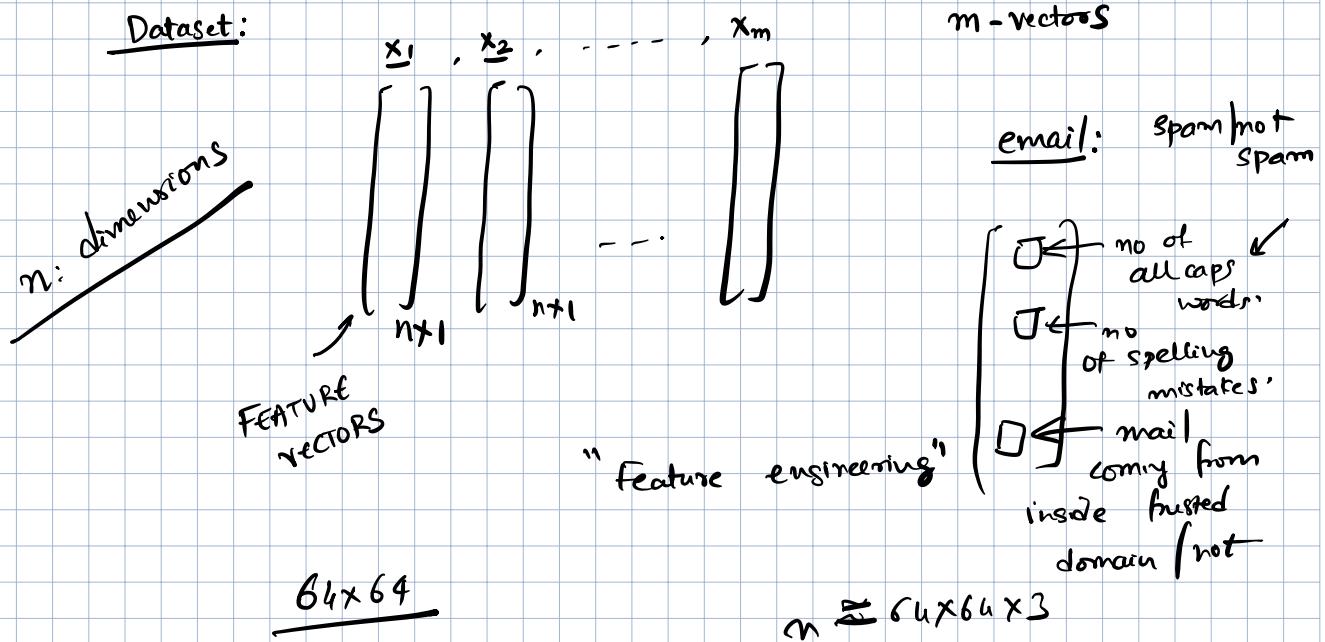


## Dimensionality Reduction

"Dimension"  $\rightarrow$  degrees of freedom.

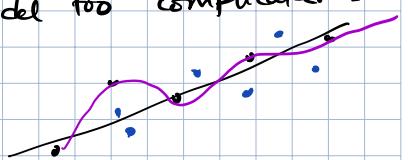
Dataset:



### Why reduce dimensions?

1) Reduced space requirement. Reduce the training complexity, complexity for classification.

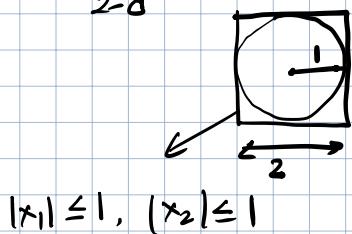
2) Avoiding overfitting  $\rightarrow$  model has too many parameters  
model too complicated -



3) Lower dimensions may have better interpretation / visualization.

4) High dimensional geometry is non-intuitive.

2-d



What fraction of space inside the box  
is taken up by the circle

$$= \frac{\text{area of circle}}{\text{area of box}} = \frac{\pi}{4}$$

3-d

$$\frac{\text{volume of Sphere}}{\text{volume of box}} = \frac{4\pi}{3 \times 8} = \frac{\pi}{6}$$

wikipedia  
volume of  
n-dimensional ball.

n-d

$$\frac{\text{volume of hypersphere}}{\text{volume of hypercube}} = \frac{\text{set of all points with distance 1 from origin}}{\text{set of all points whose co-ordinates are between -1 and 1}} = \frac{\pi^{n/2}}{\left(\frac{n}{2}\right)! 2^n}$$

Integration:

$$(0.9)^n$$

$$\frac{1.8}{2} \approx 0.9$$

Most of the volume in

the hypercube is present on  
the corners, not the center.

