

MATH 312

LECTURE 23

: PCA / MARKOV Chains.

III.1 Principal Components Analysis.

Now we will have a bunch of data points, which we write in a matrix:

$$X = \begin{bmatrix} \vec{x}_1 & \dots & \vec{x}_n \end{bmatrix}, \quad \vec{x}_i = \begin{bmatrix} x_{i1} \\ \vdots \\ x_{im} \end{bmatrix} \in \mathbb{R}^m \rightarrow \text{Each row corresponds to a variable.}$$

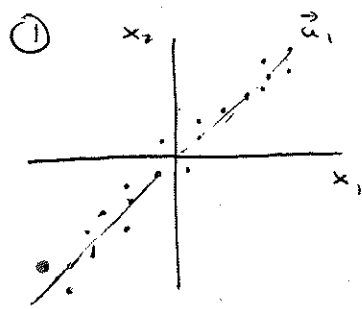
↑ ↗
points, samples, measurements

(High dimensional data)

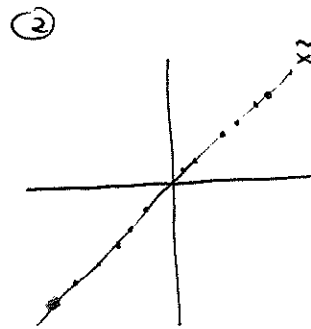
Examples: The cols. can be all the students at Penn., while the rows might be their height, age, major, ...

- The goal is to be able to distinguish or classify the points with much less than m variables, without losing much information.

Imagine we have points in \mathbb{R}^2 ($m=2$) that look as follows:



project these
points onto
 \vec{u}_1



In ①, each point is characterised by two variables (x_1, x_2).

In ②, we lose some information (red and green points cannot be distinguished in ②), however, most of them are still split apart.

→ Most variance among data remains in ②.

→ In ②, we only need one variable (z) = the distance along \vec{u}_1 .

↳ each original point \vec{x} (cols. of X) are now identified by their projection onto \vec{u}_1 .

• What is the relation with SVD?

$$\text{Write } X = \begin{bmatrix} \vec{x}_1 & \dots & \vec{x}_n \end{bmatrix} = \sigma_1 \vec{u}_1 \vec{v}_1^T + \dots + \sigma_r \vec{u}_r \vec{v}_r^T$$

$$\text{Notice that: } A \vec{v}_i = \sigma_i \vec{u}_i \Rightarrow \underbrace{A^T A \vec{v}_i}_{= \sigma_i^2 \vec{v}_i} = \sigma_i A^T \vec{u}_i \Rightarrow \| \sigma_i \vec{v}_i^T = \vec{u}_i^T A \|$$

so $X = \vec{a}_1 (\vec{a}_1^T X) + \dots + \vec{a}_r (\vec{a}_r^T X)$, where

$$\vec{a}_i (\vec{a}_i^T X) = \begin{bmatrix} \vec{a}_i \vec{a}_i^T \vec{x}_1 & \dots & \vec{a}_i \vec{a}_i^T \vec{x}_n \end{bmatrix}$$

↑
|| projection of the data points $\vec{x}_1, \dots, \vec{x}_n$ onto \vec{a}_i . ||

In our 2d example,

$$X = \vec{a}_1 (\vec{a}_1^T X) + \vec{a}_2 (\vec{a}_2^T X)$$

this gives picture ② (each column gives where the original points go onto the line spanned by \vec{a}_1).

Remark: \vec{a}_1 gives the line ~~of best~~ through the origin of best fit.

\vec{a}_1, \vec{a}_2 gives the plane through the origin of best fit. (*)

⋮

→ The \vec{a} 's are called principal components. They provide a new basis (new variables) to represent our data.

(*) We have to subtract the mean to each row to obtain good results.

Example: Database of 216 patients, 121 with ovarian cancer and 95 without it.

Each patient identified through 4000 genes.

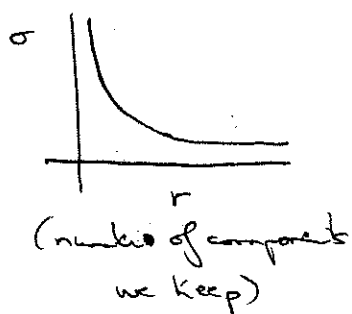
$$X = \left[\begin{array}{c} \vec{x}_1, \dots, \vec{x}_{216} \end{array} \right] \left. \vphantom{\begin{array}{c} \vec{x}_1, \dots, \vec{x}_{216} \end{array}} \right\} \text{4000 rows}$$

We want to see if there ^{are} some common factors ~~and~~ that predict the ovarian cancer.

It seems likely that many of the genes are correlated.

We perform the SVD \rightarrow we might obtain a lot of \vec{u} 's.
(up to 4000!)

