

Review

- Gibbs sampling

- ▶ MH with proposal

- ▶ $Q(\mathbf{X} | \mathbf{X}') = P(\mathbf{X}_{B(i)} | \mathbf{X}_{\neg B(i)}) I(\mathbf{X}_{\neg B(i)} = \mathbf{X}'_{\neg B(i)}) / \#B$

- ▶ failure mode: “lock-down” $X_{\neg B(i)}$ determines $X_{B(i)}$

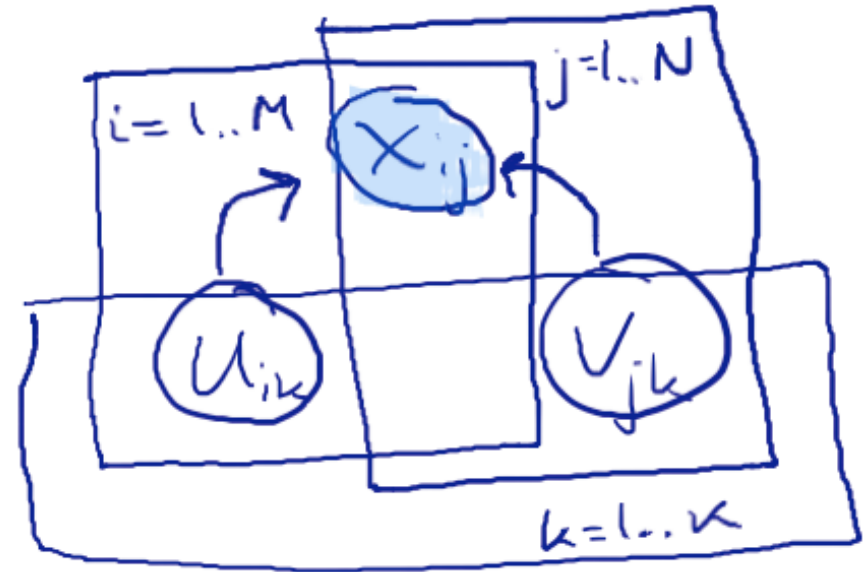
- Relational learning (properties of **sets** of entities)

- ▶ document clustering, recommender systems, eigenfaces

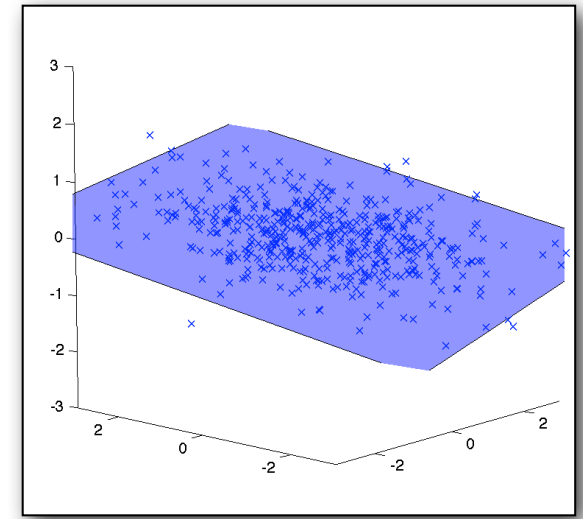
$$I(P) = \begin{cases} 1 & \text{if } P \\ 0 & \text{otherwise} \end{cases}$$

Review

- Latent-variable models
- PCA, pPCA, Bayesian PCA
 - ▶ everything Gaussian
 - ▶ $E(\underline{X} | U, V) = \underline{UV}^T$
 - ▶ MLE: use SVD
- Mean subtraction, example weights