on 2-D : How many vectors can I pick such that

they are all

to each other

perpendicular  $\rightarrow 90^\circ$

: How many vectors can I pick such that  
the angles between any two vectors is between  
 $89^\circ$  and  $91^\circ$ . 2

In n-d, there are exponentially many vectors that can be picked while satisfying this).  $\{\text{const}\}^n$

$$\|x_1\|_2 = \|x_2\|_2 = \dots = \|x_m\|_2 = 1$$

$$x_1 \quad x_2 \\ \bullet \quad \bullet \\ \vdots \quad \quad x_3$$

K-NN

given a new data point find the point in the dataset which is closest to the new data point

$$\|x - x_1\|_2^2 = 1 + 1$$

$$- 2 \underbrace{x^T}_{\downarrow} \underline{x_1}$$

In high dim spaces, there may be too many (exponential) points at roughly the same distance -

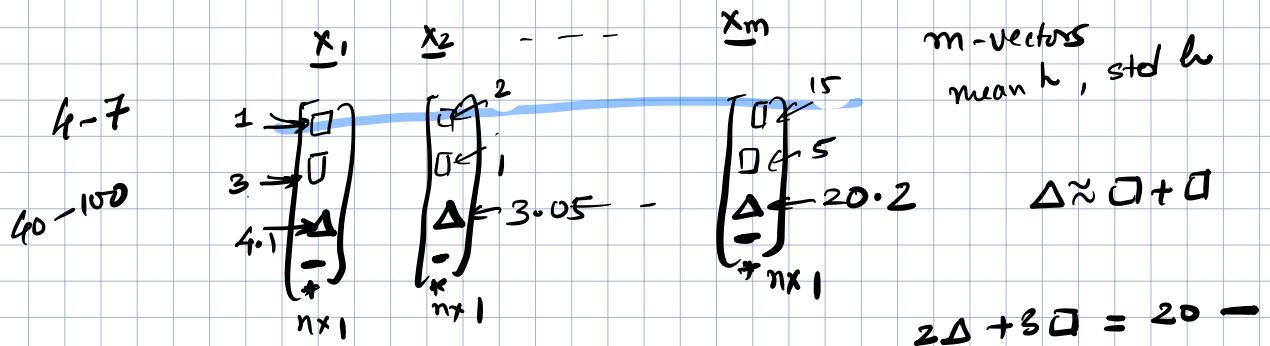
$$= 2 - 2 \cos(\text{angle between } x \text{ and } x_1)$$

$2^\circ$

$1.5 - 2.5^\circ$

s) Understand the dependencies in the engineered features.

Unsupervised learning : PCA principal component analysis ↓ dimensionality reduction technique.

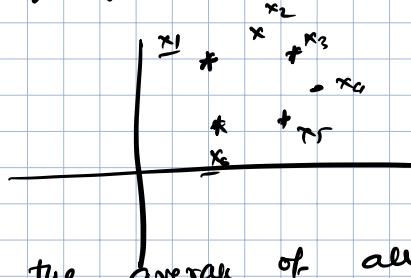


Linear dependence among features.

