16B  Prof ANANT SAHAI
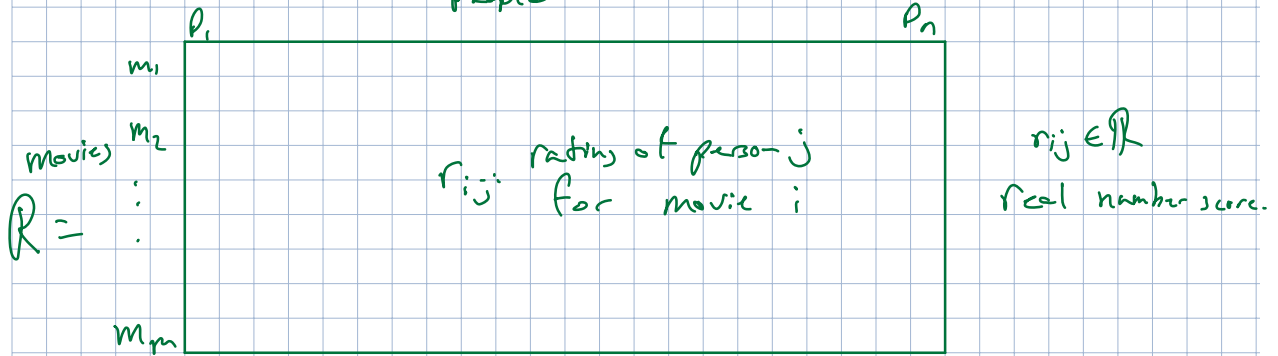
Today: Finish Movie Model
   "Checking your work" — validating with other data
   Dimensionality Reduction for classification ← happening in lab!
   Justifying this approach
   Nonlinear Models & Local Approximation
   Control around equilibria

Illustrative conceptual example: Movie Ratings,

People

$P_1$ ............................................ $P_n$

$m_1$
$m_2$
movies) $\vdots$

$R =$    $\vdots$

$m_m$

$r_{ij}$ : ratings of person-$j$ for movie $i$

$r_{ij} \in \mathbb{R}$
real number score.

16-philosophy: Make a model,
Learn the model from data,
Use the model to make predictions.

Can Modify our model : Each movie has a vector $\vec{g}$
   to reflect what is possible      of qualities associated

Each person has a vector $\vec{s}$
   of sensitivities.

SVD Gave Us: (Outer-Product Form)

$$R = \sum_{k=1}^{r} \sigma_k \vec{u}_k \vec{v}_k^T \qquad r(\vec{g}, \vec{s}) = \sum_{k=1}^{4} \sigma_k g_k s_k \quad \text{numbers}$$

Truncate to 4 terms: $\hat{R} = \sum_{k=1}^{4} \sigma_k \vec{u}_k \vec{v}_k^T$

↖ pull this out instead of
   trying to allocate it.

Now our learned model has $\vec{g}_i = \begin{bmatrix} (\vec{u}_1)_i \\ (\vec{u}_2)_i \\ (\vec{u}_3)_i \\ (\vec{u}_4)_i \end{bmatrix}$

Recall $\vec{u}$'s are $m$ long,

$\vec{v}$'s are $n$ long.

$\vec{s}_j = \begin{bmatrix} (\vec{v}_1)_j \\ (\vec{v}_2)_j \\ (\vec{v}_3)_j \\ \end{bmatrix}$

How to use this for prediction?

e.g. Given a new person, how can we extract their sensitivities $\vec{s}$?

Assume we're given $\vec{r}$ : how does this person _rate_ the m movies?

We need to go from $\vec{r}$ to $\hat{s}$

$\nearrow$        $\leftarrow$ estimated
ratings on       sensitivities
known movies      for this person--

Can do this by least squares!

<u>Aside</u> let's look at a person we already know.

