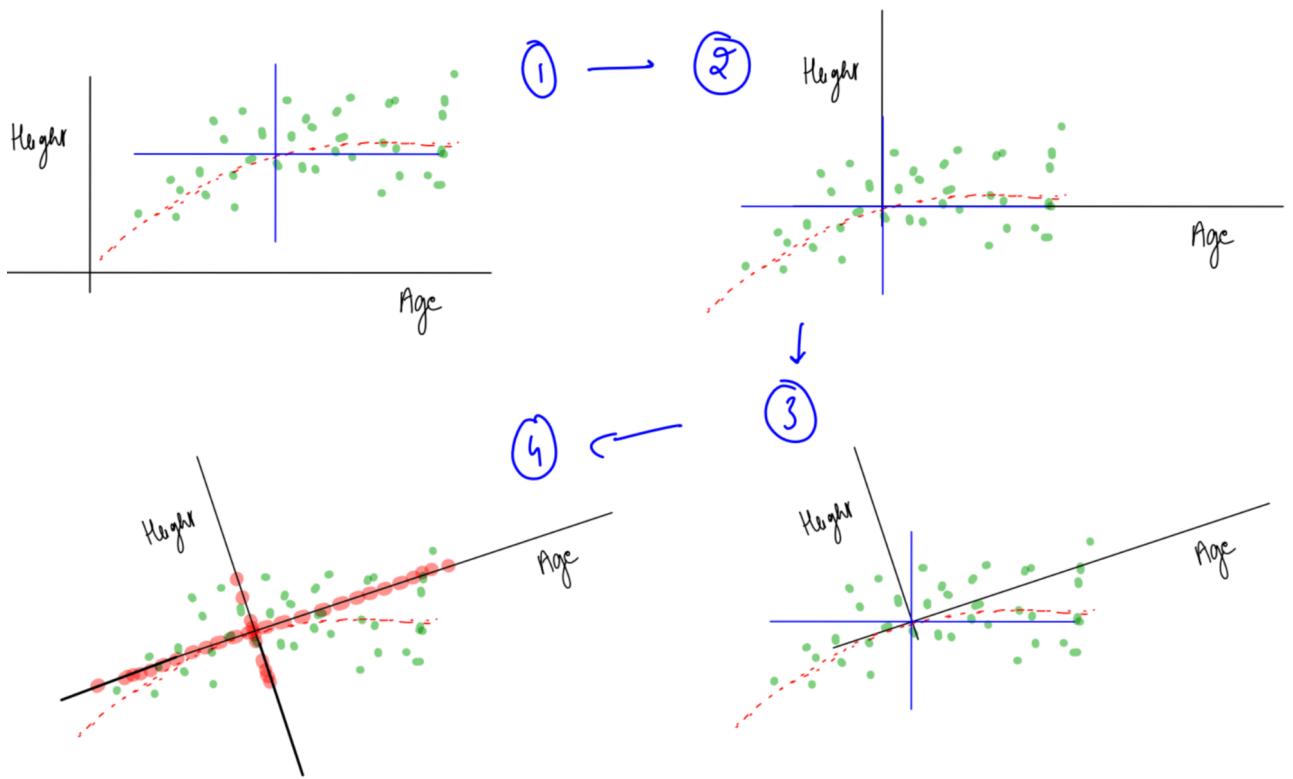


When we have 2 variables which are not strongly correlated then in that case the variances across the 2 new axis will come out to be similar hence we won't be able to choose the axis that preserves maximum information hence PCA algorithm does not work very well in this case

#### Instructors note:

- Below is the case where the features are not correlated in the 4th step demonstrate that the how the variances are almost similar hence it is difficult to

choose which axis preserves most of the information



## ▼ Implementing PCA

### Step1: Column standardization

- For every column or feature we subtract the data point value from the mean of that feature and divide it by the standard deviation of that feature and we do this for all the feature present in the dataset

