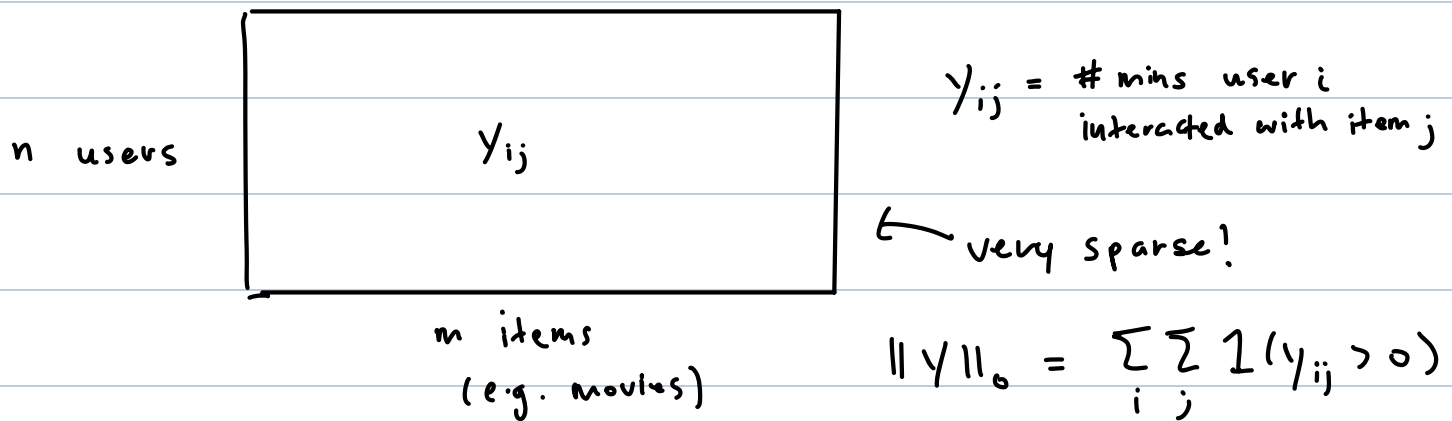


Consider the "recommender system" setting (e.g. Netflix)



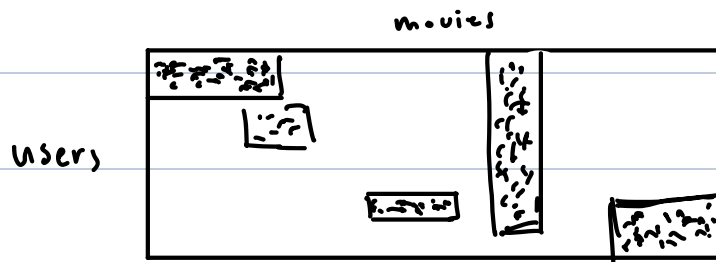
Goal: recommend items to user that they will like.

	Harry Potter movies					movies about Boston		
	HP1	...	HP8	Good Will Hunting	The Departed	The Town
user 1	1000	...	10000		0	0	0	0
user 2	2000	..	8000	..	0	0	100	700
user 3	2000	..	1000		0	0	200	300

- user 3 has seen all the Boston movies (big Boston fan)
- user 1 has seen none
- user 2 has seen many Boston movies (probably a Boston fan)
↳ recommend The Departed to user 2

In practice, we won't know all of the overlapping clusters of movies and preferences of users ahead of time. Jointly modelling them is "collaborative filtering?"

Intuition: Block structure in matrix



Model:

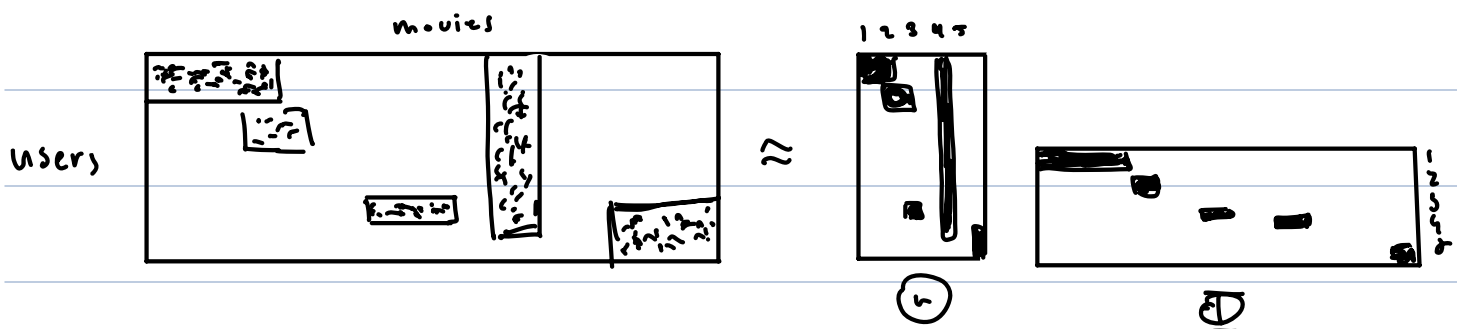
$$\mathbb{E}[y_{ij}] = \sum_k \theta_{ik} \phi_{kj}$$

