

Review

- Models that use SVD or eigen-analysis

- ▶ PageRank: eigen-analysis of **random surfer** transition matrix

- ▶ usually uses only *first* eigenvector

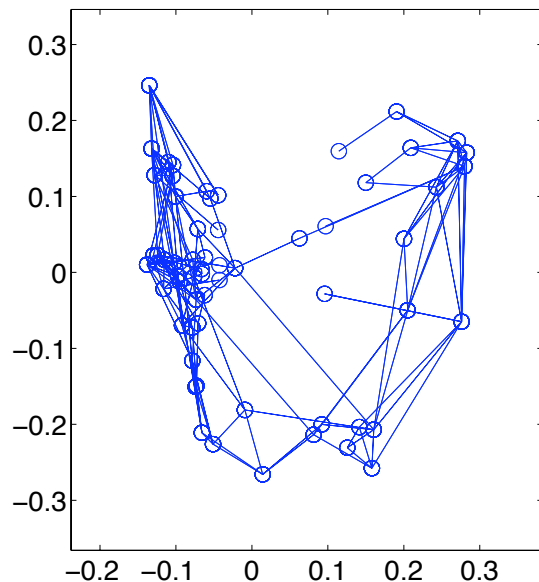
- ▶ Spectral embedding: eigen-analysis (or equivalently SVD) of random surfer model in **symmetric** graph

- ▶ usually uses 2nd–Kth EVs (small K)

- ▶ first EV is boring

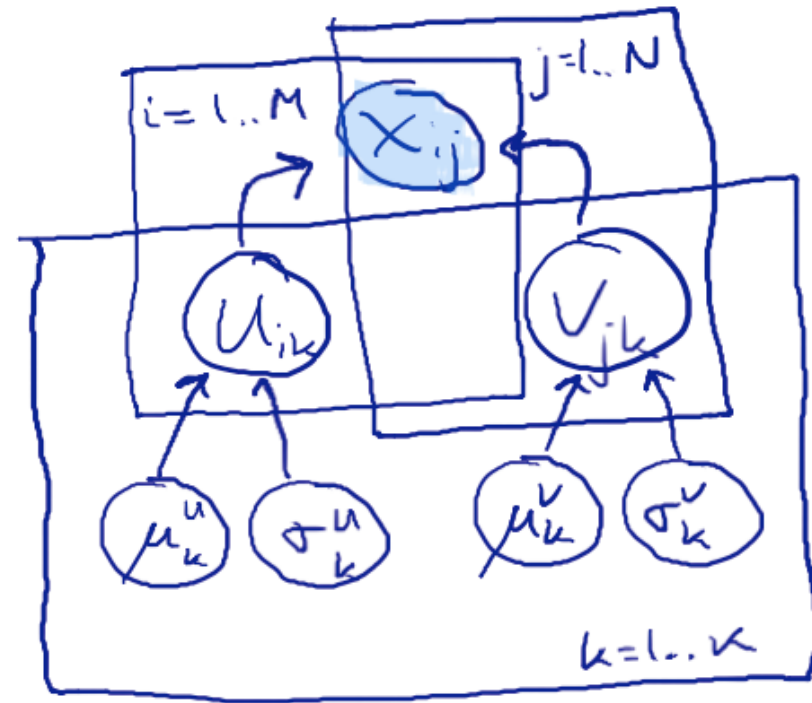
- ▶ Spectral clustering = spectral embedding followed by clustering

dolphin
friendships



Review: PCA

- 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
 - ▶ Partial answer: **hierarchical PCA**

