

S&DS 365 / 565  
Data Mining and Machine Learning

# **Unsupervised Learning: Principal Component Analysis (PCA)**

Yale

Principal component analysis (PCA) is one of the most used unsupervised learning methods.

We will explore five different interpretations of PCA

Principal component analysis (PCA) is one of the most used unsupervised learning methods.

We will explore five different interpretations of PCA

The first is a data variance/dimensionality reduction approach

Principal component analysis (PCA) is one of the most used unsupervised learning methods.

We will explore five different interpretations of PCA

The first is a data variance/dimensionality reduction approach

The second is a data representation approach

Principal component analysis (PCA) is one of the most used unsupervised learning methods.

We will explore five different interpretations of PCA

The first is a data variance/dimensionality reduction approach

The second is a data representation approach

The third is a data denoising approach

Principal component analysis (PCA) is one of the most used unsupervised learning methods.

We will explore five different interpretations of PCA

The first is a data variance/dimensionality reduction approach

The second is a data representation approach

The third is a data denoising approach

The fourth is a low-rank approximation approach

Principal component analysis (PCA) is one of the most used unsupervised learning methods.

5

We will explore five different interpretations of PCA

The first is a data variance/dimensionality reduction approach

The second is a data representation approach

The third is a data denoising approach

The fourth is a low-rank approximation approach  
*recommendation system* <sup>matrix</sup>

The fifth is a Gaussian covariance based approach  
*probabilistic interpretation*

Before we begin we need to understand what a “direction” of data means

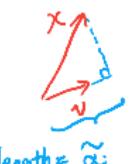
We have said before that PCA finds “interesting **directions**” of data

PCA finds components that are principle

We have data  $x_i$  with  $i \in [n]$        $x_i \in \mathbb{R}^p$

A **direction** is simply a unit vector  $v \in \mathbb{R}^p$  i.e.  $\|v\|_2 = 1$ .  
单位向量

We then project the data on that direction and have projected data.

$$\tilde{\alpha}_i = \langle x_i, v \rangle \in \mathbb{R}$$


length =  $\tilde{\alpha}_i$       vector is  $\tilde{x}$

In this way we have **dimensionality reduction**

Our original data was in  $p$  dimensions, now 1 dimension.

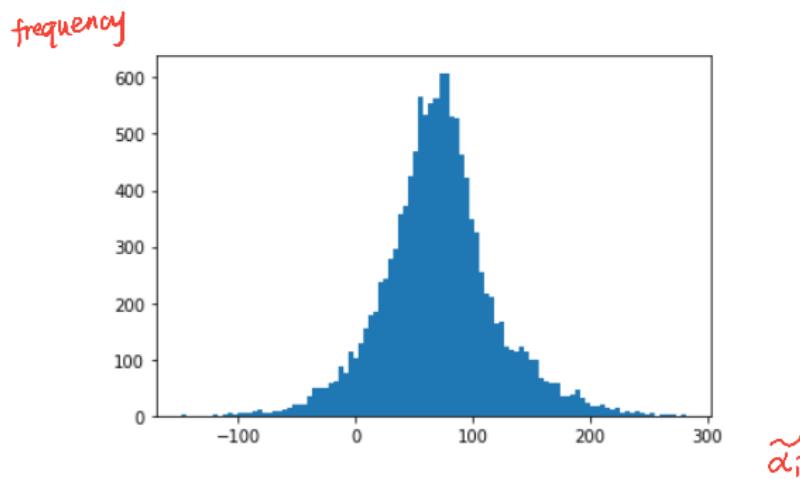
only need  $\tilde{\alpha}_i$  to represent  $x_i$

What does this projection mean? Let's take a look.

Let's remember the digits example

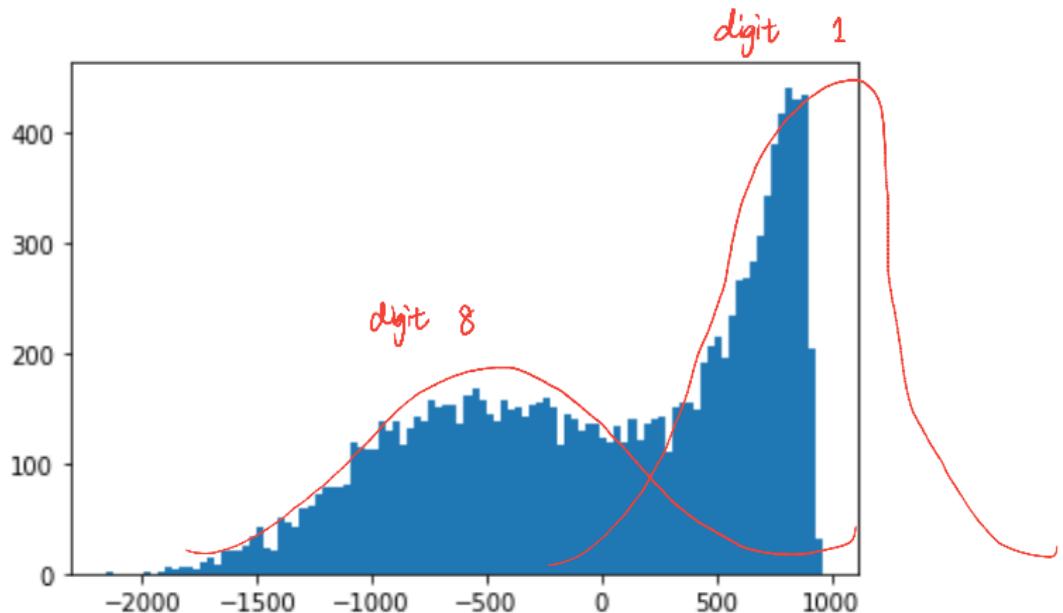
data  $12593 \times 784$

Let's take a random unit vector  $v$  and plot a histogram of  $\tilde{\alpha}_i$ .



Not very informative. Just looks pretty Gaussian.

Let's pick a direction provided by PCA.



Way more interesting.

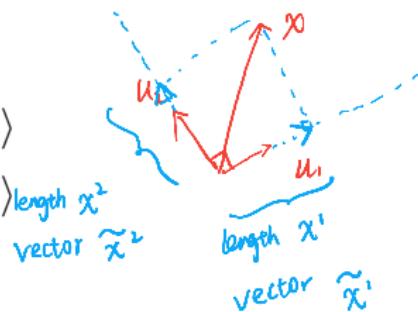
PCA assume unit vectors are orthogonal. For other <sup>no linear</sup> methods, don't need orthogonality  
 make computation easier no meaning for that

We can do this in higher dimensions. Suppose we have **two directions** now  $u_1$  and  $u_2$ . We assume that the two <sup>unit</sup> vectors are orthogonal  $\langle u_1, u_2 \rangle = 0$  and unit norm (orthonormal).

