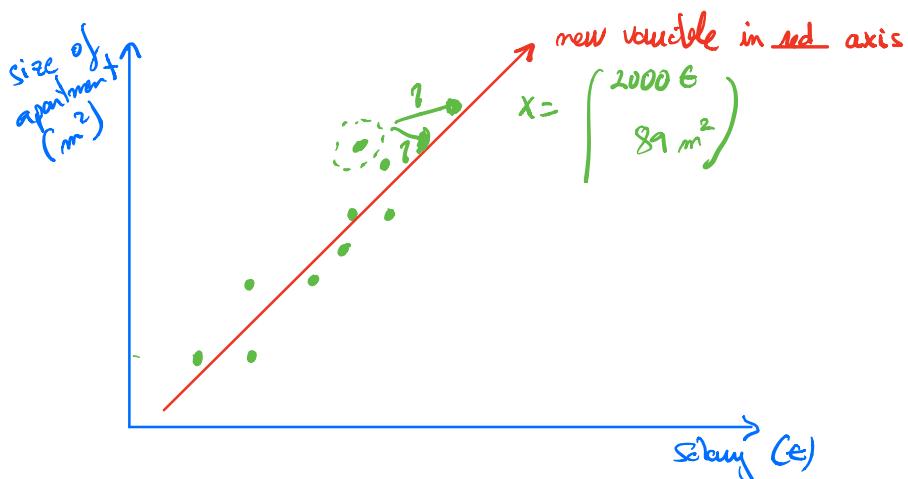


## Introduction

What does dimensionality reduction means?

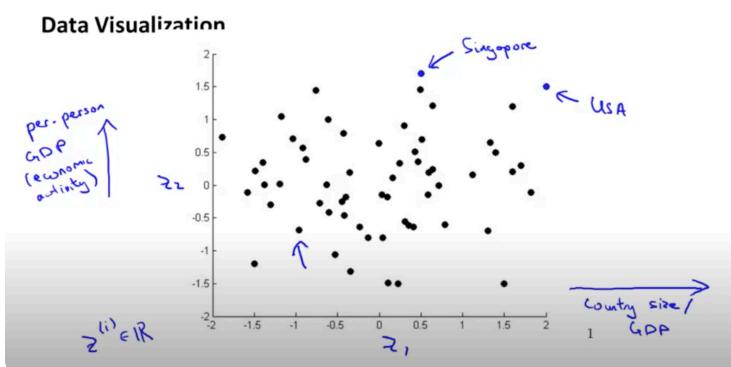
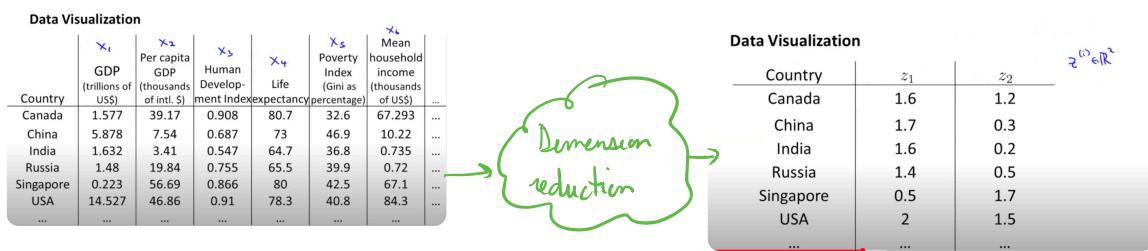
Reduce data from  $N$ -dimensional space to a lower space  
↓  
features



Which is the motivation to apply dimensionality reduction?

