

## PCA :

- Form the basis of many algs.
- Technique for dimensionality reduction.
- Compared to the "JL-Lemma" (random projection):
  - Randomly picked vectors  $r_1, \dots, r_k$
  - projected  $X$  on  $r : \langle x_1, r_1 \rangle, \dots, \langle x_1, r_k \rangle$
  - $r_i$ 's were not meaningful wrt  $S$ ; picked randomly.
  - preserved the Euclidean dist btwn points.
  - In practice,  $k > 100$  for random projections to work.
  - For PCA, we can still choose  $k=2$ .
  - PCA looks at  $S$  to come up with new representation
- High level goal of PCA is to find vectors  $v_1, \dots, v_k$  st.  $\forall x \in S, x \approx \sum_{j=1}^k a_j v_j$  approx.
- Note about preprocessing of  $S$ : Normalize each feature.
  - subtract the mean of each feature from each data point;
  - normalize by the standard deviation of each feature, bc want features to be similar range,

To do so:  $\forall$  feature  $j$ , compute:

$$\sigma_j (\text{STD of the } j^{\text{th}} \text{ feature}) = \sqrt{\frac{1}{m} \sum_{i=1}^m (x_j^{(i)})^2}$$

Then Divide all the  $j^{\text{th}}$  feature by  $\sigma_j$ .

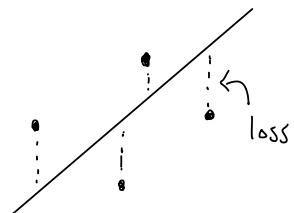
Summary :  $x_j^{(i)} \leftarrow \frac{x_j^{(i)} - \mu_j}{\sigma_j}, \quad \mu_j = \frac{1}{m} \sum_{i=1}^m x_j^{(i)}$   
 $\sigma_j^2 = \frac{1}{m} (x_j^{(i)} - \mu_j)^2$

## How to begin :

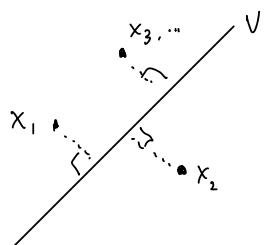
1) Find  $\vec{v}_1$  (first vector)

Look for a vector that minimizes square-distance:

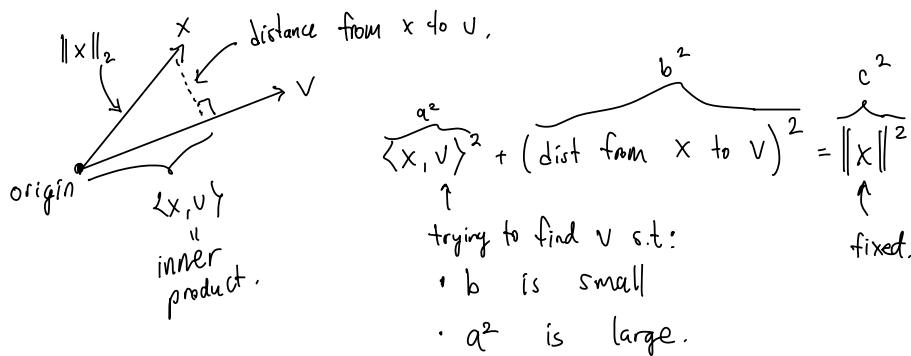
$$\min_{\mathbf{v}, \|\mathbf{v}\|_2 = 1} \frac{1}{m} \sum_{j=1}^m \left( \text{distance btwn } x_j \text{ and } \mathbf{v} \right)^2$$