$$u_1 \perp u_2$$

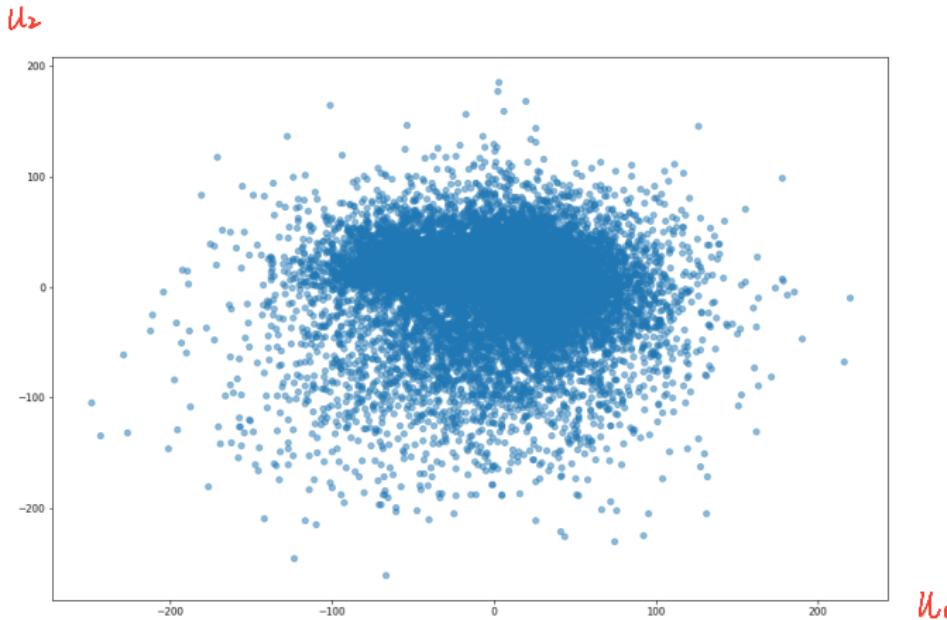
Now take

$$\begin{array}{l} (x_i^1, x_i^2) \text{ is coordinate of } x \\ \text{if base is } u_1, u_2 \end{array} \quad \left. \begin{array}{l} \text{a linear function} \\ x_i^1 = \langle x_i, u_1 \rangle \\ x_i^2 = \langle x_i, u_2 \rangle \end{array} \right\}$$

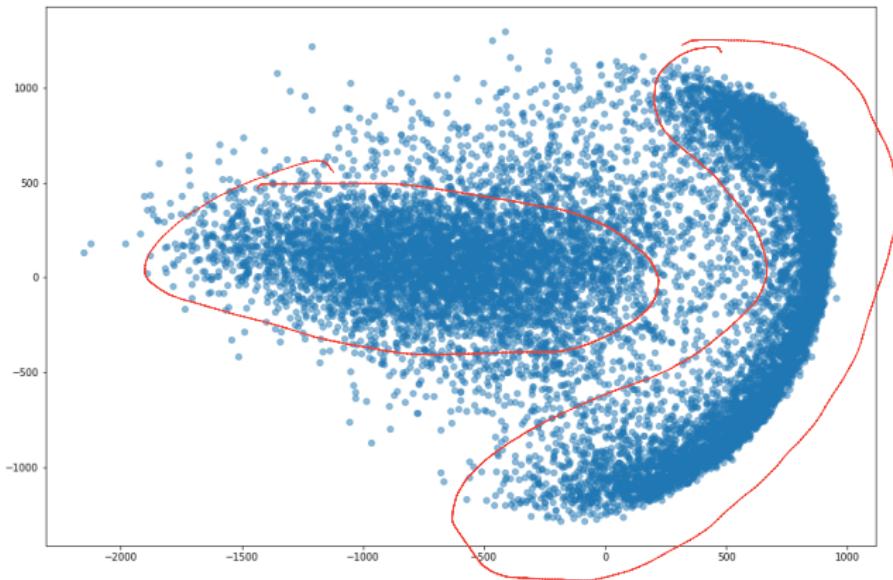


$$\begin{aligned} x &\in \mathbb{R}^P & \tilde{x}' &\in \mathbb{R}^P \\ u_1, u_2 &\in \mathbb{R}^P & \tilde{x}^2 &\in \mathbb{R}^P \\ && \tilde{x}' + \tilde{x}^2 &= x \end{aligned}$$

We can plot those coordinates now as well for two random orthonormal vectors.



Use the PCA components instead



Definitely see more structure.

## Dimensionality reduction for visualization

Example: Senate voting

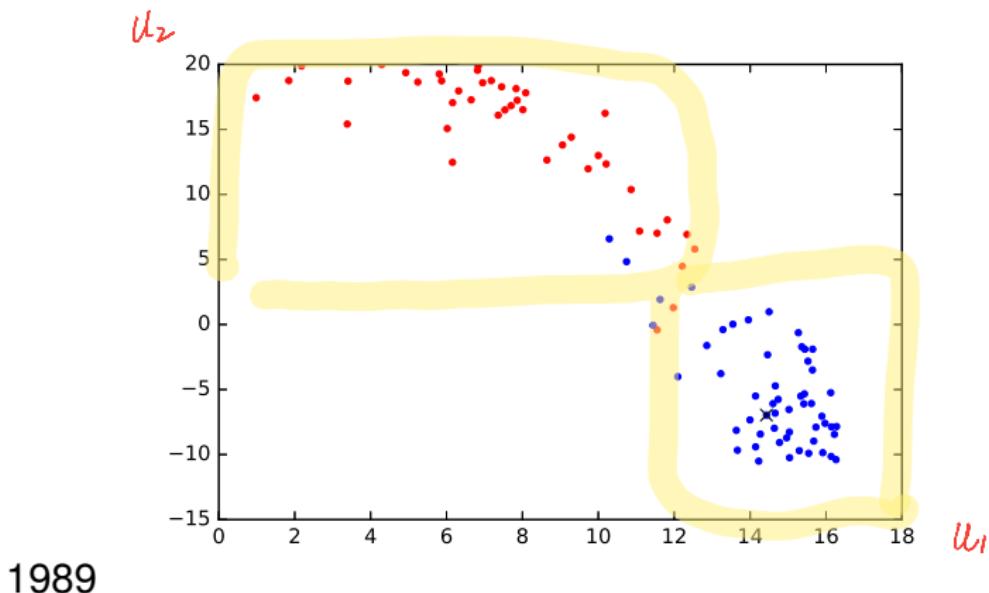
$x_i \in \{-1, 0, +1\}^p$ , voting record for Senator  $i$

↓      ↑      ↗  
disapprove    no vote    approve bill       $p$  bills       $x \in \mathbb{R}^p$