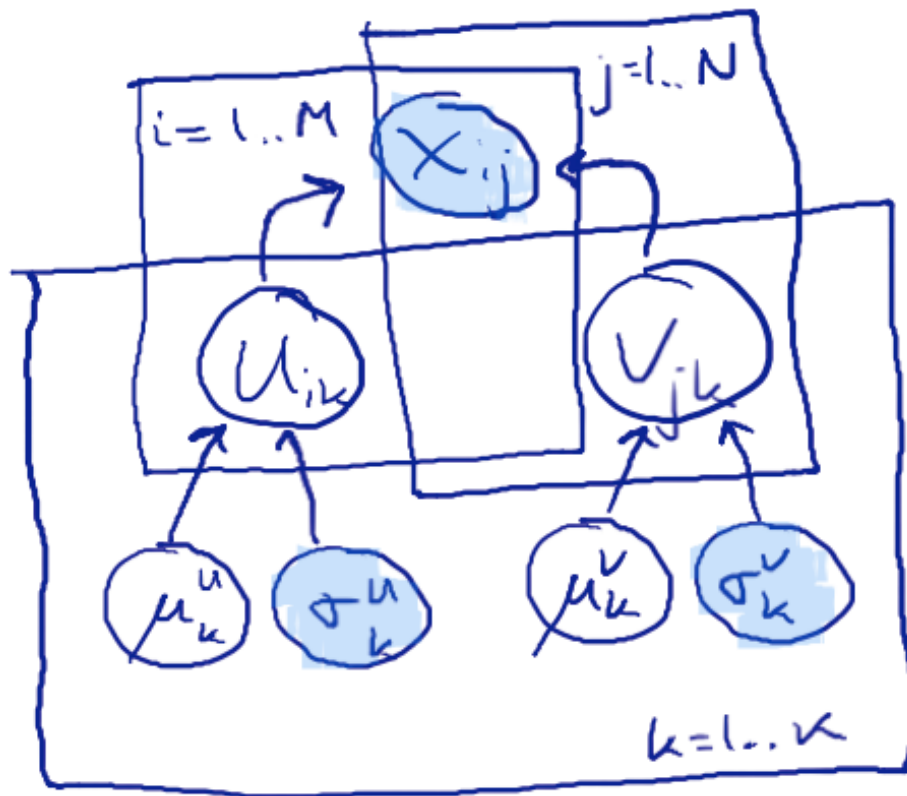


What about the second rating for a new user?

- MLE/MAP of U_i from one rating: *overfit*
 - ▶ knowing μ^U : *helps — but doesn't fix*
 - ▶ result: *confidently wrong*
- How should we fix? *be Bayesian*
- Note: often have only a few ratings per user

so, significant posterior uncertainty \Rightarrow Bayes is important no matter how much data we have

MCMC for PCA



Need:

$$\begin{aligned} &\rightarrow P(U_{ik} | V, \mu_k^u, \text{rest of } U_{i\cdot}) \\ &\quad P(V_{jk} | U, \mu_k^v, \text{rest of } V_{j\cdot}) \\ &\rightarrow P(\mu_k^u | U_{\cdot k}) \\ &\quad P(\mu_k^v | V_{\cdot k}) \end{aligned}$$

- Can do Bayesian inference by Gibbs sampling—for simplicity, assume σ s known

Recognizing a Gaussian

- Suppose $X \sim N(X \mid \mu, \sigma^2)$
- $L = -\log P(X=x \mid \mu, \sigma^2) = \log \sqrt{2\pi} \sigma + \frac{1}{2\sigma^2} (x - \mu)^2$
 - ▶ $dL/dx = \frac{1}{\sigma^2} (x - \underline{\mu})$
 - ▶ $d^2L/dx^2 = \underline{1/\sigma^2}$
- So: if we see $d^2L/dx^2 = a$, $dL/dx = a(x - b)$
 - ▶ $\mu = b$ $\sigma^2 = 1/a$

$$\frac{dL}{du_{ik}} = p^2 u_{ik} -$$

$$\sum_j \frac{v_{jk}}{\sigma^2} E_{ij} - \frac{\mu_k^u}{\sigma_k^u{}^2}$$

$$E_{ij} = x_{ij} - \sum_{k' \neq k} u_{ik'} v_{jk'}$$

Gibbs step for an

element of ~~μ~~ u_{ik}

(*) ↑

$$\bullet L = \text{const.} + \sum_{ij} (x_{ij} - \sum_{k'} u_{ik'} v_{jk'})^2 / 2\sigma^2 + \sum_{ik} (u_{ik} - \mu_k^u)^2 / 2(\sigma_k^u)^2 + \sum_{jk} (v_{jk} - \mu_k^v)^2 / 2(\sigma_k^v)^2$$