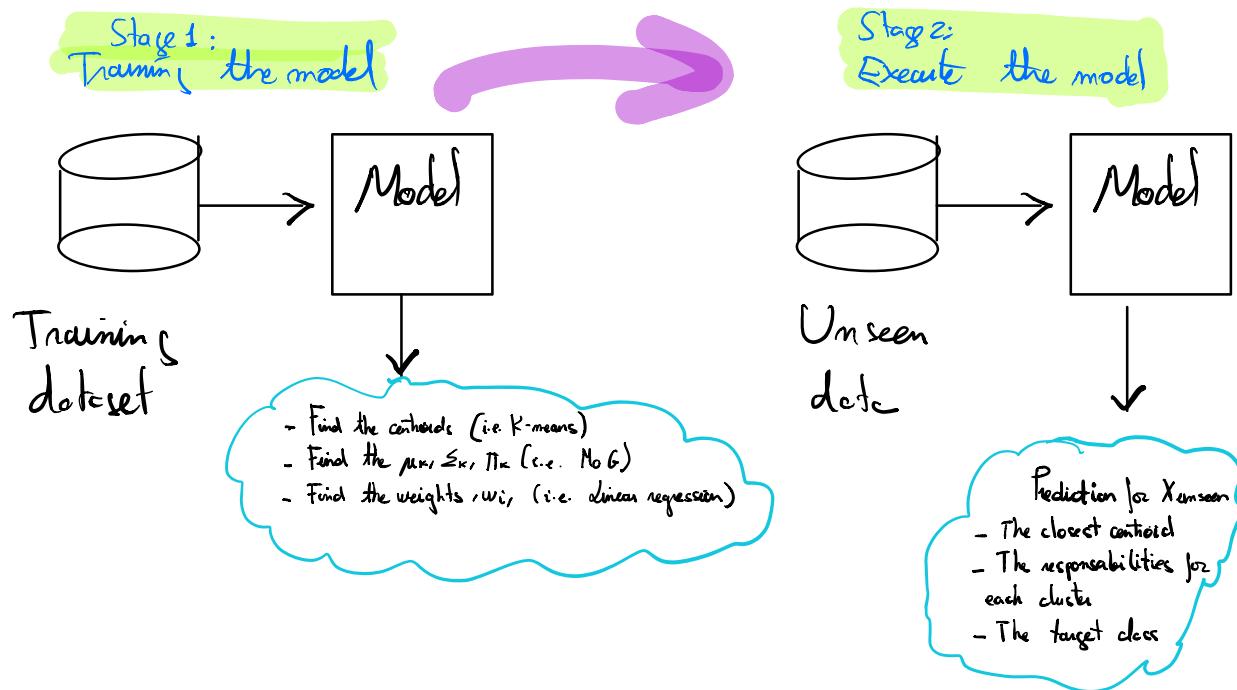
- (1) Compression of data → e.g. for video
- (2) Reduce noise → e.g. for images and video
- (3) Reduce the memory to process data
- (4) Reduce the storage
- (5) Accelerate analytics → models are executed faster
- (6) Recommendation engines
- (7) Dot visualization of high dimensional dataset

## Example of data visualization



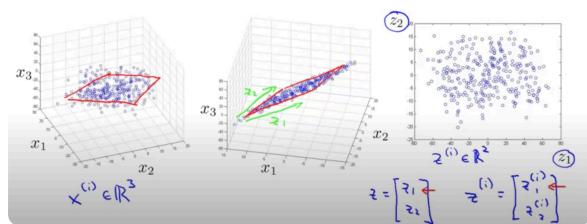
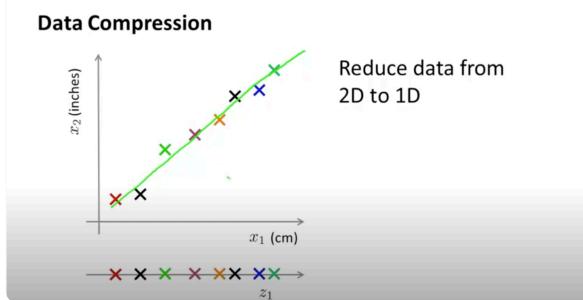
It's easier to understand the data !!

let's come back to our goal in Machine learning....