$$E(X_{ij} | U, V) = \sum_k U_{ik} V_{jk}$$

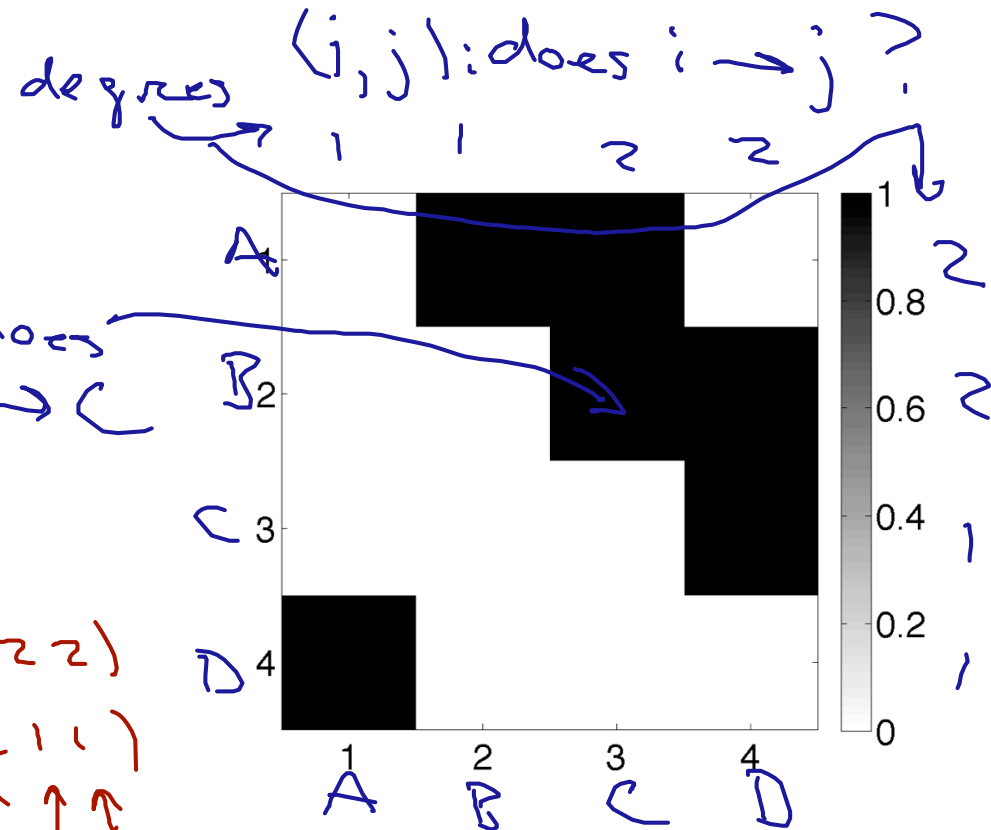
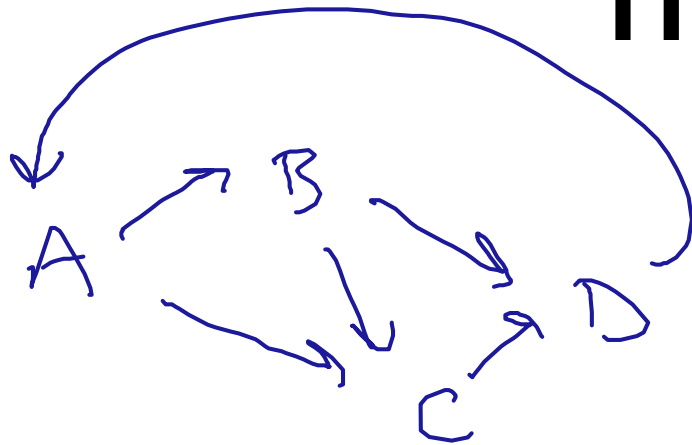


PageRank

Eigen-decomposition

- ~~SVD~~ is pretty useful: turns out to be main computational step in other models too
- A famous one: PageRank
 - ▶ Given: web graph (V, E)
 - ▶ Predict: which pages are important

PageRank: adjacency matrix



in-degrees = col sums = (1 1 2 2)
 out-degrees = row sums = (2 2 1 1)
 ↑ ↑ ↑ ↑
 for A B C D

A

Random surfer model

$$\alpha = 0.15$$

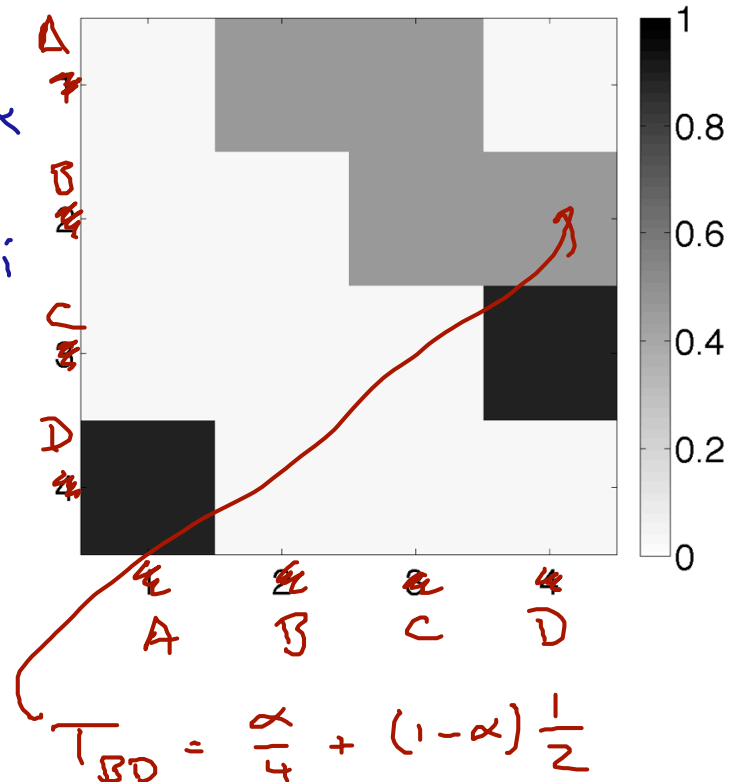
T matrix

- ▶ W.p. α : jumps uniformly to new page
- ▶ W.p. $(1-\alpha)$: follow random link