## Dimensionality reduction for visualization

Example: Senate voting *see Notebook*

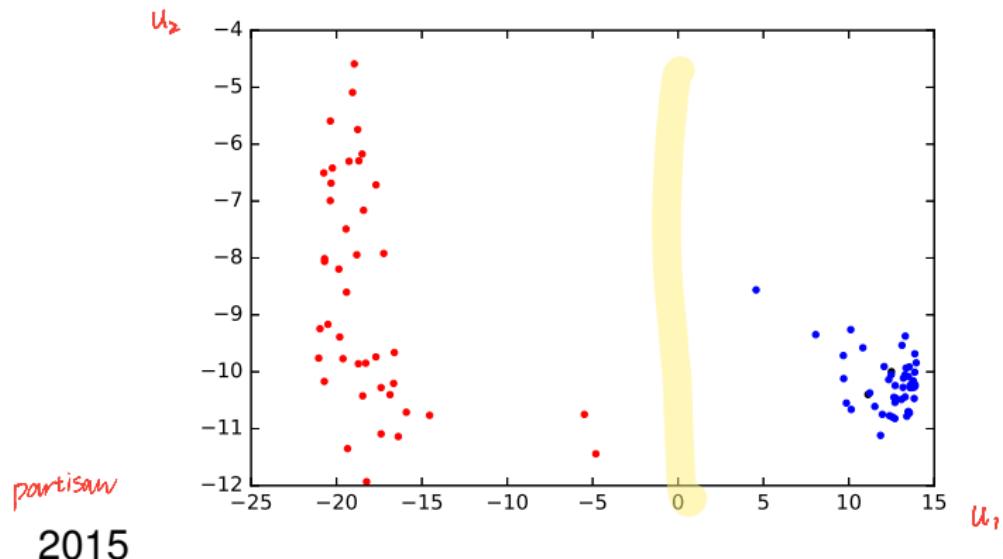
*2 principle components*



## Dimensionality reduction for visualization

Example: Senate voting

$(a_1, a_2)$  2-D coordinate of data points



For general base matrix  $A$  (not orthogonal)  $A \in \mathbb{R}^{p \times k}$

$$\text{rank}(A) = k \quad p \geq k$$

$$\tilde{x} = \arg \min_{\alpha} \|A\alpha - x\|_2^2$$

least square solution

So far  $1-D$  or  $2-D$  projections

For dictionary learning

let  $k \gg p \Rightarrow \alpha$  to be sparse

When  $A$  is orthogonal  
best representative of  $x$

solution is  $\tilde{x} = A^T x$

$K$ -dimensional projection:

$$U = \left[ \begin{array}{c|c|c|c} u_1 & u_2 & \cdots & u_k \end{array} \right] \in \mathbb{R}^{p \times k}$$

$U$  is orthogonal matrix

That is let  $\underbrace{U \in \mathbb{R}^{p \times k}}_k$  be a matrix such that  $U^T U = I$ . Then

$$\mathbb{R}^k \ni \tilde{x}_i = U^T x_i$$

$$k \times 1 = k \times p \quad p \times 1$$

What motivates this?

# First Interpretation

preserve variance (size) of  $x_i$   
energy of  $x_i$

Norm-preserving projections

$$x = \begin{bmatrix} x_1 & x_2 & \dots & x_p \end{bmatrix} \in \mathbb{R}^{p \times n}$$

Let  $X \in \mathbb{R}^{p \times n}$  with columns  $x_i$ . (Note this is different than regression where this would be  $X^T$ )

Find an orthogonal matrix  $U$  maximize (What's the best space  $K$  that we should focus on so  $x_i$  has

longest energy on that space)

$$U = \arg \max_{A|A^T A = I} \sum_{i=1}^n \|A^T x_i\|_2^2$$

rewrite as

$$\sum_{i=1}^n \|A^T x_i\|_2^2 = \|A^T X\|_F^2$$

(Matrix Norm) L<sub>2</sub>-norm square of matrix

Frobenius norm:  $\|B\|_F^2 = \sum_{ij} B_{(ij)}^2 = \text{trace}(BB^T) = \langle B, B \rangle$  is the Euclidean norm for matrices.

see before

The optimization is non-convex, but has a closed form solution involving the singular value decomposition (SVD).

Basic facts about a orthonormal collections of vectors  $U \in \mathbb{R}^{p \times k}$  and orthonormal projections

orthonormal matrix

列空间

$K = \text{colspace}(U)$  is a  $k$ -dimensional subspace

$$U = \begin{bmatrix} | & | & \dots & | \\ u_1 & u_2 & \dots & u_k \\ | & | & \dots & | \end{bmatrix}$$

$k$ 个列向量  $u_i \in \mathbb{R}^p$

Columns  $u_i$  of  $U$  form basis vectors for the space

Take  $v \in K$

vector  $v \in \mathbb{R}^k$  can be written as a linear combination of  $u_i$

$$v = \sum_{i=1}^k \alpha_i u_i = U\alpha$$

For any vector  $v \in K$ , can find unique vector  $\alpha \in \mathbb{R}^k$

Nice property:  $\alpha = U^T v$

$\alpha \in \mathbb{R}^k \Rightarrow$  coefficients of the basis vectors  $\Rightarrow$  coordinates in  $K$

$e_1, e_2$  are standard base vectors

$$e_i = \begin{bmatrix} 0 \\ 0 \\ \vdots \\ i \\ 0 \end{bmatrix} \leftarrow i\text{th coordinate}$$

Example:  $u_1 = e_1 \in \mathbb{R}^p$ ,  $u_2 = e_2 \in \mathbb{R}^p$

For any  $v \in K$ ,  $v$  can be written as

$$v = (\alpha_1, \alpha_2, 0, 0, 0, \dots, 0)^T$$

$$v = \begin{bmatrix} e_1 & e_2 \\ | & | \end{bmatrix}$$

$$v = \alpha_1 e_1 + \alpha_2 e_2 \in K$$

$$\begin{bmatrix} \alpha_1 \\ \alpha_2 \end{bmatrix} = 2 \times 1 = \begin{bmatrix} -e_1 - \\ -e_2 - \\ \vdots \\ -e_p - \end{bmatrix}_{2 \times p}^T \begin{bmatrix} \alpha_1 \\ \alpha_2 \\ 0 \\ \vdots \\ 0 \end{bmatrix}_{p \times 1}$$
$$\alpha = U^T v$$

$\alpha \Rightarrow$  coordinates  $(\alpha_1, \alpha_2)$

$$\begin{cases} \alpha_1 = e_1^T v \\ \alpha_2 = e_2^T v \end{cases}$$

Orthonormal basis need not be unique

take  $G$  is rotation matrix or any orthonormal transformation

$G \in \mathbb{R}^{k \times k}$  with  $G^T G = I = G G^T$

$G$  is orthonormal matrix

Let  $\tilde{U} = UG$ ,  $\tilde{U}_i$  is orthonormal (why?)

In general

$\tilde{U} \neq U$ , but  $\text{colspace}(U) = \text{colspace}(\tilde{U})$

recall previous eg

$$u_1 = e_1 \quad u_2 = e_2$$

$$\tilde{u}_1 = \frac{e_1 + e_2}{\sqrt{2}} \quad \tilde{u}_2 = \frac{e_1 - e_2}{\sqrt{2}}$$

$$\text{Span}\{e_1, e_2\} = \text{Span}\{\tilde{u}_1, \tilde{u}_2\}$$

$$\text{but } e_1 \neq \tilde{u}_1, e_2 \neq \tilde{u}_2$$

column vectors

Show  $U$  is orthonormal  $\Leftrightarrow$  Show column  
of  $U$  is orthonormal

$$\tilde{U}^T \tilde{U} = I_{k \times k}$$

(ii)

Because we are projecting data  $\in \mathbb{R}^p$  to  $\mathbb{R}^k$   $p > k$   
so we loose some info, But we want to know these info.  
so we have  $\Downarrow$   
正交互补空间

## Orthogonal complement of $K$

Take  $w \in \mathbb{R}^p$  with  $U^T w = 0_k \in \mathbb{R}^k$  mean  $w$  is orthogonal to every column of matrix  $U$   
then  $w$  is orthogonal to every linear combination