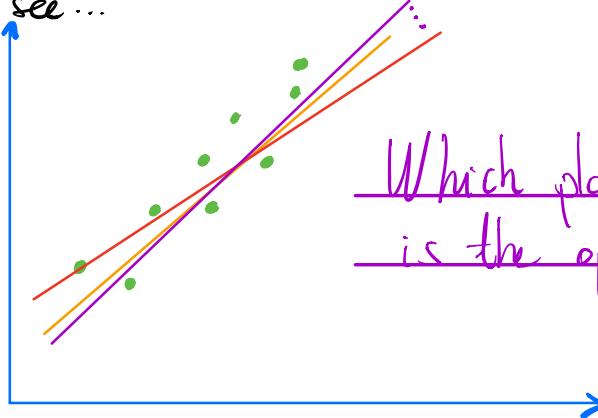
⇒ Therefore, in dimensionality reduction, our goal in the training stage is to find a new dimensional space. In other words... we need to find a mathematical expression to move from  $N$  dimensional space to an optimal new  $P$  dimensional space where  $P \ll N$

Data compression  
Reduce data from 3D to 2D



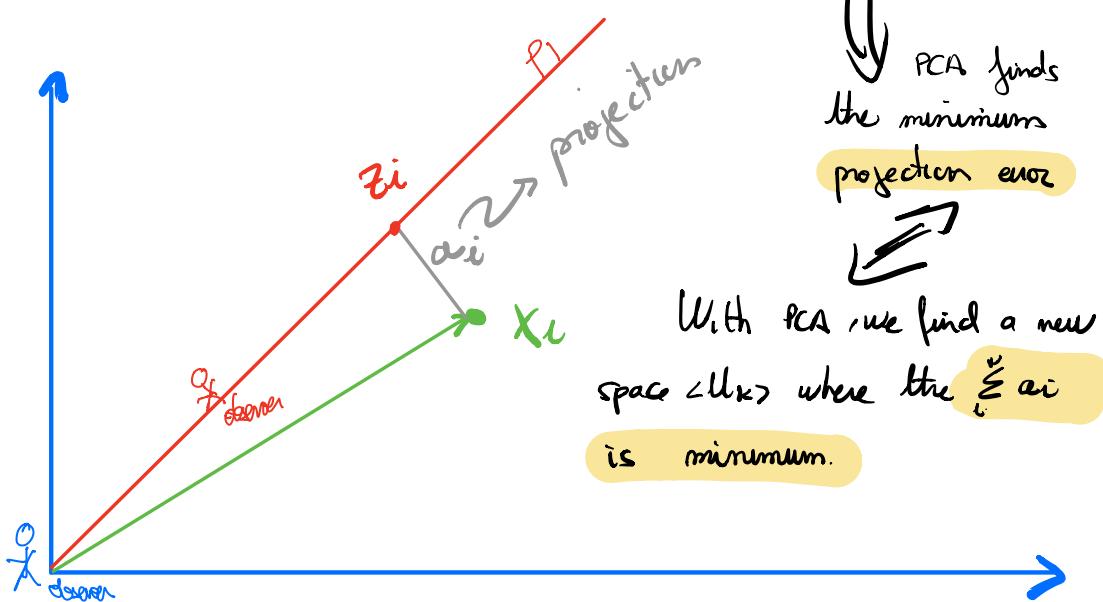
How to find the optimal new space?

let's see...



