

Low-Rank Decompositions of Matrices

Objectives

1

- Define low-rank decomposition
(matrix factorization)
- Explore applications

Matrices represent many types of information 2

1) Features in classification or modeling

2) User ratings

3) Collections of documents

Bag of words model



"Apple trees
blossom in
May."

word frequency Term-document matrix

word	frequency	Term-document matrix
agency	0	[]
apple	1	[]
blossom	1	[]
car	0	[]
currency	0	[]
happy	0	[]
May	1	[]
politics	0	[]
road	0	[]
tree	1	[]

terms ↓

documents →

Low-rank decompositions emphasize patterns 3

$$\underset{N \times M}{\boxed{\underline{A}}} \approx \underset{N \times P}{\boxed{\underline{T}}} \underset{P \times M}{\boxed{\underline{W}^T}} \quad P < N, P < M \quad \text{rank } \underline{T} = \text{rank } \underline{W} = P \\ \Rightarrow \text{rank}(\underline{T} \underline{W}^T) = P$$

$\underset{N \times P}{\boxed{\underline{T}}} \quad \text{Let } \underline{T} = [\underline{t}_1 \ \underline{t}_2 \dots \underline{t}_P], \underline{w} = [\underline{w}_1 \dots \underline{w}_P]$

$$\underline{T} \underline{W}^T = \left[\begin{array}{c|c|c|c} \hline & \underline{t}_1 & \underline{t}_2 & \cdots & \underline{t}_P \\ \hline \vdots & \vdots & \vdots & \ddots & \vdots \\ \hline & \underline{t}_1 & \underline{t}_2 & \cdots & \underline{t}_P \end{array} \right] \left[\begin{array}{c|c|c|c} \hline & -\underline{w}_1^T & - & - & - \\ \hline & -\underline{w}_2^T & - & - & - \\ \hline & \vdots & \vdots & \ddots & \vdots \\ \hline & -\underline{w}_P^T & - & - & - \end{array} \right] = \sum_{i=1}^P \underbrace{\underline{t}_i \underline{w}_i^T}_{\text{rank-one patterns}} = \sum_{i=1}^P \left| \begin{array}{c|c|c|c} \hline & \underline{t}_1 & \underline{t}_2 & \cdots & \underline{t}_P \\ \hline & \vdots & \vdots & \ddots & \vdots \\ \hline & \underline{t}_i & \underline{t}_2 & \cdots & \underline{t}_P \\ \hline & \vdots & \vdots & \ddots & \vdots \\ \hline & \underline{t}_1 & \underline{t}_2 & \cdots & \underline{t}_P \end{array} \right| = \sum_{i=1}^P \underset{N \times M}{\boxed{\underline{N}}}$$

Finding patterns -

$$1) \min_{\underline{T}, \underline{W}} \|\underline{A} - \underline{T} \underline{W}^T\|$$

singular value decomposition

2) $\underline{A} \approx \underline{T} \underline{W}^T, \underline{T}, \underline{W} \geq 0$
 non-negative matrix factorization

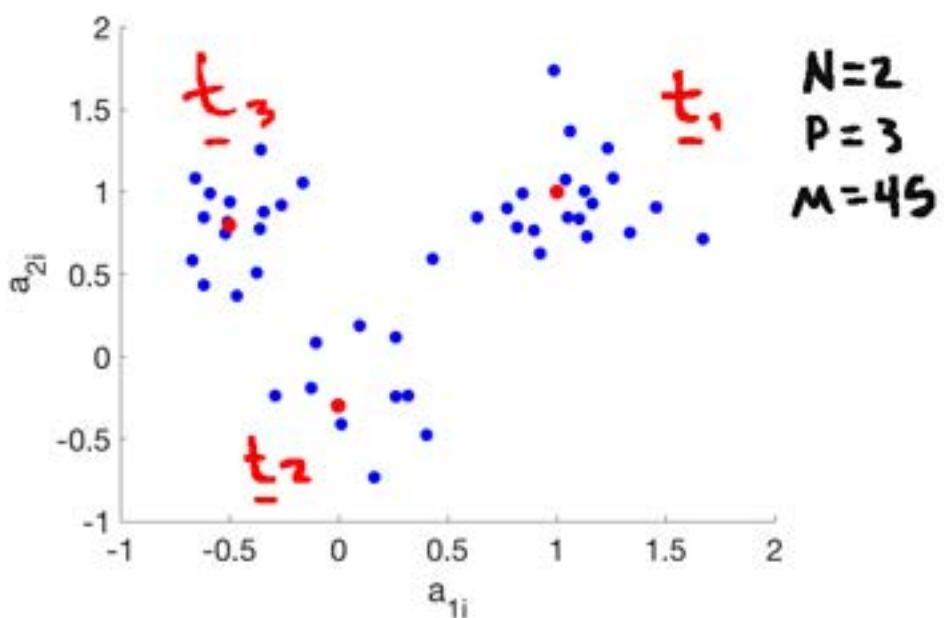
3) $\underline{A} \approx \underline{T} \underline{W}^T$
 each col \underline{W}^T all 0 w.
 single 1
 clustering

Clustering groups similar columns

4

$$\begin{bmatrix} \underline{a}_1 & \underline{a}_2 & \underline{a}_3 & \underline{a}_4 & \dots & \underline{a}_m \end{bmatrix} \approx \begin{bmatrix} \underline{t}_1 & \underline{t}_2 & \underline{t}_3 \end{bmatrix} \begin{bmatrix} 0 & 1 & 1 & 0 & & 0 \\ 1 & 0 & 0 & 0 & \dots & 1 \\ 0 & 0 & 0 & 1 & & 0 \end{bmatrix}$$

$$\Rightarrow \underline{a}_1 \approx \underline{t}_2, \underline{a}_m \approx \underline{t}_2, \underline{a}_2 \approx \underline{t}_1, \underline{a}_3 \approx \underline{t}_1, \underline{a}_4 \approx \underline{t}_3 \dots$$



Group similar documents, customers, products, etc

Many algorithms –
k-means

Low rank models "complete" missing data

5

Jill	
Star Trek	8
Pride + Prejudice	3
The Martian	7
Sense + Sensibility	4
Empire Strikes Back	?

Suppose $\underline{a} = \begin{bmatrix} 1 \\ 1 \\ -1 \\ 1 \end{bmatrix} w_1 + \begin{bmatrix} 1 \\ -1 \\ -1 \\ 1 \end{bmatrix} w_2$

Use known ratings to solve w_1, w_2

$$\begin{bmatrix} 1 & 1 \\ 1 & -1 \\ -1 & 1 \\ 1 & -1 \end{bmatrix} \begin{bmatrix} w_1 \\ w_2 \end{bmatrix} = \begin{bmatrix} 8 \\ 3 \\ 7 \\ 4 \end{bmatrix} \Rightarrow \begin{bmatrix} w_1 \\ w_2 \end{bmatrix} = \begin{bmatrix} 5.5 \\ 2 \end{bmatrix}$$

Predict ratings using
 w_1, w_2

$$\hat{\underline{a}} = \begin{bmatrix} 1 \\ 1 \\ -1 \\ 1 \end{bmatrix} 5.5 + \begin{bmatrix} 1 \\ -1 \\ -1 \\ 1 \end{bmatrix} 2 = \begin{bmatrix} 7.5 \\ 3.5 \\ 7.5 \\ 3.5 \\ 7.5 \end{bmatrix}$$

Use of "patterns" can mitigate noise 6

Noisy data $\underline{A}_m = \underline{A}_t + \underline{\epsilon}$

strong patterns no dominant patterns

Low-rank model $\hat{\underline{A}}_m = \underline{T} \underline{W}^T$ can be closer to \underline{A}_t than \underline{A}_m

Low rank classifier / model fit

$$\hat{\underline{A}}_m \underline{w} = \underline{d}$$
$$\left[\begin{array}{|c|} \hline \underline{T} \\ \hline \end{array} \right] \left[\begin{array}{|c|} \hline \underline{W}^T \\ \hline \end{array} \right] \left| \begin{array}{c} \underline{w} \\ \hline \end{array} \right. = \left| \begin{array}{c} \underline{d} \\ \hline \end{array} \right.$$

$\underline{w}' \ p \times 1$

$\underline{T} \underline{w}' = \underline{d}$
transformed features \underline{x}_i^T

$$\underline{x}_i^T = \underline{x}_i^T \underline{W} (\underline{W}^T \underline{W})^{-1}$$

New feature \underline{x}^T :

$$1) \quad \underline{x}^T = \underline{x}^T \underline{W} (\underline{W}^T \underline{W})^{-1} \quad 2) \quad \hat{d} = \text{sign}(\underline{x}^T \underline{w}')$$

**Copyright 2019
Barry Van Veen**

Clustering Data with the K-means Algorithm

Objectives -

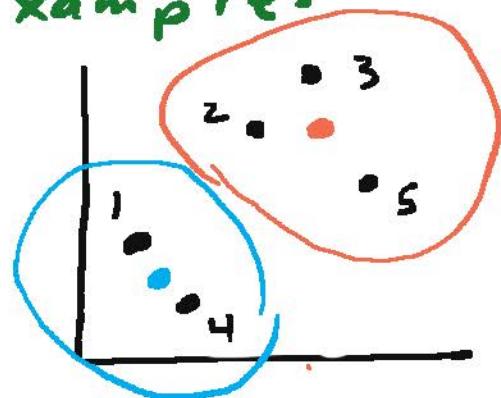
- Introduce the K-means algorithm
- Illustrate k-means with an example
- Identify considerations

Clustering: Organizing data in groups

2

given $\underline{a}_i \in \mathbb{R}^N$, $i = 1, 2, \dots, M$, find centroids $\underline{\mu}_j$, $j = 1, \dots, k$
and clusters $S_j = \{ i \mid \underline{a}_i \text{ belongs to cluster } j \}$

Example:



$$\underline{A} = [\underline{a}_1 \ \underline{a}_2 \ \underline{a}_3 \ \underline{a}_4 \ \underline{a}_5] = \begin{bmatrix} 1 & 3 & 4 & 2 & 5 \\ 2 & 4 & 5 & 1 & 3 \end{bmatrix}$$

$$S_1 = \{2, 3, 5\}, \underline{\mu}_1 = \begin{bmatrix} 4 \\ 4 \end{bmatrix}; S_2 = \{1, 4\}, \underline{\mu}_2 = \begin{bmatrix} 1.5 \\ 1.5 \end{bmatrix}$$

Unsupervised learning: data w/o labels, #clusters unknown

Matrix factorization: $\underline{A} \approx \underline{I} \ \underline{W}^T$, $\underline{I} = [\underline{\mu}_1 \ \underline{\mu}_2 \ \dots \ \underline{\mu}_k]$

$$\underline{I} = [\underline{\mu}_1 \ \underline{\mu}_2], \underline{W}^T = \begin{bmatrix} 0 & 1 & 1 & 0 & 1 \\ 1 & 0 & 0 & 1 & 0 \end{bmatrix}$$

$$[\underline{W}^T]_{km} = \begin{cases} 1, & m \in S_k \\ 0, & m \notin S_k \end{cases}$$

The K-Means Algorithm

(K clusters)

3

clusters: $S_j = \{\underline{a}_i \mid \underline{a}_i \in \text{cluster } j\}$, $|S_j| = \# \underline{a}_i \text{ in } S_j$

centroids: $\underline{\mu}_j = \frac{1}{|S_j|} \sum_{i \in S_j} \underline{a}_i$ coherence: $c_j = \sum_{i \in S_j} \|\underline{a}_i - \underline{\mu}_j\|_2^2$

Overall coherence $C = \sum_{j=1}^k c_j = \sum_{j=1}^k \sum_{i \in S_j} \|\underline{a}_i - \underline{\mu}_j\|_2^2 = \|\underline{A} - \underline{T} \underline{W}^T\|_F^2$

1) Initialize: choose $\underline{\mu}_j^0$, $j=1, 2, \dots, k$ randomly from \underline{a}_i ; set $l=0$

2) Assignment: put $\underline{a}_i \in S_j^l$ if \underline{a}_i is closest to $\underline{\mu}_j^l$