$$\frac{dL}{du_{ik}} = \sum_j (x_{ij} - \sum_{k'} u_{ik'} v_{jk'}) / \sigma^2 (-v_{jk}) + (u_{ik} - \mu_k^u) / (\sigma_k^u)^2 = p^2 u_{ik} - \sum_j \frac{v_{jk}}{\sigma^2} E_{ij} - \mu_k^u / \sigma_k^u{}^2$$

see (*)

$$\frac{d^2 L}{du_{ik}^2} = v_{jk}^2 / \sigma^2 + 1 / (\sigma_k^u)^2 = p^2$$

post. var. (of u_{ik}) = $1/p^2$

post. mean

$$\frac{\sum_j \frac{v_{jk}}{\sigma^2} E_{ij} + \frac{1}{\sigma_k^u{}^2} \mu_k^u}{p^2}$$

Gibbs: element of ~~U~~ μ_k^u

- $L = \text{const.} + \sum_{ij} (x_{ij} - \sum_{k'} u_{ik'} v_{jk'}) / 2\sigma^2$
 $+ \sum_{ik} (u_{ik} - \mu_k^u)^2 / 2(\sigma_k^u)^2 + \sum_{jk} (v_{jk} - \mu_k^v)^2 / 2(\sigma_k^v)^2$

- ~~$dL/d\mu_{ik}^u$~~ $dL/d\mu_k^u = \sum_i (u_{ik} - \mu_k^u) / \sigma_k^u^2 (-1)$
 $d^2L/d\mu_k^u^2 = 1/\sigma_k^u^2$