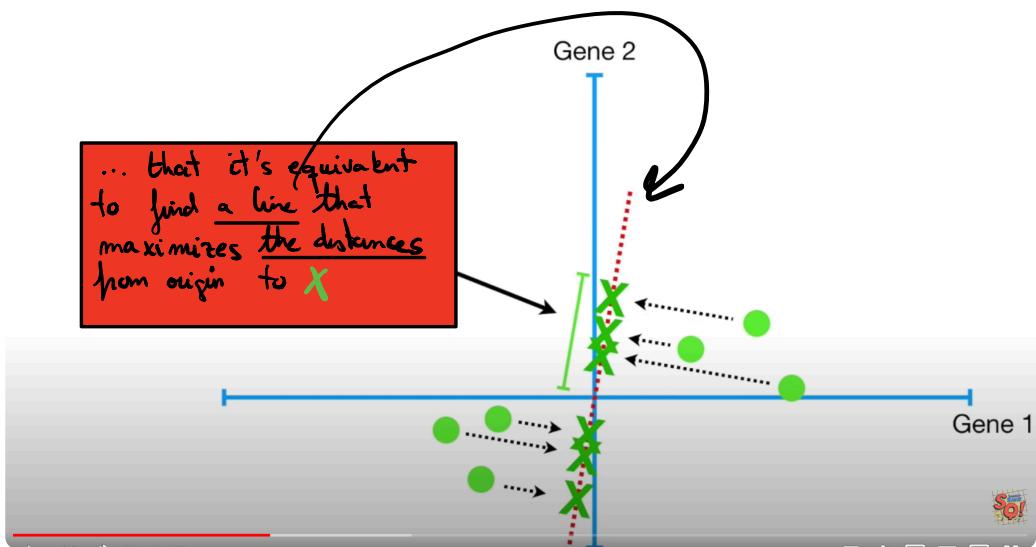
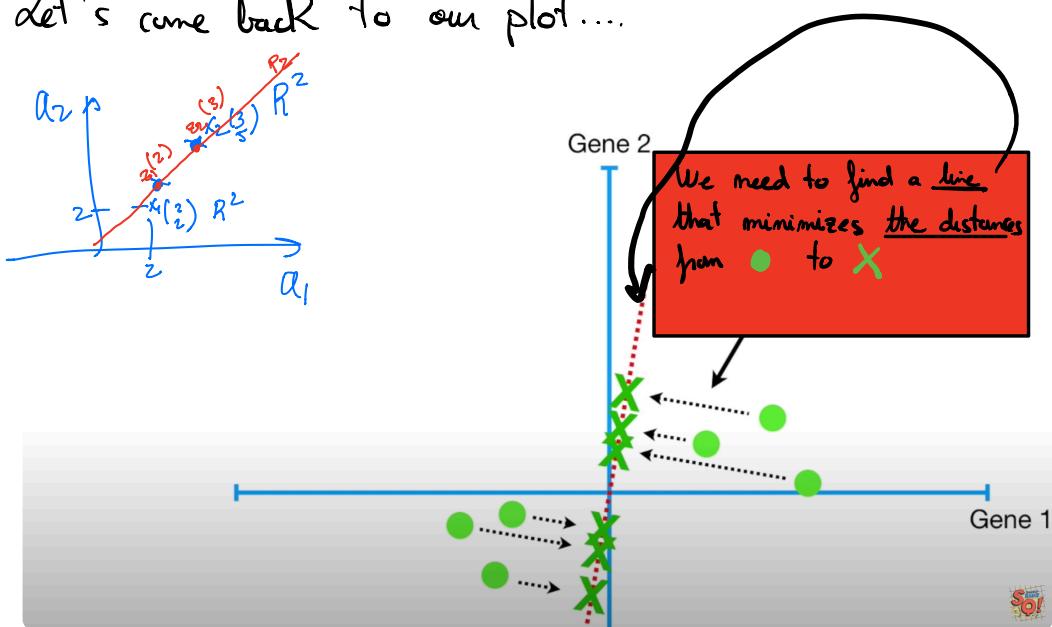
$$\begin{array}{ccc} x^{(1)} \in \mathbb{R}^2 & \longrightarrow & z^{(1)} \in \mathbb{R}^1 \\ x^{(2)} \in \mathbb{R}^2 & \longrightarrow & z^{(2)} \in \mathbb{R}^1 \\ \vdots & & \vdots \\ x^{(n)} \in \mathbb{R}^2 & \longrightarrow & z^{(n)} \in \mathbb{R}^1 \end{array}$$