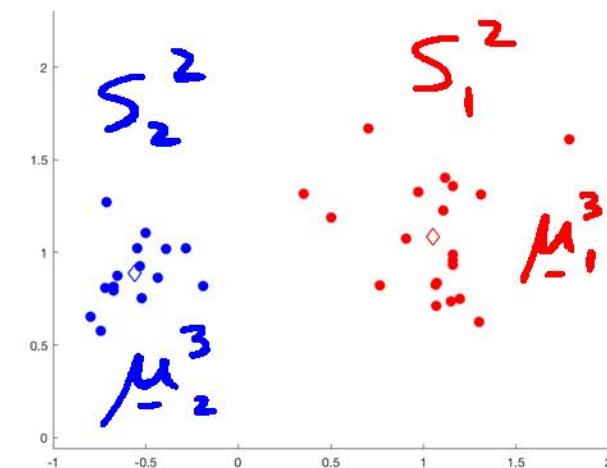
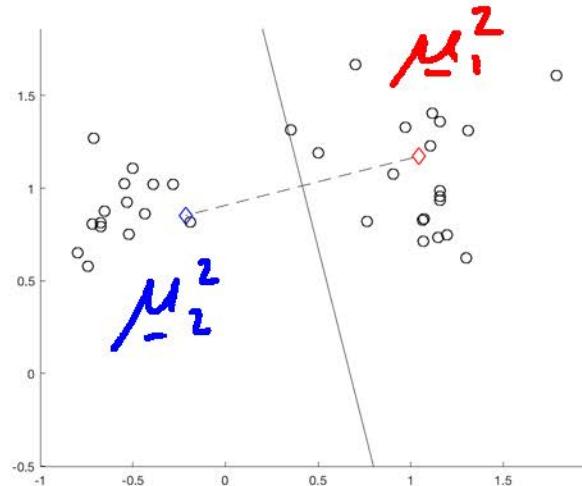
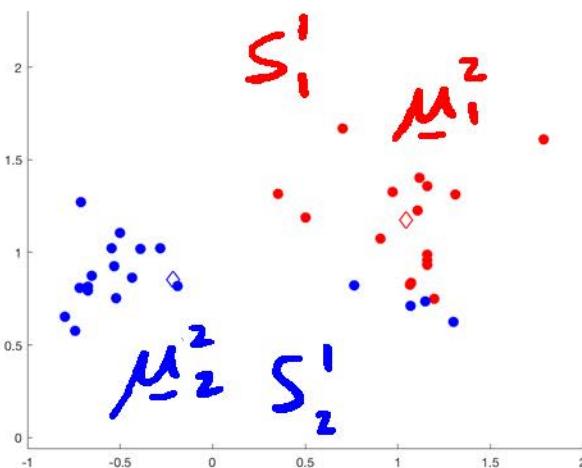
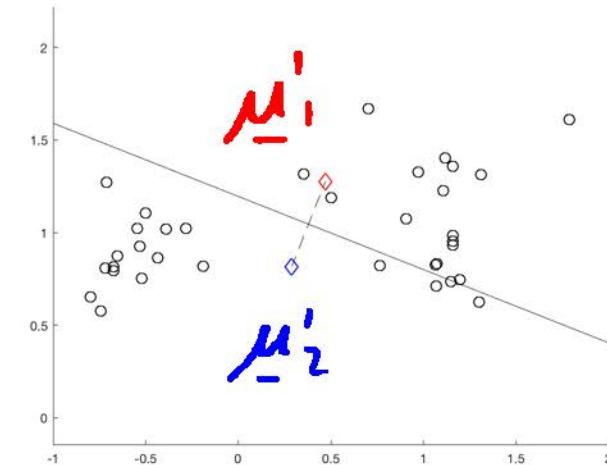
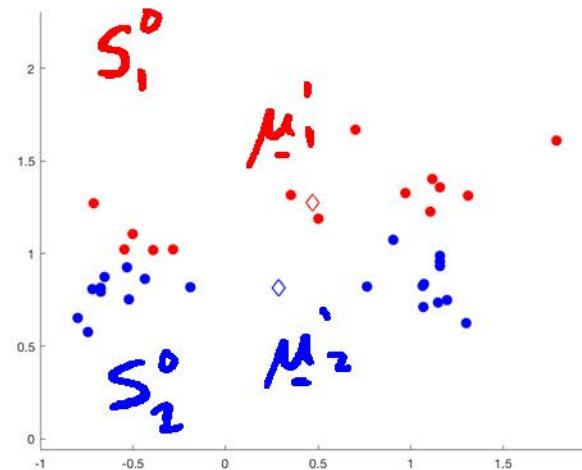
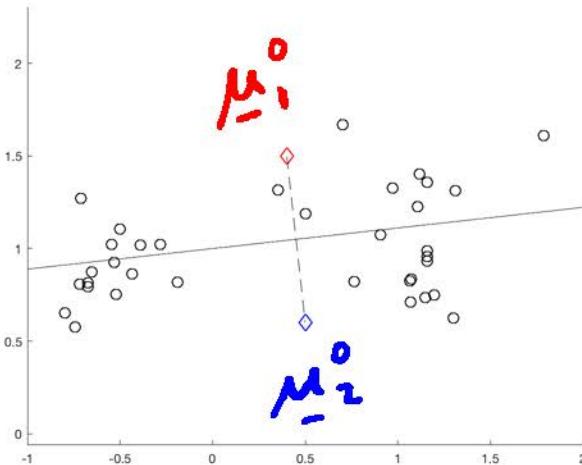
3) Update Centroids: $\underline{\mu}_j^{l+1} = \frac{1}{|S_j^l|} \sum_{i \in S_j^l} \underline{a}_i$

4) If converged \rightarrow stop

else $\rightarrow l=l+1$, go to 2)

K-Means Algorithm Example

4



K-Means Algorithm Options

5

- Initialization: many variations
- Termination: change in clusters or overall coherence or fix iterations
- Use different norms to assign clusters

Challenges

- Convergence to local minima
repeat for multiple initializations
- Unknown K
try multiple values

Copyright 2019
Barry Van Veen

The Singular Value Decomposition (SVD)

Objectives

- Define singular value decomposition (SVD)
- Express skinny SVD
- Write SVD as sum of outer products
- Use SVD to find best low-rank approximation
- Interpret matrix as an operator

SVD

- matrix decomposition that leads to good low-rank approximations
- vast range of applications

Definition:

Any $N \times M$ matrix A can be written as

$$\underline{A} = \underline{U} \Sigma \underline{V}^T$$

- U : $N \times N$, orthonormal columns
- V : $M \times M$, orthonormal columns
- Σ : $N \times M$, diagonal, $\Sigma_{ii} \geq 0$

$$N > M$$

$$\begin{bmatrix} \sigma_1 & & & \\ & \sigma_2 & & \\ & & \ddots & \\ & & & \sigma_m \\ 0 & & & \\ & 0 & & \\ & & \ddots & \\ & & & 0 \end{bmatrix}$$

$$M > N$$

$$\begin{bmatrix} \sigma_1 & & & \\ & \sigma_2 & & \\ & & \ddots & \\ & & & 0 \\ 0 & & & \\ & 0 & & \\ & & \ddots & \\ & & & \sigma_N \end{bmatrix}$$

$$\sigma_1 \geq \sigma_2 \geq \dots \geq \sigma_{N/M} \geq 0$$

SVD Dimensions

$$\underline{A} \underset{N \times M}{=} \begin{matrix} \underline{U} \\ N \times N \end{matrix} \begin{matrix} \Sigma \\ N \times N \end{matrix} \begin{matrix} \underline{V}^T \\ M \times M \end{matrix}$$

Skinny SVD

$$= \begin{matrix} \underline{U} \\ N \times M \end{matrix} \begin{matrix} \Sigma \\ M \times M \end{matrix} \begin{matrix} \underline{V}^T \\ M \times M \end{matrix}$$

$$\underline{A} \underset{N \times M}{=} \begin{matrix} \underline{U} \\ N \times N \end{matrix} \begin{matrix} \Sigma \\ N \times M \end{matrix} \begin{matrix} \underline{V}^T \\ M \times M \end{matrix}$$

$$= \begin{matrix} \underline{U} \\ N \times N \end{matrix} \begin{matrix} \Sigma \\ N \times N \end{matrix} \begin{matrix} \underline{V}^T \\ N \times M \end{matrix}$$

Sum of Outer Products Form:

"rank 1"

$$\underline{A} = \left[\begin{matrix} 1 & 1 & 1 \\ \underline{u}_1 & \underline{u}_2 & \dots & \underline{u}_M \end{matrix} \right] \left[\begin{matrix} \sigma_1 & & & \\ & \sigma_2 & & \\ & & \ddots & \\ & & & \sigma_M \end{matrix} \right] \left[\begin{matrix} -\underline{v}_1^T \\ -\underline{v}_2^T \\ \vdots \\ -\underline{v}_M^T \end{matrix} \right] = \sum_{i=1}^M \sigma_i \underline{u}_i \underline{v}_i^T = \sum_{i=1}^M \sigma_i \underline{u}_i \underline{v}_i^T \underset{N \times M}{\boxed{\Sigma}}$$

SVD gives the "best" low-rank approximation 4

$$\text{Frobenius norm } \|\underline{A}\|_F^2 = \sum_{i=1}^N \sum_{j=1}^M (\underline{A}_{i,j})^2 = \|\text{vec}(\underline{A})\|_2^2$$

Eckart-Young Theorem (1936) Let $\text{rank}(\underline{A}) = r$

$$\text{and } k < r: \min_{\text{rank}(\underline{B}) \leq k} \|\underline{A} - \underline{B}\|_F^2 = \sum_{i=k+1}^r \sigma_i^2 \text{ for } \underline{B} = \sum_{i=1}^k \sigma_i \underline{u}_i \underline{v}_i^T$$

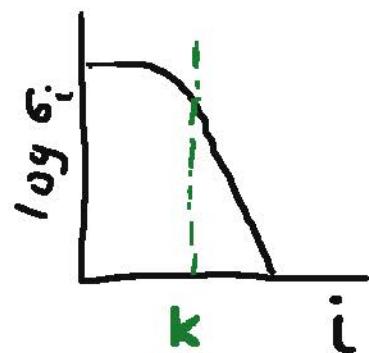
where $\underline{A} = \underline{U} \Sigma \underline{V}^T$ is the SVD.

$$\underline{A} \approx \sigma_1 \underline{u}_1 \underline{v}_1^T + \sigma_2 \underline{u}_2 \underline{v}_2^T + \dots + \sigma_k \underline{u}_k \underline{v}_k^T$$

patterus: most important 2nd most k^{th} most

σ_i provide ordered ranking of components

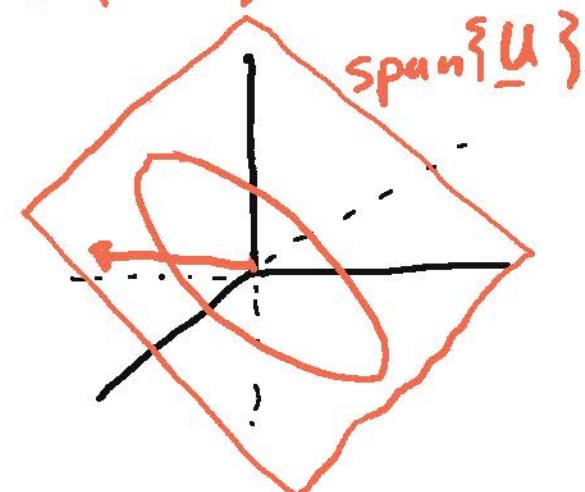
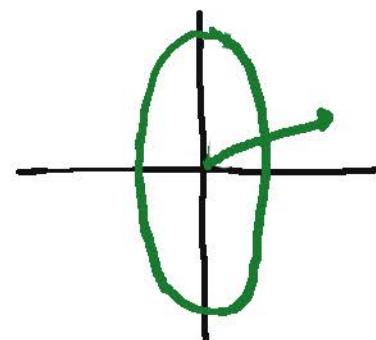
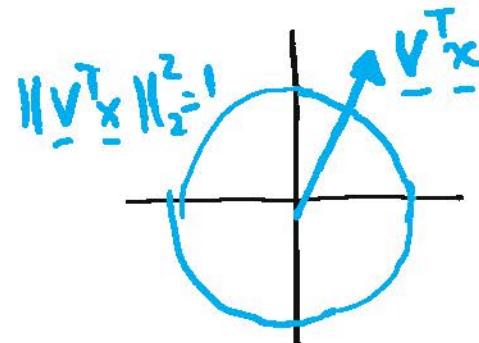
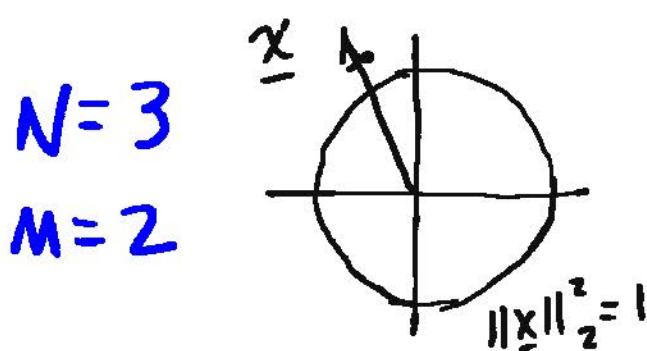
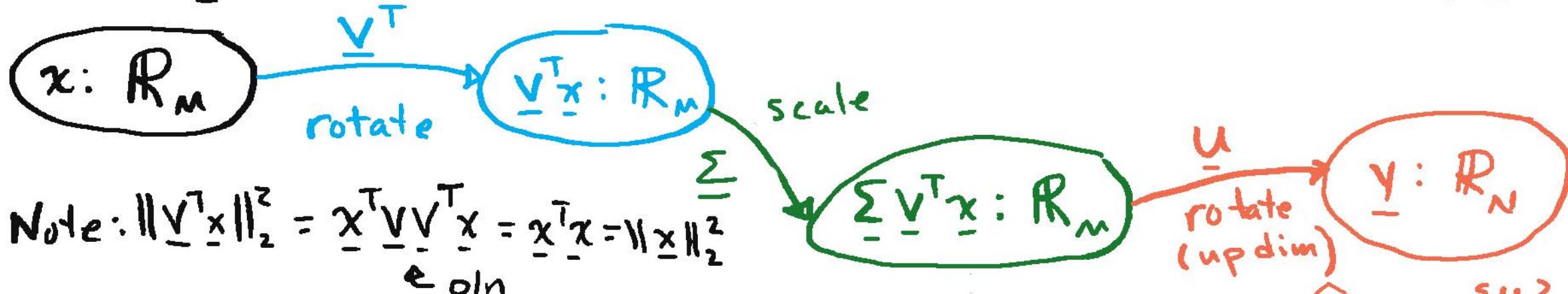
cols: scaled \underline{u}_i
rows: scaled \underline{v}_i^T