$$T_{ij} = P(j|i) = \alpha \frac{1}{N} + (1-\alpha) \frac{A_{ij}}{d_i}$$

$$d_i = \sum_j A_{ij}$$

- ▶ Intuition: page is important if a random surfer is likely to land there



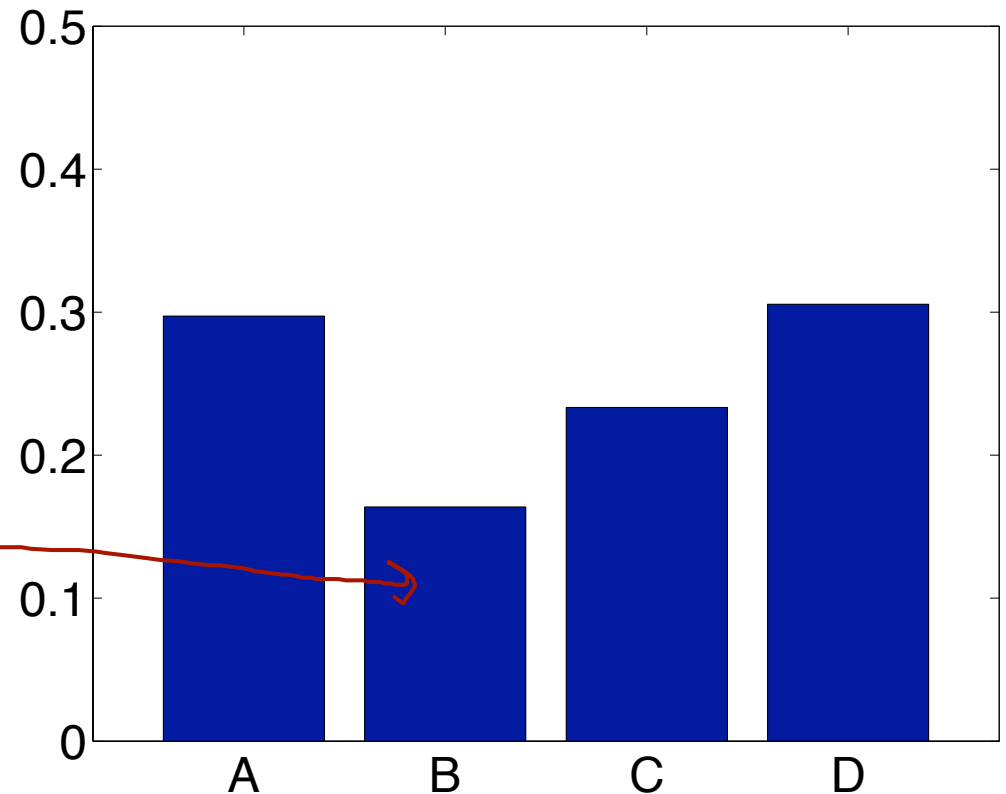
Stationary distribution

$$\pi_j = \sum_i T_{ij} \pi_i$$

$$\pi^T T = \pi^T$$

$\nearrow \pi$ is left eigenvector of T
w/ eigenvalue 1

random
surfer
visits B least
 \Rightarrow B is least
important



Thought experiment

- What if A is symmetric?
 - ▶ note: we're going to stop distinguishing A, A'

$$\begin{aligned}
 T_{ij} &= A_{ij} / d_i & d_i &= \sum_j A_{ij} & A_{ij} &= A_{ji} \\
 D &= \text{diag}(d) & T &= D^{-1} A \\
 \pi^T D^{-1} A &= \pi^T \\
 \text{suppose } \pi &= d \Rightarrow \pi^T D^{-1} = 1^T \\
 &\Rightarrow \pi^T D^{-1} A = 1^T A = d^T = \pi^T
 \end{aligned}$$

- So, stationary dist'n for symmetric A is: d
- What do people do instead?

2nd eigenvector of T

boring!!

Spectral embedding

- Another famous model: spectral embedding (and its cousin, spectral clustering)
- Embedding: assign low-D coordinates to vertices (e.g., web pages) so that similar nodes in graph \Rightarrow nearby coordinates
 - ▶ A, B similar = random surfer tends to reach the same places when starting from A or B

Where does random surfer reach?

- Given graph: adjacency A $T = D^{-1}A$ degree = d $D = \text{diag}(d)$
- Start from distribution π
 - after 1 step: $P(k \mid \pi, 1\text{-step}) = \sum_i \pi(i) P(k \mid i)$
 - after 2 steps: $P(k \mid \pi, 2\text{-step}) = \pi^T T T = \sum_i \pi(i) T_{ik} = \pi^T T$
 - after t steps: $\pi^T T^t$

SVD of normalized adjacency / transition matrix $\bar{A} = U \Sigma U^T$

Σ diagonal
 U orthogonal

Similarity

$$U^T U = I$$

$$D^{-1/2} T D^{-1/2} = D^{-1/2} D^{-1} A D^{-1/2}$$

$$= D^{-1/2} A D^{-1/2} = \bar{A}$$

symmetric

since we assumed A symmetric

- A, B similar = random surfer tends to reach the same places when starting from A or B