+ 30 \*

$n=2$

### Feature scaling

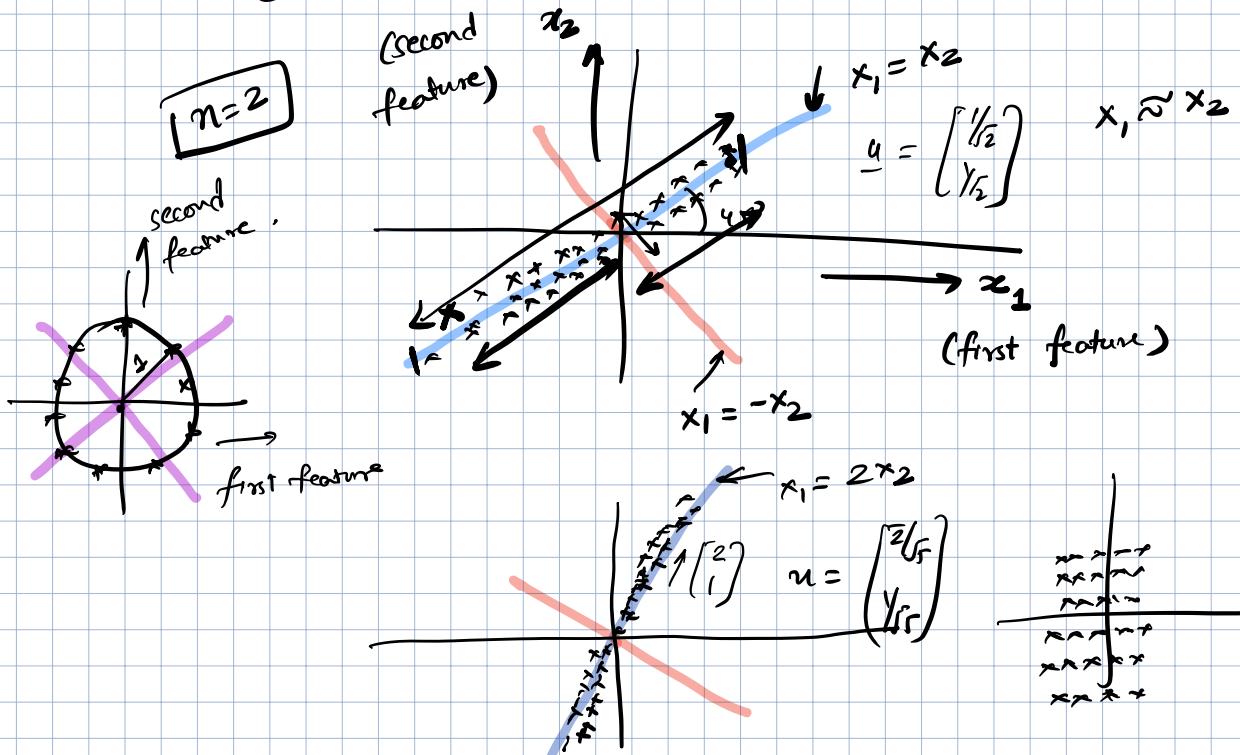


→ ensure that the average of all features is 0  
and variance is 1

$$\text{height} \rightarrow \frac{\text{height - mean}}{\text{std dev}} \quad \text{z-scores}$$

normalized set  $\mathcal{D}$   
mean = 1  
std =  $\sqrt{2}$

$$\begin{bmatrix} -1 \end{bmatrix} \quad \begin{bmatrix} 0 \end{bmatrix} \quad \begin{bmatrix} 2 \end{bmatrix} \quad \begin{bmatrix} 1 \end{bmatrix}$$



Suppose I want to reduce dimensions from  $n$  to 1.

$$\underline{x}_1 \quad \underline{x}_2 \quad \underline{x}_3 \quad \dots \quad \underline{x}_m$$

$$\begin{pmatrix} x_{11} \\ x_{12} \\ \vdots \\ x_{1n} \end{pmatrix} \quad \begin{pmatrix} x_{21} \\ x_{22} \\ \vdots \\ x_{2n} \end{pmatrix} \quad \dots$$

$n \times 1$

give a single number for each data point.

the dimensionality-reduced dataset is given by

$$\underline{x}_1 \downarrow \quad \underline{x}_2 \downarrow$$

$$\underbrace{u_1 x_{11} + u_2 x_{12} + \dots + u_n x_{1n}}_{\underline{u}^T \underline{x}_1} \quad \underbrace{u_1 x_{21} + u_2 x_{22} + \dots + u_n x_{2n}}_{\underline{u}^T \underline{x}_2} \quad \dots$$

$$\underline{u}^T \underline{x}_3 \dots$$

Pick  $\underline{u}$  such that

$$\rightarrow \underline{u}^T \underline{x}_1, \underline{u}^T \underline{x}_2, \dots, \underline{u}^T \underline{x}_m$$

$\curvearrowleft$  have the highest variance possible

$$\frac{1}{m} \left( \underline{u}^T \underline{x}_1 + \underline{u}^T \underline{x}_2 + \dots + \underline{u}^T \underline{x}_m \right)$$

$$= \underline{u}^T \left( \frac{\underline{x}_1 + \underline{x}_2 + \dots + \underline{x}_m}{m} \right)$$

$$\text{Variance}(\underline{u}) = \frac{1}{m} \left[ (\underline{u}^T \underline{x}_1)^2 + (\underline{u}^T \underline{x}_2)^2 + \dots + (\underline{u}^T \underline{x}_m)^2 \right] = 0$$

$$\frac{1}{m} \left[ \underline{u}^T (\underline{x}_1 \underline{x}_1^T) \underline{u} + \underline{u}^T \underline{x}_2 \underline{x}_2^T \underline{u} + \dots + \underline{u}^T \underline{x}_m \underline{x}_m^T \underline{u} \right]$$

$\underline{x}_i : n \times 1$

$n \times n$  matrix

$$= u^T \left[ \frac{x_1 x_1^T + x_2 x_2^T + \dots + x_m x_m^T}{m} \right] u$$

$u$   
 $n \times 1$

sum of squares of all heights  
covariance between different features  
 $n \times n$   
sum of squares of all weights

(COVARIANCE MATRIX)

In the  $n \times n$  matrix  $\Sigma$

- 1) All diagonal entries are 1
- 2) The off diagonal entries are covariances between different features.