## Implementation PCA

~~Step 1~~

## Column Standardisation

$$\frac{\bar{x}_i - \bar{\mu}}{\sigma}$$

$f_1$	$f_2$	$f_3$	$f_4$	$f_5$	$f_6$	$f_7$
• • • • • • • • • •		• • • • •				

$\bar{x}, \sigma_x$

$\bar{x}_i = \frac{1}{n} \sum_{i=1}^n x_i^o$

$\bar{y}, \sigma_y$

## Step2

- Now when we were looking at projecting the data points onto the new axis essentially we were looking for a vector  $u$  on which we can project all the data points and wherever that is maximum we'll consider that
  - Thus this step is about finding a vector  $\bar{u}$  such that the variance of  $\bar{x}$  (our original datapoints) along the vector  $\bar{u}$  is maximized

Q. How do we find the projection of a vector  $\bar{x}_i$  onto new vector  $\bar{u}$ ?

- By performing the dot product between the  $\bar{x}_i$  and  $\bar{u}$  and dividing it by the norm of the  $\bar{u}$
  - $\frac{\bar{x}_i \bar{u}}{||\bar{u}||}$  or we can also write as  $x_i \cdot \hat{u}$
  - So after projecting on the new vector we'll be getting the distance of that point from the mean point or the new origin that we shifted to the mean point

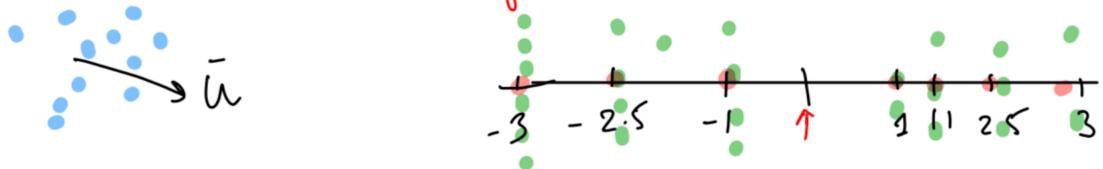
Q. Next we want to find the variance along the new axis how can we do that? let's denote the new projected points as  $x'$  thus  $x' = \frac{\bar{x}_i \bar{u}}{\|u\|}$

- Formula for variance is  $\text{Var} = \frac{1}{n} \sum (x_i - \mu)^2$
  - here  $x_i$  is  $x'$  thus our formula will become

$$\circ \text{ Var} = \frac{1}{n} \sum \left( \frac{\bar{x}_i \bar{u}}{\|u\|} - \mu \right)^2$$

- But do you agree that  $\mu$  would be 0 because we have shifted our origin to the mean point and this is called as standardizing the data

Step 2 Find a vector  $\bar{v}$  such that the variance of  $\bar{x}$  along  $\bar{v}$  is maximised.



$$\text{Projection of } \bar{x}_i \text{ on } \bar{u} = \frac{\bar{x}_i \cdot \bar{u}}{\|\bar{u}\|} : \bar{x}_i \hat{u}$$

$$\text{Var}(x_i') = \frac{1}{n} \sum_{i=1}^n \left( \frac{\bar{x}_i \cdot \bar{u}}{\|\bar{u}\|} - \mu \right)^2$$

Now we convert this into a constraint problem and then using lagrange multiplier to convert that into an unconstraint problem

Note: that if  $\|u\| = 1$  then  $\|u\|^2 = 1$

$$\text{Var}(x_i') = \frac{1}{n} \sum_{i=1}^n \left( \frac{\bar{x}_i \cdot \bar{u}}{\|\bar{u}\|} \right)^2$$

↓ Lagrange Multiplier

$$\bar{u}^* = \underset{\bar{u}}{\text{argmax}} \frac{1}{n} \sum_{i=1}^n (\bar{x}_i \cdot \bar{u})^2$$

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

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

Now we can either solve this optimization problem using gradient ascent or we have a mathematical approach to solve this optimization problem as it is more efficient