## Regression



PCA



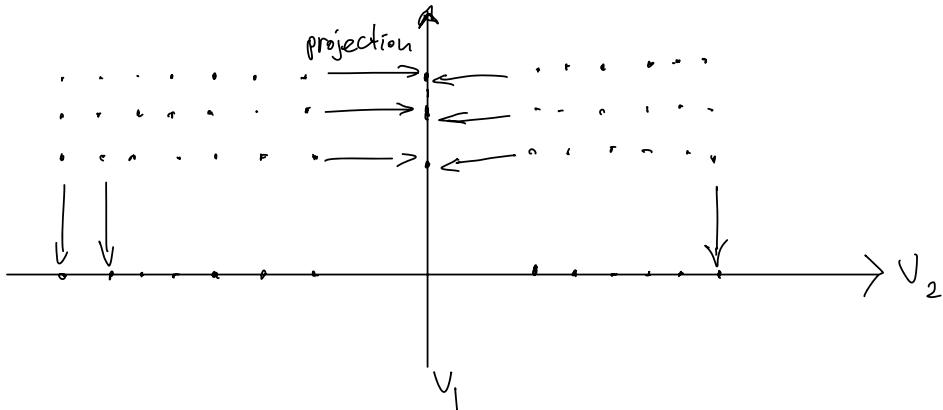
equivalently: want  $v$  that maximizes  $\langle x, v \rangle^2$ ,

Why? :  $b^2 = c^2 - a^2$ .  $c^2$  is fixed, so if  $\max a^2 = \langle x, v \rangle^2 \Rightarrow \min b^2$ .

$$\text{Find: } \underset{\|v\|_2 = 1}{V} \max \frac{1}{m} \sum_{i=1}^m \underbrace{\langle x_i, v \rangle}_\text{dot product}^2 \quad \leftarrow \text{direction of maximal variance}$$

Remember: Mean of the training set is 0, bc we subtracted the mean out of it.

$$\text{var}[x] = E[x^2] - \overline{[E[x]]}^2$$



which projection preserves the structure of the training set?  
 $v_2!$

The above was for 1 vector. What about for  $k$ -vectors?

Max subspaces  $S$  of dimension  $K$   $\left[ \frac{1}{m} \sum_{j=1}^m (\text{length of } x^j \text{ projected onto } S)^2 \right]$

A really nice/preferable basis would be an orthonormal basis,  $v_1, \dots, v_K$ .

Note: Orthonormal basis = vectors are unit + orthogonal to one another.

(distance from  $x$  to  $S$ ) $^2 = \|x\|^2 - \underbrace{\langle x_1, v_1 \rangle^2}_{c^2} - \underbrace{\dots + \langle x_k, v_k \rangle^2}_{a^2}$

PCA Objective:

$$\boxed{\begin{aligned} \text{MAX}_{\substack{v_1, \dots, v_k \\ \text{orthogonal}}} \quad & \frac{1}{m} \sum_{j=1}^m \sum_{i=1}^k \langle x^j, v_i \rangle^2 \\ & = \langle x^1 \cdot v_1 \rangle^2 + \langle x^1 \cdot v_2 \rangle^2 + \dots + \langle x^1 \cdot v_k \rangle^2 \\ & + \\ & \langle x^m \cdot v_1 \rangle^2 + \langle x^m \cdot v_2 \rangle^2 + \dots + \langle x^m \cdot v_k \rangle^2 \end{aligned}}$$

want orthogonal basis that:

- Maximizes the variance along these directions that are orthogonal.

Assume we have  $v_1, \dots, v_k$ , how to rewrite  $x$ ?

$$x = \langle x, v_1 \rangle \cdot v_1 + \langle x, v_2 \rangle \cdot v_2 + \dots + \langle x, v_k \rangle \cdot v_k$$
$$= \sum \text{projection onto the basis vectors.}$$

$\Rightarrow x$  can be written as a vector in  $R^k$  corresponding to these projections.

### Application 1: Understanding genomes:

Took 1400 people from Europe.

Each person was represented according to 200,000 genome markers in their genome.

Each person = vector of 200,000 features.

This corresponds to matrix of  $1400 \times 200,000$ .

- Ran PCA on this data to find vectors  $v_1$  and  $v_2$

- Each person corresponds to 2 numbers,

They plotted these 2 numbers; color code each point according to country of origin.

### Application 2:

Image Data Compression,

Strategy for compression data:

Each data point is an image (vector of pixels)

Each image has 65,000 pixels (65,000 features),

Images of faces. Run PCA on dataset.