- ~~$d^2L/(d\mu_{ik}^u)^2$~~

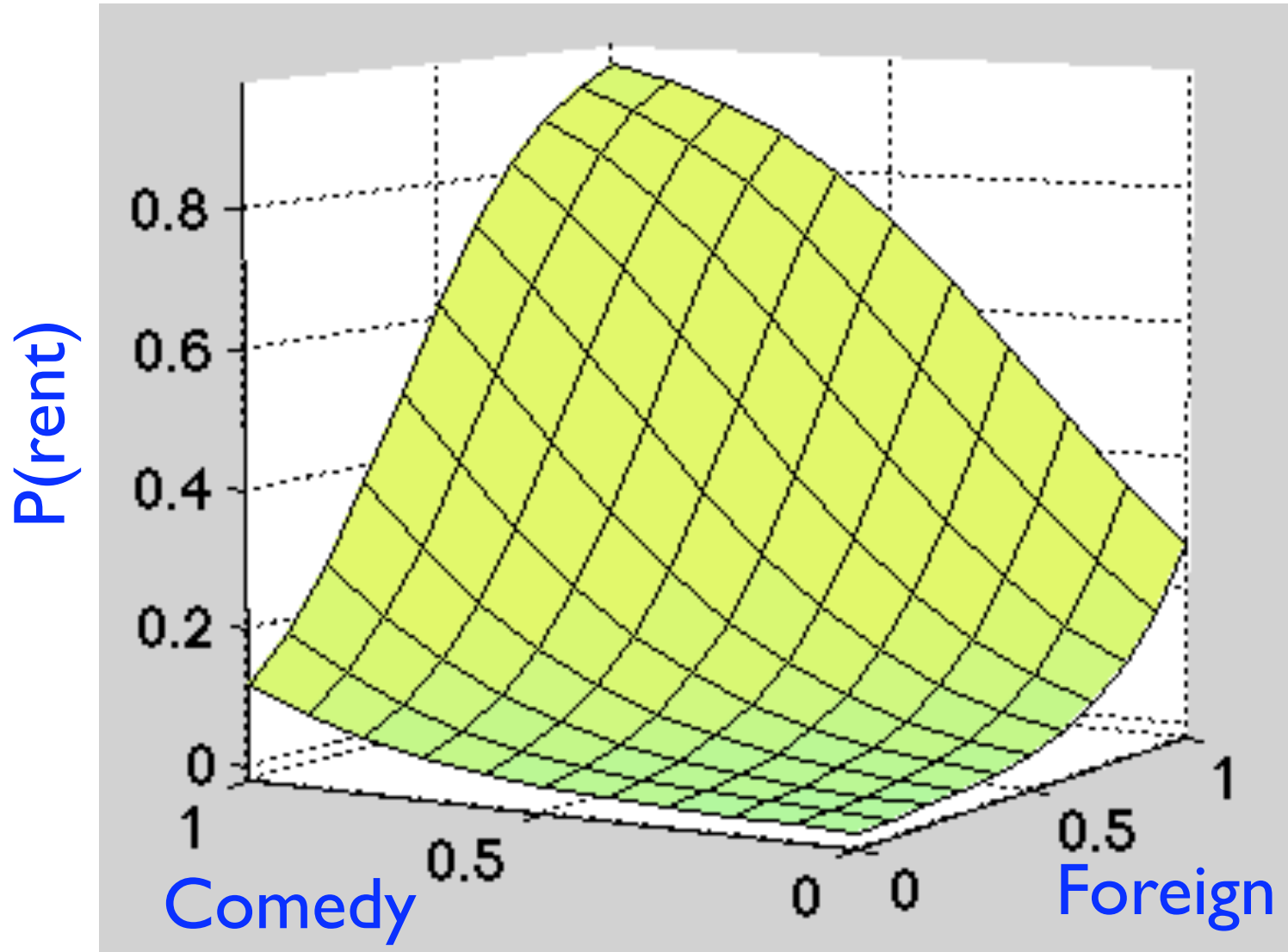
► post. mean = $\frac{\sum_i u_{ik}}{M}$