Solution = Principal Component Analysis looks for a new hyperplane (i.e. dimensional space) where the distance of the training dataset's points and their projections are minimal

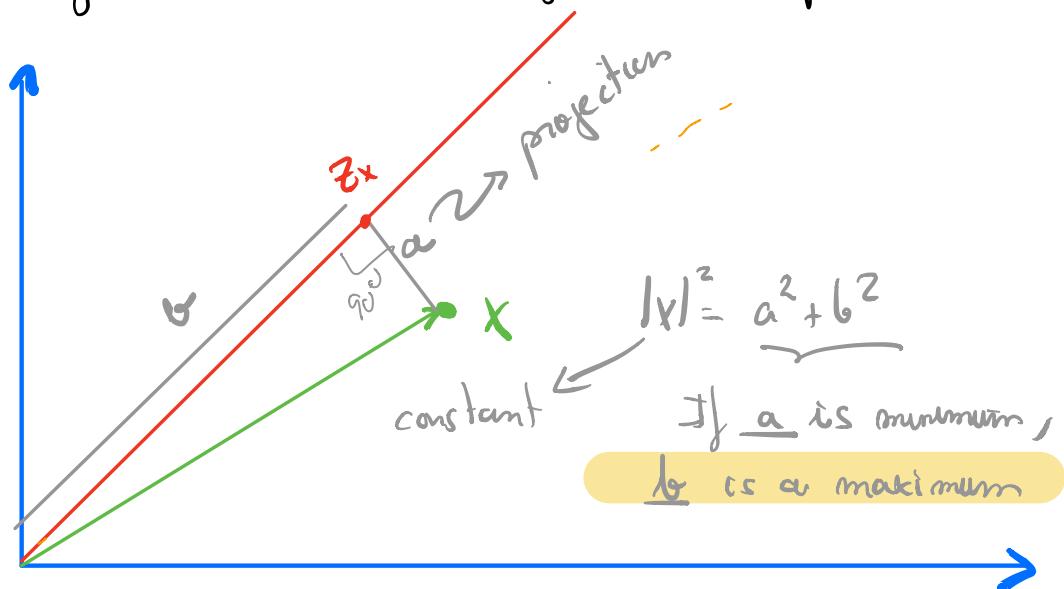


So, next question is... How to find a new space where the total projection error is a minimum?  $\leq u \geq$

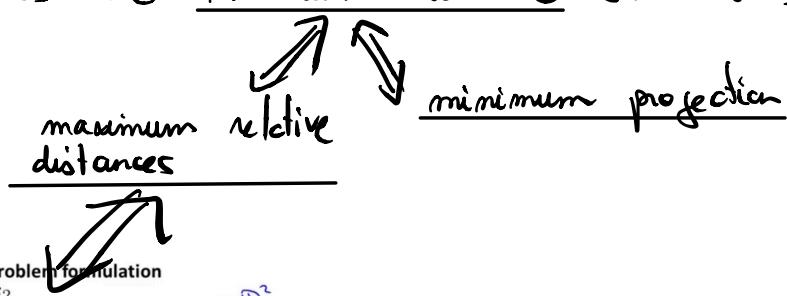
let's come back to our plot....