$K = 100$  to  $150$ ,

Image  $\approx$  linear combination of 150 vectors of length 65,000

## The Big Question:

How do we find these  $v_1, \dots, v_k$ ?

$$\max_{\substack{v_1, \dots, v_k \\ \text{orthonormal}}} \left[ \frac{1}{m} \sum_{j=1}^m \sum_{i=1}^k \langle x_j, v_i \rangle^2 \right]$$

Let  $X$  be an  $m$  by  $n$  matrix.

$v \leftarrow$  a column vector  $= \begin{bmatrix} \vdots \\ i \end{bmatrix}$

$v^T \leftarrow$  a row vector  $= [\dots \dots]$ .

$v^T v \leftarrow$  inner product (scalar):

$$[\dots \dots]^T \begin{bmatrix} 1 \\ \vdots \\ 1 \end{bmatrix}_k = \begin{bmatrix} 1 \\ \vdots \\ 1 \end{bmatrix}$$

$VV^T \leftarrow$  outer product (matrix).

$$k \begin{bmatrix} \vdots \\ \vdots \end{bmatrix} [\dots \dots]_k = k \begin{bmatrix} \vdots \\ \vdots \end{bmatrix}$$

$\frac{1}{m} X^T X$  (n by n matrix)

sample covariance matrix.

$(i, j)^{\text{th}}$  entry of  $X^T X$

corresponds to "how similar is feature  $i$  to feature  $j$ ?"

Note that  $X^T X$  is a symmetric matrix. Proof:  $X^T X$ .

$(i, j) =$  inner product of row  $i$  of  $X^T$  with column  $j$  of  $X$ .  
 $=$  " " of column  $i$  of  $X$  with column  $j$  of  $X$ .  
 $=$  " " of column  $j$  of  $X$  with " "  $i$  of  $X$ .

fact: all eigenvalues of symmetric matrices is  $\geq 0$ .

For matrix  $A$ , vector  $v$  is an

eigenvector if  $A \cdot v = \lambda v$ ,  $\lambda \in \mathbb{R}$ .  
 $\swarrow$   
eigenvalue.

Definition: An **orthogonal** matrix is one where all columns are **orthonormal**; square matrix  $A$  is orthogonal if:

$$A^T A = I \text{, and } A A^T = I$$

ex:  $\begin{bmatrix} 1 & 1 & -1 \\ 1 & 2 & 0 \\ -1 & 0 & 5 \end{bmatrix}$

Spectral theorem: Every symmetric matrix  $A$  has an

**eigen decomposition**:  $A = Q \cdot D \cdot Q^T$

$$D = \begin{bmatrix} \lambda_1 & & \\ & \lambda_2 & \\ & & \ddots & \lambda_n \end{bmatrix}$$

↑  
orthogonal  
matrix      ↑  
diagonal  
matrix

entries of  $D$  are the eigenvalues of  $A$ .

Let's try to compute  $U_1$ .

Recall  $X$  is the matrix corresponding to  $S$  (training set).

$X$  is  $m$  by  $n$ .

$$X \cdot v = \begin{bmatrix} & n \\ m & x_1 & \dots \\ & x_2 & \dots \\ & \vdots & \\ & x_m & \dots \end{bmatrix} \begin{bmatrix} v_1 \\ \vdots \\ v_n \end{bmatrix} = \begin{bmatrix} \langle x_1, v \rangle \\ \vdots \\ \langle x_m, v \rangle \end{bmatrix}$$

$$\Rightarrow (x_v)^T \cdot (x_v) = \sum_{i=1}^m \langle x_i, v \rangle^2$$

previously we stated;

$\Downarrow$

$v^T x^T x v$

↑ find  $v$  that max this.

$\Rightarrow$  Equivalent: Find a  $V$  that maximizes  $V^T \underbrace{(X^T X)}_A V$   
 "maximizing a quadratic form".

Maximize  $V, \|V\|_2 = 1 \left[ V^T A V \right]$  "max quadratic form".

Let's look at a simple case:  $A$  is diagonal,

$$A = \begin{bmatrix} \lambda_1 & & & 0 \\ & \lambda_2 & & \\ & & \lambda_3 & \\ 0 & & & \ddots \lambda_n \end{bmatrix} \quad \lambda_1 \geq \lambda_2 \geq \lambda_3 \dots \lambda_n \geq 0.$$

If want to "max quadratic form":  
 pick  $V = (1, \dots, 0)$ , since  
 $\lambda_1$  is the largest, so put all  
 the weight in it.