## Mathematical approach of optimising PCA

Let's learn about the mathematical approach

- First thing is to change how we are looking at  $(\bar{x}^T \bar{u})^2$
- We have our dataset D where  $x_i$  belongs to a D dim space and we have n data points
- Thus we'll represent each row of our dataset as a row vector ie for  $i^{th}$  data point it'll be  $x_i^T$  we'll represent this dataset a matrix of dimension  $n * d$  and this matrix is called as a **data matrix** and let's call this as X
- Next we'll take our data matrix and multiply it with our vector u which is a column vector
  - Note that  $\bar{u}$  is a d dimensional
- So when we multiply X with u matrix the resultant matrix that would have a shape of  $n * 1$ , let's call this as vector v
  - First row entry in the new vector v would be nothing but the dot product between first row of X matrix and the u vector.
  - Now when we compute the norm of the v vector we'll essentially be computing  $\sum_{i=1}^n (\bar{x}^T \bar{u})^2$  which is the first term of our optimization problem

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

- change - !  
how we are  
seeing this -

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

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

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

This matrix is called the data matrix.  
or  $X$ .

Let's redefine our Optimization problem in terms of  $X$ .  $\bar{u}$

- Our new optimization problem is
  - $\bar{u}^* = \operatorname{argmax}_{\bar{u}} \frac{1}{n} (\bar{v}^\top \bar{v}) + \lambda (||\bar{u}||^2 - 1)$ 
    - Here we have represented  $||v||^2$  as  $\bar{v}^\top \bar{v}$  it means the same thing
  - If we simply  $\bar{v}^\top \bar{v}$  in terms of  $X$  and  $\bar{u}$  then
    - $\bar{v}^\top \bar{v} = (X\bar{u})^\top \cdot X\bar{u}$
    - We can rewrite this as  $\bar{v}^\top \bar{v} = \bar{u}^\top X^\top \cdot X\bar{u}$

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

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

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

Why did we do all of this?

- So that computing gradients becomes easier

Gradient of the function  $f(\bar{u}) = \bar{u}^T S \bar{u}$  where  $S$  is any matrix then  $\nabla_{\bar{u}} f(\bar{u}) = (S + S^T) \bar{u}$

- we can also write this as  $\nabla_{\bar{u}} f(\bar{u}) = \bar{u}^T (S + S^T)$
- Note that we are taking this as a fact and not looking at the proof of this as proving this will consume a lot of time.

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

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

Now we can call the entire equation after the *argmax* as

$$f(\bar{u}) = \frac{1}{n} \bar{u}^T X^T X \bar{u} + \lambda (||\bar{u}||^2 - 1)$$

- We can represent  $||\bar{u}||^2$  as  $\bar{u}^T \bar{u}$
- So here  $X^T X$  would be our matrix  $S$  from fact 1
- So using Fact 1 let's compute the gradient of  $f(\bar{u})$  would be
  - $\nabla_{\bar{u}} f(\bar{u}) = \frac{1}{n} ((X^T X) + (X^T X)^T) \bar{u} + \lambda(2\bar{u})$
- Now  $(A \cdot B)^T$  is equal to  $A^T \cdot B^T$  using this we can rewrite  $(X^T X)^T$  as  $X^T X$
- Thus the final equation becomes  $\nabla_{\bar{u}} f(\bar{u}) = \frac{2}{n} (X^T X) \bar{u} + \lambda(2\bar{u})$

$$f(\bar{u}) = \frac{1}{n} \underbrace{\bar{u}^T (X^T X) \bar{u}}_{\text{Fact 1: } \downarrow} + \lambda (\underbrace{\bar{u}^T \bar{u} - 1}_{\text{Fact 2: } \downarrow})$$

