

# PCA

- DIMENSIONALITY REDUCTION
- COMPARE TO "RANDOM PROJECTION" "JL-LEMMA"

RANDOMLY PICKED VECTORS  $r_1, \dots, r_k$

PROJECTED  $x$   $\langle x, r_1 \rangle, \dots, \langle x, r_k \rangle$

- $r_i$ 's were ~~not~~ NOT MEANINGFUL WITH RESPECT TO  $S$
- PRESERVE EUCLIDEAN DISTANCE BETWEEN POINTS
- IN PRACTICE  $k > 100$  FOR RANDOM PROJECTION TO WORK.
- FOR PCA WE CAN CHOOSE  $k = 2$
- PCA LOOKS AT  $S$  TO COME UP WITH A NEW REPRESENTATION

HIGH LEVEL GOAL OF PCA IS TO FIND VECTORS

$$v_1, \dots, v_k \text{ s.t. } \forall x \in S \quad x \approx \sum_{j=1}^k a_j v_j$$

NOTE ABOUT PRE-PROCESSING OF S

- SUBTRACT THE MEAN OR CENTER OF MASS FROM EACH DATA POINT.
- NORMALIZE THE STANDARD DEVIATION OF EACH FEATURE.  
BY

$$V_i: \text{ COMPUTE } \sqrt{\frac{1}{n} \sum_{j=1}^n (x_i^j)^2} = \sigma_i$$

DIVIDE ALL THE  $i^{\text{th}}$  FEATURES BY  $\sigma_i$

How to BEGIN? FIND  $V_1$

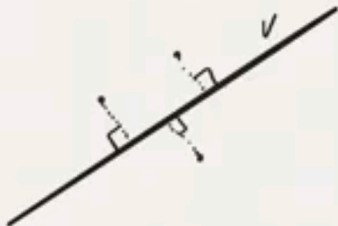
LOOK FOR A VECTOR THAT MINIMIZES SQ-DISTANCE

$$\min_{V, \|V\|_2=1} \frac{1}{m} \sum_{i=1}^m (\text{DISTANCE BTWN } x^i \text{ AND } V)^2$$

PICTURE

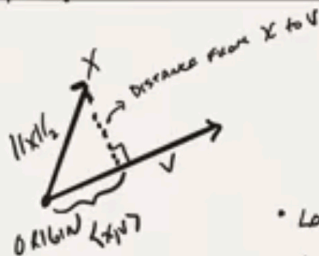


REGRESSION



PCA

$$\min_{V, \|V\|_2=1} \frac{1}{m} \sum_{i=1}^m (\text{DISTANCE FROM } X^i \text{ AND } V)^2$$



$$\underbrace{\langle x, v \rangle^2}_{\text{small}} + (\text{DIST FROM } X \text{ TO } V)^2 = \underbrace{\|x\|^2}_{\text{FIXED}}$$

- LOOKING FOR  $V$  SMALL
- LOOK FOR A  $V$  TO MAKE  $\langle x, v \rangle^2$  LARGE.

EQUIVALENTLY, WE WANT TO FIND A  $V$  THAT MAXIMIZES  $\langle x, v \rangle^2$

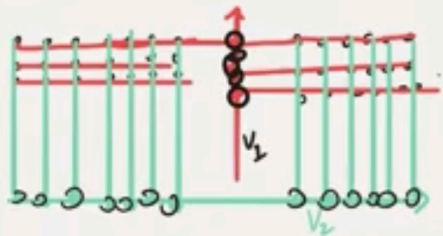
FIND  $V$

$$\|V\|_2=1$$

$$\max \frac{1}{m} \sum_{i=1}^m \langle x^i, v \rangle^2$$

← DIRECTION OF  
MAXIMAL VARIANCE.

$$\text{VAR } \underline{E[X]} - (E[V])^2$$



MAIN IDEA  
VISUALLY

THIS WAS FOR 1 VECTOR WHAT ABOUT FINDING  $K$ -VECTORS  
 $K$ -COMPONENTS

MAX SUBSPACES  
 $S$  OF DIMENSION  $K$   $\frac{1}{m} \sum_{j=1}^m (\text{LENGTH OF } X^j \text{ PROJECTED ONTO } S)^2$

A REALLY NICE/PREFERRABLE BASIS WOULD BE AN ORTHONORMAL BASIS  
 $v_1, \dots, v_K$

$$(\text{DISTANCE FROM } X \text{ TO } S)^2 = \|x\|^2 - [\langle x, v_1 \rangle^2 + \dots + \langle x, v_K \rangle^2]$$

$$\max_{\substack{v_1, \dots, v_K \\ \text{ORTHOGONAL}}} \frac{1}{m} \sum_{j=1}^m \sum_{i=1}^K \langle x_j^i, v_i \rangle^2$$

PCA OBJECTIVE

ASSUME WE HAVE  $v_1, \dots, v_k$

$$\underline{X = \langle x, v_1 \rangle \cdot v_1 + \langle x, v_2 \rangle \cdot v_2 + \dots + \langle x, v_k \rangle \cdot v_k}$$

$X$  CAN BE WRITTEN AS A VECTOR IN  $\mathbb{R}^k$  CORRESPONDING TO THESE PROJECTIONS.