Then for all  $v \in K$ ,  $\langle w, v \rangle = 0$  

of column of  $U$   
 $\downarrow$   
vector  $v$

$w$  lives in orthogonal complement of  $K$  denoted  $K^\perp$

i.e.  $w$  lives in Null Space (kernel) of  $U^T$

for  $v \in K$

$$UU^T(w + v) = v \text{ (why?)}$$

Another nice property: Energy can be decomposed between  $k$  and  $k^\perp$

For  $v \in K$ ,  $\|v\|_2^2 = \|U\alpha\|_2^2 = \alpha^T \underbrace{U^T U}_{I} \alpha = \|\alpha\|_2^2 = \|U^T v\|_2^2$ .  
energy of  $v$

General vector  $v \in \mathbb{R}^p \Rightarrow v = UU^T v + w$  where  $w = v - UU^T v \in K^\perp$

Show  $w \in k^\perp$

$$U^T w = U^T(v - UU^T v) = 0 \quad = U^T v - \underbrace{U^T U U^T v}_I$$

$$\|v\|_2^2 = \underbrace{\|UU^T v\|_2^2}_{\text{part of } v \text{ that lives in } K} + \underbrace{\|(I - UU^T)v\|_2^2}_{\text{part of } v \text{ that lives in } K^\perp}$$

## Orthogonal projections

$K$  is subspace of  $\mathbb{R}^k$   
 $V \in \mathbb{R}^{p \times k}$

$K = \text{colspace } U \subset \mathbb{R}^p$  is  $k$ -dimensional

Orthogonal projection of vector  $v \in \mathbb{R}^p$   
onto space  $K$  is vector  $g \in \mathbb{R}^p$  that  
minimize distance between  $v$  and  $g$

exists  $\exists \alpha \in \mathbb{R}^k$  with  $g = U\alpha$  since  $g \in K$   
change from finding  $g$  to finding  $\alpha$

$$P_K(v) = U(\arg \min_{\alpha \in \mathbb{R}^k} \|v - U\alpha\|_2^2)$$

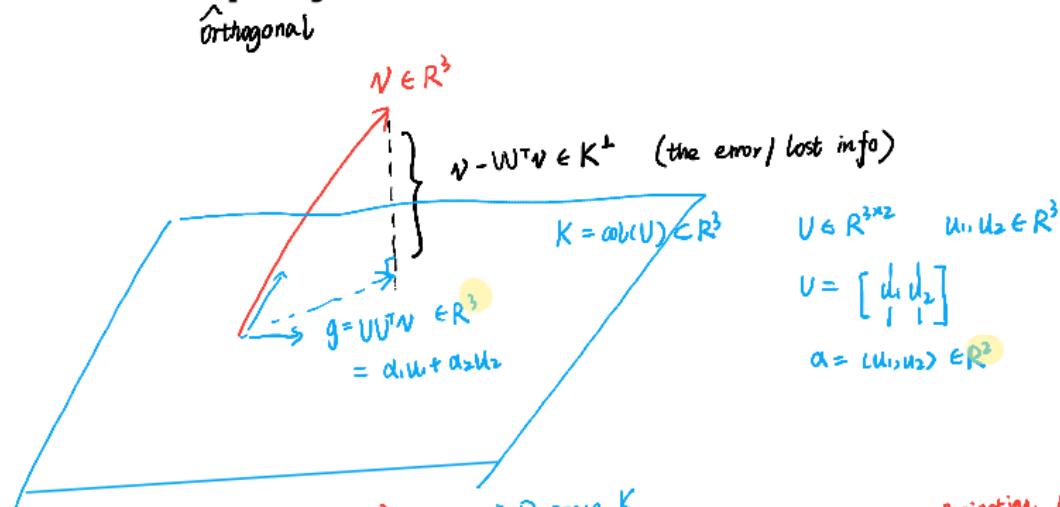
Linear regression again  $\hat{\alpha} = (\underbrace{U^T U}_I)^{-1} U^T v$

$$U^T U = I \implies \hat{\alpha} = U^T v \implies P_K(v) = UU^T v.$$

$UU^T$  is an **orthogonal projection matrix** onto the space  $K$ .

Note that  $v \neq UU^T v$  in general unless  $v \in K$ .

# Illustration of projections



Orthogonal project  $v \in \mathbb{R}^3$  onto 2-D space  $K$   
is try to minimize  $l_2$ -norm of error vector  $(v - UV^T v) \in K^\perp \Rightarrow$  projection matrix is

$$\text{we can show } I - UU^T = P_{K^\perp}$$

The orthogonal projection matrix of  
projection onto orthogonal complement of  $K$  is  $I - UU^T$   
it minimize the distance between  $v$  and projected  $v$

$$(I - UU^T)v$$

Column vector of orthonormal matrix A



The orthonormal basis vectors are each directions. Now we return to PCA and the SVD.

$$U = \arg \max_{A|A^T A = I} \|A^T X\|_F^2$$

How to solve it?

Recall  $X \in \mathbb{R}^{p \times n}$ . Assume  $n \geq p$ .

effectively find those coefficients  
that represent the projection of  $X$   
onto space that spanned by  $A$

For any matrix  $X \in \mathbb{R}^{p \times n}$

**Singular value decomposition (SVD):**  $X = USV^T$  where

$U \in \mathbb{R}^{p \times p}$  is an orthonormal matrix,  $S \in \mathbb{R}^{p \times p}$  is diagonal with positive entries, and  $V^T \in \mathbb{R}^{p \times n}$  is orthonormal.

the rows of  $V^T$  is orthonormal

the columns of  $V$  ~

Equivalently:  $U \in \mathbb{R}^{p \times p}$  and  $V \in \mathbb{R}^{n \times n}$  with  $S \in \mathbb{R}^{p \times n}$  where  $S_{(ij)} \geq 0$  and  $S_{(ij)} = 0$  for  $i \neq j$ . *diagonal matrix with positive entries*

$$X = USV^T = \begin{bmatrix} \text{rank}(X) \\ p \times p \end{bmatrix} \begin{bmatrix} \text{diag} \\ p \times n \end{bmatrix} \begin{bmatrix} 0 \\ n \times n \end{bmatrix}$$

$\text{rank}(X) = \text{nnz}(S)$  (number of non-zeros)

$$X = \sum_{i=1}^{\text{rank}(X)} S_{(ii)} u_i v_i^T$$

scalar  $\swarrow$  column vectors

Convention:  $S_{(ii)} \geq S_{(jj)}$  for  $i \leq j$ .  $S = \begin{bmatrix} \text{diag} & 0 \\ 0 & 0 \end{bmatrix}$

$u_i$  are the **left singular vectors**

$v_i$  are the **right singular vectors**

Diagonal entry of  $S$

$\sigma_i = S_{(ii)}$  are the **singular values**

$$S_{(ii)} \sim S_{(kk)} \rightarrow \begin{cases} \text{left : } u_i \sim u_k \\ \text{right : } v_i \sim v_k \end{cases}$$

We say the top- $k$  singular vectors are those corresponding to the  $k$  largest singular values.