$j^{th}$ person : $\vec{r_j} = \sum_{k=1}^{4} \sigma_k \; \vec{u_k} \; (\vec{v_k}^T)_j$  $\leftarrow$ This is a scalar: $j^{th}$ position in vector $\vec{v_k}$

How can I get $\vec{s_j}$ from $\vec{r_j}$ ?

$$\left( \vec{s_j} \right)_k = \frac{1}{\sigma_k} \vec{u_k}^T \vec{r}$$ : I can get the $k^{th}$ sensitivity by projecting $\vec{r}$ onto the $k^{th}$ principal component $[\vec{u_k}]$
<u>and scaling</u> by $\frac{1}{\sigma_k}$ to fit our model

Similarly : We can get $\vec{g_i}$ from $\vec{q_i}^T = [r_{i,1} \; r_{i,2} \; \cdots \; r_{i,n}]$

$$B_y \quad \left( \vec{g_i} \right)_k = \frac{1}{\sigma_k} \vec{v_k}^T \vec{q_i}$$

$\Longrightarrow$ So for a new person with ratings $\vec{r}$,

Can get $\hat{s} = \begin{bmatrix} \frac{1}{\sigma_1} \vec{u_1}^T \vec{r} \\ \frac{1}{\sigma_2} \vec{u_2}^T \vec{r} \\ \vdots \\ \frac{1}{\sigma_4} \vec{u_4}^T \vec{r} \end{bmatrix}$

Note:
Dividing by $\sigma_i$ is a warning flag.
$\Rightarrow$ Beware of small $\sigma_i$.

And similarly for a new movie. $\hat{g} = \begin{bmatrix} \frac{1}{\sigma_1} \vec{v_1}^T \hat{q} \\ \vdots \\ \frac{1}{\sigma_4} \vec{v_4}^T \hat{q} \end{bmatrix}$

Can mess things up.

So. Given a new person (who comes with ratings $\vec{r}$ on the m known movies)

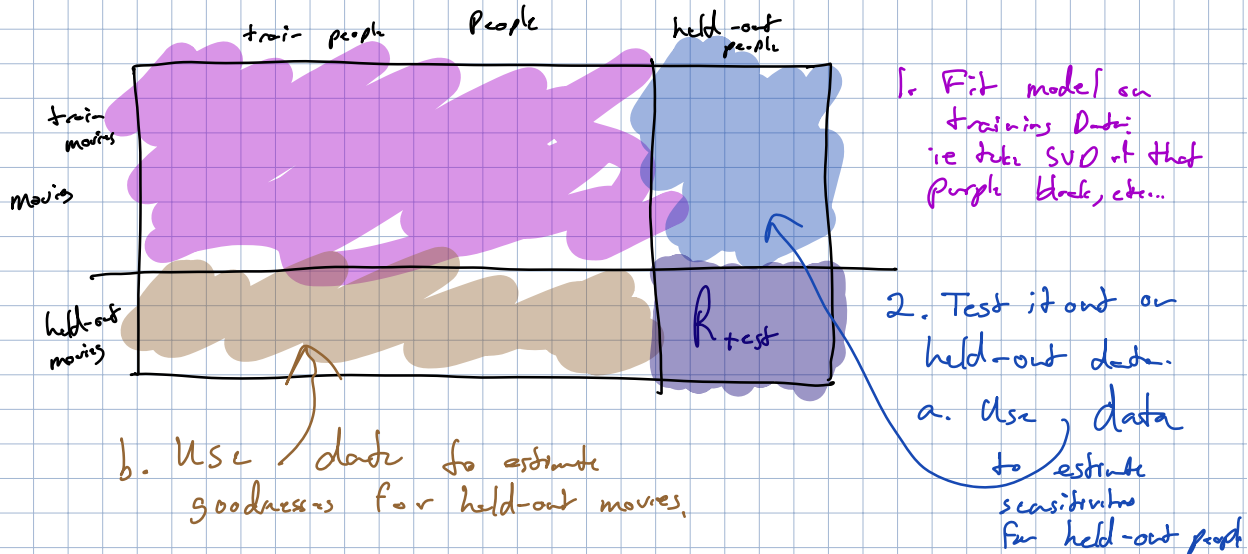and a new movie (which comes with ratings $\bar{q}$ by the n known people)

we can predict the rating $r = \sum_{k=1}^{4} \sigma_k \hat{g}_k \hat{s}_k$

↑ ↑

Fit as above.

Question: How do we "check our work"

e.g. How do we justify that 4 is the right number of components to use? ← 4 qualities are important.