$\theta_{ik} > 0$: user i 's preference for "genre" k

$$\phi_{kj} = P(j | k) \rightarrow \sum_j \phi_{kj} = 1$$

Matrix factorization

eg. $k=5$



Likelihood

$$y_{ij} \sim \text{Pois} \left(\underbrace{\sum_k \theta_{ik} \phi_{kj}}_{= \mathbb{E}[y_{ij}] = \mu_{ij}} \right)$$

MLE:

— We could try to fit using MLE

$$\hat{\theta}, \hat{\phi} \leftarrow \underset{\theta, \phi}{\operatorname{argmax}} \log \prod_i \prod_j \operatorname{Pois}(y_{ij}; \mu_{ij})$$

$$= \sum_i \sum_j \log \frac{\mu_{ij}^{y_{ij}}}{y_{ij}!} e^{-\mu_{ij}}$$

$$\propto \sum_i \sum_j y_{ij} \log \mu_{ij} - \sum_i \sum_j \mu_{ij}$$

$$= \sum_i \sum_j y_{ij} \log(\phi_i \theta_j^T \phi_j) - \sum_i \sum_k \phi_{ik}$$

only need to compute at the non-zeros.

This is the (negative) "I-divergence loss" aka "generalized KL".

Various algorithms for efficient MLE, including EM.

"Non-negative matrix factorization" (NMF) [Lee and Seung 2000]

Using the fitted model

- $\hat{\phi}_k \in \Delta_m$ for $k \in [K]$: K learned "genres" of movies
 \uparrow simplex
- $\hat{\theta}_i \in \Delta_K$: learned user preferences

e.g. recommend $\underset{j: y_{ij}=0}{\operatorname{argmax}} \sum_k \hat{\theta}_{ik} \hat{\phi}_{kj}$

Bayesian model averaging

MLE only learns one solution, but there might be many ways to represent overlapping clusters which are each good in different contexts.

Say we have samples from a posterior distribution :

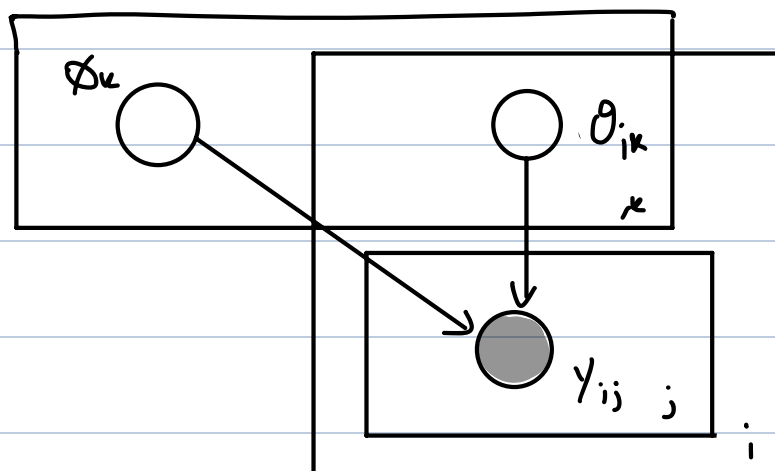
$$(\theta_i^s, \phi_1^s, \dots, \phi_k^s) \sim P(\theta_i, \phi_1, \dots, \phi_k | Y)$$

then recommend $\text{argmax}_{j: y_{ij} > 0} \frac{1}{S} \sum_s \sum_k \theta_{ik}^s \phi_{kj}^s$

Priors

$$\theta_{ik} \sim \text{Gamma}(a, b)$$

$$\phi_k \sim \text{Dir}(\alpha_1, \dots, \alpha_k)$$



Data augmentation and auxiliary variables

$$P(\theta_i | \dots) \propto \prod_{ik} \text{Gamma}(\theta_{ik} | \dots) \prod_j \text{Pois}(y_{ij} | x_i \theta_i^T \phi_j)$$

Not conjugate (even conditional on Φ)

However we can re-express the generative process by introducing auxiliary variables which will provide conditional conjugacy.