### APPLICATION 1: UNDERSTANDING GENOMES

TOOK 1400 PEOPLE FROM EUROPE

EACH PERSON WAS REPRESENTED ACCORDING TO 200,000 GENETIC MARKERS IN THEIR GENOME.

CORRESPONDING TO MATRIX OF DIM  $1400 \times 200,000$

• 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).

ANOTHER APPLICATION:

IMAGE DATA COMPRESSION

STRATEGY FOR COMPRESSING 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 DATA SET

$K \approx 100-150$

IMAGE  $\approx$  LINEAR COMBINATION OF 150 VECTORS OF LENGTH 65,000  
IS AN IMAGE!

# BIG QUESTION:

How DO WE FIND THESE  $V_3, \dots, V_k$ 's?

$$\max_{V_3, \dots, V_k} \frac{1}{m} \sum_{j=1}^m \sum_{i=1}^k \langle x_j^i, V_i \rangle^2 \text{ IS MAXIMIZED}$$

ORTHOGONAL.

Let  $X$  be an  $m$  by  $n$  matrix  $\underbrace{\frac{1}{m} X^T X}_{\text{SAMPLE COVARIANCE MATRIX}} \quad (n \text{ by } n \text{ MATRIX})$

$V \leftarrow$  a column vector

$V^T \leftarrow$  a row vector

$V^T V \leftarrow$  INNER PRODUCT (SCALAR)

$V V^T \leftarrow$  OUTER PRODUCT (MATRIX)

SAMPLE COVARIANCE  
MATRIX

$(i, j)$  entry of  $X^T X$

CORRESPONDS TO "How SIMILAR IS  
FEATURE  $i$  TO FEATURE  $j$ ?"



NOTE THAT  $X^T X$  IS A SYMMETRIC MATRIX

$$X^T X \begin{pmatrix} i \\ j \end{pmatrix} = \begin{array}{l} \text{inner product of row } i \text{ of } X^T \text{ with column } j \text{ of } X \\ \text{" column } i \text{ of } X \text{ with column } j \text{ of } X \\ \text{= column } j \text{ of } X \text{ with " } i \text{ of } X \end{array}$$

FACT: ALL EIGENVALUES OF SYMMETRIC MATRICES  $\geq 0$ .

FOR A MATRIX  $A$ ,  $v$  IS AN EIGENVECTOR IF  $A \cdot v = \lambda \cdot v$   $\lambda \in \mathbb{R}$   
EIGENVALUE.

DEFINITION: AN ORTHOGONAL MATRIX IS ONE WHERE ALL COLUMNS ARE ORTHONORMAL.  $\Leftrightarrow$

$$A^T A = I \quad (A A^T = I)$$

SPECTRAL THM:

EVERY <sup>SYMMETRIC</sup> MATRIX  $A$  HAS AN  
EIGENDECOMPOSITION:

$$A = Q \cdot D \cdot Q^T$$

ORTHOGONAL  
MATRIX

DIAGONAL MATRIX

ENTRIES OF  $D$  ARE THE  
EIGENVALUES OF  $A$ .

Let's TRY TO COMPUTE  $V_1$ .

$X$  IS MATRIX CORRESPONDING TO  $S$ .

$X$  IS  $n \times n$

$X$  is  $m$  by  $n$        $X \cdot v = \begin{bmatrix} \langle x_1, v \rangle \\ \vdots \\ \langle x_m, v \rangle \end{bmatrix}$  ← rows of  $X$   $0 \leq n \leq m$

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

---

FIND A  $v$  THAT MAXIMIZES  
THAT MAXIMIZES  $v^T \underbrace{(X^T X)}_A \cdot v$

FIND  $v$

"MAXIMIZING A QUADRATIC FORM"

MAXIMIZE

$$V, \|V\|_2 = 1$$

$$V^T A \cdot V$$

"MAX QUADRATIC FORM"

Let's look at a simple case:  $A$  is DIAGONAL.

$$A = \begin{pmatrix} \lambda_1 & & \\ & \lambda_2 & \\ & & \ddots \\ & & & \lambda_n \end{pmatrix}$$

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

PICK  $V = (1, \dots, 0)$

$$A = X^T X$$

$$V^T A V = v_1, \dots, v_n \cdot \begin{pmatrix} \lambda_1 v_1 \\ \vdots \\ \lambda_n v_n \end{pmatrix} = \sum_{i=1}^n v_i^2 \cdot \lambda_i$$

WE DON'T KNOW IF  $A$  IS DIAGONAL IN GENERAL  
 $A$  IS "ALMOST" DIAGONAL.