$$\frac{1}{\sqrt{d_i}} u_i (\Sigma^t U^T D^{1/2})$$

$$\bullet P(k \mid \pi, t\text{-step}) =$$

$$\pi^T T^t = \pi^T (D^{-1} A)^t$$

$$= \pi^T D^{-1} A D^{-1} A \dots$$

$$\uparrow \quad \uparrow \quad \uparrow \quad \uparrow$$

$$D^{-1/2} D^{-1/2} D^{-1/2} D^{-1/2}$$

$$= \pi^T D^{-1/2} (D^{-1/2} A D^{-1/2}) \dots$$

$$(D^{-1/2} A D^{-1/2}) \dots$$

$$= \pi^T D^{-1/2} \bar{A}^t D^{1/2}$$

$$= \pi^T D^{-1/2} U \Sigma^t U^T D^{1/2}$$

$$= \pi^T D^{-1/2} U \Sigma^t U^T D^{1/2}$$

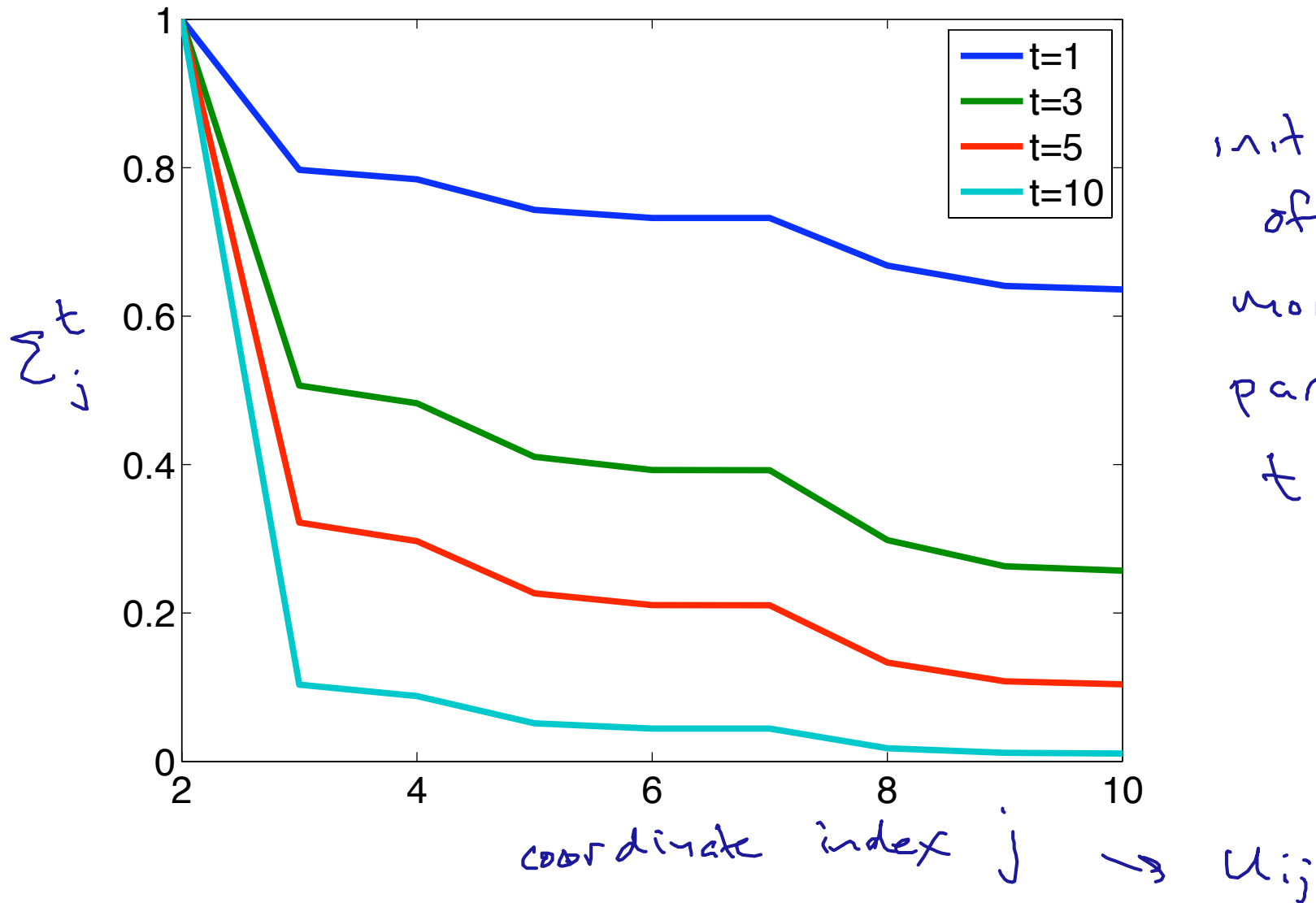
► If π has all mass on i :

► Compare i & j :