Def'n: Poisson additivity:

if: $y_k \stackrel{\text{ind.}}{\sim} \text{Pois}(\mu_k)$,

$$y_{\cdot} \triangleq \sum_k y_k$$

then:

$$y_{\cdot} \sim \text{Pois}(\mu_{\cdot})$$

$$\text{where } \mu_{\cdot} \triangleq \sum_k \mu_k$$

$$\parallel y_{ij} \sim \text{Pois}\left(\sum_k \vartheta_{ik} \varnothing_{kj}\right)$$

\Downarrow

$$\parallel y_{ijk} \stackrel{\text{ind.}}{\sim} \text{Pois}(\vartheta_{ik} \varnothing_{kj})$$

$$\parallel y_{ij} = \sum_k y_{ijk}$$

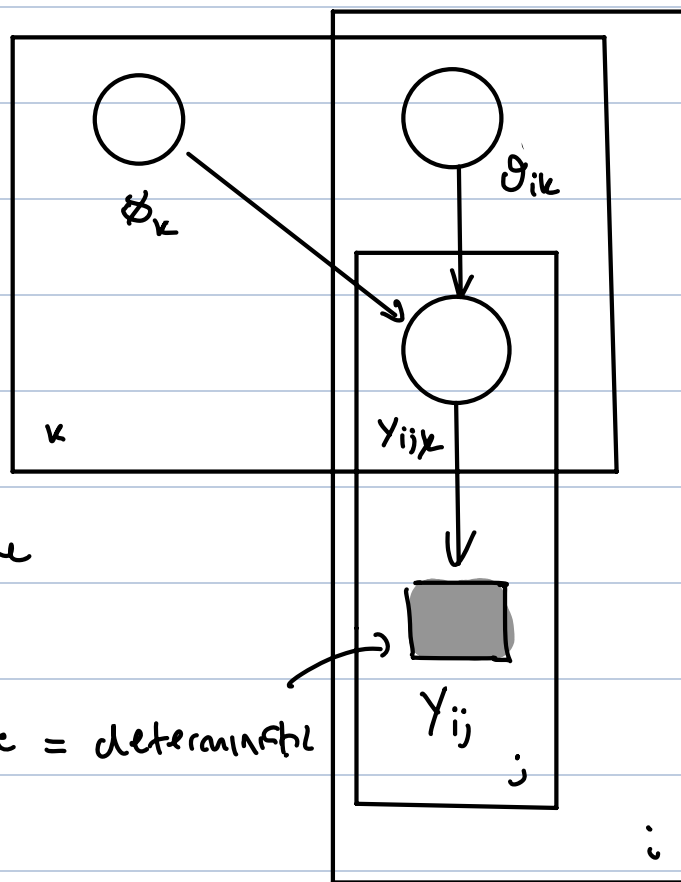
\uparrow latent sub-counts

Alternative model

with auxiliary variables



The key property is
that both this and
the original model have the
same marginal $P(y_{ij})$.



If y_{ijk} are observed, then:

$$P(\vartheta_{ik} | \dots) \propto_{\text{g}} \text{Gam}(\vartheta_{ik} | \dots) \prod_j \prod_k \text{Pois}(y_{ijk} ; \vartheta_{ik} \varnothing_{kj})$$

$$\propto \text{Gam}(\vartheta_{ik} ; a + \sum_j y_{ijk}, b + \underbrace{\sum_j \varnothing_{kj}}_{=1})$$

Conditional conjugacy after data augmentation.

Note that $P(\sigma_{ik} | \dots)$ above is the same as it would be in the following model:

$$\sigma_{ik} \sim \text{Gamma}(a, b)$$

$$y_{ik} \sim \text{Pois}(\sigma_{ik})$$

This is not a coincidence.

Yet another way to write the model:

$$\sigma_{ik} \sim \text{Gamma}(a, b)$$

$$y_{ik} \sim \text{Pois}(\sigma_{ik})$$

$$(y_{ik1} \dots y_{ikm}) \sim P(y_{ik1} \dots y_{ikm} | y_{ik0}, \sigma_{ik}, \phi_k)$$

$$y_{ij} = \sum_k y_{ikj}$$

↑ what is this dist?