Diagonal entry of  $S$  descending order

sample correlation matrix of  $X$

$$\text{if } E(X) = 0 \quad \hat{\Sigma} = \frac{1}{n} X X^T$$

Diagonal Matrix property : ① eigen value decomposition

② eigen value  $\lambda \in \mathbb{R}$

$$\text{Eigen value of } \hat{\Sigma} \geq 0 \Rightarrow \hat{\Sigma} = U D V^T$$

means The variation of  $\hat{\Sigma}$  is captured by singular

vectors

values

means

captures

variation

of

$\hat{\Sigma}$

is

captured

by

singular

vectors

values

means

captures

variation

of

$\hat{\Sigma}$

is

captured

by

singular

vectors

values

means

captures

variation

of

$\hat{\Sigma}$

is

captured

by

singular

vectors

values

means

captures

variation

of

$\hat{\Sigma}$

is

captured

by

singular

vectors

values

means

captures

variation

of

$\hat{\Sigma}$

is

captured

by

singular

vectors

values

means

captures

variation

of

$\hat{\Sigma}$

is

captured

by

singular

vectors

values

means

captures

variation

of

$\hat{\Sigma}$

is

captured

by

singular

vectors

values

means

captures

variation

of

$\hat{\Sigma}$

is

captured

by

singular

vectors

values

means

captures

variation

of

$\hat{\Sigma}$

is

captured

by

singular

vectors

values

means

captures

variation

of

$\hat{\Sigma}$

is

captured

by

singular

vectors

values

means

captures

variation

of

$\hat{\Sigma}$

is

captured

by

singular

vectors

values

means

captures

variation

of

$\hat{\Sigma}$

is

captured

by

singular

vectors

values

means

captures

variation

of

$\hat{\Sigma}$

is

captured

by

singular

vectors

values

means

captures

variation

of

$\hat{\Sigma}$

is

captured

by

singular

vectors

values

means

captures

variation

of

$\hat{\Sigma}$

is

captured

by

singular

vectors

values

means

captures

variation

of

$\hat{\Sigma}$

is

captured

by

singular

vectors

values

means

captures

variation

of

$\hat{\Sigma}$

is

captured

by

singular

vectors

values

means

captures

variation

of

$\hat{\Sigma}$

is

captured

by

singular

vectors

values

means

captures

variation

of

$\hat{\Sigma}$

is

captured

by

singular

vectors

values

means

captures

variation

of

$\hat{\Sigma}$

is

captured

by

singular

vectors

values

means

captures

variation

of

$\hat{\Sigma}$

is

captured

by

singular

vectors

values

means

captures

variation

of

$\hat{\Sigma}$

is

captured

by

singular

vectors

values

means

captures

variation

of

$\hat{\Sigma}$

is

captured

by

singular

vectors

values

means

captures

variation

of

$\hat{\Sigma}$

is

captured

by

singular

vectors

values

means

captures

variation

of

$\hat{\Sigma}$

is

captured

by

singular

vectors

values

means

captures

variation

of

$\hat{\Sigma}$

is

captured

by

singular

vectors

values

means

captures

variation

of

$\hat{\Sigma}$

is

captured

by

singular

How to solve  $U = \arg \max_{A|A^T A=I} \|A^T X\|_F^2$

How to solve  $U = \arg \max_{A|A^T A=I} \|A^T X\|_F^2$

Plug in  $X = USV^T$   
 $U \in \mathbb{R}^{p \times p}$   $V \in \mathbb{R}^{n \times n}$   $S \in \mathbb{R}^{p \times n}$

Solution

$$\|A^T USV^T\|_F^2 = \|A^T US\|_F^2$$

Multiply a matrix  $V$  by right or left whose rows are orthonormal  
The Frobenius-norm square remain the same ①

Guess  $A = U_k$

i<sup>th</sup> column of  $U_k$

$U_k \in \mathbb{R}^{p \times k}$ ,  $[U_k]_{:,i} = u_i$  (matrix whose columns are the top  $k$  left singular vectors of the matrix  $X$ )

Then,

plug in  $A = U_k$

This means the energy of  $X$  is captured by singular value of  $X$ :  $\sigma_i$

$$\|U_k^T US\|_F^2 = \sum_{i=1}^k \sigma_i^2 \quad \text{Proof in graph Notes}$$

Also turns out:

$$\|X\|_F^2 = \|USV^T\|_F^2 = \|S\|_F^2 = \sum_{i=1}^{\min(p,n)} \sigma_i^2$$

both are orthonormal matrix according to ①

The left over energy =  $\sum_{i=k+1}^{\text{rank}(X)} \sigma_i^2$

sum of the least ( $\text{rank}(X) - k$ ) squared singular values

Show our guess is right There is not any other measure that maximize  
 On the other hand  $\|A^T U S\|_F^2$

$$\begin{aligned}\|A^T U S\|_F^2 &= \|\tilde{U} S\|_F^2 \quad \text{Define } A^T U = \tilde{U} \\ &= \sum_i^k \sigma_i^2 \|\tilde{u}_i\|_2^2\end{aligned}$$

where we take  $\tilde{u}_i = A^T u_i$  where  $u_i$  is the  $i^{th}$  column of  $U$ . Clearly,  
 $0 \leq \|\tilde{u}_i\|_2^2 \leq 1$ . Furthermore,  $\sum_i^k \|\tilde{u}_i\|_2^2 = k$ . Therefore,

as  $A$  is an orthonormal matrix whose column vector is orthonormal  
 multiplying by an unit vector will not expand any energy  
 the  $L_2$ -norm square of  $\tilde{u}_i$  is up-bounded by 1

↑ prof:  $\|A^T U S\|_F^2 = \|\tilde{U} S\|_F^2 = \|A^T U\|_F^2 = \|A\|_F^2 = \sum_{i=1}^k \|A_{i,:}\|_2^2 = k$

↑  
 columns of  $A$  are all orthonormal  
 $U$  is a square matrix with orthogonal columns  
 ↑  
 $L_2$ -norm of column vectors of  $A$

$$\begin{aligned}\|A^T U S\|_F^2 &= \sum_{i=1}^k \sigma_i^2 \|\tilde{u}_i\|_2^2 \\ &\stackrel{p>k}{\leq} \sum_{i=1}^k \sigma_i^2 \\ &\leq \sum_{i=1}^k \sigma_i^2\end{aligned}$$

Thus, the optimization is upper-bounded by the same quantity achieved by  $A = U_k$ , so that is the solution. Note that if for any  $i, j \leq k$  ( $i \neq j$ ) then the solution need not be unique.

To recap. Any matrix  $X \in \mathbb{R}^{p \times n}$  can be written as  $X = USV^T$

$$U_k = \arg \max_{A \in \mathbb{R}^{p \times k}} \|A^T X\|_F^2$$

where  $U_k$  is the matrix whose columns are the respective top  $k$  left singular vectors of  $X$ .

$$\tilde{\alpha}_i = U_k^T x_i \quad x_i \in \mathbb{R}^p \quad \tilde{\alpha}_i \in \mathbb{R}^k$$

is the reduced dimension representation of  $x$

Put another way,  $\tilde{\alpha}_i$  make up the coordinates to represent  $x$  in the space spanned by  $U_k$ .

## Second Interpretation

representative data of  
we want to find data live in  $P-D$  that live in  $k-D$  ( $k < p$ )

Represent the data in a **low-rank** space

This interpretation is directly connected to the previous one.

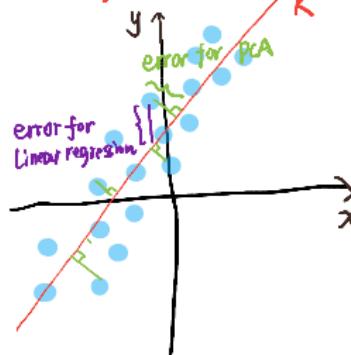
We denote  $\tilde{x}_i = U_k \tilde{\alpha}_i = U_k U_k^T x_i$ . Thus,  $\tilde{x}_i$  is just the projection of  $x$  onto the space spanned by  $U_k$ .