How to pick a unit vector  $u$  such that  $u^T \Sigma u$  is the largest possible?

$$u = \left[ \frac{1}{\sqrt{2}}, \frac{1}{\sqrt{2}} \right]$$

6 pm

What if  $\Sigma$  were a diagonal matrix?

$n=2$

Suppose,  $\Sigma = \begin{bmatrix} 2 & 0 \\ 0 & 5 \end{bmatrix}$   $\underline{u} = \begin{bmatrix} u_1 \\ u_2 \end{bmatrix}$

$$\underline{u}^T \Sigma \underline{u} = [u_1 \ u_2] \begin{bmatrix} 2 & 0 \\ 0 & 5 \end{bmatrix} \begin{bmatrix} u_1 \\ u_2 \end{bmatrix}$$

$$= [u_1 \ u_2] \begin{bmatrix} 2u_1 \\ 5u_2 \end{bmatrix} = 2u_1^2 + 5u_2^2$$

$$u_1^2 + u_2^2 = 1$$

we should pick

$$\underline{u} = \begin{bmatrix} 0 \\ 1 \end{bmatrix} \quad u_2 = 1, u_1 = 0.$$

Suppose  $\Sigma$  is a diagonal matrix with diagonal entries  $d_1, d_2, d_3, \dots$

then  $\underline{u} = \begin{bmatrix} u_1 \\ u_2 \\ u_3 \\ \vdots \end{bmatrix}$

$$\underline{u}^T \Sigma \underline{u} = d_1 u_1^2 + d_2 u_2^2 + d_3 u_3^2 + \dots$$

$$u_1^2 + u_2^2 + u_3^2 + \dots = 1$$

$$\begin{bmatrix} 2 & 0 & 0 \\ 0 & 10 & 0 \\ 0 & 0 & 10 \end{bmatrix}$$

$$\underline{u} = \begin{bmatrix} 0 \\ 0 \\ 1 \end{bmatrix}$$

what if the matrix  $\Sigma$  is not diagonal?

$$\Sigma = V D V^T$$

symmetric       $\nwarrow$  orthogonal  
 diagonal

$$\underline{u}^T \Sigma \underline{u}$$

$$\underline{u}^T V D V^T \underline{u}$$

$$\leftarrow \boxed{V} \leftarrow \boxed{D} \leftarrow \boxed{V^T} \leftarrow$$

$$\underline{w}^T \underline{w}$$

$$\underline{w} = \underline{V}^T \underline{u}$$

$\underline{w}$ :  $w$ -ordinates  
after rotation by  
 $\underline{V}$ .

$$\underline{w}^T = \underline{u}^T \underline{V}$$

$$\sqrt{\underline{w}^T \underline{w}} =$$

$$\underline{w} = \begin{bmatrix} w_1 \\ w_2 \\ \vdots \\ w_n \end{bmatrix}$$

$$\begin{aligned} w_1^2 + w_2^2 + \dots + w_n^2 \\ = u_1^2 + u_2^2 + \dots + u_n^2 \\ = 1 \end{aligned}$$

$$\text{Variance}(\underline{u}) = \underline{u}^T \Sigma \underline{u}$$

$$= \underline{w}^T D \underline{w}$$

$$\boxed{\underline{w} = \underline{V}^T \underline{u}} \rightarrow \boxed{\underline{u} = \underline{V}^{-1} \underline{w}}$$

$\underline{w}$  is also a unit vector.

For simplicity, assume that the Diagonal entries in  $D$

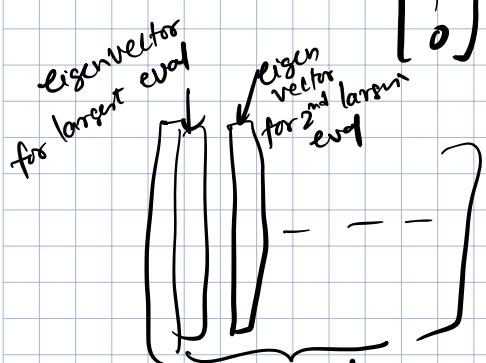
are in sorted order

$$\begin{bmatrix} 10 \\ 5 \\ 3 \\ 2 \end{bmatrix}$$

$$\underline{w} = \begin{bmatrix} 1 \\ 0 \\ \vdots \\ 0 \end{bmatrix}$$

$$\underline{u} = \underbrace{\begin{bmatrix} \cos\theta & \sin\theta \\ -\sin\theta & \cos\theta \end{bmatrix}}_{\sqrt{}} \begin{bmatrix} 1 \\ 0 \end{bmatrix}$$