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

$$\rightarrow \text{Since } (X^T X)^T = X^T X \quad \cdot + \lambda (2 \bar{u})$$

$$= \frac{2}{n} X^T X \bar{u} + \lambda (2 \bar{u})$$

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

$$\frac{\partial \bar{u}^T \bar{u}}{\partial \bar{u}} = 2 \bar{u}$$

So our gradient equation is

- $\nabla_{\bar{u}} f(\bar{u}) = 2 \left( \frac{1}{n} X^T X \bar{u} + \lambda \bar{u} \right)$
- to find the minima or maxima we'll equate the gradient to 0
- Thus  $2 \left( \frac{1}{n} X^T X \bar{u} + \lambda \bar{u} \right) = 0$ 
  - Here 2 will go on the other side cancel out and shifting  $\lambda \bar{u}$  and  $n$  to other side our equation will be
    - $X^T X \bar{u} = -n \lambda \bar{u}$
    - We can represent  $-n \lambda$  as  $\lambda'$ 
      - $(X^T X) \bar{u} = \lambda' \bar{u}$
      - This equation represents the eigen value eigen vector expression

So this equation tells us that the vector  $\bar{u}$  which maximizes  $f(\bar{u})$  is the one which satisfies equation 1

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

$$\Rightarrow \frac{1}{n} X^T \times \bar{u} = -\lambda \bar{u}$$

$$\Rightarrow X^T \times \bar{u} = \underbrace{-(n\lambda)}_{\lambda'} \bar{u}$$

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

↑

Eigenvalue, eigenvector expression!!

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

By observing equation 1 we notice that  $\bar{u}$  is nothing but eigen vector of the matrix  $X^T X$

## What is eigen vector

For any square matrix  $S$  if we find an  $\bar{x}$  such that  $S\bar{x} = \lambda\bar{x}$  then  $\bar{x}$  is said to be an **eigen vector**

- Here square matrix means number of rows equals number of columns
- Word eigen comes from the german word same or identity

What is an eigenvector?

"eigen" → German word  
= identity.

For any square matrix  $S$ ,

↳ no. of rows = no. of cols. same direction

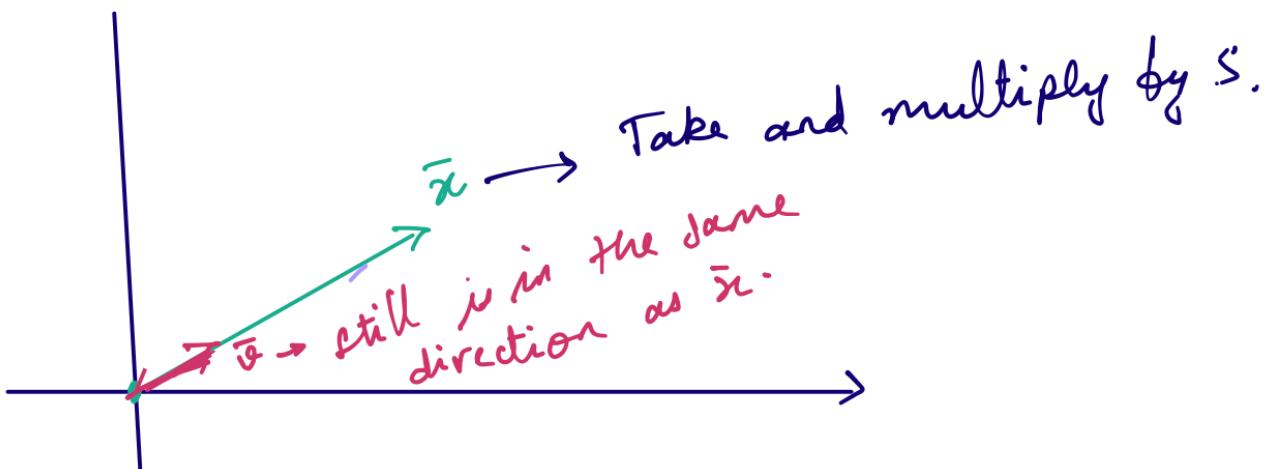
if we find an  $\bar{x}$  such that

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

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

Let's say in our 2D space we have a vector  $\bar{x}$  and we multiply it with a square matrix  $S$  which is of size  $d * d$  the result is another vector  $\bar{v}$  of size  $d * 1$