# Illustration

We believe true signal live in low dimensional space

Noise are just added noise - sense

PCA: 2-D data project to 1-D



VS. Linear regression: minimize  $\frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2$   
minimizing prediction error

$$\text{PCA: } \max_{\mathbf{A}^T \mathbf{A} = \mathbf{I}} \sum_{i=1}^n \|\mathbf{A} \mathbf{A}^T \mathbf{x}_i\|_2^2 = \min_{K \in \mathbb{R}^{n \times k}} \sum_{i=1}^n \|\tilde{\mathbf{x}}_i - \mathbf{x}_i\|^2$$

minimizing reconstruction error

We can also set this interpretation as an optimization

Find  $A$  that minimize reconstruction error

$$k = \underset{K, \tilde{x}_i \in K}{\operatorname{argmin}} \sum_{i=1}^n \|\tilde{x}_i - x_i\|_2^2 \text{ equiv } U_k = \arg \min_{A \in \mathbb{R}^{p \times k}} \sum_{i=1}^n \|AA^T x_i - x_i\|_2^2$$

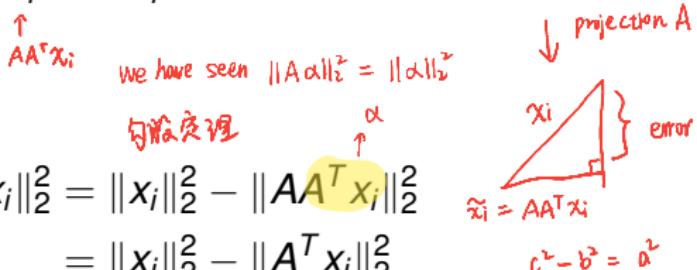
$A^T A = I$

reconstruction error

So we wish to find the orthogonal  $k$  dimensional projection that minimizes the error between  $\tilde{x}_i$  and  $x_i$ .

Can verify that

$$\begin{aligned} \arg \min & \|AA^T x_i - x_i\|_2^2 = \|x_i\|_2^2 - \|AA^T x_i\|_2^2 \\ & = \|x_i\|_2^2 - \|A^T x_i\|_2^2 \end{aligned}$$



Thus, the minimization above is exactly equivalent to the maximization we looked at in the previous interpretation.

$$\arg \max_A \|A^T x_i\|_2^2$$

## Third Interpretation

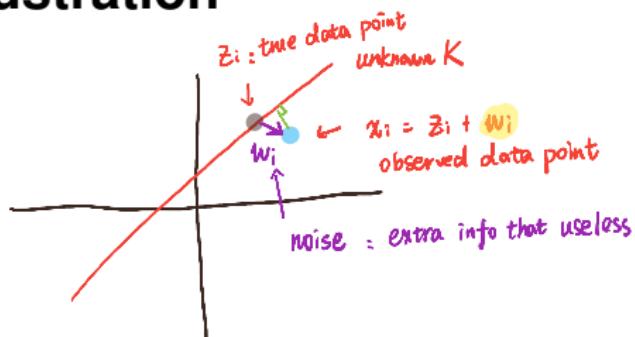
This interpretation is a **data-denoising**. It is also very connected to the previous interpretation

A Probabilistic Model

Gaussian noise

Suppose that our data  $x_i = z_i + w_i$  where  $z_i \in K$  for some **unknown**  $k$ -dimensional space  $K$  and  $w_i \sim N(0, \sigma^2 I)$ .

# Illustration



We can show the maximum likelihood of estimate z<sub>i</sub> is just projection x<sub>i</sub> onto K

But we need find K

We wish to reduce the noise and uncover the “latent” space  $K$ . We can frame this as an optimization.

The maximum likelihood estimate of  $z_1, \dots, z_n$

$$K, z_i = \arg \min_{K, z_1, z_2, \dots, z_n \in K} \sum_{i=1}^n \|z_i - x_i\|_2^2$$

this optimization is equivalent to  $\arg \min_{K, z_1, z_2, \dots, z_n} \sum_{i=1}^n \|AA^T z_i - x_i\|_2^2$   $\arg \min_{K, z_1, z_2, \dots, z_n} \sum_{i=1}^n \|\tilde{z}_i - x_i\|_2^2$

Note that we are optimizing over  $K$  and  $z_i$ . Let  $A \in \mathbb{R}^{p \times k}$  be a matrix whose columns are orthogonal and span the space  $K$ . Let  $z_i = A\alpha_i$ .

$$z_i = \arg \min_{A \in \mathbb{R}^{p \times k}, \alpha_i \in \mathbb{R}^k} \sum_{i=1}^n \|A\alpha_i - x_i\|_2^2$$

One can verify that

$$\begin{aligned} \|A\alpha_i - x_i\|_2^2 &= \|A\alpha_i - AA^T x_i - (I - AA^T)x_i\|_2^2 \\ &= \underbrace{\|A\alpha_i - AA^T x_i\|_2^2}_{0} + \underbrace{\|(I - AA^T)x_i\|_2^2}_{\text{independent of } \alpha} \end{aligned}$$

The optimization then becomes

$$A, \alpha_i = \arg \min_{A \in \mathbb{R}^{p \times k}, \alpha_i \in \mathbb{R}^k} \sum_{i=1}^n \|A\alpha_i - AA^T x_i\|_2^2 + \|(I - AA^T)x_i\|_2^2$$

But the second term is independent of  $\alpha$ . For a fixed  $A$ , the optimal choice of  $\alpha_i = A^T x_i$  can be found in closed form using linear regression/orthogonality/etc...  $\|A\alpha_i - AA^T x_i\|_2^2 = \|\alpha_i - A^T x_i\|_2^2 = \|\alpha_i - \alpha_i\|_2^2 = 0$

Thus, the first term is zero, so the optimization becomes

$$U_k = \arg \min_{A \in \mathbb{R}^{p \times k}} \sum_{i=1}^n \|(I - AA^T)x_i\|_2^2$$

This optimization yields exactly the same optimization as the second interpretation.

$$U_k = \arg \min_{A \in \mathbb{R}^{p \times k} | AA^T = I} \sum_{i=1}^n \|AA^T x_i - x_i\|_2^2$$

**Recap second and third interpretations:** If we wish to denoise data by projecting onto a low-dimensional space  $K$  then we optimize over projection  $P_K$  to minimize reconstruction error

$$P_K = \arg \min_{P_K} \sum_{i=1}^n \|P_K x_i - x_i\|_2^2$$

reconstruction error ( $L_2$ -norm square)

The optimal choice is  $P_K = U_k U_k^T$ .

Another thing need to be optimized is the size of  $k$ .

Method ①

Plot singular value of data  $X$ , 用眼看着眼睛 where the point to pick

Note: we can't use cross-validation to pick  $k$  with smallest reconstruction error  
as as  $k \uparrow$  reconstruction error will always ↓

② we have intuition on what level noise is on, then any singular value < threshold might attribute to noise, so just set those singular value to 0

## Fourth interpretation

We think true data are low-rank matrix  
other things are just noise

Low-rank matrix interpretation.

This is maybe best motivated by an example: recommendation systems.