SVD describes Matrix as an operator

5

$$A: N \times M, \begin{cases} \underline{x}: M \times 1 \\ \underline{y}: N \times 1 \end{cases} \quad \underline{y} = A \underline{x} = \underline{U} \sum \underline{V}^T \underline{x} = \underline{U} [\sum (\underline{V}^T \underline{x})]$$



Operator Norm $\|A\|_2 = \|A\|_{op} := \max_{\underline{x} \neq 0} \frac{\|A \underline{x}\|_2}{\|\underline{x}\|_2} = \sigma_1$

(proof: notes)

**Copyright 2019
Barry Van Veen**

Properties of Singular Value Decomposition

Objectives

- review orthonormality of singular vectors
- review rank and singular values
- explore singular vectors as bases
- Connect SVD and matrix inversion

Singular Value Decomposition

2

$$\underline{A} \underset{N \times M}{=} \underline{U} \underset{N \times N}{|} \begin{matrix} \diagup \\ \diagdown \end{matrix} \underset{M \times M}{\Sigma} \underset{M \times N}{|} \underline{V}^T$$

Orthonormality

$$\underline{U}^T \underline{U} = \underline{I} ; \underline{V}^T \underline{V} = \underline{I}$$

full
econ

$$\underline{U} \underline{U}^T = \underline{I}_N ; \underline{V} \underline{V}^T = \underline{I}_M$$

full only

\square Rank

$$\underline{A} \underset{N \times M}{=} \underline{U} \underset{N \times N}{|} \begin{matrix} \diagup \\ \diagdown \end{matrix} \underset{M \times M}{\Sigma} \underset{M \times N}{|} \underline{V}^T$$

left sing.
vectors

sing.
values

$$\sigma_1 \geq \sigma_2 \geq \dots \geq \sigma_{\min\{N, M\}} \geq 0$$

right sing.
vectors

$$\underline{V} = \begin{bmatrix} \underline{v}_1 & \underline{v}_2 & \dots & \underline{v}_M \end{bmatrix}$$

$$\text{rank } (\underline{A}) = p \Leftrightarrow \sigma_1 \geq \dots \geq \sigma_p > \sigma_{p+1} = \dots = \sigma_{\min\{N, M\}} = 0$$

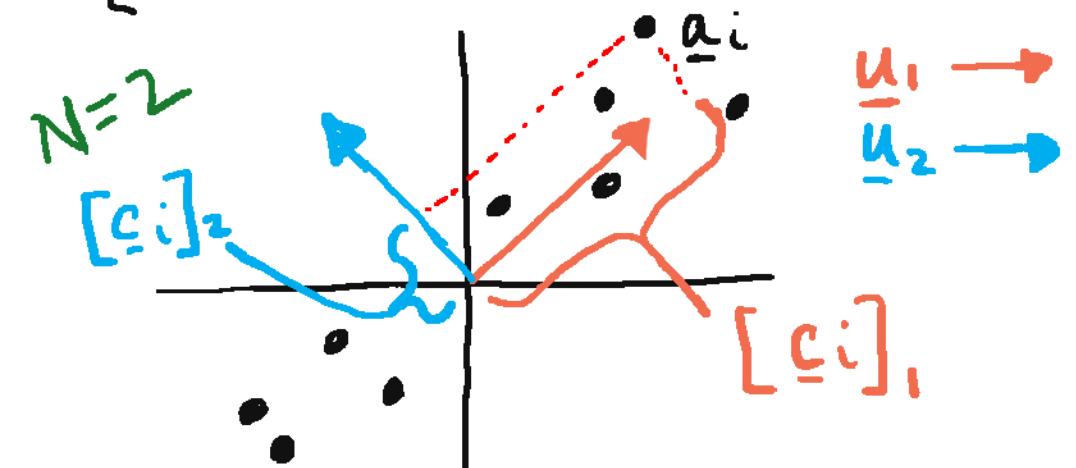
$$\underline{U} \underset{N \times N}{|} \begin{matrix} \diagup \\ \diagdown \end{matrix} \underset{M \times M}{\Sigma} \underset{M \times N}{|} \underline{V}^T$$

$$\underline{A} = \sum_{i=1}^p \sigma_i \underline{U}_i \underline{V}_i^T$$

Singular vectors are o/n bases for rows/columns 3

$$\begin{bmatrix} \frac{1}{\alpha_1}, \frac{1}{\alpha_2}, \dots, \frac{1}{\alpha_m} \end{bmatrix} = \begin{bmatrix} \frac{1}{u_1}, \frac{1}{u_2}, \dots, \frac{1}{u_p} \end{bmatrix} \underbrace{\begin{bmatrix} \frac{1}{c_1}, \frac{1}{c_2}, \dots, \frac{1}{c_n} \end{bmatrix}}_{C = \underline{V}^T} \Rightarrow \underline{a}_i = \sum_{j=1}^p u_j [c_i]_j$$

coords of \underline{a}_i
in basis \underline{u}



left sing. vec. \underline{u} : o/n basis cols \underline{A}

j^{th} coord $\sim \sigma_j$ $[c_i]_j = \sigma_j \underbrace{[\underline{v}^T]_{j,i}}_{\text{max}=1}$

$$\begin{bmatrix} -x_1^T \\ -x_2^T \\ \vdots \\ -x_N^T \end{bmatrix} = \underbrace{\begin{bmatrix} -d_1^T \\ -d_2^T \\ \vdots \\ -d_N^T \end{bmatrix}}_{D = \underline{U} \Sigma} \begin{bmatrix} -v_1^T \\ -v_2^T \\ \vdots \\ -v_p^T \end{bmatrix} \Rightarrow x_i^T = \sum_{j=1}^p v_j^T [d_i]_j$$

right sing. vec. \underline{v} : o/n basis
for rows of \underline{A}

j^{th} coord $\sim \sigma_j$ $[d_i]_j = \sigma_j [\underline{u}]_{i,j}$

SVD gives inverse of square matrices 4

$$N=M \quad \underline{A} = \underline{U} \underline{\Sigma} \underline{V}^T \quad \underline{U}, \underline{\Sigma}, \underline{V} : N \times N$$

Noninvertible (singular): $\text{rank}(\underline{A}) < N$

$$\sigma_1 \geq \sigma_2 \geq \cdots \geq \sigma_p > \sigma_{p+1} = \cdots = \sigma_N = 0$$

Invertible: $\text{rank}(\underline{A}) = N \quad \underline{A}^{-1} = \underline{V} \underline{\Sigma}^{-T} \underline{U}^T$

$$\begin{aligned} \underline{A} \cdot \underline{A}^{-1} &= \underline{U} \underline{\Sigma} \underline{V}^T \underline{V} \underline{\Sigma}^{-T} \underline{U}^T = \underline{U} \underline{\Sigma} \underline{\Sigma}^{-T} \underline{U}^T = \underline{U} \underline{\Sigma} \underline{\Sigma}^{-1} \underline{U}^T \\ &= \underline{U} \underline{U}^T = \underline{I} \quad (\text{no econ SVD for full rank square}) \end{aligned}$$

$$\underline{A} = \sum_{i=1}^N \sigma_i \underline{u}_i \underline{v}_i^T, \quad \underline{A}^{-1} = \sum_{i=1}^N \frac{1}{\sigma_i} \underline{v}_i \underline{u}_i^T$$

SVD of \underline{A} gives
SVD of \underline{A}^{-1}

**Copyright 2019
Barry Van Veen**

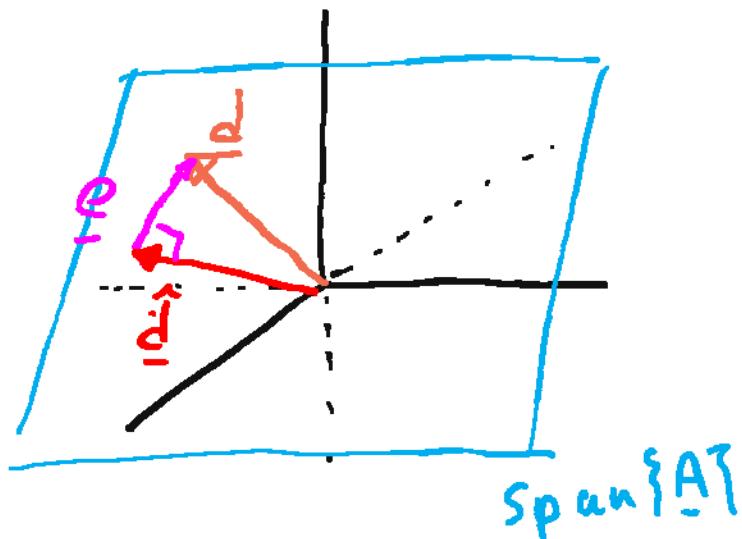
The SVD and Least-Squares Problems

Objectives

1

- Express least-squares solution in terms of SVD
- Express least-squares error in terms of SVD
- Use SVD to solve the orthobases classification problem

SVD gives insight into the least-squares problem²



$$\min_{\underline{w}} \|\underline{d} - \underline{A}\underline{w}\|_2^2$$

$$\underline{A} : N \times P \quad \text{rank}(\underline{A}) = P$$

$$\underline{d} : N \times 1$$

$$\underline{e} = \underline{d} - \underline{\hat{d}}$$

$$= (\underline{I} - \underline{P}_A) \underline{d} = \underline{P}_{A^\perp} \underline{d}$$

$$\underline{w} = (\underline{A}^T \underline{A})^{-1} \underline{A}^T \underline{d}$$

$$\begin{aligned} \underline{\hat{d}} &= \underline{A} \underline{w} = \underline{A} (\underline{A}^T \underline{A})^{-1} \underline{A}^T \underline{d} \\ &= \underline{P}_A \underline{d} \end{aligned}$$

SVD -

$$\underline{A} =$$

$$\begin{bmatrix} \underline{\tilde{U}} & \vdots & \underline{U}_L \end{bmatrix}$$

$$\begin{bmatrix} \Sigma & & \\ & \ddots & \\ & & 0 \end{bmatrix}$$

$$\underline{V}^T$$

$$\underline{U} = [\underline{\tilde{U}} : \underline{U}_L] \quad N \times N$$

$$\underline{A} = \underline{\tilde{U}} \underline{\Sigma} \underline{V}^T$$

Least-squares solution

3

$$\underline{w} = (\underline{A}^T \underline{A})^{-1} \underline{A}^T \underline{d} = (\underline{V} \underline{\Sigma}^T \underline{\tilde{U}}^T \underline{\tilde{U}} \underline{\Sigma} \underline{V}^T)^{-1} \underline{V} \underline{\Sigma}^T \underline{\tilde{U}}^T \underline{d}$$

$$= (\underline{V} \underline{\Sigma}^2 \underline{V}^T)^{-1} \underline{V} \underline{\Sigma} \underline{\tilde{U}}^T \underline{d} \quad \text{recall } (\underline{E} \underline{F} \underline{G})^{-1} = \underline{G}^{-1} \underline{F}^{-1} \underline{E}^{-1}$$

$$= \underline{V} \underline{\Sigma}^{-2} \underline{V}^T \underline{V} \underline{\Sigma} \underline{\tilde{U}}^T \underline{d} \quad \underline{V}^{-1} = \underline{V}^T$$

$$= \underline{V} \underline{\Sigma}^{-1} \underline{\tilde{U}}^T \underline{d} = \sum_{i=1}^p \frac{1}{\sigma_i} v_i (\underline{\tilde{U}}^T \underline{d})$$

pseudo-inverse of A

$$\underline{\Sigma}^{-1} \underline{\Sigma} = \begin{bmatrix} \frac{1}{\sigma_1^2} & & & \\ & \frac{1}{\sigma_2^2} & & \\ & & \ddots & \\ & & & \frac{1}{\sigma_p^2} \end{bmatrix} \begin{bmatrix} \sigma_1 & \sigma_2 & \dots & 0 \end{bmatrix}$$