If  $A$  = diagonal

$$\Rightarrow V^T \overset{\checkmark}{A} V = (V_1, \dots, V_n) \cdot \begin{pmatrix} \lambda V_1 \\ \vdots \\ \lambda V_n \end{pmatrix} = \sum_{i=1}^n V_i^2 \cdot \lambda_i$$

We don't know if  $A$  is diagonal in general.

but we do know:  $A = Q \cdot D Q^T$  (spectral theorem)

$\uparrow$  symmetric     $\uparrow$  orthogonal     $\curvearrowright$  diagonal.

$\Rightarrow A$  is "almost" diagonal!

let  $D = e_1 = (1, 0, 0, \dots, 0)$

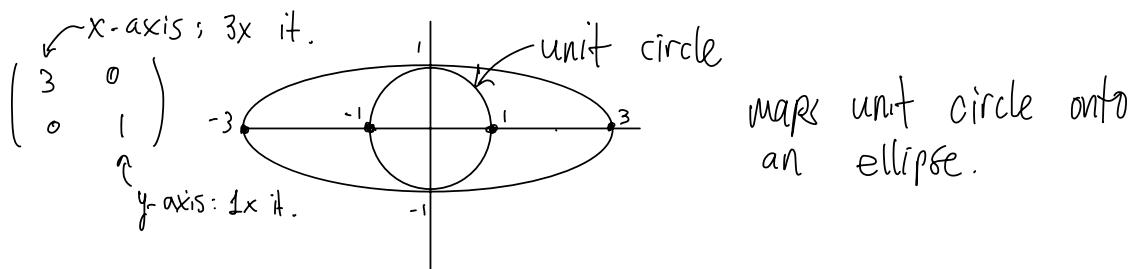
choose  $V = (Q \cdot e_1)$  ← maximizes the top eigenvector of  $A$ .

## PCA Part II:

We studied the following optimization problem:

$$\max_{v, \|v\|_2=1} [V^T A V] \quad \left( \begin{array}{l} \text{matrix } A \text{ corresponded to a} \\ \text{covariance matrix } x^T x \text{ or } \underbrace{\frac{1}{m} X^T X}_{\text{normalized}}. \end{array} \right)$$

Last time: we also discussed the "easy case" when  $A$  was a diagonal matrix.



Recall from linear algebra: rotational matrices,

These are orthogonal matrices, for example:

$$\begin{pmatrix} \cos \theta & -\sin \theta \\ \sin \theta & \cos \theta \end{pmatrix} = \text{rotate } \theta \text{ degrees counter clockwise.}$$

$$\begin{pmatrix} \cos \theta & \sin \theta \\ -\sin \theta & \cos \theta \end{pmatrix} = \text{rotate } \theta \text{ degrees clockwise,}$$

(orthogonal matrices do not change the norms).

For Example: Remember  $A$  is diagonal.

$$A = \begin{pmatrix} 3 & 0 \\ 0 & 1 \end{pmatrix} \Rightarrow \text{solution to } \max_{v, \|v\|_2=1} [V^T \begin{pmatrix} 3 & 0 \\ 0 & 1 \end{pmatrix} V] \text{ is } V = \begin{pmatrix} 1, 0 \end{pmatrix}.$$

↑ corresponds to choosing 3.

$$V^T A V = \sum_{i=1}^n v_i^2 d_i , \quad \sum v_i^2 = 1.$$

$A$  is diagonal

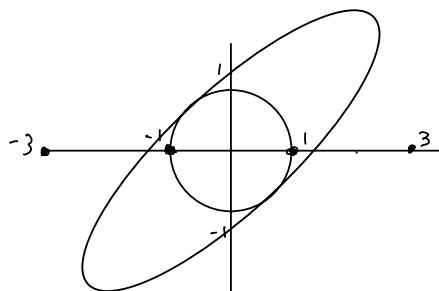
$$\lambda_1, \lambda_2, \dots, \lambda_n \geq 0.$$