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

$$A = Q \cdot D \cdot Q^T$$

CHOOSE

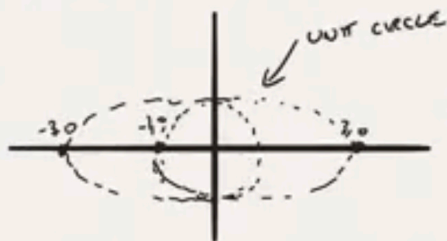
$$V = (Q \cdot e_1)$$

MAXIMIZES  
TOP  
EIGENVECTOR  
OF  $A$

DIAGONAL

LAST TIME: WE ALSO DISCUSSED THE "EASY CASE"  
WHEN  $A$  WAS A DIAGONAL MATRIX

$$\begin{pmatrix} 3 & 0 \\ 0 & 1 \end{pmatrix}$$



MAPS UNIT CIRCLE  
INTO AN ELLIPSE.

RECALL FROM LINEAR ALGEBRA: ROTATION MATRICES  
ORTHOGONAL MATRICES

FOR EXAMPLE  $\begin{pmatrix} \cos \theta & -\sin \theta \\ \sin \theta & \cos \theta \end{pmatrix}$

ROTATE  
 $\theta$  DEGREE  
COUNTERCLOCKWISE

$$\begin{pmatrix} \cos \theta & \sin \theta \\ -\sin \theta & \cos \theta \end{pmatrix}$$

ROTATE  
THE  
AXES  $\theta$  DEGREES  
CLOCKWISE

$$\begin{pmatrix} 3 & 0 \\ 0 & 1 \end{pmatrix}$$

THE solution to  $\max_{V, \|V\|_2=1} V^T \begin{pmatrix} 3 & 0 \\ 0 & 1 \end{pmatrix} V$

$$V = (1, 0)$$

$$\underline{V^T A V = \sum_{i=1}^n v_i^2 \lambda_i}$$

A is diagonal

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

$$\sum v_i^2 = 1 \quad \text{corresponds to choice 3}$$

$$V = (1, 0)$$

x

o d c

FOR EXAMPLE

$$\begin{pmatrix} 2 & 1 \\ 1 & 2 \end{pmatrix} = \underbrace{\begin{pmatrix} \frac{1}{\sqrt{2}} & -\frac{1}{\sqrt{2}} \\ \frac{1}{\sqrt{2}} & \frac{1}{\sqrt{2}} \end{pmatrix}}_{\text{ROTATE COUNTER CLOCKWISE}} \cdot \underbrace{\begin{pmatrix} 3 & 0 \\ 0 & 1 \end{pmatrix}}_{\text{STRETCHING}} \cdot \underbrace{\begin{pmatrix} \frac{1}{\sqrt{2}} & \frac{1}{\sqrt{2}} \\ -\frac{1}{\sqrt{2}} & \frac{1}{\sqrt{2}} \end{pmatrix}}_{\text{ROTATE } 45^\circ \text{ CLOCKWISE}}$$



SPECTRAL THM: 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  $A = \underline{X^T X}$  THEN ALL EIGENVALUES  $\geq 0$ .

CLAIM 1: FOR ANY  $v$ ,  $v^T A v \geq 0$

BECOME  $A = X^T X$   $v^T A v = (\underline{Xv})^T \cdot \underline{Xv} \geq 0$ .

CLAIM 2  $A$  CANNOT HAVE NEGATIVE EIGENVALUES. (PROOF IS BY CONTRADICTION)

LET'S ASSUME BY CONTRADICTION THAT  $\lambda_i < 0$  (i.e. EIGENVALUE IS NEGATIVE)

$A = QDQ^T$  LET'S CONSIDER VECTOR  $Q \cdot e_i$  ( $0 \ 0 \ 0 \ 0 \ \underset{i^{th}}{1} \ 0 \ 0 \ 0 \ 0$ )  $= e_i$

$v = Q \cdot e_i$   $v^T A v$

$$e_i^T \cancel{Q^T} Q \cancel{D} Q^T Q e_i = \underline{e_i^T D e_i} < 0$$



TO RECAP:

PCA:

- 1) SUBTRACT THE MEAN FROM YOUR DATA
- 2) NORMALIZE THE COLUMNS 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

TO  $P$   $K$  PRINCIPAL COMPONENTS.

PROVE  $i^{\text{th}}$  ROW OF  $Q^T$  IS AN EIGENVECTOR OF  $A$

$$i^{\text{th}} \text{ ROW OF } Q^T = Q \cdot e_i$$

$$\begin{aligned} \underline{A} \cdot \underline{Q \cdot e_i} &= Q D Q^T Q \cdot e_i = Q D e_i \\ &= \underbrace{\lambda_i}_{\text{eigenvalue}} \cdot \underbrace{Q \cdot e_i}_{\text{eigenvector}} \end{aligned}$$

HOW DO WE COMPUTE THIS DECOMPOSITION?

ANOTHER PROBLEM "SINGULAR VALUE DECOMPOSITION"  
SVD

KNOWN: POLYNOMIAL-TIME ALGORITHMS FOR COMPUTING SVD.