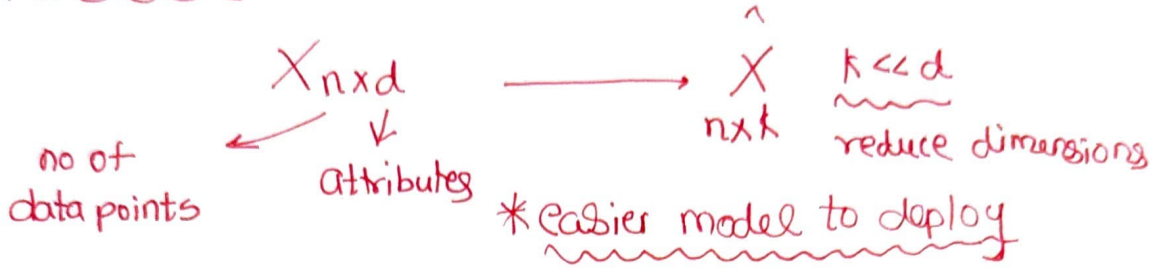


*Dimensionality Reduction

OBJECTIVES

- 1) Preserve in some sense, the original data
 - 2) Smooth out the data for better generalization
 - 3) Smaller dimension results in faster models.
- } overfitting



*Covariance matrices

attributes of dataset $\left\{ \begin{array}{l} a_1 = [a_{11} \dots a_{1n}] \\ a_2 = [a_{21} \dots a_{2n}] \end{array} \right\}$

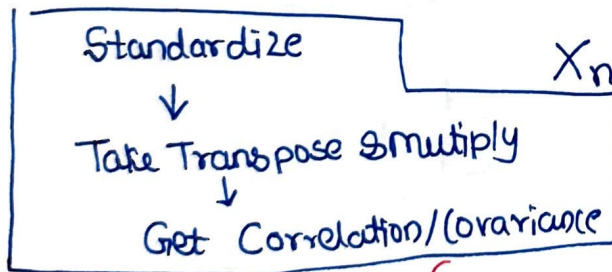
covariance between these 2 vectors $E(a_1)$ and $E(a_2)$ $\xrightarrow{\text{calculate average of both}}$

Covariance $(a_1, a_2) = \frac{((a_1 - E[a_1])(a_2 - E[a_2])^T)}{n}$

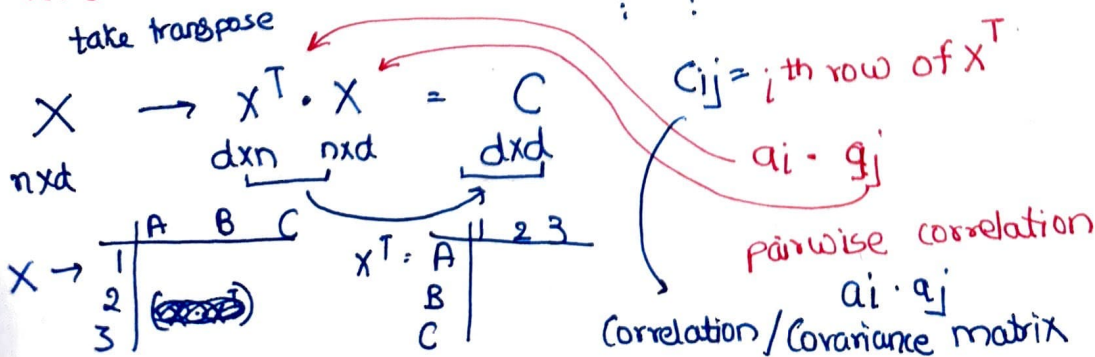
\downarrow measures correlation

$1 \times n$ $n \times 1$

$a_1 \quad a_2 \quad \dots$



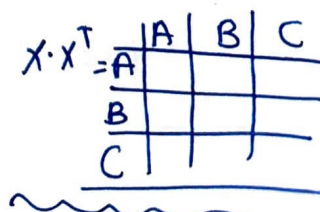
$X_{n \times d}$ $\xrightarrow{\text{standardization}}$ $\left\{ \begin{array}{l} \text{average of each attribute is zero} \\ a_1 - E[a_1] \\ a_2 - E[a_2] \\ \vdots \end{array} \right\}$ standardization