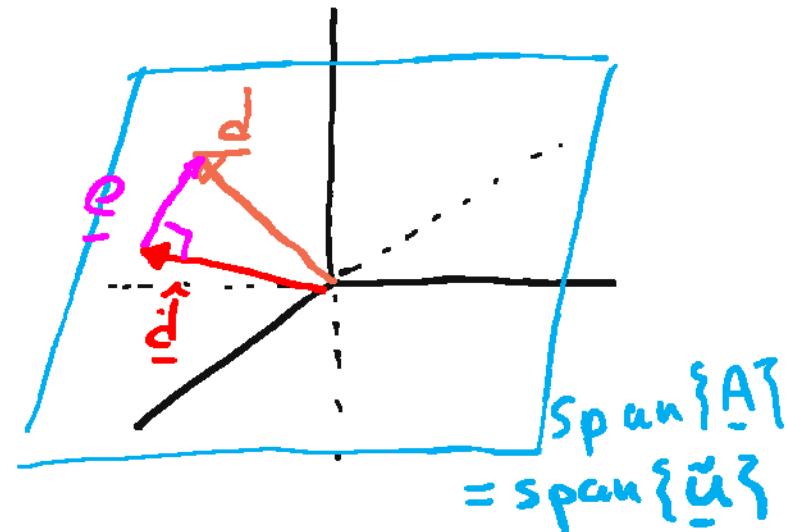
$$[(\underline{A}^T \underline{A})^{-1} \underline{A}^T] \underline{A} = (\underline{A}^T \underline{A})^{-1} \underline{A}^T \underline{A} = \underline{I}$$

$$= \begin{bmatrix} \frac{1}{\sigma_1} & \frac{1}{\sigma_2} & \dots & 0 \end{bmatrix} = \underline{\Sigma}^{-1}$$

$$\underline{V} \underline{\Sigma}^{-1} \underline{\tilde{U}}^T \underline{\tilde{U}} \underline{\Sigma} \underline{V}^T = \underline{V} \underline{\Sigma}^{-1} \underline{\Sigma} \underline{V}^T = \underline{V} \underline{V}^T = \underline{I}$$

Least-squares error and projections

4



$$\begin{aligned}\hat{\underline{d}} &= \underline{A} \left[(\underline{A}^T \underline{A})^{-1} \underline{A}^T \right] \underline{d} = \underline{P}_A \underline{d} \\ &= \underline{\tilde{U}} \underline{\Sigma} \underline{V}^T \left[\underline{V} \underline{\Sigma}^{-1} \underline{\tilde{U}}^T \right] \underline{d} = \underline{\tilde{U}} \underline{\Sigma} \underline{\Sigma}^{-1} \underline{\tilde{U}}^T \underline{d} \\ &= \underline{\tilde{U}} \underline{\tilde{U}}^T \underline{d} \Rightarrow \underline{P}_A = \underline{\tilde{U}} \underline{\tilde{U}}^T\end{aligned}$$

$$\underline{e} = \underline{d} - \hat{\underline{d}} = (\underline{I} - \underline{\tilde{U}}\underline{\tilde{U}}^T) \underline{d} = P_{A^\perp}\underline{d} \Rightarrow P_{A^\perp} = \underline{I} - \underline{\tilde{U}}\underline{\tilde{U}}^T$$

Recall $\underline{U} = [\underline{\tilde{U}} : \underline{U}_1]$ ($N \times N$) $\underline{U}^T \underline{U} = \underline{U} \underline{U}^T = \underline{I} = [\underline{\tilde{U}}^T : \underline{U}_1^T] \begin{bmatrix} \underline{\tilde{U}} \\ \underline{U}_1^T \end{bmatrix}$

$$\text{so } \underline{I} = \underline{\tilde{U}}\underline{\tilde{U}}^T + \underline{U}_1\underline{U}_1^T \Rightarrow P_{A^\perp} = \underline{U}_1\underline{U}_1^T$$

Classification using SVD O/n bases

5

Training

$$\begin{bmatrix} \underline{x}_1^T \\ \underline{x}_2^T \\ \vdots \\ \underline{x}_N^T \end{bmatrix} \underline{w} = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_N \end{bmatrix} \Rightarrow \underline{A}\underline{w} = \underline{d}$$

\Downarrow

Prediction

$$\tilde{y} = \text{sign}(\tilde{\underline{x}}^T \underline{w})$$

features

labels

Orthobases

$$\begin{aligned} \underline{A} &= \tilde{\underline{U}} \Sigma \underline{V}^T \\ \tilde{\underline{U}} \underline{w}' &= \underline{d} \quad \underline{w} = \underline{V} \Sigma^{-1} \tilde{\underline{U}}^T \underline{d} \\ \underline{w}' &= \tilde{\underline{U}}^T \underline{d} \quad = \underline{V} \Sigma^{-1} \underline{w}' \end{aligned}$$

$|w'_i|$ indicates importance
of i^{th} ortho feature

Orthobases Prediction

$$\tilde{y} = \text{sign}(\tilde{\underline{x}}^T \underline{w}) = \text{sign}(\tilde{\underline{x}}^T \underline{V} \Sigma^{-1} \underline{w}')$$

$$\tilde{\underline{x}}^T = \tilde{\underline{x}}^T \underline{V} \Sigma^{-1}$$

transformed feature

$$\tilde{y} = \text{sign}(\tilde{\underline{x}}'^T \underline{w}')$$

orthobasis classifier

Copyright 2019
Barry Van Veen

SVD and Regularization of Least-Squares Problems

Objectives

- Analyze impact of errors in least-squares problems using SVD
- Introduce truncated SVD regularization
- Analyze ridge regression using SVD

III-conditioned least-squares problems 2

$\min_{\underline{w}} \|\underline{A}\underline{w} - \underline{d}\|_2^2 \Rightarrow \underline{w} = (\underline{A}^T \underline{A})^{-1} \underline{A}^T \underline{d}$ have small singular values

$$\text{SVD: } \underline{A} = \underline{U} \Sigma \underline{V}^T \Rightarrow \underline{w} = \underline{V} \Sigma^{-1} \underline{U}^T \underline{d} = \sum_{i=1}^P \frac{1}{\sigma_i} \underline{v}_i (\underline{u}_i^T \underline{d})$$

$N \times P$, rank P

$$\Rightarrow \|\underline{w}\|_2^2 = \sum_{i=1}^P \left(\frac{1}{\sigma_i}\right)^2 (\underline{u}_i^T \underline{d})^2 \quad \text{Small } \sigma_i \Rightarrow \text{large } \|\underline{w}\|_2$$

Prediction with errors: $\tilde{\mathbf{y}} = (\tilde{\mathbf{x}} + \underline{\boldsymbol{\xi}})^T \underline{w} = \mathbf{x}^T \underline{w} + \underline{\boldsymbol{\xi}}^T \underline{w}$

$$|\underline{\boldsymbol{\xi}}^T \underline{w}|^2 = \|\underline{w}\|_2^2 \|\underline{\boldsymbol{\xi}}\|_2^2 \cos^2 \theta$$

large $\|\underline{w}\|_2^2 \Rightarrow$ sensitive

to errors

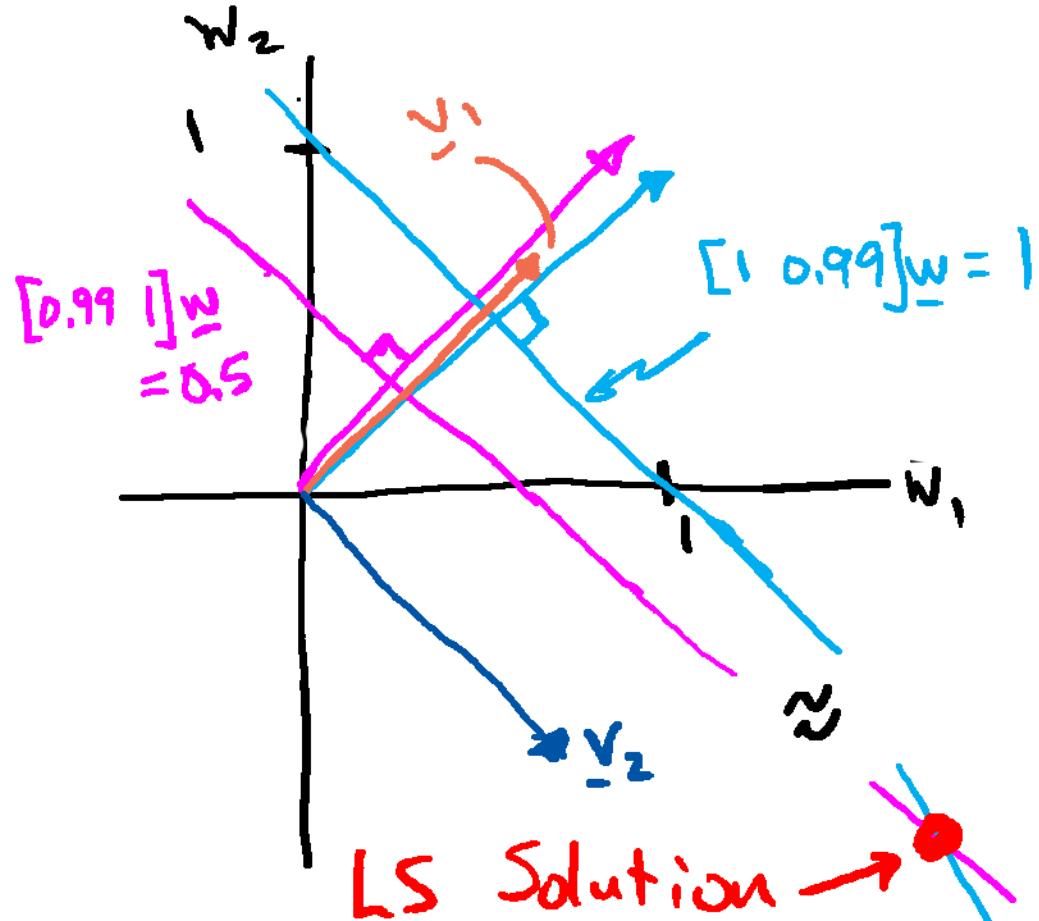
What if $\text{rank}(\underline{A}) < P$?

$$\sigma_P = 0$$

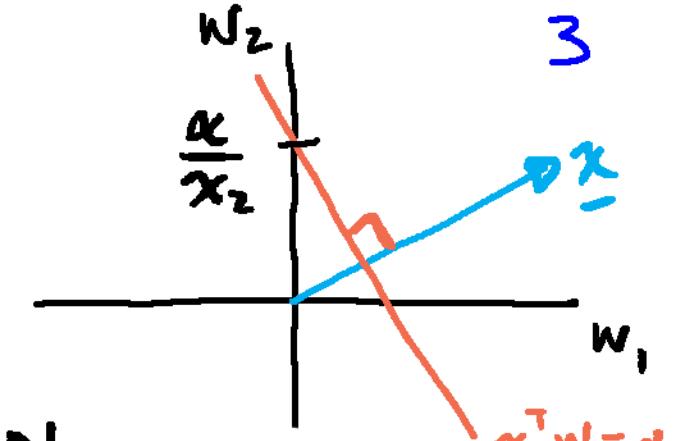
no unique solution

Example: ill-conditioned A

$$\underline{A} = \begin{bmatrix} 1 & 0.99 \\ 0.99 & 1 \end{bmatrix}, \underline{d} = \begin{bmatrix} 1 \\ 0.5 \end{bmatrix}$$



Geometry $\underline{x}^T \underline{w} = \alpha$



$$\sigma_1 = 1.99 \quad \sigma_2 = 0.01$$

$$1/\sigma_1 \approx 0.5, \quad 1/\sigma_2 = 100$$

$$\underline{v}_1 = \frac{1}{\sqrt{2}} \begin{bmatrix} 1 & -1 \end{bmatrix}^T \quad \underline{v}_2 = \frac{1}{\sqrt{2}} \begin{bmatrix} 1 & 1 \end{bmatrix}^T$$

$$\underline{u}_1^T \underline{d} = \frac{1.5}{\sqrt{2}}, \quad \underline{u}_2^T \underline{d} = \frac{0.5}{\sqrt{2}}$$

$$\underline{w} = \underline{v}_1 \frac{\underline{u}_1^T \underline{d}}{\sigma_1} + \underline{v}_2 \frac{\underline{u}_2^T \underline{d}}{\sigma_2} = \begin{bmatrix} 1 \\ -1 \end{bmatrix} \frac{3}{\sqrt{2}} + \begin{bmatrix} 1 \\ 1 \end{bmatrix} \frac{25}{\sqrt{2}}$$

$\|\underline{w}\|_2 \uparrow$ as $\sigma_2 \downarrow$ (\rightarrow)

Regularized LS via truncated SVD

4

Replace $\sum_{i=1}^p \frac{1}{\sigma_i} \underline{v}_i (\underline{u}_i^\top \underline{d})$ with $\sum_{i=1}^r \frac{1}{\sigma_i} \underline{v}_i (\underline{u}_i^\top \underline{d})$
where $r < p$.