How many are important? \rightarrow we can look the weights σ^2 .
(graph)



Let's plot the patients in the basis given by $\{\vec{u}_1, \vec{u}_2, \vec{u}_3\}$,
that, we project each \vec{x}_i into the subspace spanned by \vec{J}

$$\vec{x}_i = \begin{bmatrix} x_{i1} \\ x_{i2} \\ x_{i3} \\ \vdots \\ x_{i4000} \end{bmatrix} = \vec{u}_1 (\vec{u}_1^T \vec{x}_i) + \vec{u}_2 (\vec{u}_2^T \vec{x}_i) + \vec{u}_3 (\vec{u}_3^T \vec{x}_i) \Rightarrow$$

vector of 4000 components

$$\vec{x}_i \Big|_{\{\vec{u}_1, \vec{u}_2, \vec{u}_3\}} = \begin{bmatrix} \vec{u}_1^T \vec{x}_i \\ \vec{u}_2^T \vec{x}_i \\ \vec{u}_3^T \vec{x}_i \end{bmatrix} \quad \text{This can be plotted in a 3d graph!}$$

(see Matlab...)

Example: Eigenfaces

Database 1600 pictures \rightarrow 25 people \times 64 pictures/person.

Each picture is $192 \cdot 168$ pixels.
= 32256

$$X = \begin{bmatrix} \vec{x}_1 & \dots & \vec{x}_{1600} \end{bmatrix} \left. \vphantom{\begin{bmatrix} \vec{x}_1 & \dots & \vec{x}_{1600} \end{bmatrix}} \right\} 32256 \text{ rows (pixels)}$$

each image is stored as a column vector.

• We do the SVD of X - mean rows.

↳ The \vec{u} 's are called "eigenfaces" \rightarrow they keep most common features in human faces.

↳ They provide a basis for the space of "faces".

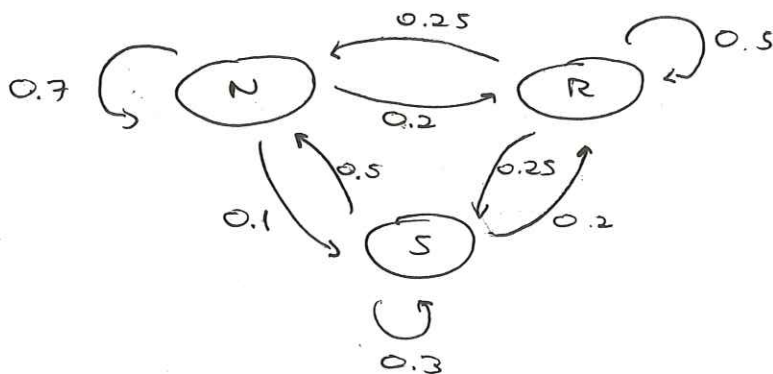
A new face ("test face") can be projected into the eigenspace given by the \vec{u} 's (building using X , the "training set").

↳ Now we can see if it is close to one of the faces in X or not \rightarrow face recognition.

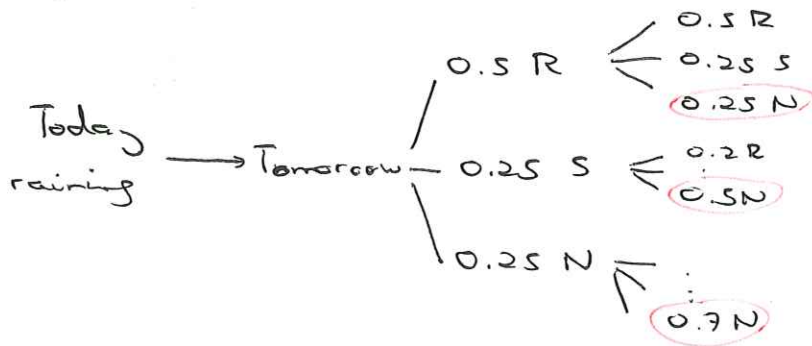
$$\text{testface} = \begin{bmatrix} \vdots \\ \vdots \\ \vdots \end{bmatrix} \rightsquigarrow \underbrace{U_r U_r^T \text{testface}}_{\substack{\text{approximation of testface} \\ \text{using first } r \text{ eigfaces.}}} \quad (U_r = [\vec{u}_1, \dots, \vec{u}_r])$$

New Chapter: Markov Chains.

Consider the following silly model for weather prediction:



If today is raining, what is the probability that the day after tomorrow is a nice day?



So after tomorrow,

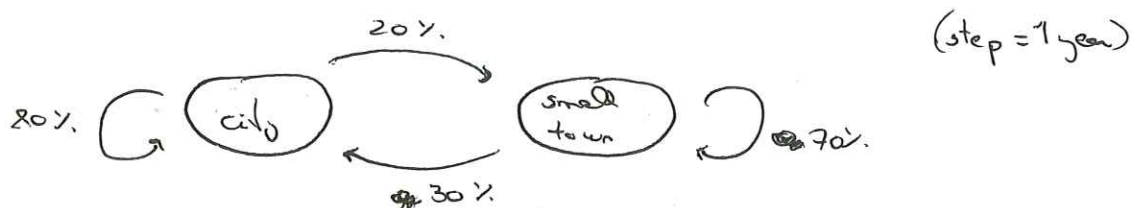
$$p = (0.5)(0.25) + (0.25)(0.5) + (0.25)(0.7) = 0.425$$

What is the probability that it will snow in a week, if today is a nice day?