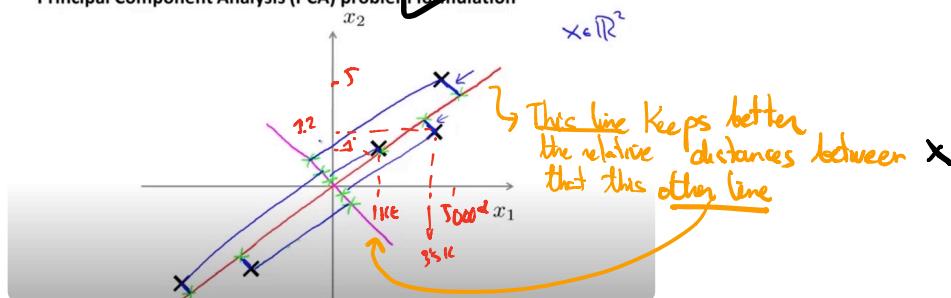
If we do zoom in just one datapoint...



Therefore, the goal in PCA is to find a new space  
<u>s</u> that captures the maximum variance (c.e. b is maximum)



Principal Component Analysis (PCA) problem formulation



In summary, PCA calculates the directions <u>s</u> of greatest variance. As eigenvectors are in the directions of maximum variance, eigenvectors are in the basis of PCA

Training Stage : How to calculate PCA of a dataset  $X$ ?

Imagine  $X$  as ....

	# cars	salary	size home	# members
Family 1	2	3000	180	5
Family 2	1	2850	95	3
...	...	...	...	...
Family N	3	4250	120	4

	$x_1$	$x_2$	$x_3$	$x_4$
$x^{(1)}$	2	3000	180	5
$x^{(2)}$	1	2850	95	3
...	...	...	...	...
$x^{(n)}$	3	4250	120	4

Stage 1 → For each feature, we calculate the mean:

$$\mu_j = \frac{1}{N} \sum_{i=1}^N x_j^{(i)} \rightarrow \text{For feature or variable (D dimensions)}$$

→ If features have different scale, it makes sense to rescale the values to have comparable range of values

Stage 2 → Remove the mean of each feature to each datapoint  $\Leftrightarrow$  It's equivalent to center the datapoints in the origin

$$x \rightarrow \tilde{x} = x - \mu$$

**Stage 3** → We calculate the "covariance matrix" of  $\bar{X}$ , i.e.  $\Sigma$

**Stage 4** → Find eigenvectors and eigenvalues of  $\Sigma$  using SVD :

$$\Sigma = U^{-1} \cdot D \cdot U$$

where  $U$  is the matrix of eigenvectors of  $\Sigma$ ; i.e.  $U = \begin{bmatrix} | & | & | \\ u_1 & u_2 & \dots & u_D \\ | & | & | \end{bmatrix}$

where  $D$  is the matrix of eigenvalues of  $\Sigma$ , i.e.  $D = \begin{bmatrix} \lambda_1 & 0 & \dots & 0 \\ 0 & \lambda_2 & \dots & 0 \\ 0 & \dots & \dots & \lambda_D \end{bmatrix}$

**Stage 5** → We order eigenvectors based on their eigenvalues :  $\lambda_1 > \lambda_2 > \dots > \lambda_D$

↑ greatest eigenvalues, greatest variance

**Stage 6** → We select the first  $k$  eigenvectors

$$\begin{array}{ccccccccc} \lambda_1 & > & \lambda_2 & > & \dots & > & \lambda_k & > & \lambda_{k+1} & > \dots & > & \lambda_N \\ \downarrow & & \downarrow & & \dots & & \downarrow & & \downarrow & & \dots & & \downarrow \\ U_{\text{reduce}} & & u_1 & & u_2 & & \dots & & u_k & & u_{k+1} & & \dots & & u_N \end{array}$$

Stage 7 → We are ready to project datapoints in the new Reduce space as follows:

$$\underline{\tilde{x}} = U_{\text{reduce}}^{-1} \cdot \underline{\tilde{X}}$$



$U_{\text{reduce}}^{-1} = U_{\text{reduce}}^T$  because  $U_{\text{reduce}}$  is orthonormal

We should take  $\leq k$  eigenvectors which explain the 80%-95% of total variance