If $a_1 = a_2$
 $a_1 = -a_2$
 or
 $a_2 = -a_1$

a_1 temp C°
 a_2 temp F°

high correlation
 high covariance



Eigen Values

* What is the action of a matrix to a vector?

$$A = \begin{bmatrix} 2 & -1 \\ -1 & 2 \end{bmatrix} \begin{matrix} \text{Input} \\ \begin{bmatrix} 1 \\ 0 \end{bmatrix} \end{matrix} = \begin{matrix} \text{Input} \\ \begin{bmatrix} 2 \\ -1 \end{bmatrix} \end{matrix}$$

w_1

$$A = \begin{bmatrix} 2 & -1 \\ -1 & 2 \end{bmatrix} \begin{matrix} \text{Input} \\ \begin{bmatrix} 1 \\ 1 \end{bmatrix} \end{matrix} = \begin{matrix} \text{Output} \\ \begin{bmatrix} 1 \\ 1 \end{bmatrix} \end{matrix} \cdot I$$

$$A = \begin{bmatrix} 2 & -1 \\ -1 & 2 \end{bmatrix} \begin{matrix} \text{no change} \\ \begin{bmatrix} 1 \\ -1 \end{bmatrix} \end{matrix} = \begin{matrix} \begin{bmatrix} 3 \\ -3 \end{bmatrix} \end{matrix} = \begin{bmatrix} 1 \\ -1 \end{bmatrix} \cdot 3$$

Input output break into

Same direction
only length changes

$\frac{2 \times 2}{\text{matrix}}$ vector $s \rightarrow$ 2 eigen vectors
2 eigen values.

* A property of eigen values

$$\begin{bmatrix} 1 & 1 \end{bmatrix} \begin{bmatrix} 1 \\ -1 \end{bmatrix} = 0 \quad \left\{ \begin{matrix} w_1^T \cdot w_2 = 0 \end{matrix} \right\} \text{orthogonal}$$

$$* \lambda_1^2 + \lambda_2^2 = 1 + 9 = 10$$

$$A = \begin{bmatrix} 2 & -1 \\ -1 & 2 \end{bmatrix}$$

Square each entry
 $= 2^2 + (-1)^2 + (-1)^2 + 2^2$
 $= 4 + 4 + 1 + 1$
 $= 10$

* Generalization

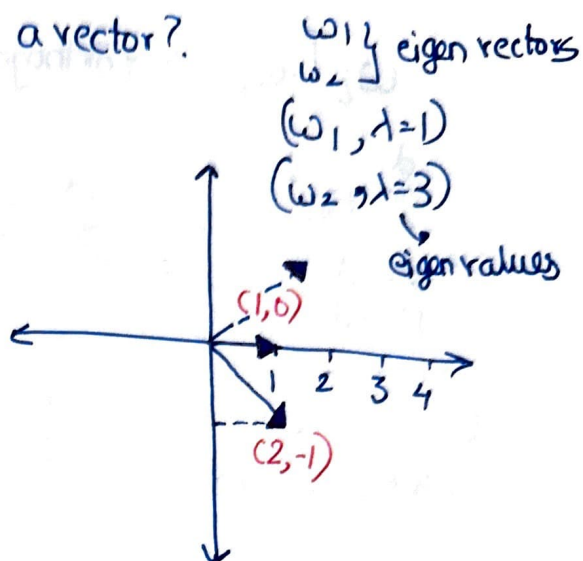
eigen vectors

A
 $d \times d$

w_1, w_2, \dots, w_d : not changing direction

$\lambda_1, \lambda_2, \dots, \lambda_d$: tell us how much vector stretch
eigen values