- Avoid inverting small/zero singular values
- Equivalent to replacing $\underline{A} = \sum_{i=1}^p \sigma_i \underline{u}_i \underline{v}_i^\top$ with the rank- r approximation $\underline{A}_r = \sum_{i=1}^r \sigma_i \underline{u}_i \underline{v}_i^\top$
- Increases $\min_w \| \underline{A} \underline{w} - \underline{d} \|_2^2$
- Can choose r using intuition or cross-validation

Regularized LS via ridge regression

$$\min_{\underline{w}} \|\underline{A}\underline{w} - \underline{d}\|_2^2 + \lambda \|\underline{w}\|_2^2 \Rightarrow \underline{w} = (\underline{A}^T \underline{A} + \lambda \underline{\underline{I}})^{-1} \underline{A}^T \underline{d}$$

controls norm!

Use SVD: $\underline{A}^T \underline{A} = \underline{V} \Sigma^2 \underline{V}^T$, $\lambda \underline{\underline{I}} = \underline{V} \lambda \underline{\underline{I}} \underline{V}^T$

$$\underline{w} = (\underline{V} (\Sigma^2 + \lambda \underline{\underline{I}}) \underline{V}^T)^{-1} \underline{V} \Sigma \underline{U}^T \underline{d} = \underline{V} (\Sigma^2 + \lambda \underline{\underline{I}})^{-1} \Sigma \underline{U}^T \underline{d}$$

$$\underline{D} = \begin{bmatrix} \frac{1}{\sigma_1^2 + \lambda} & & 0 \\ & \ddots & \\ 0 & & \frac{1}{\sigma_p^2 + \lambda} \end{bmatrix} \begin{bmatrix} \sigma_1 & & 0 \\ & \ddots & \\ 0 & & \sigma_p \end{bmatrix} = \begin{bmatrix} \sigma_1 / (\sigma_1^2 + \lambda) & & 0 \\ & \ddots & \\ 0 & & \sigma_p / (\sigma_p^2 + \lambda) \end{bmatrix}$$

Controlled!

$$\underline{w} = \sum_{i=1}^p \frac{\sigma_i}{\sigma_i^2 + \lambda} \underline{v}_i (\underline{U}^T \underline{d})$$

- as $\sigma_i \rightarrow 0$, $\frac{\sigma_i}{\sigma_i^2 + \lambda} \rightarrow \sigma_i / \lambda$
- increased value $\|\underline{A}\underline{w} - \underline{d}\|_2^2$

Copyright 2019
Barry Van Veen

Principal Component Analysis

Objectives

- Define principal components
- Relate principal components to singular vectors
- Relate geometry of data matrix to singular vectors and singular values

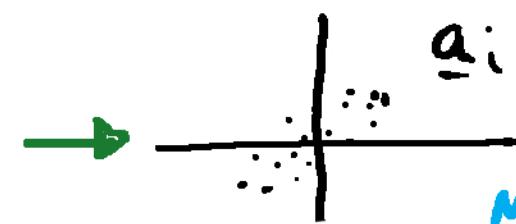
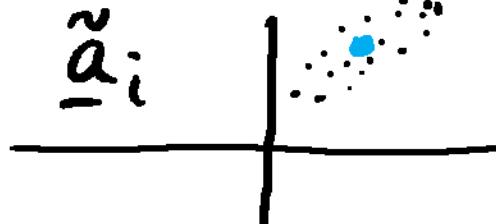
PCA represents maximum "variance"

²

Data: $\underline{a}_i, i=1, 2, \dots, N$ ($N \times 1$) vectors, $\underline{A} = [\underline{a}_1 \ \underline{a}_2 \ \dots \ \underline{a}_N]$

- PCA assumes zero mean \Rightarrow 1st step: center data

- First principal component:



direction \underline{f} accounting for

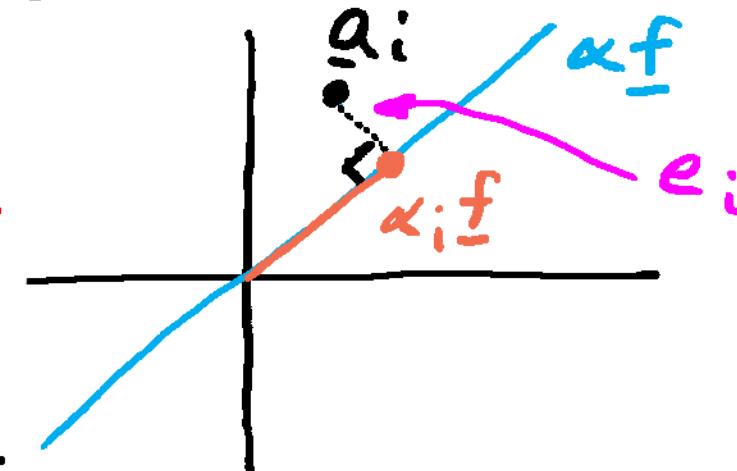
maximum variance in data, $\|\underline{f}\|_2^2 = 1$

$$\max_{\|\underline{f}\|_2^2 = 1} \left\{ \frac{1}{m} \sum_{i=1}^m \|\alpha_i \underline{f}\|_2^2 \right\} \quad \text{best line}$$

$$\min_{\alpha_i} \|\underline{a}_i - \alpha_i \underline{f}\|_2^2 \Rightarrow \max_{\|\underline{f}\|_2^2 = 1} \left\{ \frac{1}{m} \sum_{i=1}^m (\underline{f}^\top \underline{a}_i)^2 \right\}$$

$$\alpha_i = \underline{f}^\top \underline{a}_i$$

$$\|\underline{f}^\top \underline{A}\|_2^2 = \|\underline{A}^\top \underline{f}\|_2^2$$



$$\underline{a}_i = \tilde{\underline{a}}_i - \frac{1}{m} \sum_{j=1}^m \tilde{\underline{a}}_j$$

Principal Components are singular vectors

3

$$\max_{\|\underline{f}\|_2^2=1} \frac{1}{m} \underline{f}^T \underline{A} \underline{A}^T \underline{f} \Rightarrow \underline{A} = \underline{U} \underline{\Sigma} \underline{V}^T \quad \underline{A} \underline{A}^T = \underline{U} \underline{\Sigma} \underline{V}^T \underline{V} \underline{\Sigma} \underline{U}^T \\ = \underline{U} \underline{\Sigma}^2 \underline{U}^T$$

$$\max_{\|\underline{f}\|_2^2=1} \frac{1}{m} \underline{f}^T \underline{U} \underline{\Sigma}^2 \underline{U}^T \underline{f} \Rightarrow \underline{f} = \underline{u}_1 \text{ (notes) "best line"}$$

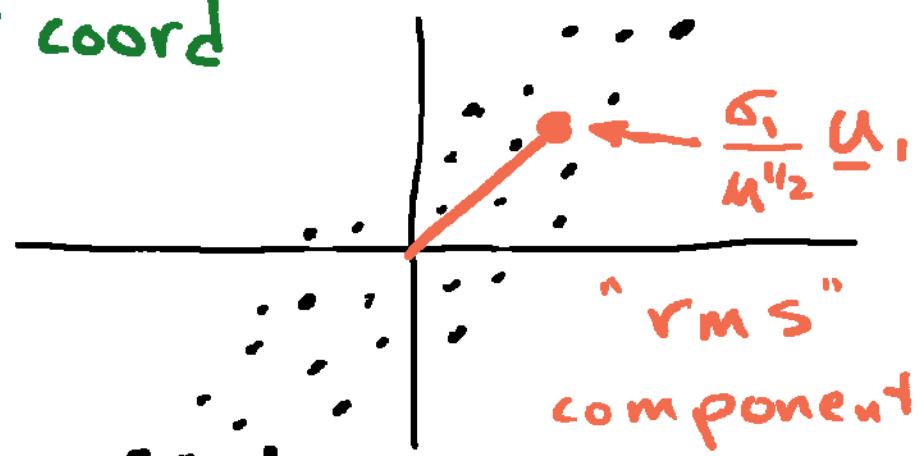
Variance associated w. 1st PC $\frac{1}{m} \underline{u}_1^T \underline{A} \underline{A}^T \underline{u}_1 = \frac{\sigma_1^2}{m}$

Coordinates of data: $\alpha_i = \underline{u}_1^T \underline{a}_i$, $\underline{\alpha}^T = [\alpha_1 \alpha_2 \dots \alpha_m]$

$\underline{\alpha}^T = \underline{u}_1^T \underline{A}$ root mean square coord

$$= \underline{u}_1^T \underline{U} \underline{\Sigma} \underline{V}^T \left(\frac{1}{m} \sum_{i=1}^m |\alpha_i|^2 \right)^{1/2}$$

$$= \sigma_1 \underline{v}_1^T \quad = \frac{1}{m^{1/2}} \|\underline{\alpha}\|_2 = \frac{\sigma_1}{m^{1/2}}$$



PC are singular vectors

4

$$2^{\text{nd}} \text{ PC: } \max_{\|\underline{g}\|_2^2 = 1, \underline{g}^T \underline{u}_1 = 0} \frac{1}{m} \sum_{i=1}^m |\underline{g}^T \underline{a}_i|^2$$

$\Rightarrow \underline{g} = \underline{u}_2$ (2nd left singular vector)

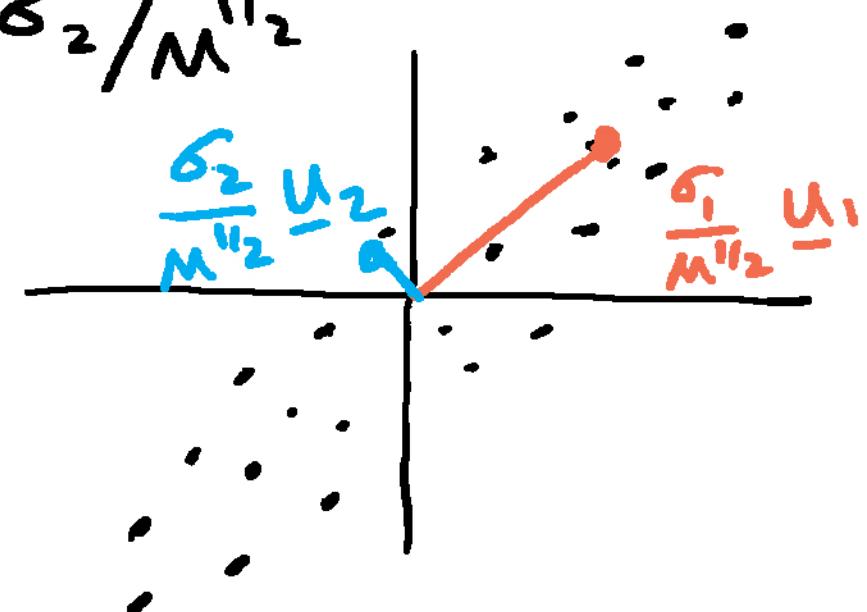
Variance associated w. 2nd PC: $\frac{1}{m} \underline{u}_2^T A A^T \underline{u}_2 = \frac{\sigma_2^2}{m}$

RMS value of 2nd PC coord: $\sigma_2 / m^{1/2}$

kth PC: \underline{u}_k

kth PC Variance: $\frac{\sigma_k^2}{m}$

kth PC coord RMS: $\frac{\sigma_k}{m^{1/2}}$



Summary

$$\underline{A} = \underline{U} \Sigma \underline{V}^T$$

5

- Left sing. vectors \leftrightarrow PC for columns of \underline{A}
- Sing. values \leftrightarrow \sim RMS value PC coords
- PC for rows of \underline{A}
 - use columns of $\underline{A}^T = \underline{V} \Sigma \underline{U}^T$
 - Right sing. vectors \underline{v}_i are PC
 - Sing. values \sim RMS value PC coords
- Eckhart - Young: SVD gives best rank r approximation to \underline{A} $\underline{A} \approx \sum_{k=1}^r \sigma_k \underline{u}_k \underline{v}_k^T$
 $r=1 \quad \underline{A} \approx \underline{u}_1 (\sigma_1 \underline{v}_1^T)$