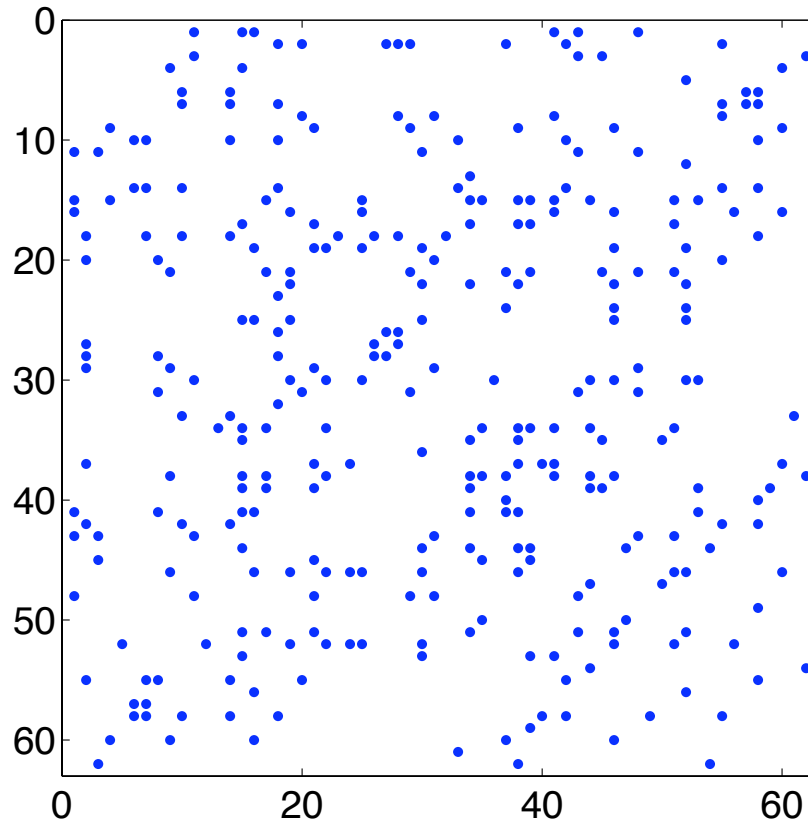
► Role of Σ^t : see next slide

nodes w/ similar u_i & u_j get to similar places

Role of Σ^t (real data)