Def'n : Joint distribution of Poissons and their sum.

$$\textcircled{1} \parallel \begin{array}{l} y_j \stackrel{\text{ind}}{\sim} \text{Pois}(\mu_j) \quad j=1 \dots m \\ y_{\cdot} = \sum_j y_j \end{array}$$

$$\textcircled{2} \parallel \begin{array}{l} y_{\cdot} \sim \text{Pois}(\mu_{\cdot}) \\ y_{1:m} \sim \text{Mult}(y_{\cdot}, \underbrace{\frac{\mu_1}{\mu_{\cdot}} \dots \frac{\mu_m}{\mu_{\cdot}}}_{\approx \tilde{\mu}_{1:m}}) \end{array}$$

$$P(y_{1:m}, y_{\cdot} | \mu_{1:m}) =$$

$$\underbrace{P(y_{1:m} | \mu_{1:m})}_{\textcircled{1}} \underbrace{P(y_{\cdot} | y_{1:m}, \mu_{1:m})}_{\textcircled{2}} = \underbrace{P(y_{\cdot} | \mu_{1:m})}_{\textcircled{1}} \underbrace{P(y_{1:m} | y_{\cdot}, \mu_{1:m})}_{\textcircled{2}}$$

$$\textcircled{1} \quad \underbrace{\prod_j \text{Pois}(y_j; \mu_j)}_{\textcircled{1}} \underbrace{\delta(y_{\cdot} = \sum_j y_j)}_{\textcircled{2}} = \underbrace{\text{Pois}(y_{\cdot}; \mu_{\cdot})}_{\textcircled{1}} \underbrace{\text{Mult}(y_{1:m} | y_{\cdot}, \tilde{\mu}_{1:m})}_{\textcircled{2}} \quad \textcircled{2}$$

To confirm:

$$\prod_j \frac{\mu_j^{y_j}}{y_j!} e^{-\mu_j} \delta(\dots) = \left[\frac{\mu_{\cdot}^{y_{\cdot}}}{y_{\cdot}!} e^{-\mu_{\cdot}} \right] \prod_j \frac{y_{\cdot}!}{y_j!} \frac{\mu_j^{y_j}}{\mu_{\cdot}^{y_j}} \delta(\dots)$$

$$\frac{y_{\cdot}!}{\prod_j y_j!} \prod_j \left(\frac{\mu_j}{\mu_{\cdot}} \right)^{y_j} \delta(\dots) = \frac{y_{\cdot}!}{\prod_j y_j!} \prod_j \hat{\mu}_j^{y_j} \delta(\dots)$$

Applying this definition to our model:

$$\begin{aligned} \theta_{ik} &\sim \text{Gamma}(a, b) \\ y_{ik} &\sim \text{Pois}(\theta_{ik}) \\ (y_{ik1} \dots y_{ikm}) &\sim p(y_{ik1} \dots y_{ikm} \mid y_{ik}, \theta_{ik}, \phi_k) \\ y_{ij} &= \sum_k y_{ikj} \\ &= \text{Mult}(y_{ik1} \dots y_{ikm}; y_{ik}, \dots) \end{aligned}$$