2014 Netflix release a dataset for public competition  
The winner uses low-rank matrix interpretation as workhorse

# Recommendations

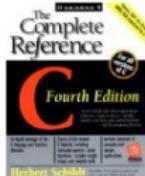
## Recommendations for You in Books



Convex Optimization  
Stephen Boyd, Lieven Vandenberghe  
Hardcover  
★★★★★ (15)  
\$90.00 \$70.64

Why recommended?

[See more recommendations](#)



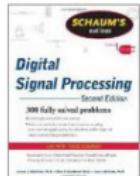
C: The Complete Reference, 4th Ed.  
Herbert Schildt  
Paperback  
★★★★★ (41)  
\$41.99 \$27.71

Why recommended?  
column = rating vector of book 1



Principles of Mathematical Analysis  
Walter Rudin  
Hardcover  
★★★★★ (114)  
\$97.89

Why recommended?



Schaums Outline of Digital Signal Processing  
M. H. Hayes  
Paperback  
★★★★★ (16)  
\$22.00 \$14.96

Why recommended?



Partial Differential Equations...  
Lawrence C. Evans  
Hardcover  
★★★★★ (22)  
\$90.68

Why recommended?

Rating { like / dislike } better  
1-5

Rating Matrix  $X \in \mathbb{R}^{d_{\text{users}} \times d_{\text{books}}}$

Users	book1	book2	book3	book4	book5	book6	book7	book8	book9	...
	3	2	*	1	*	3	*	*	3	...
	*	*	4	*	*	3	*	3	4	...
	*	3	*	3	5	3	*	*	1	...
	3	2	4	1	3	3	3	3	4	...

row : rating vector of user i

missing info  
but might redundant  
Predict them  
① less robust method:  
clustering movie  
average similar movies  
to impute

We might expect that certain books are similar as are certain users.  
Using that information we can make recommendations.

We can let our recommendation matrix  $X$  be such that  $X_{(ij)}$  is user  $i$ 's  
(probably un-observed) rating for item  $j$ .

How can we model this?

One option for modeling is using low-rank latent spaces.

Imagine for item  $j$  we have  $\text{features } m_j \in \mathbb{R}^k$ .  
*unknown*

For instance, for books features could be: fiction, fantasy, math, physics, etc...

Then for user  $i$  we might have  $\text{coefficient vector } a_i$  such that  
 $a_i^T m_j \approx X_{(ij)}$ . This is exactly the **linear model** for linear regression!

*rating data      inner product represent how similar objects are*  
Each entry in  $a_i$  represents how much a user likes a certain feature of the object. For example with books the features listed above. So,  $a_i$  is like the  $\text{user's feature vector}$ .

So, if the  $m_j$  were known, then we could just use linear regression to estimate  $a_i$  for each user based on the movies they have observed.

Recommendation is just prediction  $\hat{x}_{i\ell} = m_i^T \hat{a}_i$

Then we can just recommend which movie they might like by estimating the user  $i$ 's rating for an unobserved movie  $\ell$  by computing  $m_\ell^T a_i$ , and finding the largest predicted ratings.

What if we do not have features available?

We can just learn them. That is we can learn the latent  $k$  dimensional subspace of features.

More precisely, suppose we have  $d_1$  users and  $d_2$  movies. we can let  $A \in \mathbb{R}^{d_1 \times k}$  be the matrix such that the rows are the user coefficients and  $M \in \mathbb{R}^{d_2 \times k}$  be the matrix such that the rows are the object (books, movies) features.

Then, the rating matrix  $X \approx AM^T$ . Of course, in reality there is noise and things do not perfectly fit the model. Also, for recommendations we have missing entries.

Thus, suppose we have a model  $X = AM^T + W$  where  $W$  is some noise term with bounded Frobenius norm. often take  $W$  as Gaussian. Assume we observed all the entries (we can assume some entries are missing, in this case,  $W$  is Gaussian Matrix with deletions. people have shown this also works) Then we can consider the optimization

$$\arg \min_{A \in \mathbb{R}^{d_1 \times k}, M \in \mathbb{R}^{d_1 \times k}} \|AM^T - X\|_F^2$$

*Frobenius Norm of Least square*

Ambiguity:  $(\frac{A}{2})(2M)^T$  for identifiability

Note that  $AM^T$  is always a rank  $k$  matrix. Therefore, we know we can write it as  $AM^T$  where  $A \in \mathbb{R}^{d_1 \times k}$  and  $A^T A = I$ .

If  $A$  has orthogonal columns

Thus, the solution is exactly as the previous interpretation!

$$\arg \min_{A \in \mathbb{R}^{d_1 \times k}, M \in \mathbb{R}^{k \times d_2}} \sum_{i=1}^{d_1} \|A\alpha_i - x_i\|_2^2$$

*3rd*  
*A $\alpha_i$ : column of original data  $X$*   
*x $_i$ : rating column for each movie i*

*4th*  
*m $_i$ : feature vector of movie i*  
*A $m_i$ : feature vector of all users*

Set  $A = U_k$  and  $M = V_k S_k$ . The role of  $\alpha_i$  from the previous interpretation is played by  $m_i$  in this interpretation.

In General

Suppose we have a matrix  $X \in \mathbb{R}^{d_1 \times d_2}$ . For example user movie ratings.

We wish to find a low-rank matrix approximation of  $X$

$$\hat{X}_k = \arg \min_{Y \mid \text{rank}(Y) \leq k} \|Y - X\|_F^2$$

↓  
AM<sup>T</sup>

$$\left\{ \begin{array}{l} A = U_k \\ M = V_k S_k \end{array} \right.$$

Then,  $\hat{X}_k = U_k S_k V_k^T$

Netflix: This optimization + scaling (70%) to get 10% Recommendation Improvement

# Fifth interpretation

Notice that we have not really talked about probability or statistics in any of the interpretations.

vs. Numerical Interpretation (1~4)

Analogous to Factor Analysis

We now introduce the **Gaussian random vector** interpretation, which looks at the covariance matrix of that data.

Often before PCA, remove mean of column

To that end, suppose we have  $x \sim N(\mathbf{0}, \Sigma)$ . Suppose I wish to find the unit vector  $v$  that maximizes this direction of variance

$$\arg \max_{v \mid \|v\|=1} \mathbb{E}(v^T x)^2$$

That is,  $v$  is the direction such that the variance along its direction is maximized.