post. var. = $\frac{\sigma_k^u^2}{M}$

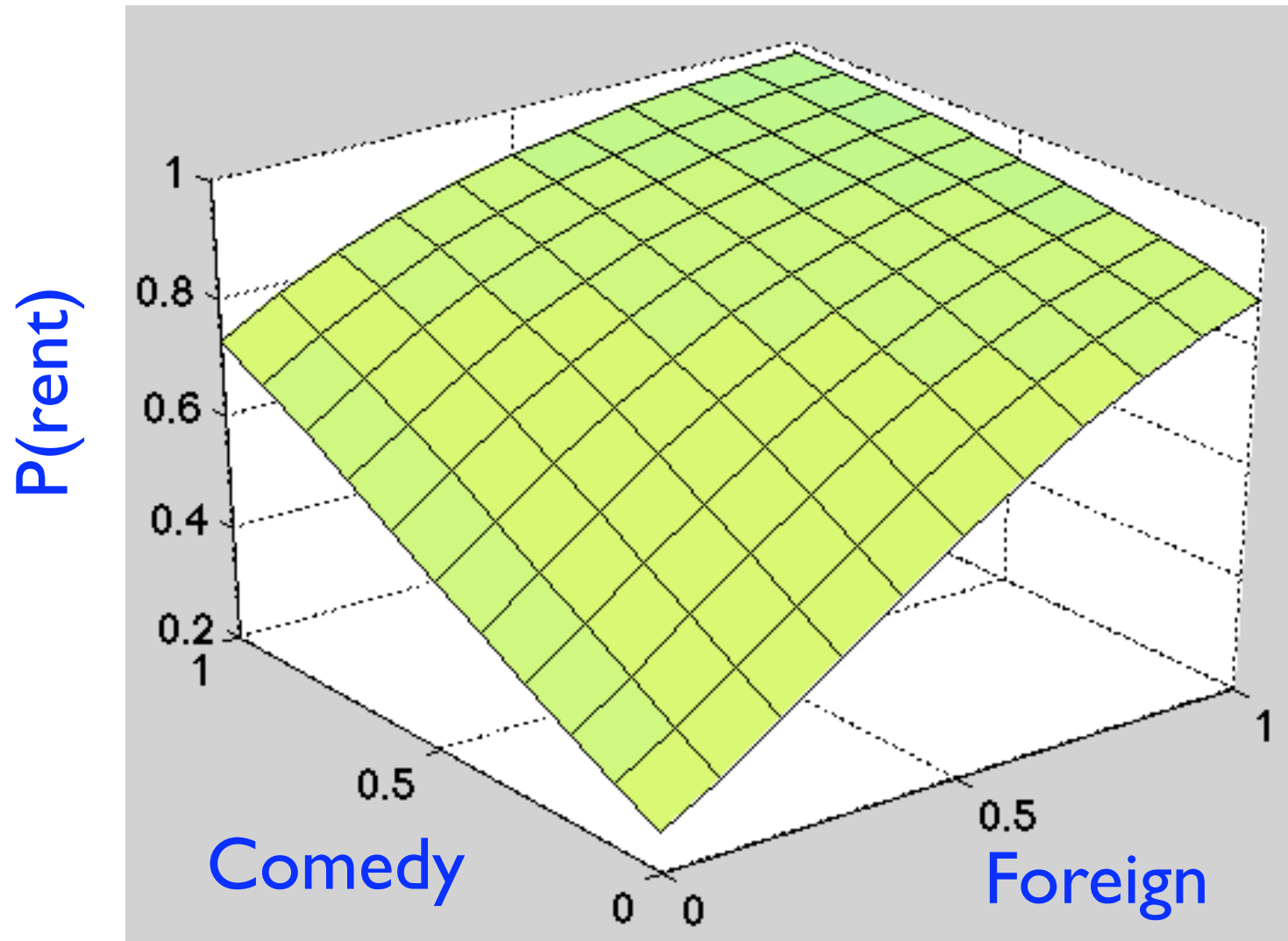
In reality

- Above, blocks are single elements of U or V
- Better: blocks are entire rows of U or V
 - ▶ take gradient, Hessian to get mean, covariance
 - ▶ formulas look a lot like linear regression (normal equations)
- And, want to fit σ^U, σ^V too
 - ▶ sample $1/\sigma^2$ from a **Gamma** (or Σ^{-1} from a **Wishart**) distribution