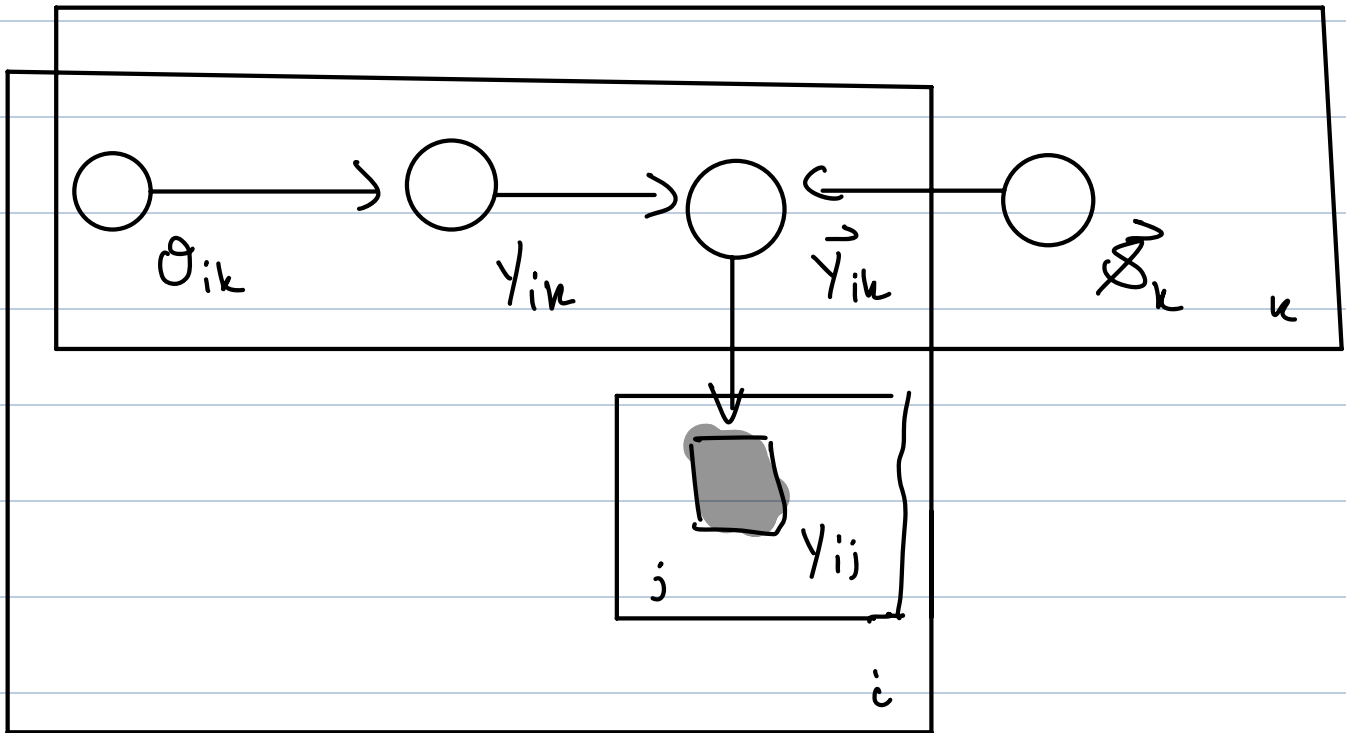
$$\left(\frac{\theta_{ik} \phi_{k1}}{\sum_j \theta_{ik} \phi_{kj}} \dots \frac{\theta_{ik} \phi_{km}}{\sum_j \theta_{ik} \phi_{kj}} \right) = (\phi_{k1} \dots \phi_{km}) \equiv \phi_k.$$

$\underbrace{\sum_j \theta_{ik} \phi_{kj}}_{=1} \quad \underbrace{\sum_j \theta_{ik} \phi_{kj}}_{=1}$

So:

$$\begin{aligned} y_{ik} &\sim \text{Pois}(\theta_{ik}) \\ y_{ik1} \dots y_{ikm} &\sim \text{Mult}(y_{ik}, \phi_k) \end{aligned}$$

The model with all auxiliary variables:



$$\theta_{ik} \sim \text{Gam}(a, b)$$

$$\vec{\phi}_k \sim \text{Dir}(\alpha_1, \dots, \alpha_m)$$

$$y_{ik} \sim \text{Pois}(\theta_{ik})$$

$$\vec{y}_{ik} \sim \text{Mult}(y_{ik}, \vec{\phi}_k)$$

$$y_{ij} = \sum_k y_{ijk}$$

conditionally
independent!

Gibbs sampler:

$$\theta_{ik}^s \sim \text{Gamma}(a + \sum_j y_{ijk}, b + 1)$$

$$\phi_k^s \sim \text{Dir}(\alpha_1 + \sum_i y_{i1k}, \dots, \alpha_m + \sum_i y_{imk})$$

$$(y_{ij1}^s \dots y_{ijn}^s) \sim \text{Mult}(y_{ij}, \frac{\theta_{i1}^s \phi_{1j}^s}{\theta_i^T \phi_j}, \dots, \frac{\theta_{in}^s \phi_{nj}^s}{\theta_i^T \phi_j})$$

Auxiliary variable MCMC

At a high level, all MCMC methods generate samples of latent variables:

$$z^s \sim p(z | y)$$

Such that:

$$\lim_{S \rightarrow \infty} \frac{1}{S} \sum_s \mathbb{1}(f(z_s) \in B) \rightarrow P(f(z) \in B | y)$$

Auxiliary variables are variables A such that

$$p(z, y, A) = \int p(z, y, A=c) da$$

Adding such variables into MCMC is always valid:

$$\text{Define: } f(z, A) = z$$

$$\frac{1}{S} \sum_s \mathbb{1}(f(z^s, A^s) \in B) \rightarrow P(z \in B | y).$$

e.g. Gibbs:

$$\parallel p(z_i | z_{-i}, y) \quad \text{intractable}$$

$$\parallel \begin{array}{l} p(z_i | z_{-i}, y, A) \quad \text{tractable} \\ p(A | z, y) \quad \text{tractable} \end{array}$$