- $\bar{v} = S \cdot \bar{x}$
- If  $S \cdot \bar{x} = \lambda \cdot \bar{x}$  where  $\lambda$  is a number then  $\bar{v} = \lambda \cdot \bar{x}$
- which means that the direction of  $\bar{v}$  will be same as  $\bar{x}$  thus  $\bar{x}$  is a eigen vector because when we take vector  $x$  and multiply it with a matrix  $s$  the new vector still points in the same direction as the vector  $x$
- Eigen vectors are the property of the matrix  $S$  where when we multiply a matrix with a vector if the direction remains the same and the magnitude of the resultant vector changes then the vector that we multiply with the matrix is the eigne vector



$S_{d \times d}$

$\bar{x}_{d \times 1}$

=

$\bar{v}_{d \times 1}$

Why did we start this whole thing?

- the problem that we wanted to solve was that we wanted to find a vector  $\bar{u}$  such that the variance along the  $\bar{x}$  is maximized
- So by doing all this math we have proven that this vector  $\bar{u}$  is nothing but the eigen vector

So to solve the problem of finding a vector  $\bar{u}$  such that the variance is maximized we'll first have to find an eigen vector

- To find this we'll use numpy we have a function which gives us eigen vectors for a given matrix

Now each matrix has many eigen vectors and each eigen vector has an associated eigen value so we pick that eigen vector corresponding to the maximum eigen value for our  $\bar{u}$

So once we have the eigen values and the eigen vector of the matrix we simply have to sort the eigen values in descending order and select the eigen vector with the highest eigen value.

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

How to find this? → Numpy

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

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

We got our vector  $\bar{u}$  which maximizes our variance now how to find other directions?

- So if we want to reduce the dimensions from  $d$  to  $d'$  we simply pick the first  $d'$  eigen vector with the largest corresponding eigen value
- So let's say we have a  $3 \times 3$  matrix it'll have  $(x_1, \lambda_1), (x_2, \lambda_2), (x_3, \lambda_3)$ 
  - Because for any square matrix the maximum number of eigen vector will be equal to the dimension of the matrix
  - Here  $x_1, x_2, x_3$  are eigen vectors and  $\lambda_1, \lambda_2, \lambda_3$  are eigen values
- So we'll arrange this eigen values in descending order and let say that we want to reduce it to 2 dimensions so we'll pick 2 eigen vector  $x_1$  and  $x_2$
- And we'll use this as our new axis on which we'll project our points to find the new coordinates of the points.

### Step 3: other directions?

If we want to reduce

$$d \rightarrow d'$$

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

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

<sup>2</sup> Pick  $(\bar{x}_1, \bar{x}_2)$ ,  $(0.5, 0.3)$ ,  
These are the new ones.

Next let's see how to project the old points which are present in  $d$  dimensional space to  $d'$  dimensional space

- we'll compute the projection of each data point along the new axis ie vectors  $\bar{u} \in d'$  ie  $\bar{u}_1, \bar{u}_2 \dots \bar{u}_{d'}$  where vector

The last thing that we want to compute is how much of information have we captured or retained from the original dataset present in d dimension

- to do that  $\frac{\lambda_1 + \lambda_2 + \dots + \lambda_d}{\sum_{i=1}^d \lambda_i}$
- So based on the previous example where we saw that we were picking 2 eigen vector with eigen values of 0.5 and 0.3 then that means that 80% of information we have retained and 20% of it was lost when we transformed from 3D to 2D
- Thus to choose the d' we'll always look at what percentage of information lost or retained is acceptable based on that we'll take that many eigen vector such that the information retained is within the acceptable value

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

$$\underbrace{\frac{\bar{x}_i^\top \bar{u}_1}{\|\bar{u}_1\|}}_{\text{coordinate along } \bar{u}_1} + \underbrace{\frac{\bar{x}_i^\top \bar{u}_2}{\|\bar{u}_2\|}}_{\text{coordinate along } \bar{u}_2} \dots + \underbrace{\frac{\bar{x}_i^\top \bar{u}_{d'}}{\|\bar{u}_{d'}\|}}_{\text{coordinate along } \bar{u}_{d'}}$$

Step 5: How to calculate the information loss?