$$A \underline{w_i} = \lambda_i \underline{w_i}$$



$w_i^T \cdot w_j = 0$ } orthogonal for all $i \neq j$

$$\sum_{\substack{i=1 \dots d \\ j=1 \dots d}} A_{ij}^2$$

Sum of all
entries Square

$$= \sum_{j=1}^d \lambda_j^2$$

Sum of squares
of eigen values

} Frobenius
norm
of matrix

Dimensionality Reduction - the notions

Principal Component

vectors

$$V_1 = PC1$$

$$V_2 = PC2$$

$$V_3 = PC3$$

\times
 $n \times 3$

V_1 & V_2 orthogonal

$V_1^T V_2 = 0$, $V_2^T V_3 = 0$, $V_1^T V_3 = 0$
mutually orthogonal vectors.

$$C_1 V_1 + C_2 V_2 + C_3 V_3$$

$W =$

can be
in our
dataset
 \times

any
point w /vector
represented
as
sum of these PC

Say w has less variance
with V_3

$$\text{So } C_1 V_1 + C_2 V_2 + C_3 V_3$$

\downarrow
Small
(noise)

transformation

$$\begin{matrix} \times & \rightsquigarrow & \hat{\times} \\ n \times 3 & & n \times 2 \end{matrix}$$

every point
 w

point in new dataset

$$[x_1, y_1, z_1]$$

$$[C_1, C_2]$$

take every
point w

transform

and only keep 2 components (high variance components)

$$\begin{matrix} \hat{\times} \\ n \times k \end{matrix} = \begin{matrix} \times & \cdot & W \\ n \times d & & d \times k \end{matrix} \left. \begin{array}{l} \text{Find how to} \\ \text{calculate } w \\ \text{linear} \\ \text{dimensionality} \\ \text{reduction} \end{array} \right\}$$

C_1 & $C_2 \rightarrow$ capture the variance
information most

(kept)