Stage 8 → Come back to original coordinate basis:

$$\underline{\tilde{X}}_{\text{reprojected}} = U_{\text{reduce}} \cdot \underline{\tilde{x}}$$

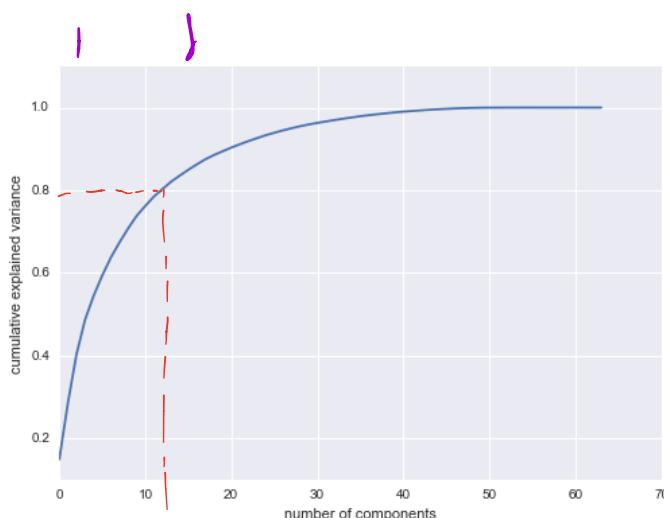
$$\underline{x}_{\text{reprojected}} = \underline{\tilde{X}}_{\text{reprojected}} + \mu$$

To evaluate the performance, we calculate the

$k = \text{selected } k \text{ components from PCA}$

$$\frac{\text{Proportion of variance}}{= \frac{\sum_{i=1}^k \lambda_i}{\sum_{i=1}^D \lambda_i}}$$

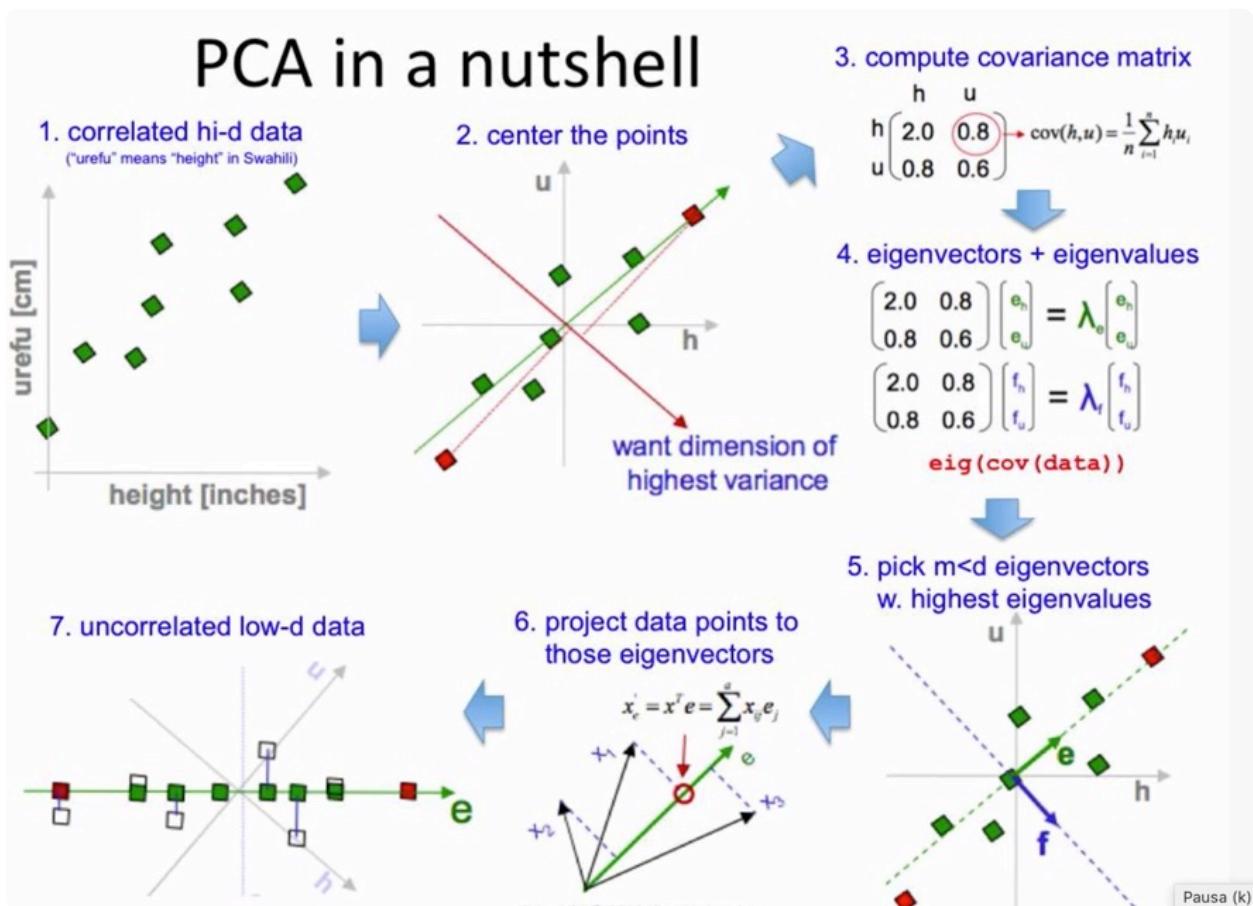
$D = \text{dimension of original dataset}$