From Model-order-selection problem on the HW, we know we need to use some extra data.



train-people    People    held-out people

train-movies

movies

held-out movies

$R_{test}$

1. Fit model on training Data: ie take SVD of that purple block, etc...

2. Test it out on held-out data.
   a. Use data to estimate sensitivities for held-out people

b. Use data to estimate goodnesses for held-out movies.

c. Use $R_{test}$ to check quality of predictions.
   Compare predictions to actual ratings. (Use Mean-Squared Error)

This procedure, sweeping through candidate #s of components can justify model order choice.

Recall: This is the "silver standard" (like a "unit test" for learning from data)

The Gold Standard is always an "integration test" or "functional test." See if it works in practice.

<u>Step back to classification.</u>

Wanted a way to reduce the dimensionality of data.

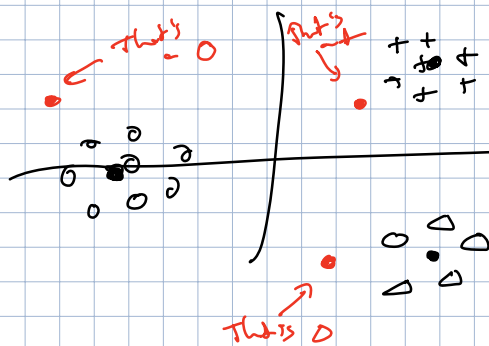$\Longrightarrow$ Want to keep most important parts.

Can use the PCA approach: Project onto $\vec{u}_i$ if data is columns

or $\vec{v}_i$ if data is rows.

<u>Heuristic</u> $\leftarrow$ Based on <u>hope</u> that accurately capturing the pattern in the data helps with classification.

By reducing the dimensionality, we get an easier problem.

In particular, in 2d or 3d can visualize.

e.g. Three Categories $+$, $\triangle$, $\circ$



That's ≈ O

Puts O

Given.

How would we categorize a new point?

This O

First & Most Basic Approach: Which group is the new point closest to?

An approach: Compute the centers of each group, and compare distances to centers.

$\angle AB$ will do this approach.

How can we justify our approach to learning the pattern by SVD.

We are viewing $\hat{R} = \sum_{i=1}^{\ell} \sigma_i \vec{u}_i \vec{v}_i^\sigma$ as an approximation to $R$.

In what sense is this a good approximation?!

$\longrightarrow$ $\hat{R}$ certainly has its columns all in an $\ell$-dim subspace

So $\hat{R}$ is "simple" $\leftarrow$ low-dimensional.

Is $\hat{R}$ "close" to $R$?

We will our favorite squared-error.

Want $\|R - \hat{R}\|_F^2$ to be as small as possible.

$\underbrace{\qquad\qquad\qquad}$ Frobenius Norm

Treat $R$ and $\hat{R}$ like big vectors.

Underlying model: $r_{ij} = \left[ \displaystyle\sum_{k=1}^{\ell} \sigma_k (\vec{g}_k)_i (\vec{s}_k)_j \right] + \underset{\uparrow}{w_{ij}}$

Hopefully small.

Key Question: Why does the SVD help??

$$\|R - \hat{R}\|_F = \|\overbrace{U\Sigma V^T}^{\text{Full SVD}} - \hat{R}\|_F$$

$\begin{array}{l}\text{Premultiply}\\ \text{by}\\ U^T\end{array} = \|\Sigma V^T - U^T\hat{R}\|_F$

$\text{transpose} = \|V\Sigma^T - \hat{R}^T U\|_F$

$\begin{array}{l}\text{Premultiply}\\ \text{by } V^T\end{array} = \|\Sigma^T - V^T\hat{R}^T U\|_F$