Nonlinearity: conjunctive features



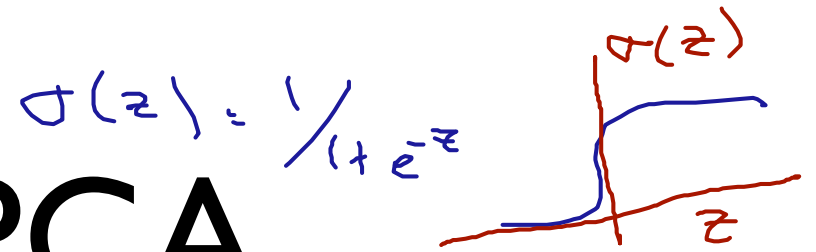
Disjunctive features



Non-Gaussian

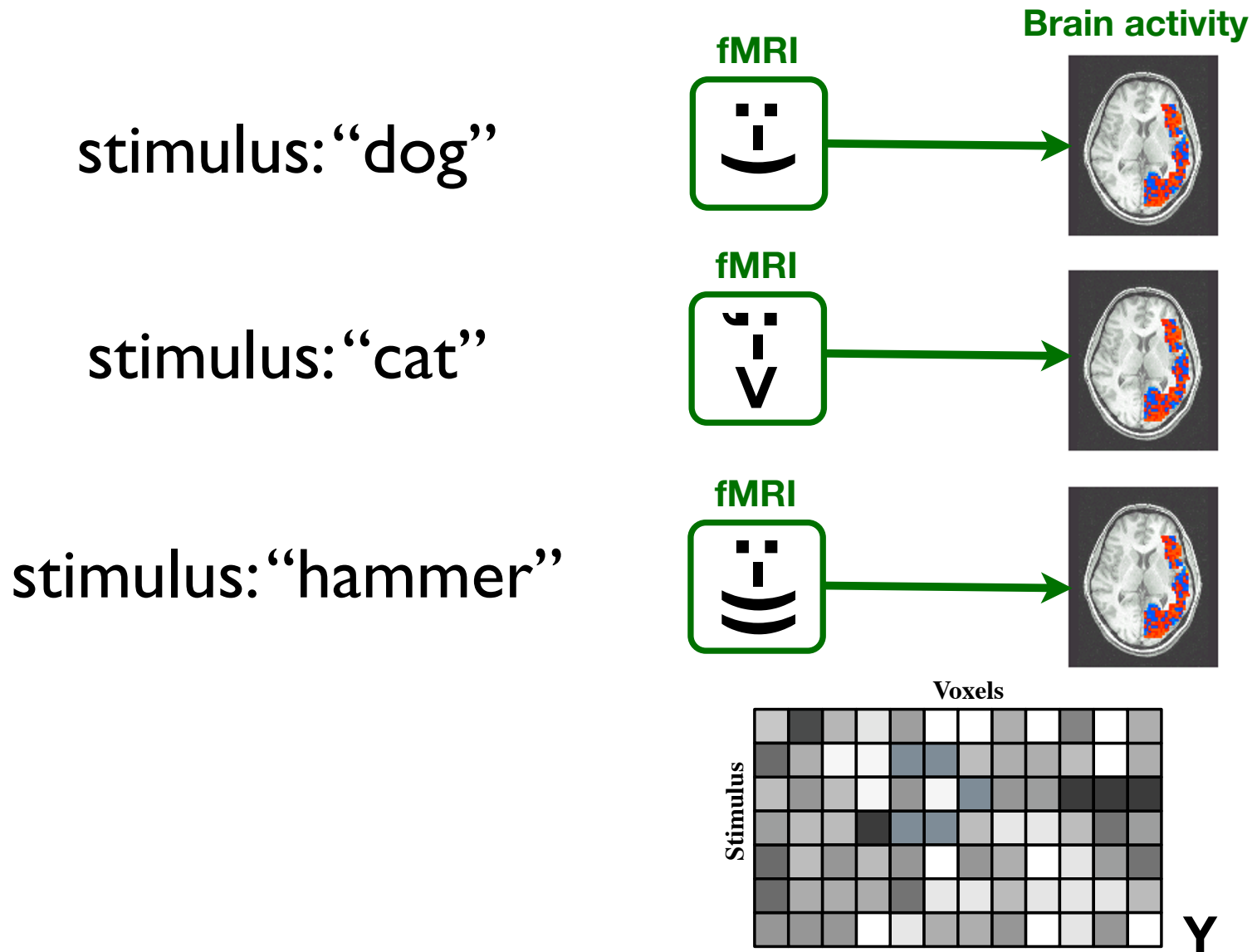
- X , U , and V could each be non-Gaussian
 - ▶ e.g., binary!
 - ▶ $\text{rents}(U, M)$, $\text{comedy}(M)$, $\text{female}(U)$
- For X : predicting -0.1 instead of 0 is only as bad as predicting $+0.1$ instead of 0
- For U, V : might infer -17% comedy or 32% female