Note: can write any covariance matrix (form:  $X^T X$ ), as diagonal matrix . rotational matrix,

for example,

$$\begin{pmatrix} 2 & 1 \\ 1 & 2 \end{pmatrix} = \underbrace{\begin{pmatrix} \frac{1}{\sqrt{2}} & -\frac{1}{\sqrt{2}} \\ \frac{1}{\sqrt{2}} & \frac{1}{\sqrt{2}} \end{pmatrix}}_{\text{rotates counter-clockwise}} \cdot \underbrace{\begin{pmatrix} 3 & 0 \\ 0 & 1 \end{pmatrix}}_{\text{diagonal matrix. stretching in } X\text{-axis.}} \cdot \underbrace{\begin{pmatrix} \frac{1}{\sqrt{2}} & \frac{1}{\sqrt{2}} \\ -\frac{1}{\sqrt{2}} & \frac{1}{\sqrt{2}} \end{pmatrix}}_{\text{rotates } 45^\circ \text{ clockwise.}}$$

$$A = X^T X = Q D Q^T$$



## Spectral Theorem Revisited: X axis

Any symmetric matrix can be written as  $QDQ^T$  where  $Q$  is orthogonal and  $D$  is diagonal with real values on the diagonal eigenvalues.

Furthermore, If  $A = X^T X$  then all eigenvalues  $\geq 0$ .

Claim 1: For any  $V$ ,  $V^T A V \geq 0$ .

$$\begin{aligned} \text{Because } A &= X^T X, V^T A V = V^T (X^T X) V = V^T X^T X V \\ &= (XV)^T \cdot XV \geq 0, \\ &\text{inner product of self} \geq 0. \end{aligned}$$

Claim 2:  $A$  cannot have negative eigenvalues (proof by contradiction).

Let's assume by contradiction that  $\lambda_i < 0$ , ( $i$ th eigenvalue is negative).

Rewrite  $\underbrace{A}_{\text{symmetric}} = \underbrace{QDQ^T}_{\text{spectral theorem}}$ . Let's consider vector  $Q \cdot e_i$ ,

$$\text{where } e_i = (0 \ 0 \ 0 \ 0 \ 0 \underset{i\text{th}}{\overset{1}{\dots}} 0 \ 0),$$

$$V = Q \cdot e_i$$

$$\text{Consider } V^T A V = (\overbrace{Qe_i}^V)^T \underbrace{A}_{(\overbrace{QDQ^T})} \overbrace{(Qe_i)}^V$$

$$= e_i^T \underbrace{Q^T Q}_{\text{Identity since orthogonal.}} \underbrace{D}_{\substack{\text{vector of 0} \\ \text{except 1 entry which is negative,} \\ \text{in the } i\text{th position.}}} D e_i = e_i^T D e_i < 0, \text{ which contradicts Claim 1.}$$

e.g.  $[0 \ 0 \ 0 \ -1 \ 0 \ 0]$

## Recap of PCA:

- 1) subtract the mean from your data,
- 2) Normalize the columns/features of your data,
- 3) compute eigenvalue/eigenvector decomposition of your matrix  $QDQ^T$
- 4) The first  $k$  rows of  $Q^T$  are the  $k$  eigenvectors you're looking for.  
Top  $K$  principal components.

Prove that the  $i^{th}$  row of  $Q^T$  is an eigenvector of  $A$ :

$$i^{th} \text{ row of } Q^T = Q \cdot e_i$$

$$\underbrace{A \cdot Q \cdot e_i}_{\text{i}^{th} \text{ row of } Q^T} = \underbrace{\overbrace{Q}^A \overbrace{D}^{\lambda} \overbrace{Q^T}^I}_{\text{I}} \underbrace{Q \cdot e_i}_{\substack{\uparrow \\ \text{QDe}_i}} = \underbrace{Q D e_i}_{\substack{\uparrow \\ \text{diagonal matrix}}} = \underbrace{\lambda_i}_{\text{eigenvalue}} \underbrace{Q e_i}_{\text{eigenvector}}$$

How do we compute this eigenvector/eigenvalue decomposition?

One method: "singular value decomposition" (SVD),

- Can use it to decompose, but expensive.
- Known: polynomial-time algorithm for computing SVD.

Another method: "power method".