Copyright 2019
Barry Van Veen

Bias-Variance Tradeoff in Low-Rank Approximations

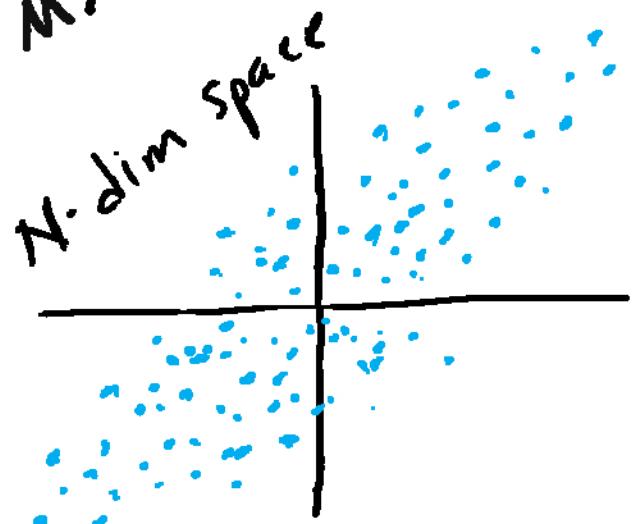
Objectives

- Introduce concept of noisy data
- Consider impact of noise on SVD
- Define bias and variance
- Use low-rank models to trade bias for variance

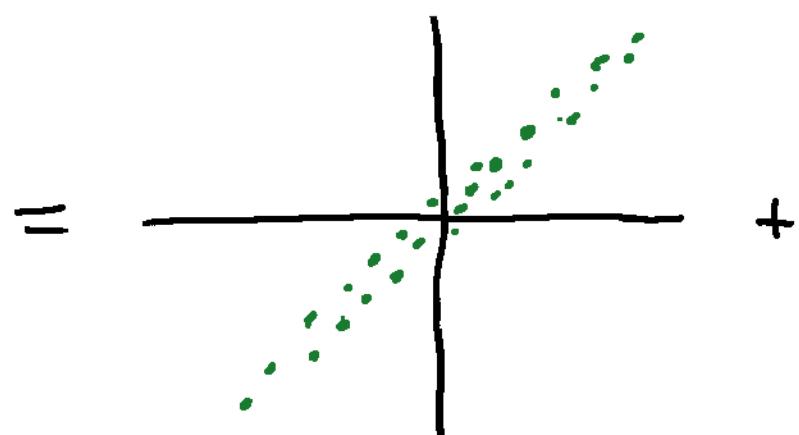
Data is often contaminated by noise 2

$$N \times M \text{ measured } \underline{A} = \underline{S}_{\text{clean}} + \underline{G}_{\text{noise}}$$

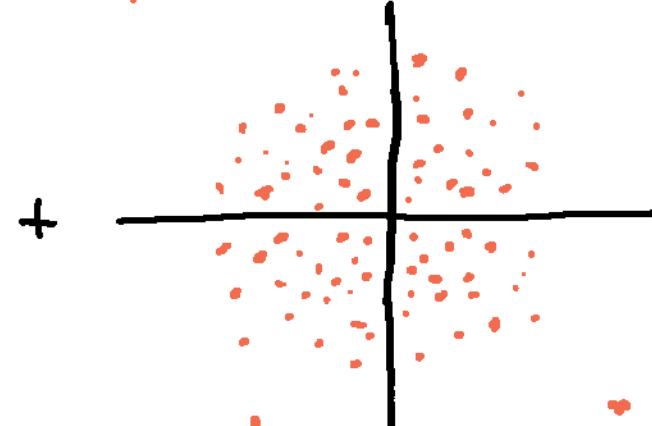
electronics in sensing systems
environmental static
limited precision in computers



diffuse structure



very structured



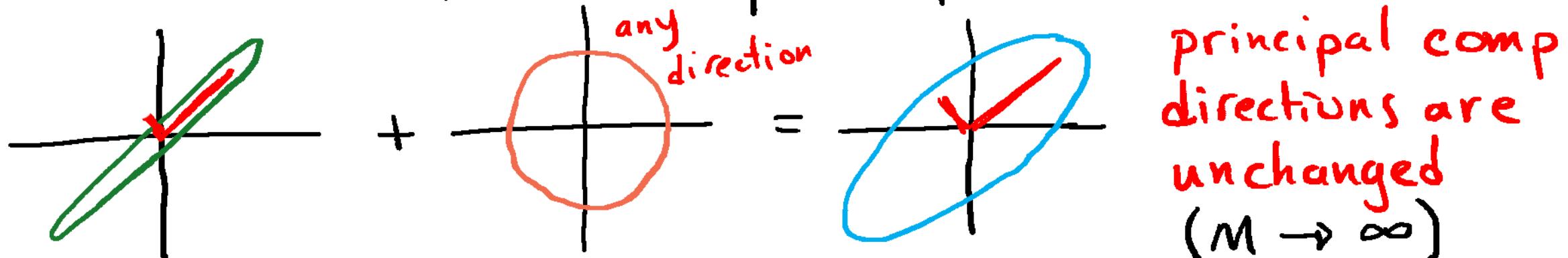
isotropic or "white"
noise - no preferred
direction

Sum of squared errors: $\|\underline{G}\|_F^2$

$$\|\underline{G}\|_F^2 = \sum_{i=1}^N M \left(\frac{1}{M} \sum_{j=1}^n g_{ij}^2 \right) = M \sum_{i=1}^N \text{var}_i \stackrel{\text{isotropic}}{\approx} MN \sigma_g^2$$

Singular vectors are invariant (approx) to isotropic noise 3

- Proof uses probability concepts



variance along each component (singvals) changes

$$\underline{\Sigma}_A \underline{V}^T \approx \underline{\Sigma}_S \underline{V}^T + \underline{\Sigma}_G \underline{V}^T$$

$$\sigma_{A_i} \approx \sigma_{S_i} + M^{1/2} \sigma_g$$

$$\Sigma_A = \begin{bmatrix} \sigma_{A_1} & & \\ & \ddots & \\ & & \sigma_{A_N} \end{bmatrix}, \quad \Sigma_S = \begin{bmatrix} \sigma_{S_1} & & \\ & \ddots & \\ & & \sigma_{S_N} \end{bmatrix}, \quad \Sigma_G \approx \begin{bmatrix} M^{1/2} \sigma_g & & \\ & M^{1/2} \sigma_g & \\ & & \ddots \\ & & & M^{1/2} \sigma_g \end{bmatrix}$$

(isotropic)

$\sigma_{g_i} = M^{1/2} \cdot \text{RMS}$

Low-rank models trade bias for variance 4

Original: $\underline{A} = \underline{S} + \underline{G}$ Error: $\|\underline{A} - \underline{S}\|_F^2 \approx NM\sigma_g^2$

Low rank: $\hat{\underline{A}}_r = \sum_{i=1}^r \sigma_{A_i} \underline{u}_i \underline{v}_i^\top \approx \hat{\underline{S}}_r + \hat{\underline{G}}_r$

$$\hat{\underline{S}}_r = \sum_{i=1}^r \sigma_{S_i} \underline{u}_i \underline{v}_i^\top \quad \hat{\underline{G}}_r \approx \sum_{i=r+1}^N M^{1/2} \sigma_g \underline{u}_i \underline{v}_i^\top$$

Bias²: $b^2(r) = \|\underline{S} - \hat{\underline{S}}_r\|_F^2$

$$b^2(r) = \left\| \sum_{i=r+1}^N \sigma_{S_i} \underline{u}_i \underline{v}_i^\top \right\|_F^2 \\ = \sum_{i=r+1}^N \sigma_{S_i}^2 \quad (\text{notes})$$

sum of squared "tail"
singular values

Variance: $v(r) = \|\hat{\underline{G}}_r\|_F^2$

$$v(r) = \left\| \sum_{i=1}^r M^{1/2} \sigma_g \underline{u}_i \underline{v}_i^\top \right\|_F^2 \\ = rM\sigma_g^2$$

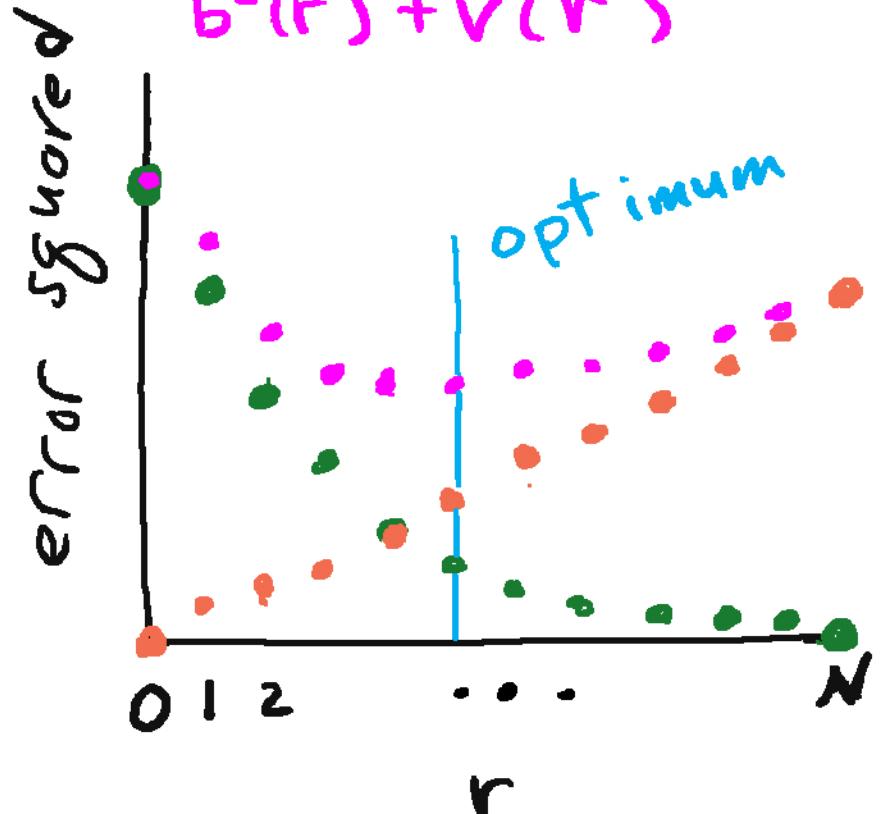
dimensions x variance
dimension

Trading bias for variance

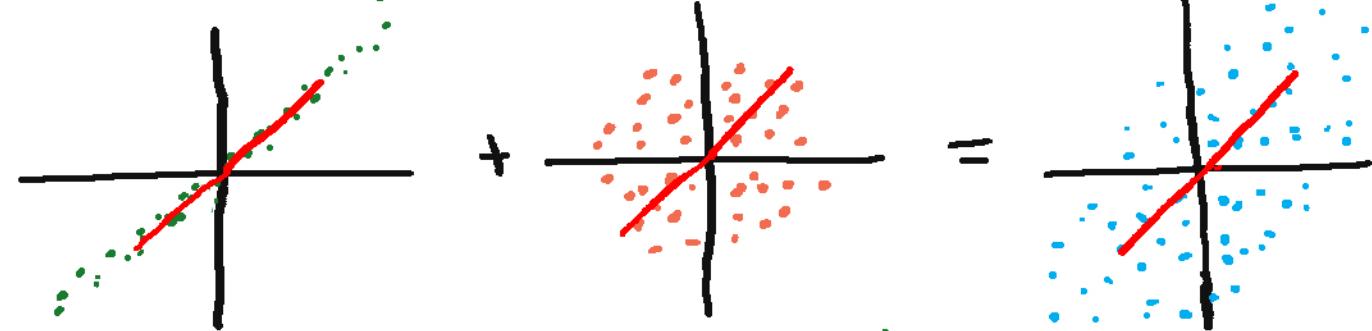
$b^2(r) = \sum_{i=r+1}^N \sigma_{S_i}^2$ decreases as r increases

$V(r) = rM\sigma_g^2$ increases as r increases

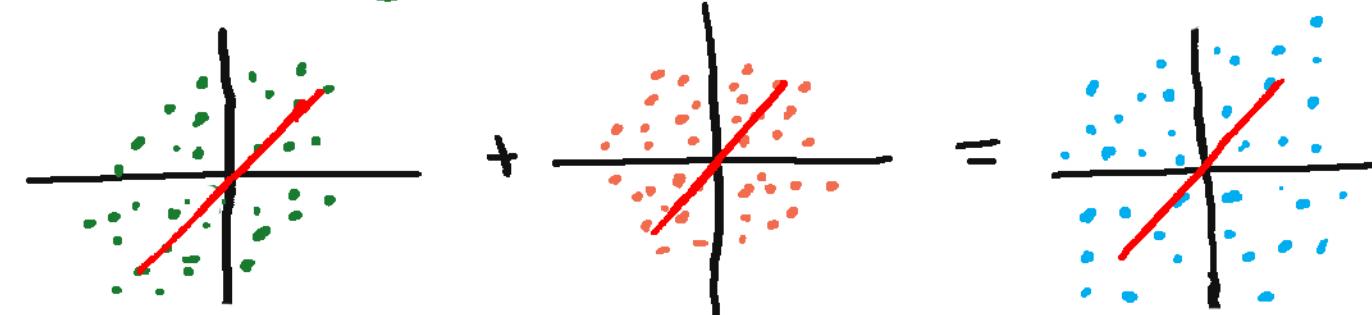
$$b^2(r) + V(r)$$



S highly structured $\sigma_1 \gg \sigma_2$



S weakly structured $\sigma_1 \approx \sigma_2$



Copyright 2019
Barry Van Veen

Eigendecomposition, SVD, and Power Iterations

Objectives

- Define eigenvectors and eigenvalues
- Relate the eigen decomposition to SVD
- Power iterations for computing eigenvector with largest eigenvalue