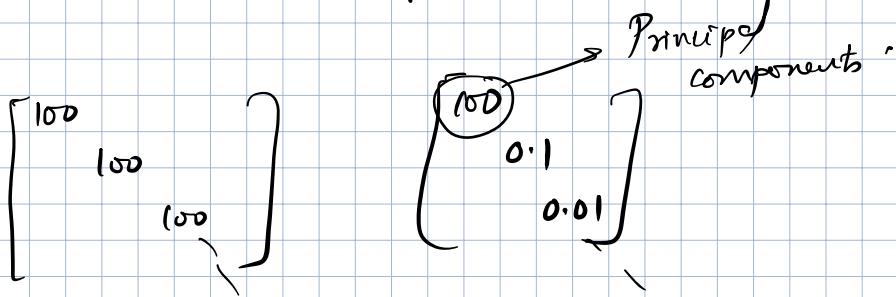
$$= \begin{bmatrix} \cos\theta \\ -\sin\theta \end{bmatrix} \quad \text{first column of the matrix } \underline{V}.$$



We should pick  $\underline{u}$  to be the first column of the (sorted) eigen vector matrix  $\underline{V}$ .

$\underline{u}$  = eigenvector corresponding to largest eigenvalue of  $\Sigma$ .

How good will this dimensionality reduction be?



Principal components are the directions along which the data is "most expressive" or "highly informative" or "most varying".

→ obtained by taking eigenvectors corresponding to the top eigenvalues of the covariance matrix,

### Algorithm (PCA: Principal component analysis)

Input: Total  $m$  vectors,  $\underline{x}_1, \underline{x}_2, \dots, \underline{x}_m$  with each vector having  $n$ -coordinates.

$K$  : target dimension

( $K \ll n$ )

① feature scaling

② find the covariance matrix  $\sum = \frac{1}{m} \underline{x}_1 \underline{x}_1^T + \underline{x}_2 \underline{x}_2^T + \dots$

③ find the eigen vectors for the  $K$ -largest eigenvalues of  $\sum$ .

$\underline{u}_1 \underline{u}_2, \dots, \underline{u}_K$

④ the output (dimensionality reduced dataset) is given by

$x_1 \quad n \times 1 \quad x_2 \quad \dots \quad x_m$

$$Kx_1 = \begin{bmatrix} u_1^T x_1 \\ u_2^T x_1 \\ u_3^T x_1 \\ \vdots \\ u_k^T x_1 \end{bmatrix}$$

$$= \begin{bmatrix} u_1^T x_2 \\ u_2^T x_2 \\ u_3^T x_2 \\ \vdots \\ u_k^T x_2 \end{bmatrix}$$

$$\vdots$$

$$= \begin{bmatrix} u_1^T x_m \\ u_2^T x_m \\ u_3^T x_m \\ \vdots \\ u_k^T x_m \end{bmatrix}$$

Sort all evals  
 $\lambda_1 > \lambda_2 > \dots > \lambda_n$

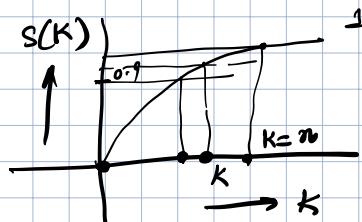
Construct

$$s(k) = \frac{\sum_{i=1}^k \lambda_i}{\sum_{i=1}^n \lambda_i}$$

$\uparrow$  all evals.

$$\begin{bmatrix} 100 \\ 50 \\ 10 \\ 0.001 \\ 0.0001 \end{bmatrix}$$

$s(k)$  to be close to 1.



$$\begin{bmatrix} 10 \\ 10 \\ 1.0 \\ 1.0 \end{bmatrix}$$

① Sthem Alpaydin

Intro to machine learning

(entire chapter on dimensionality reduction)

② Bishop ,  
PRML

Murphy  
 Probabilistic perspective  
 Machine learning ,

Convex  
optimization : KKT conditions

--

$f(x)$

Duality  
Gradient  
descent

$$u \max_{\underline{u}} \underbrace{\underline{u}^\top \Sigma \underline{u}}_{f(\underline{u})}$$

$$\sum_{i=1}^n (y_i - w^\top x_i)^2$$

Textbook:

Convex Optimization

Stephen Boyd

email:

stachitza@ee.  
iith.ac.in