Logistic PCA

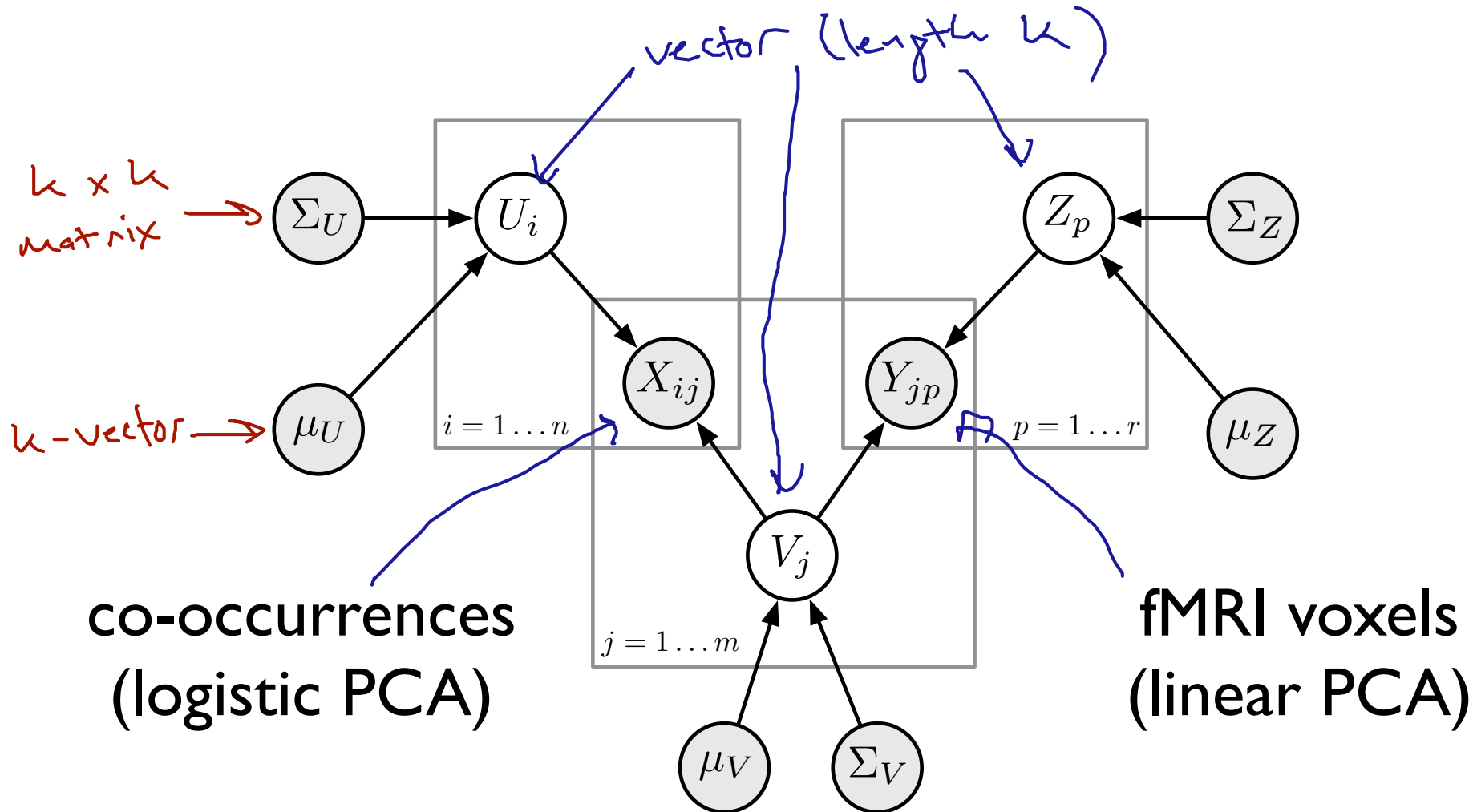


- Regular PCA: $X_{ij} \sim N(U_i \cdot V_j, \sigma^2)$
- Logistic PCA: $P(X_{ij} = 1 | u, v) = \sigma(u_i + v_j)$
- Might expect learning, inference to be hard
 - ▶ but, MH works well, using $dL/d\theta$, $d^2L/d\theta^2$
- Generalization: **exponential family PCA**
 - ▶ w/ optional hierarchy, Bayesianism