Eigendecomposition applies to square matrices ²

Eigenvector \underline{e}_i , eigenvalue λ_i , \underline{B} ($k \times k$)

$$\underline{B} \underline{e}_i = \lambda_i \underline{e}_i \quad \text{matrix mult} \Leftrightarrow \text{scalar mult}$$

$$\underline{e}_i \xrightarrow{\underline{B}} \lambda_i \underline{e}_i \quad i=1, 2, \dots, k$$

- K eigenvalues, possibly complex valued
- Distinct $\lambda_i \Rightarrow$ linearly independent \underline{e}_i
- Symmetric $\underline{B} \Rightarrow K$ orthonormal $\underline{e}_i \quad \underline{E} \underline{E}^T = \underline{E}^T \underline{E} = \underline{I}$

$$\underline{B} \underline{e}_i = \lambda_i \underline{e}_i \Rightarrow \underline{B} [\underline{e}_1 \underline{e}_2 \dots \underline{e}_k] = [\underline{e}_1 \underline{e}_2 \dots \underline{e}_k] \begin{bmatrix} \lambda_1 & & & \\ & \lambda_2 & & \\ & & \ddots & \\ & & & \lambda_k \end{bmatrix}$$
$$\underline{B} \underline{E} = \underline{E} \underline{\Lambda} \Rightarrow \underline{B} = \underline{E} \underline{\Lambda} \underline{E}^T = \sum_{i=1}^k \lambda_i \underline{e}_i \underline{e}_i^T$$

Symmetric PSD matrices and SVD

3

$$\underline{A} = [\underline{a}_1 \ \underline{a}_2 \cdots \underline{a}_M] = \begin{bmatrix} \underline{x}_1^T \\ \underline{x}_2^T \\ \vdots \\ \underline{x}_N^T \end{bmatrix} = \underline{U} \underline{\Sigma} \underline{V}^T \quad (N \times M, N > M)$$

1) $\underline{B} = \underline{A} \underline{A}^T = \sum_{i=1}^M \underline{a}_i \underline{a}_i^T$ full SVD

$$\underline{B} = \underline{U} \underline{\Sigma} \underline{V}^T \underline{V} \underline{\Sigma}^T \underline{U}^T = \underline{U} \underline{\Sigma} \underline{\Sigma}^T \underline{U}^T = \underline{U} \begin{bmatrix} \sigma_1^2 & & & \\ & \ddots & & 0 \\ & & \sigma_n^2 & 0 \\ 0 & & & \ddots 0 \end{bmatrix} \underline{U}^T$$

left SV of $\underline{A} \Leftrightarrow$ eigenvectors \underline{B}

$$\lambda_i = \begin{cases} \sigma_i^2 & i=1, 2, \dots, M \\ 0 & i=M+1, \dots, N \end{cases}$$

2) $\underline{B} = \underline{A}^T \underline{A} = \sum_{i=1}^N \underline{x}_i \underline{x}_i^T$
 $= \underline{V} \underline{\Sigma}^T \underline{U}^T \underline{U} \underline{\Sigma} \underline{V}^T = \underline{V} \underline{\Sigma}^T \underline{\Sigma} \underline{V}^T = \underline{V} \begin{bmatrix} \sigma_1^2 & & & \\ & \ddots & & 0 \\ & & \sigma_n^2 & 0 \\ 0 & & & \ddots 0 \end{bmatrix} \underline{V}^T$

right SV of $\underline{A} \Leftrightarrow$ eigenvectors \underline{B} , $\lambda_i = \sigma_i^2, i=1, 2, \dots, N$

Power iteration for computing 1st principal component

A: $N \times M$, $N \gg M$ want \underline{v}_1 , 1st principal component
 right SV of A, eigenvector of $\underline{B} = \underline{A}^T \underline{A}$ ($M \times M$)

Power Iteration

pick \underline{c}_0 (random)
 for $k=1, 2, \dots$ to converge

$$\underline{c}_k = \underline{B} \underline{c}_{k-1} / \|\underline{B} \underline{c}_{k-1}\|_2$$

end

$$\underline{v}_1 = \underline{c}_{\text{end}}$$

$$\underline{B} \underline{c}_{k-1} = \underbrace{\underline{B} \cdot \underline{B} \cdots \underline{B}}_{k \text{ times}} \underline{c}_0 = \underline{B}^k \underline{c}_0$$

$$\begin{aligned} \underline{B}^k &= \underbrace{\underline{V} \underline{\Lambda} \underline{V}^T \underline{V} \underline{\Lambda} \underline{V}^T \cdots \underline{V} \underline{\Lambda} \underline{V}^T}_{k \text{ times}} \\ &= \underline{V} \underline{\Lambda}^k \underline{V}^T \end{aligned}$$

$$\text{Let } \underline{c}_0 = \underline{V} \underline{g} = \underline{V} \begin{bmatrix} g_1 \\ g_2 \\ \vdots \\ g_M \end{bmatrix}$$

$$\begin{aligned} \underline{B}^k \underline{c}_0 &= \underline{V} \underline{\Lambda}^k \underline{V}^T \underline{V} \underline{g} \\ &= \underline{V} \underline{\Lambda}^k \underline{g} \end{aligned}$$

Power iteration ...

$$c_k = \underline{B} c_{k-1} / \| \underline{B} c_{k-1} \|_2 = \underline{\lambda^k g} / \| \underline{\lambda^k g} \|_2$$

$$\underline{\lambda^k g} = [v_1 \ v_2 \ \dots \ v_m] \begin{bmatrix} \lambda_1^k & 0 & & \\ 0 & \lambda_2^k & \dots & 0 \\ & & \ddots & \\ 0 & 0 & \dots & \lambda_m^k \end{bmatrix} \begin{bmatrix} g_1 \\ g_2 \\ \vdots \\ g_m \end{bmatrix}$$

but $\frac{\lambda_i}{\lambda_1} < 1$

$$= \lambda_1^k g_1 \underline{v} \begin{bmatrix} 1 & & & \\ (\frac{\lambda_2}{\lambda_1})^k & 0 & & \\ & \ddots & \ddots & 0 \\ 0 & & & (\frac{\lambda_m}{\lambda_1})^k \end{bmatrix} \begin{bmatrix} 1 \\ g_2/g_1 \\ \vdots \\ g_m/g_1 \end{bmatrix}$$

so $(\frac{\lambda_i}{\lambda_1})^k \rightarrow 0$

$$\underline{\lambda^k g} \rightarrow \lambda_1^k g_1 \underline{v} \begin{bmatrix} 1 & & & \\ 0 & 0 & & \\ 0 & 0 & \ddots & 0 \\ 0 & 0 & & 0 \end{bmatrix} \begin{bmatrix} 1 \\ \vdots \\ \vdots \end{bmatrix} = \lambda_1^k g_1 \underline{v}_1$$

$$c_k \rightarrow \frac{\lambda_1^k g_1 \underline{v}_1}{\| \lambda_1^k g_1 \underline{v}_1 \|_2} = \frac{\underline{v}_1}{\| \underline{v}_1 \|_2} = \underline{v}_1$$

**Copyright 2019
Barry Van Veen**

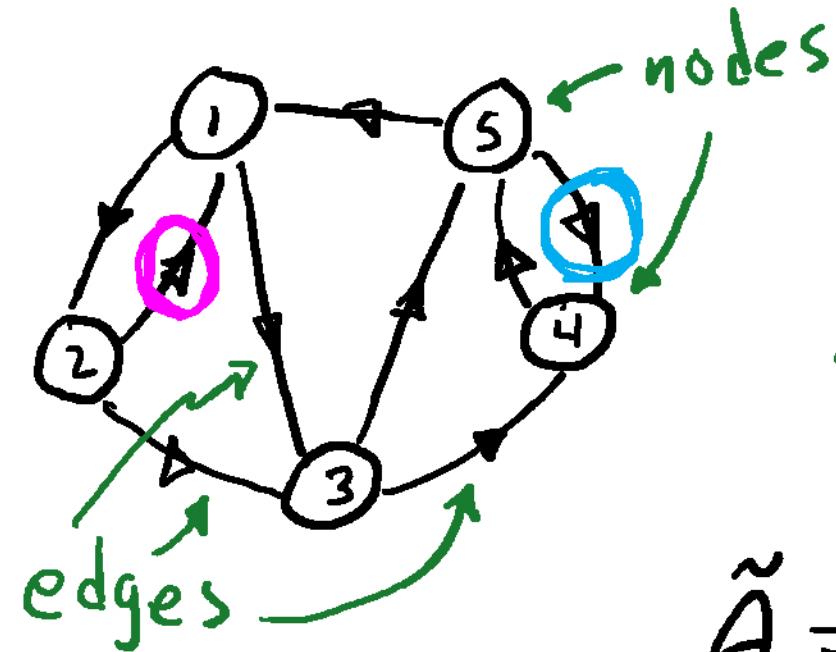
Network Graphs and the PageRank Algorithm

Objectives

- Introduce matrix representations for network graphs
- Define transition probability matrix and paths on graph
- Illustrate PageRank algorithm concepts

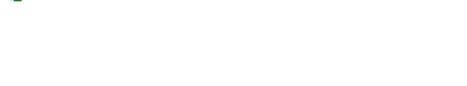
Matrices represent network graphs

2



Examples: webpages/links,
cities/roads, routers/wires

Adjacency matrix; connection
topology

edges 

$$\tilde{A} = \begin{bmatrix} 0 & 1/2 & 0 & 0 & 1/2 \\ 1/4 & 0 & 0 & 0 & 0 \\ 3/4 & 1/2 & 0 & 0 & 0 \\ 0 & 0 & 7/8 & 0 & 1/2 \\ 0 & 0 & 1/8 & 1 & 0 \end{bmatrix}$$

from

to

edge from node 1 to node 5

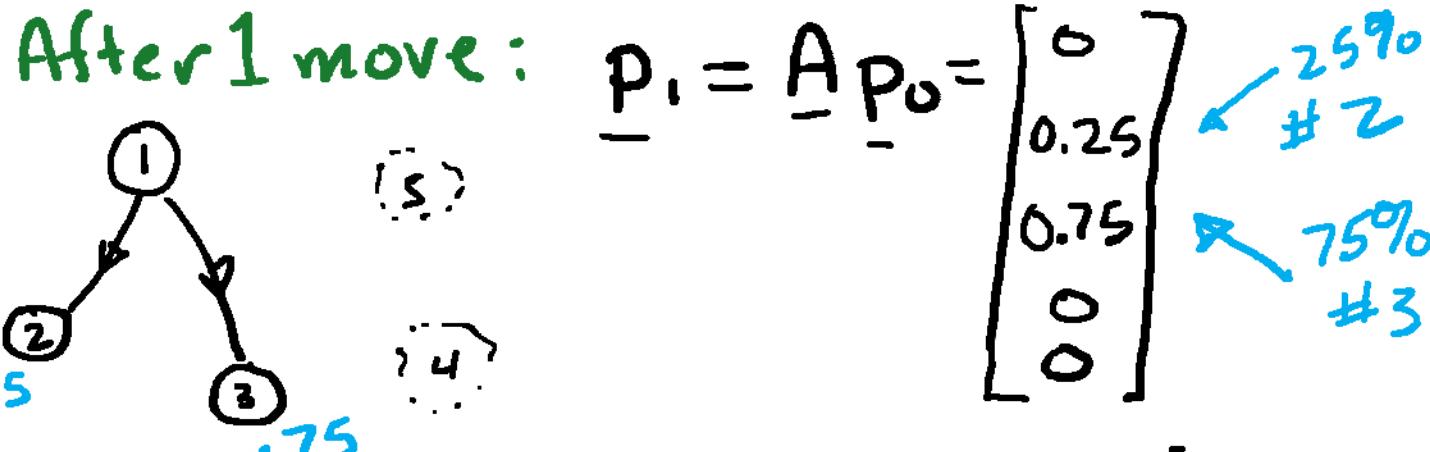
edge from node 4 to node 5

transition probability matrix - columns sum to 1

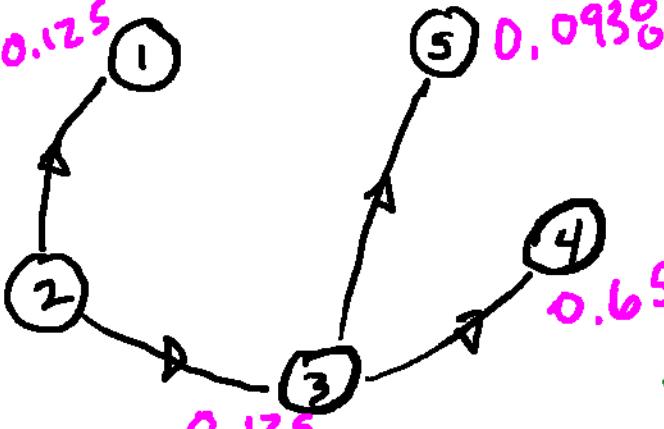
Transition probability matrix predicts "paths" 3

$$\underline{A} = \begin{bmatrix} 0 & \frac{1}{2} & 0 & 0 & \frac{1}{2} \\ \frac{1}{4} & 0 & 0 & 0 & 0 \\ \frac{3}{4} & \frac{1}{2} & 0 & 0 & 0 \\ 0 & 0 & \frac{7}{8} & 0 & \frac{1}{2} \\ 0 & 0 & \frac{1}{8} & 1 & 0 \end{bmatrix}$$