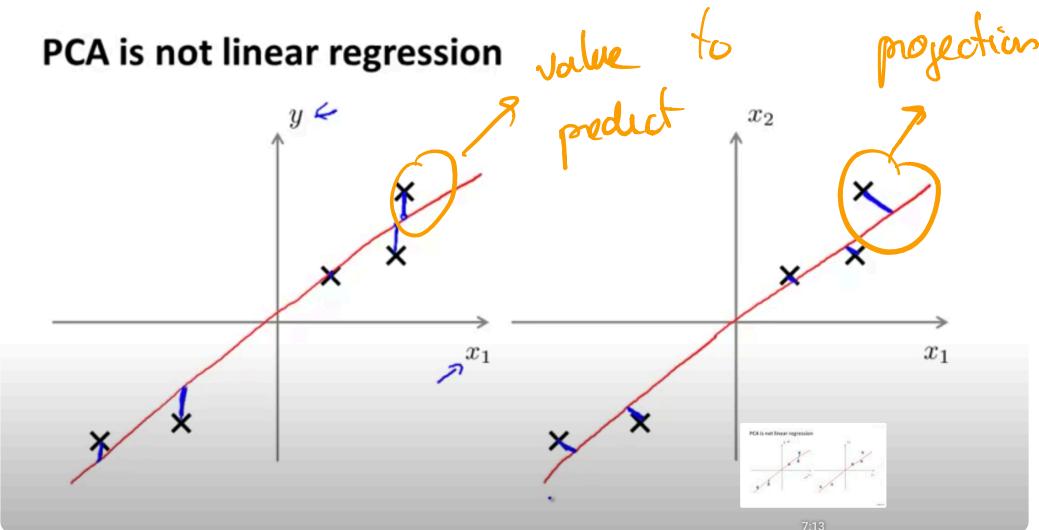
## Summary

### Review / Definition of PCA

- systematic way to transform input features into principal components
- use principal components as new features
- PCs are directions in data that maximize variance (minimize information loss) when you project/compress down onto them
- more variance of data along a PC, higher that PC is ranked
- most variance/most information → first PC  
second - most variance (without overlapping w/ first PC) → second PC
- max no. of PCs = no. of input features



Please, no confuse PCA with linear regression



An interpretation in images:

