Covariance /correlation matrix

⟶ standardize data

⟶ create cov matrix.

Covariance: Dataset with $m$ samples, $n$ features

$$cov(x_a, x_b) = \frac{1}{m} \sum_{i=1}^{m} (z_{i,a} - \bar{q}_a)(z_{i,b} - \bar{q}_b)$$

10 features ⟹ cov matrix: $10 \times 10$

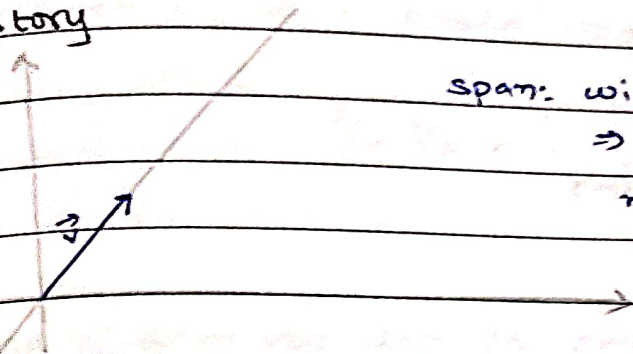$X X^T$ ⟹ No. of samples × No. of samples.

$$X = \begin{bmatrix} & x_1 & x_2 & x_3 & \cdots & x_n \\ s_1 & z_{11} & z_{12} & z_{13} & & z_{1n} \\ s_2 & z_{21} & & & & \\ s_3 & & & & & \\ \vdots & & & & & \\ s_m & z_{m1} & & & & z_{mn} \end{bmatrix}$$

$$C = \begin{bmatrix} cov(x_1, x_1) & cov(x_1, x_2) & \cdots & cov(x_1, x_n) \\ cov(x_2, x_1) & cov(x_2 x_2) & \cdots & cov(x_2, x_n) \\ \vdots & & & \\ cov(x_n, x_1) & cov(x_n, x_2) & \cdots & cov(x_n, x_n) \end{bmatrix} = \frac{1}{m} X^T X$$

n-dim Cov Matrix

Assume the
matrix is cente...

# Eigen - Story

span: will be on the same line
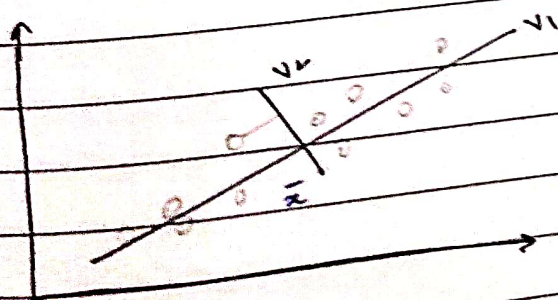⇒ same dir"

multiply with any number.

What happens with a matrix operation on the vector?

$$A\vec{v} = \begin{bmatrix} 3 & 2 \\ 1 & 3 \end{bmatrix} \begin{bmatrix} 2 \\ 1 \end{bmatrix} = \begin{bmatrix} 8 \\ 5 \end{bmatrix}$$

- The vectors change their directions most of the time [the span changes]
- we have special cases where the vectors stay on the span [only getting scaled].

$$A\vec{v} = \begin{bmatrix} 1 & 3 \\ 3 & 1 \end{bmatrix} \begin{bmatrix} 1 \\ 1 \end{bmatrix} = 4 \begin{bmatrix} 1 \\ 1 \end{bmatrix}$$

- The vector here is the eigenvector & the scaling 4 is the eigenvalue.

## DEMO - PCA

$v_1$ → Most distribution and spread

→ consider the variation along a $dir^n$ v among all the points

$$var(v) = \frac{1}{n} \sum_{x \ points} |(x - \bar{x})^T \cdot v|^2$$

NOTE: The E. vectors of the cov matrix give you the direction that maximizes the variance. The direction of the green line is where the variance is max. Compare that with the projection on the $V_2$: the spread is small

→ Unit vector $v$ maximizes var:

$$V_1 = max_v \{var(v)\}$$

PCA :-

large E-value capturing alot of variance.

most of the sample variance.

How many PC's
* -A dataset with $m$ samples & $n$ features will give
  to a $n \times n$ cov matrix.

- The $n \times n$ covariance matrix will have $n$ e-vectors,
  so $n$ PC's.

  $n$ features $\longrightarrow$ $n$ Principal components.

Take dataset $\rightarrow$ Standardize features
numpy commands E.vectors/value $\rightarrow$ chose how much value/  im
plot scree plot $\rightarrow$ 90% data $\Rightarrow$ 7/8 P.C's.    variance $\rightarrow$

where does dim reduction comes from?
$\hookrightarrow$ can always ignore the components of lesser significance.

$\rightarrow$ we do lose some information, but if the E. values
  are small

Step by Step Computation      numpy commands
         e.values
         e. vectors

steps (1) standardization of the data

(2) Compute the Covariance matrix

(3) Calculate the E.values & vectors of cov matrix

(4) Compute the principal components by selecting the first D E.vectors.

(5) Reduces the dimensions of the dataset.

The features in which data is most spread out.

PCA → UNSUPERVISED. LEARNING

final prediction result ⇒ we only concentrate on " FEATURES "

PCA works best when → good amount of spread

→ relation among features

# SINGULAR VALUE DECOMPOSITION

→ SVD gives the decomposition for any arbitrary matrix,

$$M = U \Lambda V^T$$

$$M_{m \times n} = U_{m \times r} \; \Lambda_{r \times r} \; V_{r \times n}^T$$

$$\downarrow$$

diagonal matrix = Root of +ve e.values of M [$x^T x$ or $x x^T$]

- U & V → orthogonal matrices, $U^T U = 1$ ; $V^T V = 1$
- U consists of orthonormal eigen vectors of M [$x x^T$].
- V consists of orthonormal eigenvectors of $M^T$ [$x^T x$].

$$Cov(M) = X$$

→ The SVD of the data matrix, $X = U \Lambda V^T$

→ After standardization, the cov matrix of the data matrix, $\Sigma = \frac{1}{m} X^T X$.

$$\Sigma = \frac{1}{m} X^T X = \frac{1}{m} (U \Lambda V^T)^T (U \Lambda V^T) = \frac{1}{m} (V \Lambda^T U^T)(U \Lambda V^T)$$

$$= \frac{1}{m} \left( V \Lambda^T \Lambda V^T \right) = \frac{1}{m} \left( V (\Lambda)^2 V^T \right)$$

→ $(\Lambda)^2$ is a diagonal matrix whose entries are $\Lambda_{ii} = \lambda_i^2$, the squares of the E·values of the SVD of X.

→ we can run SVD on X without ever

→ Both X and $X^T X$ share the same e·vectors in their SVD.

For ex: Given an image of a woman; where given no. of features = 200

at PC = 0       Just black & white lines
PC = 10       blurry
    ⋮
PC = 50       Almost Similar

PC → capturing most of the variances

## SUMMARY

→ PCA allows us to find the highest variance (lowest square distance) direction to project to.

→ E-values gives an indication of the number of dimensions to choose.

→ Can be computed in multiple ways (SVD is popular)

→ It is an unsupervised algo.

→ Ensures pre-processing for effectiveness

→ Is used in a variety of applications.