Start at node 1: $\underline{p}_0 = [1 \ 0 \ 0 \ 0 \ 0]^T$



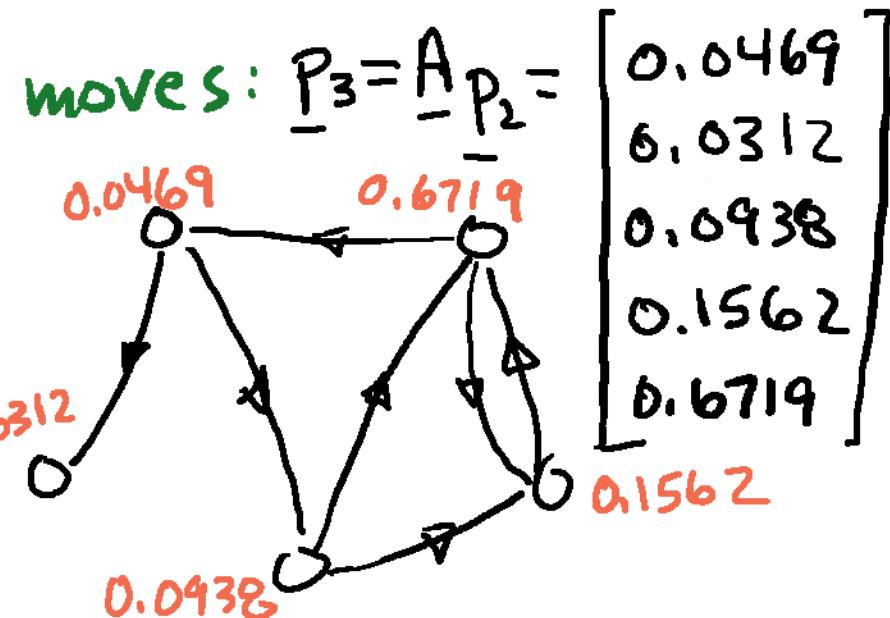
After 2 moves: $\underline{p}_2 = \underline{A} \underline{p}_1 =$



$$\begin{bmatrix} 0.125 \\ 0 \\ 0.125 \\ 0.6562 \\ 0.0938 \end{bmatrix}$$

Markov chain: next state depends only on current state

After 3 moves: $\underline{p}_3 = \underline{A} \underline{p}_2 =$



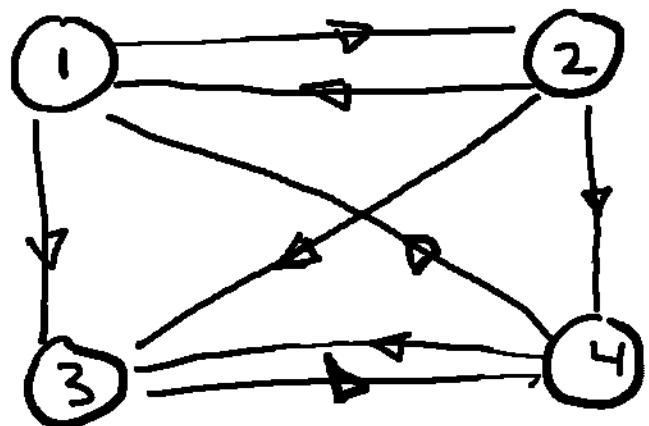
$$\begin{bmatrix} 0.0469 \\ 0.6719 \\ 0.0312 \\ 0.0938 \\ 0.1562 \\ 0.6719 \end{bmatrix}$$

0.1562

PageRank algorithm ranks web pages

4

where will I visit most?



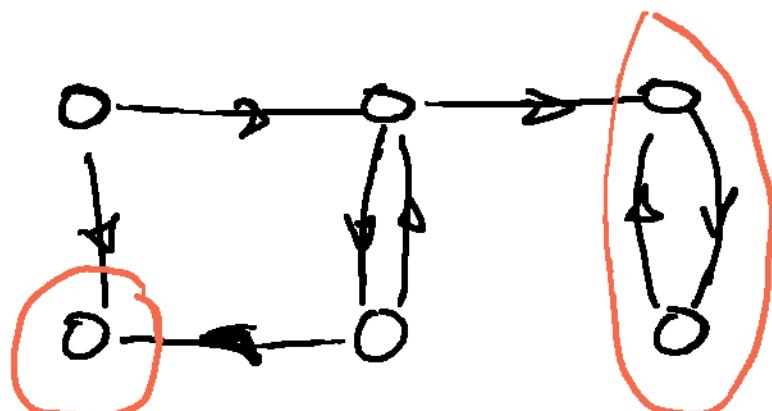
1) Adjacency matrix

$$\tilde{A} = \begin{bmatrix} 0 & 1 & 0 & 1 \\ 1 & 0 & 0 & 0 \\ 1 & 1 & 0 & 1 \\ 0 & 1 & 1 & 0 \end{bmatrix}$$

2) Normalize columns
(equal prob. outlinks)

$$A = \begin{bmatrix} 0 & \frac{1}{3} & 0 & \frac{1}{2} \\ \frac{1}{2} & 0 & 0 & 0 \\ \frac{1}{2} & \frac{1}{3} & 0 & \frac{1}{2} \\ 0 & \frac{1}{3} & 1 & 0 \end{bmatrix}$$

3) Eliminate traps



introduce small probability
to go from any node to
any other node

Transition matrix: $\underline{Q} = (1-\alpha)\underline{A} + \frac{\alpha}{N} \begin{bmatrix} 1 & 1 & \dots & 1 \\ 1 & 1 & \dots & 1 \\ \vdots & \vdots & \ddots & \vdots \\ 1 & 1 & \dots & 1 \end{bmatrix}$

$\frac{\alpha}{N}$: random jump probability

Eigenvector of \underline{Q} ranks pages

5

\underline{Q} is irreducible (no traps) column stochastic (cols sum to 1), with non negative entries \Rightarrow (Perron - Frobenius) Largest eval is 1, evect $\underline{P} = [p_1 \cdots p_N]^T$ satisfies $p_i > 0$, $\sum_i p_i = 1$

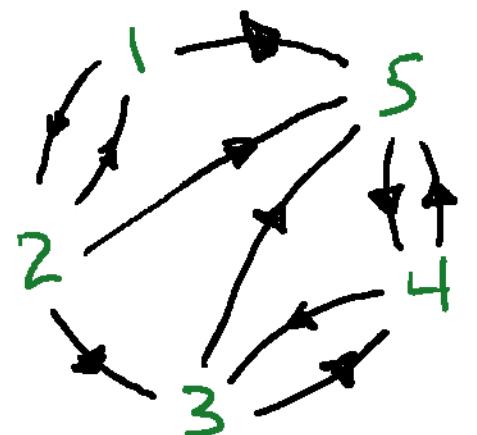
$$\underline{u} = \frac{1}{N} [1 \ 1 \ \dots \ 1]^T \quad \lim_{k \rightarrow \infty} \underline{Q}^k \underline{u} = \underline{P}$$

Steady-state
Distribution

$$\underline{Q} \underline{P} = \underline{P}$$

\underline{P} ranks importance of pages

Example:



$$\tilde{A} = \begin{bmatrix} 0 & 1 & 0 & 0 & 0 \\ 1 & 0 & 0 & 0 & 0 \\ 0 & 1 & 0 & 1 & 0 \\ 0 & 0 & 1 & 0 & 1 \\ 1 & 1 & 1 & 1 & 0 \end{bmatrix}$$

$$\underline{Q} = \begin{bmatrix} -0.002 & .332 & .002 & .002 & .002 \\ .497 & .002 & .002 & .002 & .002 \\ .002 & .332 & .002 & .497 & .002 \\ .002 & .002 & .497 & -.062 & .992 \\ -.497 & .332 & .497 & .497 & .002 \end{bmatrix} \quad (\alpha = 0.01)$$

$$\underline{u} = \begin{bmatrix} 0.2 \\ 0.2 \\ 0.2 \\ 0.2 \\ 0.2 \end{bmatrix}, \quad \underline{Qu} = \begin{bmatrix} .07 \\ .10 \\ .17 \\ .30 \\ .36 \end{bmatrix}, \quad \underline{Q^2 u} = \begin{bmatrix} .003 \\ .004 \\ -.18 \\ -.45 \\ -.30 \end{bmatrix}, \quad \dots \quad \underline{Q^{10} u} = \begin{bmatrix} .003 \\ .004 \\ .22 \\ -.44 \\ -.33 \end{bmatrix}$$

1 2 3 4 5

$$\underline{P} = \begin{bmatrix} 0.0032 & 0.0036 & 0.2211 & 0.4401 & 0.3320 \end{bmatrix}$$

Copyright 2019
Barry Van Veen

Matrix Completion

Objectives

- define the matrix completion problem
- approach missing data using low-rank models
- introduce iterative singular value thresholding

Use "patterns" to fill in missing entries 2

Ratings

matrix

$$\underline{X} \in \mathbb{R}^{N \times m}$$

	5	4	9	1	\times	N movies ↓
9	6	\times	\times	7		
\times	10	\times	2	4		
3	7	3	\times	\times		
8	\times	\times	6	2		
m users →						

Can we predict the missing entries?

Model: assume \underline{X} is well approximated with a small number of patterns

$$\underline{X} \approx \sum_{i=1}^r \underline{t}_i \underline{s}_i^T = \underline{T} \underline{S}$$

genres, actors, director...
hobbies, age, address...

Matrix completion: use known data to find
patterns and predict missing entries 3

$\Omega = \{(i,j) : \underline{x}_{ij} \text{ given}\}$ indices of known values

1) Rank minimization $\underline{X} = \underset{\underline{M}}{\operatorname{arg\,min}} \operatorname{rank}(\underline{M})$ s.t. $\underline{M}_{ij} = \underline{x}_{ij} \quad \forall i,j \in \Omega$
minimum number of patterns matching given values

Intractable!

$$\operatorname{rank} \underline{M} = \#\{l : \sigma_l > 0\}$$

2) Nuclear norm minimization

$\underline{X} = \underset{\underline{M}}{\operatorname{arg\,min}} \|\underline{M}\|_* \text{ s.t. } \underline{M}_{ij} = \underline{x}_{ij} \quad \forall i,j \in \Omega$ Nuclear / trace norm

$$\|\underline{M}\|_* = \sum_k \sigma_k$$

Computationally tractable

Iterative Singular Value Thresholding

4

is one possible algorithm

Initialize

$$\underline{M}^{(0)} = \underline{0}$$

Set threshold or r

Iterate

for $k = 1, 2, 3, \dots$

$$\underline{M}^{(k)} = \underline{M}^{(k-1)}$$

$\underline{M}_{\Sigma}^{(k)} = \underline{X}_{\Sigma}$ (fill in known values)

$$[\underline{U}, \underline{\Sigma}, \underline{V}] = \text{svd}(\underline{M}^{(k)})$$

$$\hat{\Sigma}_{ii} = \Sigma_{ii} \cdot \begin{cases} 1 & \Sigma_{ii} > \text{threshold} \\ 0 & \Sigma_{ii} \leq \text{threshold} \end{cases}$$

- or -

$$\hat{\Sigma}_{ii} = \begin{cases} \Sigma_{ii}, & i \leq r \\ 0, & i \geq r+1 \end{cases}$$

$$\underline{M}^{(k)} = \underline{U} \sum \underline{V}^T$$

$$\text{if } \|\underline{M}^{(k)} - \underline{M}^{(k-1)}\|_F < \varepsilon$$

stop

else

next k

Matrix completion is an open problem

- choosing r or threshold in ISVT
- multiple algorithms:
 - convergence
 - complexity
 - noise
- results depend on distribution of missing entries
- applications include missing pixels in images, position from partial distance info, ...

Copyright 2019
Barry Van Veen