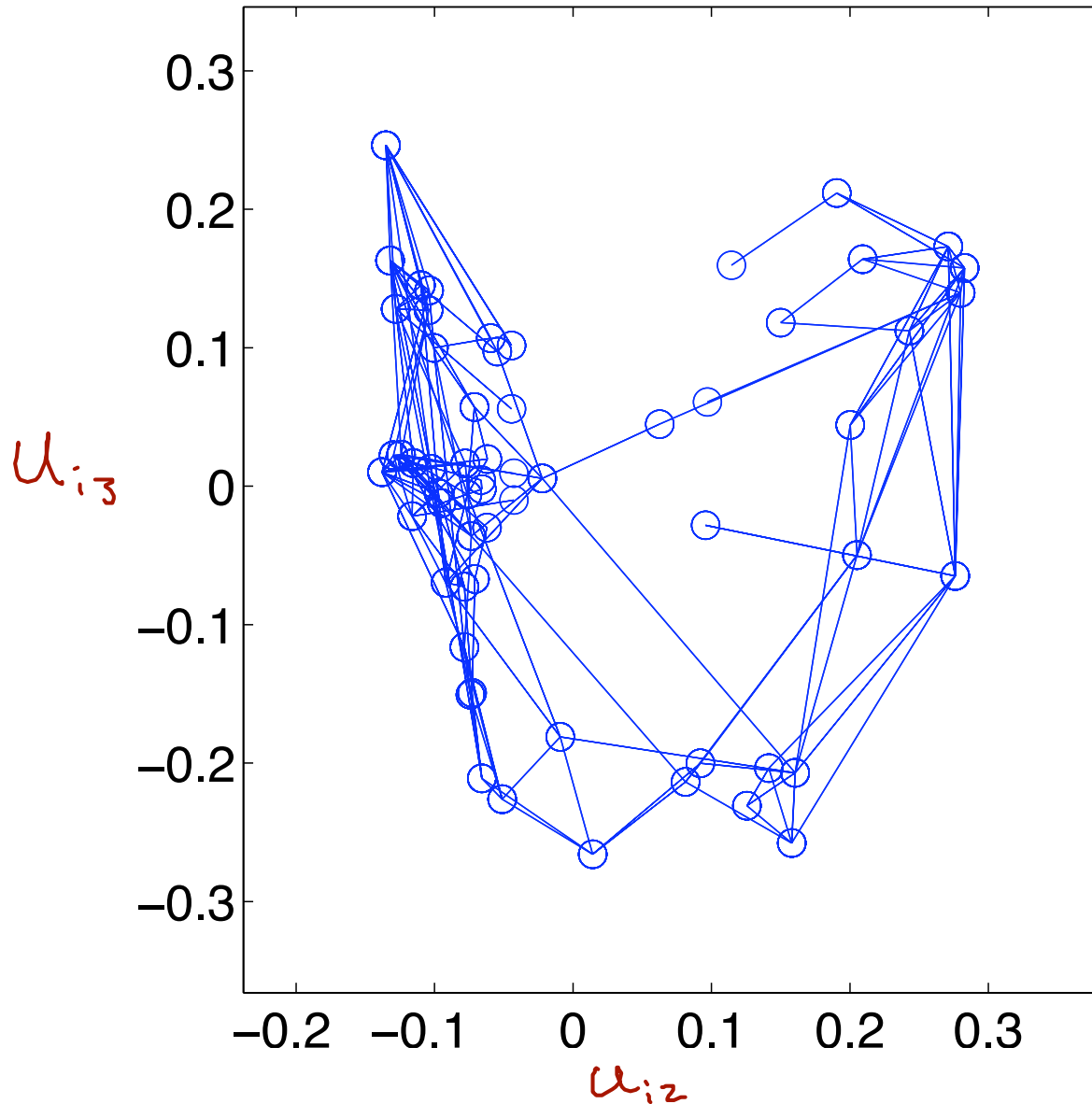


Example: dolphins

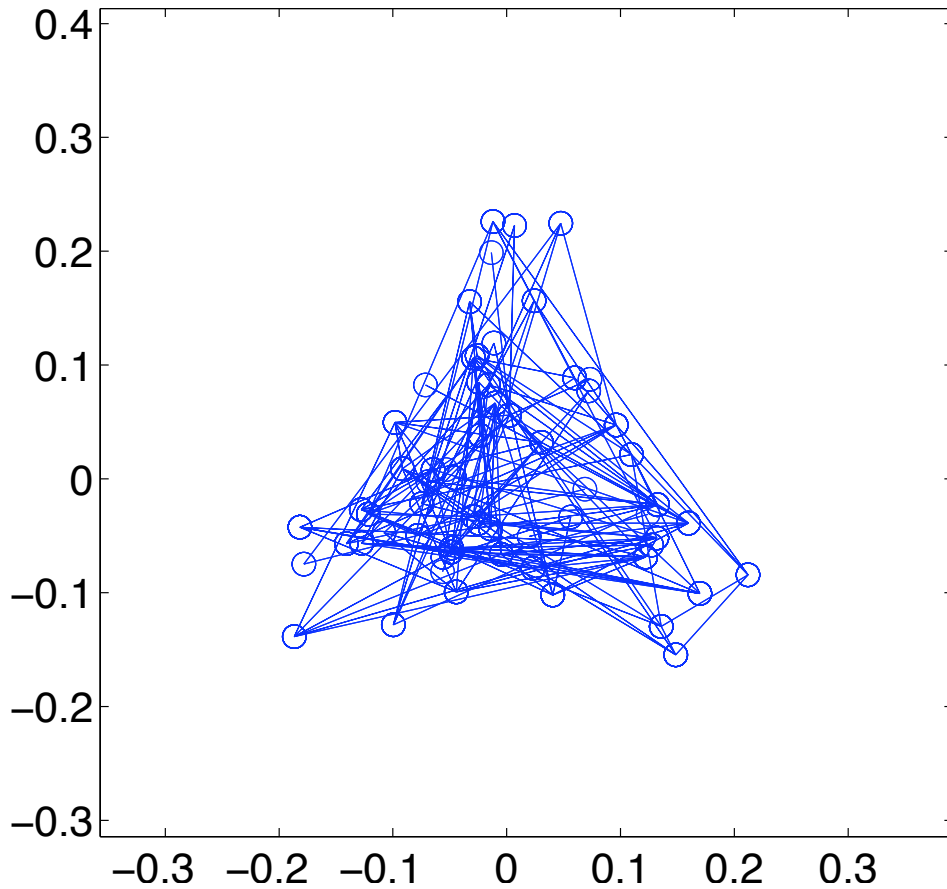


- 62-dolphin social network near Doubtful Sound, New Zealand
 - ▶ $A_{ij} = 1$ if dolphin i friends dolphin j