Application: fMRI



2-matrix model

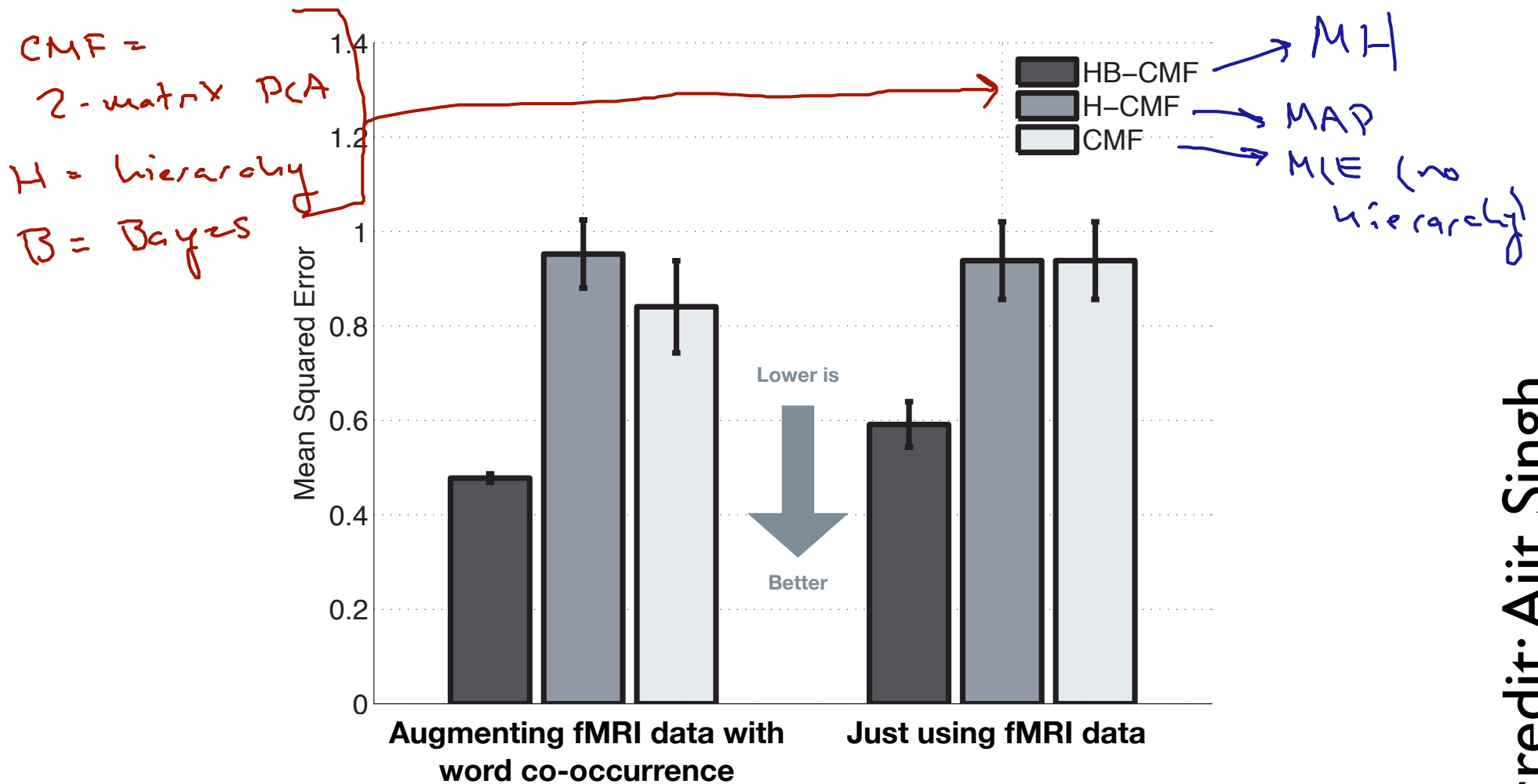


Results (logistic PCA)

linear PCA &

linear PCA side of model

Y (fMRI data): Fold-in



credit: Ajit Singh