$\text{transposing} = \|\Sigma - \underbrace{U^T\hat{R}V}\|_F$
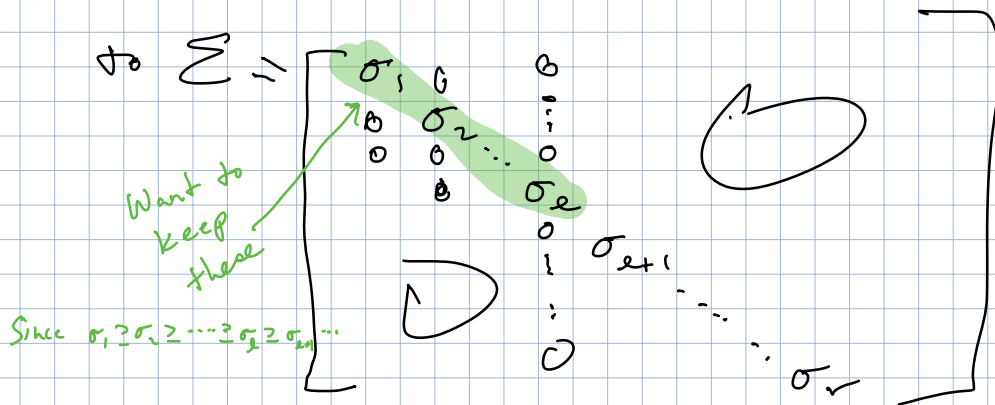
still a rank-$\ell$ matrix.
or lower

$\ell$

$$\text{since } \hat{R} = \sum_{k=1}^{r} \sigma_k (\vec{g_k}) \vec{s_k}^T$$

$$U^T \hat{R} V = \sum_{k=1}^{\ell} \sigma_k (U^T \vec{g_k})(\vec{s_k}^T V)$$

$\underbrace{\qquad}$ just a vector.

$\rightarrow$ There are l.t them.

$\Rightarrow U^T \hat{R} V$ has rank at most $\ell$.

So, $U^T \hat{R} V$ should be the best rank $- \ell$ approximation

to $\Sigma =$



Want to keep these

Since $\sigma_1 \geq \sigma_2 \geq \cdots \geq \sigma_\ell \geq \sigma_{\ell+1} \cdots$

Intuitively clear! should just keep the $\ell$ biggest #s.

Want $U^T \vec{g_k}$ to be the k-th column of the identity. So $\vec{g_k} = \vec{u_k}$. Similarly we want $\vec{s_k}^T V$ to be the $k^{th}$ row of the identity. So $\vec{s_k} = \vec{v_k}$.

The full theorem is called the Eckart-Young-Mirsky Thm.

$\Rightarrow$ The best rank $-\ell$ approximation to a matrix

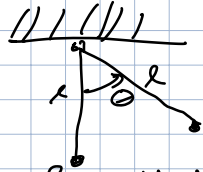is given by keeping the first $\ell$ terms in the SVD.

---

Step Way Back: All of our 16AB models have been linear

except Switches / comparators $\in$ nonlinear.

Real World is largely nonlinear.



$\leftarrow$ Actually Smoothly nonlinear.

Almost all mechanical systems are nonlinear.

Consider a pendulum:  e.g. Get sines & cosines in equations governing motion

Need a tool & philosophy for attacking nonlinear problems

Tool: Local Approximation. Just pretend things are linear.

Taylor Series: $f(x) = f(x_0) + f'(x_0)(x-x_0)$

$$+ \frac{f''(x_0)(x-x_0)^2}{2}$$

$$+ \cdots$$

When $x$ is close to $x_0$, the higher-order terms contribute very little.

$\Longrightarrow$ lump all dropped terms into the disturbance

Consider: $f(x, u) = 32 + x^3 u^2$

expand around $(x_0, u_0)$

$$f(x_0 + \delta x, u_0 + \delta u) = 32 + (x_0 + \delta x)^3 (u_0 + \delta u)^2$$

$$= 32 + (x_0^3 + 3x_0^2 \delta x + 3x_0 (\delta x)^2 + \delta x^3)$$
$$\cdot (u_0^2 + 2u_0 \delta u + (\delta u)^2)$$

If $\delta x$ and $\delta u$ are both tiny, then $(\delta x)(\delta u)$

$(\delta x)^2$, etc.

All higher powers are extremely small.

$$f(x_0 + \delta x, u_0 + \delta u) \approx 32 + x_0^3 u_0^2 + \boxed{3x_0^2 u_0^2} \delta x$$

$$+ \boxed{2 x_0^3 u_0} \delta u$$

$$+ \frac{\text{super}}{\text{tiny}} \leftarrow \text{lump this into disturbance}$$

This gives an approximate "linear" model in the neighborhood of $(x_0, u_0)$.

<span style="color:magenta">↑Actually, "affine" since there is the constant term in front</span>