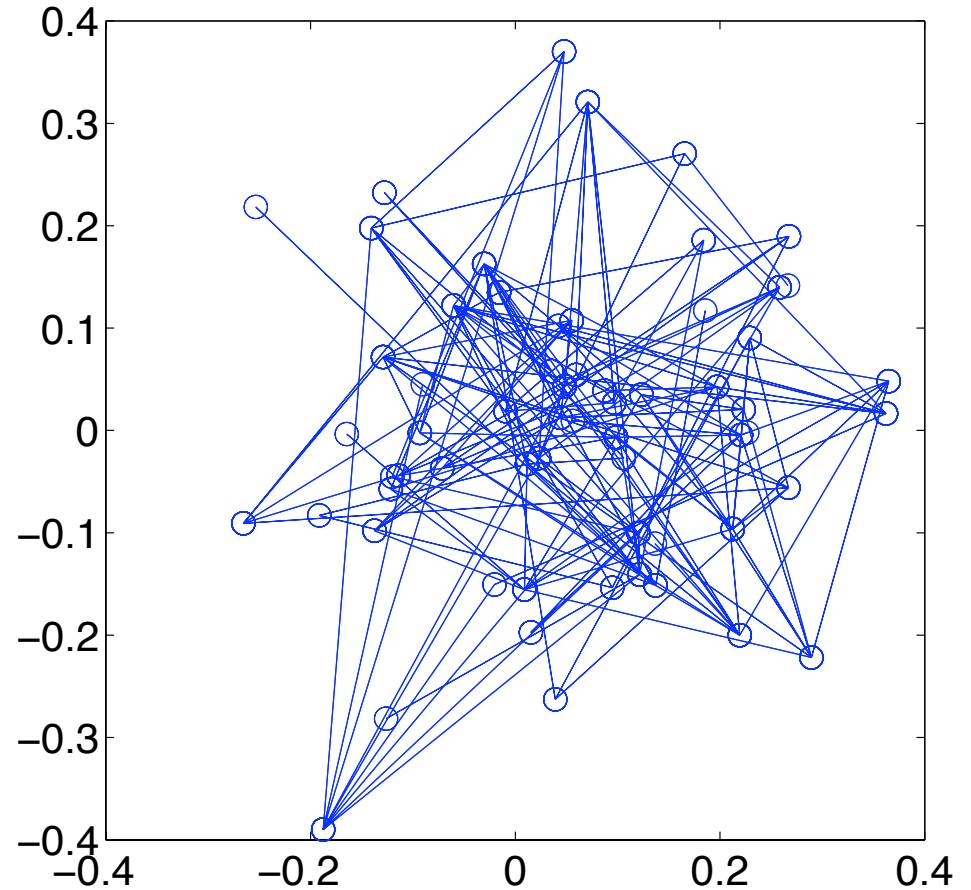
Dolphin network



Comparisons

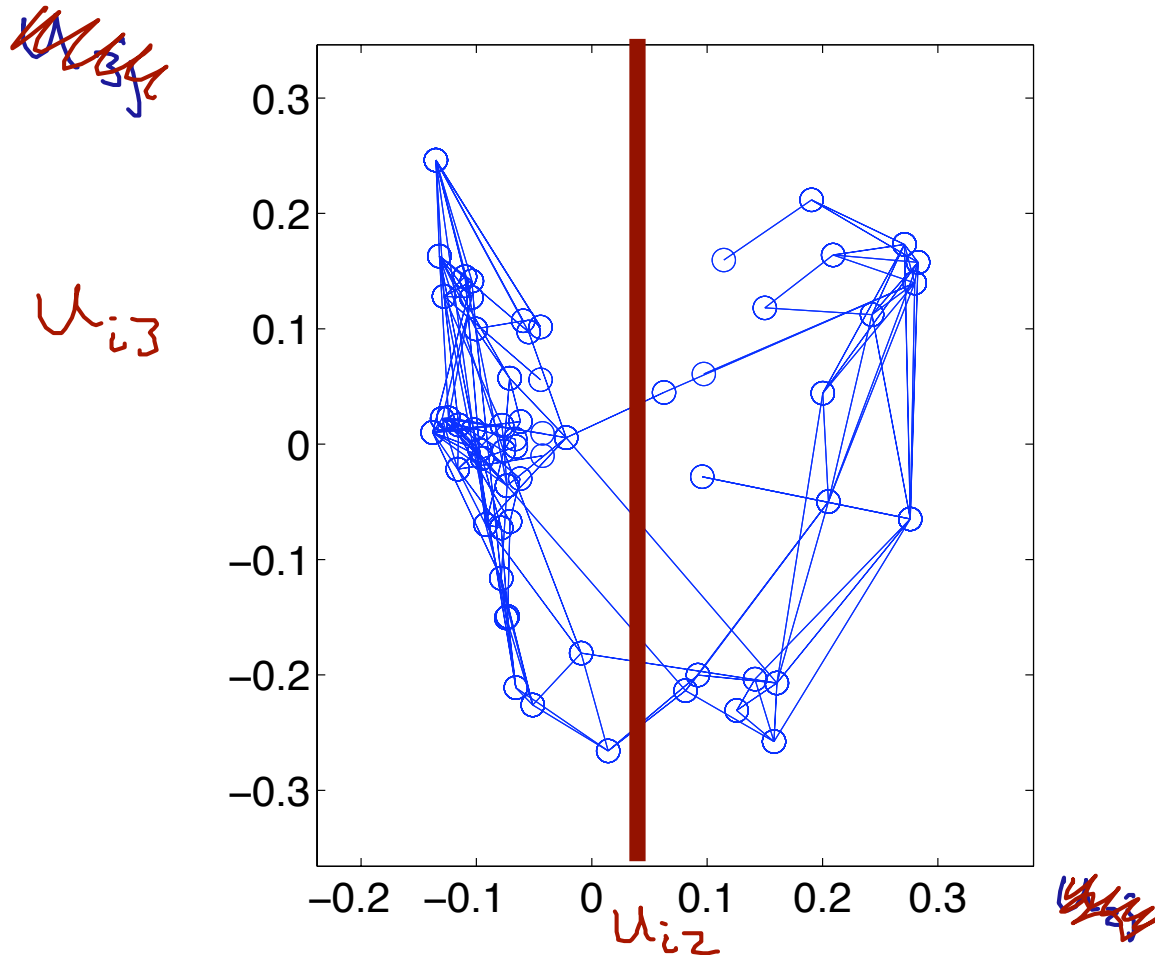


spectral embedding of
random data



random embedding of
dolphin data

Spectral clustering



- Use your favorite clustering algorithm on coordinates from spectral embedding

PCA: the good, the bad, and the ugly

- The good: simple, successful
 - The bad: linear, Gaussian
 - ▶ $E(X) = UV^T$
 - ▶ $X, U, V \sim \text{Gaussian}$
 - The ugly: failure to generalize to new entities
- fold-in problem*

Consistency

(convergence is w.p. 1)

consistent = as $|data| \rightarrow \infty$

est. params \rightarrow true params.

- Linear & logistic regression are **consistent**

- What would consistency mean for PCA?

▶ forget about row/col means for now $E(x|u,v) = uv^T$

- Consistency:

all increasing

▶ #users, #movies, #ratings (= nnz(W)) \leftarrow data

▶ numel(U), numel(V) \leftarrow parameters

▶ consistency = want $\hat{U} \rightarrow U_{true}$ \leftarrow infinite size

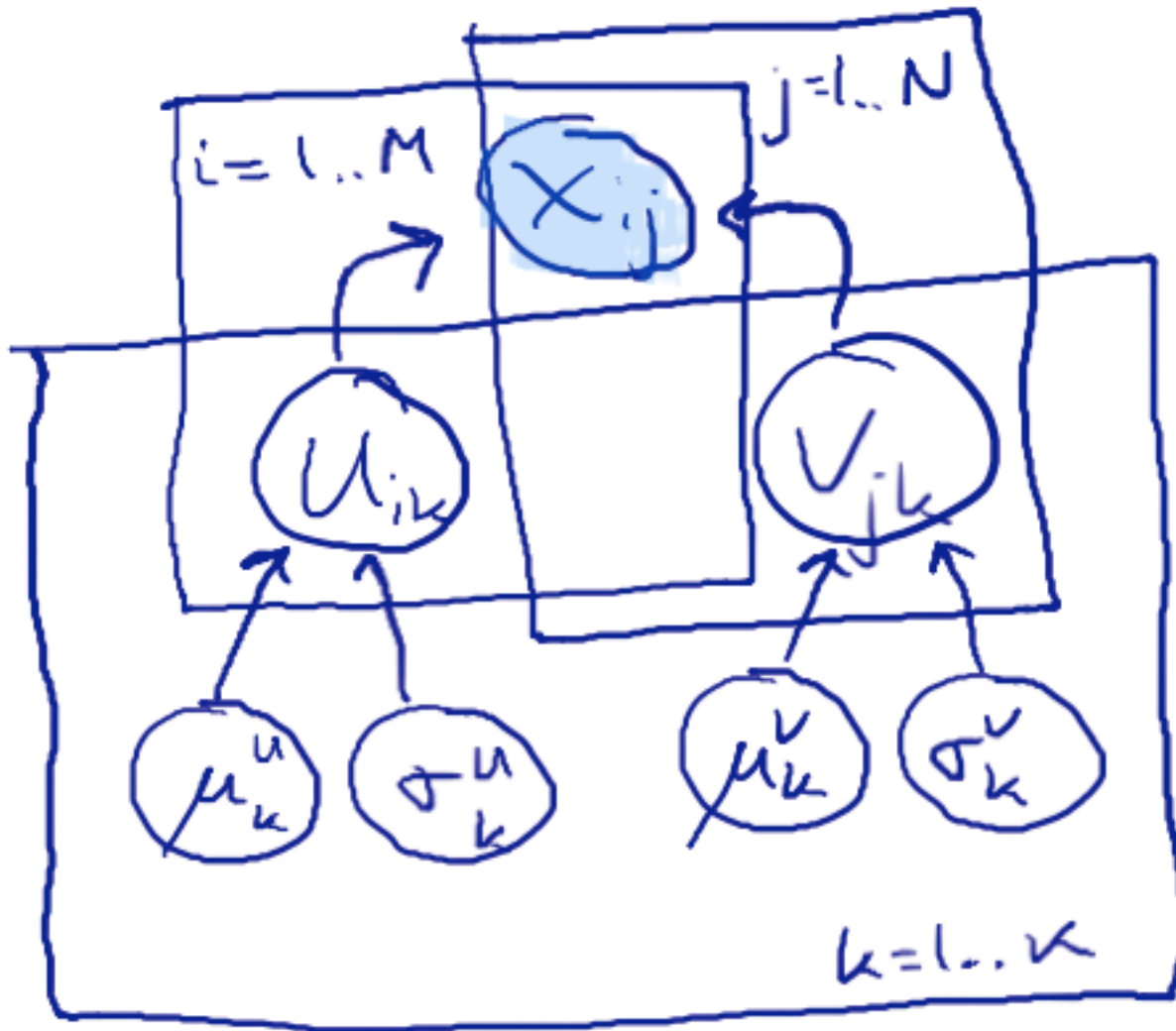
can never happen!!

(some entries U_{ij} not seen yet)

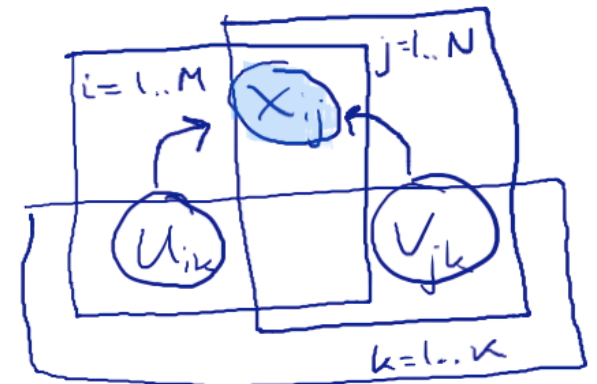
Failure to generalize

- What does this mean for generalization?
 - ▶ new user's rating of movie_j: only info is *nothing!*
 - ▶ new movie rated by user_i: only info is *nothing!*
 - ▶ all our carefully-learned factors give us: *nothing!*
- Generalization is:
only to new entries in existing rows/cols

Hierarchical model



old, non-hierarchical
model



Benefit of hierarchy

- Now: only K μ^U latents, K μ^V latents (and corresponding σ s)
 - ▶ can get consistency for these ^{only} if we observe more and more X_{ij}

- For a new user or movie:

μ^U tells us corresponding thing
 μ^V tells us how existing users will probably rate movie
i.e., how user will probably rate existing movies

Mean subtraction

- Can now see that mean subtraction is a special case of our hierarchical model
 - ▶ Fix $V_{j1} = 1$ for all j ; then $U_{i1} = \text{row mean}$
 - ▶ Fix $U_{i2} = 1$ for all i ; then $V_{j2} = \text{col mean}$
 - ▶ global mean: unnecessary \rightarrow include in row or col mean