$X$  is return     $v$ : How related to every stock    Elevation can be negative  
Risk Analysis in Economics : find direction that maximize variance (vulnerability) (related to Risk in Stocks)  
Prior - Factual

# Illustration

Gaussian Matrix with i.i.d

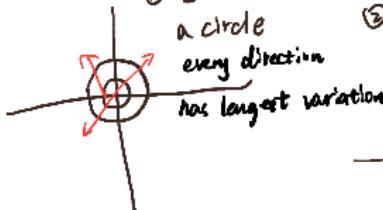
$$\text{Pdf } \frac{1}{2\pi \det(\Sigma)^{\frac{1}{2}}} \exp\left(-\frac{x^T \Sigma^{-1} x}{2}\right)$$

①  $\Sigma = I$

a circle

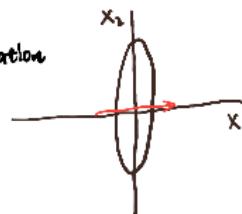
every direction

has longest variation



② general eg.  $\Sigma = \begin{pmatrix} 2 & 0 \\ 0 & 1 \end{pmatrix}$

Ellipse



Now variance change fastest in direction  $x_1$

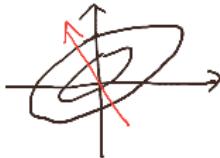
Ellipses align with axis

③ Ellipse don't align with axes

$$\text{Rotate data } \tilde{X} = UX \quad \Sigma = \tilde{X} \tilde{X}^T$$

$$E(\tilde{X} \tilde{X}^T) = U \Sigma' U^T \rightarrow \text{a diagonal matrix}$$

Eigenvalue decomposition of  $\Sigma$  is  $U \Sigma' U^T$



$$E[v^T x^2] = E[(v^T x)(x^T v)] = v^T E(x x^T) v$$

scalar

$$= v^T E[x x^T] v$$

*v is a constant, a linear combination ( $\sum \dots$ )*

$$\arg \max_{v \mid \|v\|=1} \mathbb{E}(v^T x)^2 = \arg \max_{v \mid \|v\|_2=1} v^T \Sigma v$$

We recall that  $\Sigma$  is a covariance matrix. Therefore, it is symmetric and positive semi-definite.  $\Rightarrow$  all eigen values of  $\Sigma \geq 0$  *always*

*why :  $\Sigma$  is symmetric*

Thus, it has an eigen-decomposition  $\Sigma = U D U^T$  where  $D$  is diagonal with positive entries, and  $U$  is a square orthogonal matrix, so  $U^T U = U U^T = I$ .

*For a symmetric and 半正定 matrix SVD = ED*

Note that this is exactly the SVD of  $\Sigma$  as well. We will show that the optimal choice of  $v$  is  $u_1$  column of  $U$

Case ②

$$\begin{aligned}\arg \max_{\nu} \nu^T \Sigma \nu &= \nu^T \left( \sum_{i=1}^p D_{(ii)} u_i u_i^T \right) \nu \\ &= \sum_{i=1}^p D_{(ii)} (\nu^T u_i)^2\end{aligned}$$

Both,  $\nu_i$  and  $u_i$  are unit vectors, so  $0 \leq (\nu^T u_i)^2 \leq 1$ . Also  
 $\sum_i (\nu^T u_i)^2 = 1$

The optimal choice of  $\nu$  is therefore  $\nu = u_1$ , the first eigenvector of  $\Sigma$ !

Put more weight onto component with longest  $D_{(ii)}$  (eigen value)

$P_k$  is projection to  $k$ -D space with highest Gaussian variance

As before we can generalize to  $k$ -dimensional projections.

$$K = \arg \max_K \mathbb{E} \|P_K \overset{x}{X}\|_2^2 = \arg \max_K \text{tr}(P_K \Sigma P_K^T) \quad \downarrow P_K = AA^T$$

$$\|P_K x\|^2 = \underset{P_K}{\overbrace{x^T P_K^T P_K x}} = \text{tr}(P_K^T P_K \overset{x^T}{x}) = \arg \max_{A \in \mathbb{R}^{p \times k} | A^T A = I} \text{tr}(AA^T \Sigma A^T A)$$

$$= \arg \max_{A \in \mathbb{R}^{p \times k} | A^T A = I} \text{tr}(AA^T U D^{1/2} D^{1/2} U^T A^T A)$$

$$\begin{aligned} &= \arg \max_{A \in \mathbb{R}^{p \times k} | A^T A = I} \text{tr}(\underbrace{\|AA^T UD^{1/2}\|_F^2}_F) \\ &= \arg \max_{A \in \mathbb{R}^{p \times k} | A^T A = I} \text{tr}(\|A^T UD^{1/2}\|_F^2) \quad \text{if } A \text{ is orthogonal matrix} \\ &= \arg \max_{A \in \mathbb{R}^{p \times k} | A^T A = I} \text{tr}(\|A^T UD^{1/2}\|_F^2) \end{aligned}$$

This is exactly the same optimization as in Interpretation 1! So,  
 $A = U_k$ , the top  $k$  eigenvectors of  $\Sigma$ !

Suppose we do not know what  $\Sigma$  is but we still wish to estimate  $v$ ? Suppose we have samples  $x_i \sim N(0, \Sigma)$ . Just use empirical averages instead! Rather than  $E(v^T x_i)^2$

$$\arg \max_v \frac{1}{n} \sum_{i=1}^n (v^T x_i)^2$$

empirical variance

But this is exactly the first interpretation of PCA!

Actually, we can rewrite it as

$$\frac{1}{n} \sum_{i=1}^n v^T x_i x_i^T v = v^T \hat{\Sigma} v$$

where  $\hat{\Sigma}$  is the empirical covariance matrix of the data!

often big data, if we just want to know left singular vectors, we only need to compute  $\Sigma$  covariance matrix

More generally can consider  $K$ -dimensional orthogonal projections

$$K = \arg \max_K \frac{1}{n} \sum_{i=1}^n \|P_K x_i\|_2^2$$

Recap: PCA for data generated according to a Gaussian model simply finds the top  $k$  eigenvectors of the sample covariance matrix.

△

Σ

That's one motivation. Of course, the data need not be Gaussian.

So another way to solve PCA is to compute the sample covariance matrix of the data (with or without re-centering) and compute the top eigenvectors.

# Constructing the SVD from eigenvectors

The above shows how we can compute the PCA using eigenvectors from the sample covariance matrix of the data.

We can use that motivation to show how one can construct the SVD.

Given a matrix  $X \in \mathbb{R}^{p \times n}$  let  $Y = XX^T$ . Clearly,  $Y$  is both positive semi-definite and symmetric. Thus, we can write  $Y = UDU^T$  for its eigen-decomposition.

We will show that  $X = UD^{1/2}V^T$  for an appropriate orthogonal matrix  $V$ .

We deal with the case  $p \leq n$  and  $X$  has rank  $p$ . The general case is a simple extension using projections.

$Y = XX^T \implies D^{-1/2}U^TXX^TUD^{-1/2} = I$ . Possible since  $X$  has rank  $p$ , so  $Y$  has rank  $p$  as well, so  $D_{(ii)} > 0$  for all  $i$ .

Thus,  $X^TUD^{-1/2} = (D^{-1/2}U^TX)^{-1} = (D^{-1/2}U^TX)^T$

Therefore,  $X^TUD^{1/2} = V$  an orthogonal matrix.

Multiplying yields the desired result.

# Conclusion

PCA is a powerful tool for unsupervised learning.

It can be used for visualization, dimensionality reduction, data-denoising, and low-rank matrix approximations.

Given data  $X \in \mathbb{R}^{p \times n}$  where we have  $n$  data points in  $p$  dimensions the top  $k$  principal components of  $X$  are denoted as  $U_k$  which are exactly the top  $k$  left singular vectors of  $X$ .

If we do dimensionality reduction and project  $U_k^T X$ , then we are left with

$$\begin{aligned} U_k^T X &= U_k^T U S V^T \\ &= S_k V_k^T \end{aligned}$$

So the lower-dimensional coordinates  $\hat{\alpha}_i$  of example  $x_i$  are exactly  $[\hat{\alpha}_i]_{(j)} = \sigma_j [V_k]_{(ij)}$ , i.e.  $\sigma_j$  times the  $i^{th}$  component of the  $j^{th}$  right singular vector.