$C_3 \rightarrow$ dropped after transformation

Remember

→ Correlation

gives correlation / variance between features



cannot directly use it

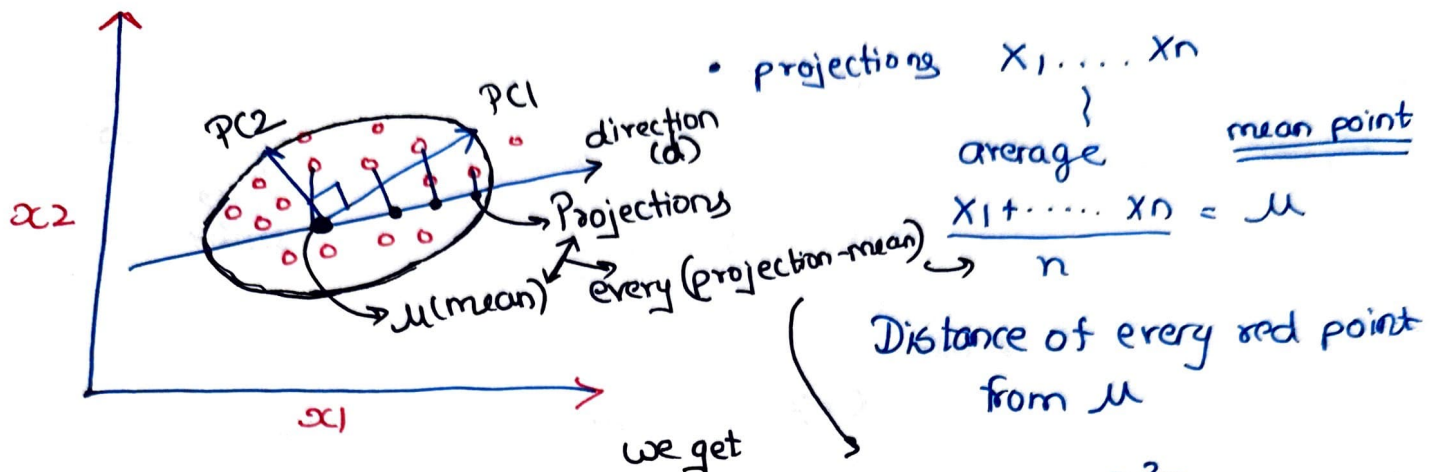
Convert/
decompose

we use eigen vector / values → multiply with data points



find how much scaling
the features need / do
to the data.

* Dimensionality Reduction - Math Intuition :-



How to compute PCA

Variance $= \frac{(x_2 - \mu)^2}{n}$

Scalar value

$CW = \lambda W$

$(C - \lambda I)W = 0$

$C - \lambda I = 0$

Solve for $\lambda \rightarrow$ then substitute back

get W, λ both

Look for direction that maximized variance

max α

Var $d = PC1$

Look for direction

How to compute PCA

standardize

$X \sim C = X^T X$

$n \times d$ $d \times d$ correlation

$Cw_i = \lambda_i w_i$

$\begin{bmatrix} w_{i1}, \lambda_{i1} \\ \vdots \\ w_{id}, \lambda_{id} \end{bmatrix}$

$\lambda_1 > \lambda_2 > \lambda_3 \dots > \lambda_d$

$\lambda_1^2 \rightarrow$ variance of our data

$\hat{X} = X \cdot W = X \cdot \begin{pmatrix} w_1 | w_2 | \dots | w_k \end{pmatrix}$

$d \times k$

What dimension to reduce to? k -eigen vectors $d \times k$

\rightarrow Sum of first k eigen values

Total Variance $= \|C\|_F = \sum C_{ii}^2 = \lambda_1^2 + \lambda_2^2 + \dots + \lambda_d^2$

Stop at λ_k^2

90% of Variance \rightarrow Drop other eigen vector, value

Kernel PCA

(non linear dimensionality reduction)

Covariance matrix

$$\begin{aligned} C_{d \times d} \omega_i &= \lambda_i \cdot \omega_i \Rightarrow X^T \cdot X \cdot \omega_i = \lambda_i \omega_i \quad \text{multiply by } X \\ &\Rightarrow X \cdot X^T \cdot X \omega_i = \lambda \cdot X \cdot \omega_i \quad \left. \begin{array}{l} \text{matrix} \\ n \times n \end{array} \right\} \begin{array}{l} \text{not} \\ \text{efficient} \end{array} \end{aligned}$$

\downarrow
eigen values

$$X \cdot X^T \cdot Z_i = \lambda \cdot Z_i \quad \left. \begin{array}{l} \lambda_1 \quad Z_1 \\ \lambda_2 \quad Z_2 \\ \vdots \quad \vdots \end{array} \right\} Z_i = X \omega_i$$

$$\hat{X} = [Z_1 | Z_2 | \dots | Z_k] : \text{identical } X \cdot \omega$$

Introducing the kernel

$$K = X \cdot X^T \quad \left| \quad K_{ij} = \underbrace{X_i \cdot X_j}_{\substack{\text{i-th row} \quad \text{j-th row}}} \right.$$

Similar Correlation matrix to

$$K_{i,j} = \underbrace{\langle X_i, X_j \rangle}_{\text{linear}} : \text{Similarity}$$

$$K_{ij} = \exp(-\|X_i - X_j\|_2) \quad \text{RBF kernel} \quad \left. \begin{array}{l} \text{not explicitly} \\ \text{high dimension} \\ \text{similarity} \end{array} \right\}$$

non linear reduction \rightarrow linear d... reduction

\rightarrow Similarities measured in some space

kernels differentiates features based on
similarity / by calculating
distances in 2D, 3D... spaces

* Kernel Trick