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

## ✓ PCA coding implementation

```
import numpy as np
import pandas as pd
from matplotlib import pyplot as plt
```

```
!gdown 16elpyaPZhS63ibzIdetEjZhC04FpVSZu
```

Downloading...

From: <https://drive.google.com/uc?id=16elpyaPZhS63ibzIdetEjZhC04FpVSZu>

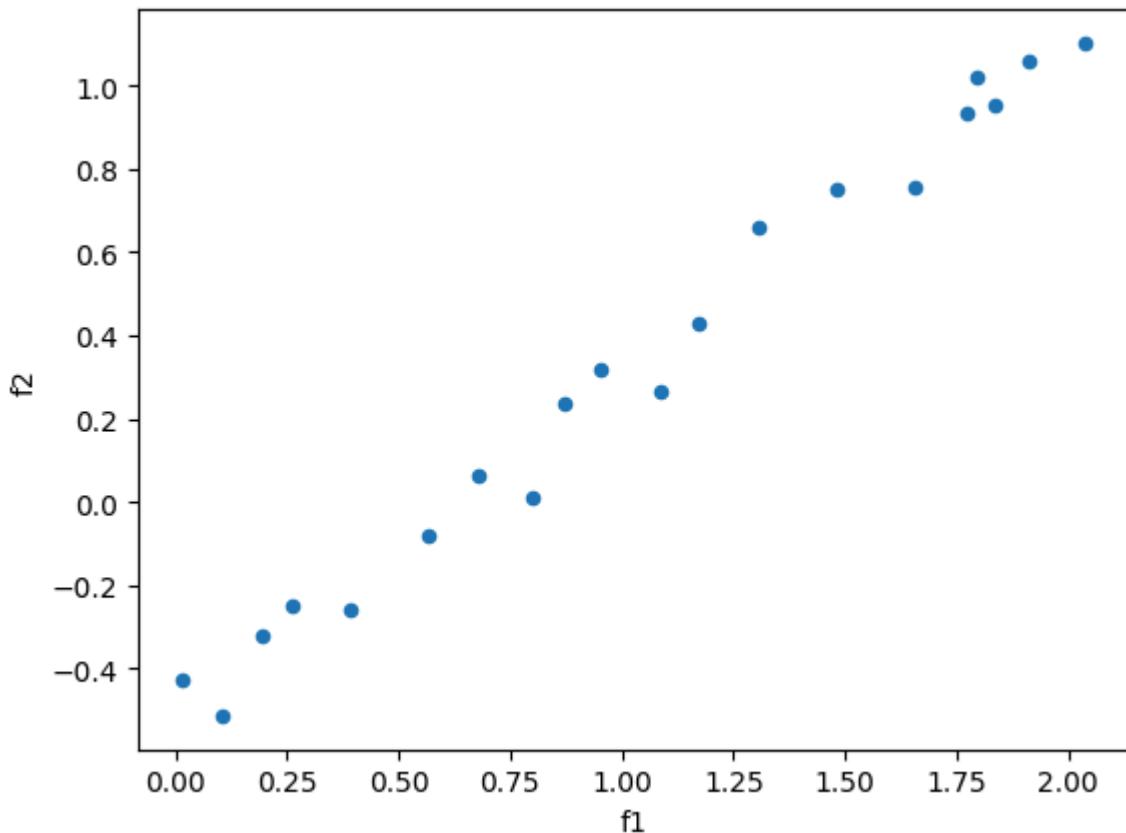
To: /content/pca\_2d\_data.csv

100% 778/778 [00:00<00:00, 4.02MB/s]

```
data = pd.read_csv('/content/pca_2d_data.csv')
data.head()
```

	f1	f2
0	2.038772	1.102760
1	0.869536	0.238030
2	1.795888	1.020364
3	0.103510	-0.513661
4	1.479554	0.750702

```
data.plot(x='f1', y='f2', kind='scatter')
plt.show()
```



```
## Implement PCA!!
# Step 1: Column standardization

mu = data.mean(axis = 0)

std = data.std(axis = 0)

X_st = (data - mu)/std
```