Moreover, along the years, which fraction of days are nice, raining or snowy?

→ The answers are not straight forward, but we will see they are easy to compute.

• Let's use a different simpler example:



If today ~~50%~~ ^{50%} of people live in cities, what is the distribution next year?

$$\text{In cities: } (0.5) \cdot (0.8) + (0.5) \cdot (0.3) = 0.55 \rightarrow 55\%$$

$$\text{In towns: } (0.5) \cdot (0.7) + (0.5) \cdot (0.2) = 0.45 \rightarrow 45\%$$

And next one:

(same numbers)

$$\text{In cities: } (0.55) \cdot (0.8) + (0.45) \cdot (0.3) = 0.575 \rightarrow 57.5\%$$

$$\text{So in towns } \rightarrow 42.5\%$$

Create the "state" vector $\vec{x}(k) = \begin{bmatrix} x_1(k) \\ x_2(k) \end{bmatrix}$, where

$x_1(k) = \%$ in cities in year k
 $x_2(k) = \%$ in town in year k
 (of course, $x_1(k) + x_2(k) = 1$ $\forall k$).

How to find $\vec{x}(k+1)$?

$$\begin{aligned} x_1(k+1) &= x_1(k) \cdot (0.8) + x_2(k) \cdot (0.3) \\ x_2(k+1) &= x_1(k) \cdot (0.2) + x_2(k) \cdot (0.7) \end{aligned} \Rightarrow \vec{x}(k+1) = \begin{bmatrix} 0.8 & 0.3 \\ 0.2 & 0.7 \end{bmatrix} \vec{x}(k)$$

transition matrix

Notice: columns add up to 1.

Q: What is the transition matrix for the weather model?

State vector: $\vec{x} = \begin{bmatrix} x_S \\ x_N \\ x_R \end{bmatrix}$

$$\vec{x}(k+1) = P \vec{x}(k) \rightarrow P = \begin{bmatrix} 0.3 & 0.1 & 0.25 \\ 0.5 & 0.7 & 0.25 \\ 0.2 & 0.2 & 0.5 \end{bmatrix}$$

P_{ij} = probability of going in one step from i to j

Check previous result:

$$p = P^2 \begin{bmatrix} 0 \\ 0 \\ 1 \end{bmatrix} = P \begin{bmatrix} 0.25 \\ 0.25 \\ 0.5 \end{bmatrix} = \begin{bmatrix} (0.3)(0.25) + (0.1)(0.25) + (0.5)(0.25) \\ (0.5)(0.25) + (0.7)(0.25) + (0.5)(0.25) \\ \dots \end{bmatrix}$$

↑ probability after two days
 ↑ today raining
 ↑ probability of nice day after 2 days.

• Def. Markov (or stochastic) matrix

An $n \times n$ matrix is Markov if:

1) All columns add up to 1.

2) All entries are ≥ 0 .

• The sequence of states $\{\vec{x}(1), \vec{x}(2), \dots, \vec{x}(k), \dots\}$ defined by a Markov matrix is called a Markov chain (discrete)

Remark: key feature of this model: Every state only depends on the previous step: $\vec{x}(k+1) = A \vec{x}(k)$.

Example: 1) The evolution of a particle, defined through its position and velocity, is a Markov process (by Newton's Law, if we know position and velocity now, we know the future).

\leftarrow at t, T considering classical mechanics.

\Rightarrow Poker is not a Markov process: you use (or should use) the knowledge about which cards have been appearing from the start.

Some questions:

1) Each state $\vec{x}(k)$ describes the distribution at that moment

(% in cities, % in towns, and so on).

So $\vec{x}(k+1)$ needs to be a % distribution too, i.e., the components should add up to 1.

Is it true? That is, for a Markov matrix, does $A\vec{x}$ add up to 1 if \vec{x} does?

Let's see that yes:

Notice that the addition of the components of a vector is given

by

$$[1 \dots 1] \vec{x} = [1 \dots 1] \begin{bmatrix} x_1 \\ \vdots \\ x_n \end{bmatrix} = x_1 + \dots + x_n = 1$$

↑
assumption (%...)

But then,

$$[1 \dots 1] A \vec{x} = [1 \dots 1] \begin{bmatrix} \vec{a}_1 & | & \dots & | & \vec{a}_n \end{bmatrix} \vec{x} = [1 \dots 1] \vec{x} = x_1 + \dots + x_n = 1$$

↑

Markov matrix \Rightarrow $\begin{cases} [1 \dots 1] \vec{a}_1 = \text{addition components first column} = 1 \\ [1 \dots 1] \vec{a}_2 = \text{" " second column} = 1 \\ \vdots \\ \text{All cols. of } A \text{ add up to } 1! \end{cases}$