

# Probabilistic Graphical Models

10-708

More on learning fully observed BNs, exponential families, and generalized linear models

Eric Xing

Lecture 10, Oct 12, 2005

Reading: MJ-Chap. 7,8



## Exponential family

- For a numeric random variable  $\mathcal{X}$

$$\begin{aligned} p(\mathcal{X} | \eta) &= h(\mathcal{X}) \exp\{\eta^T T(\mathcal{X}) - A(\eta)\} \\ &= \frac{1}{Z(\eta)} h(\mathcal{X}) \exp\{\eta^T T(\mathcal{X})\} \end{aligned}$$

is an exponential family distribution with natural (canonical) parameter  $\eta$

- Function  $T(\mathcal{X})$  is a *sufficient statistic*.
- Function  $A(\eta) = \log Z(\eta)$  is the log normalizer.
- Examples: Bernoulli, multinomial, Gaussian, Poisson, gamma,...
- A distribution  $p(\mathcal{X})$  has finite sufficient statistics (independent of number of data cases) iff it is in the exponential family.



# Multivariate Gaussian Distribution



- For a continuous vector random variable  $\mathbf{X} \in \mathbb{R}^k$ :

$$p(\mathbf{x}|\mu, \Sigma) = \frac{1}{(2\pi)^{k/2} |\Sigma|^{1/2}} \exp\left\{-\frac{1}{2}(\mathbf{x} - \mu)^T \Sigma^{-1} (\mathbf{x} - \mu)\right\}$$

$$= \frac{1}{(2\pi)^{k/2}} \exp\left\{-\frac{1}{2} \text{tr}(\Sigma^{-1} \mathbf{x} \mathbf{x}^T) + \mu^T \Sigma^{-1} \mathbf{x} - \frac{1}{2} \mu^T \Sigma^{-1} \mu - \log|\Sigma|\right\}$$

Moment parameter

- Exponential family representation

$$\eta = \left[\Sigma^{-1} \mu; -\frac{1}{2} \text{vec}(\Sigma^{-1})\right] = [\eta_1, \text{vec}(\eta_2)], \quad \eta_1 = \Sigma^{-1} \mu \text{ and } \eta_2 = -\frac{1}{2} \Sigma^{-1}$$

Natural parameter

$$T(\mathbf{x}) = [\mathbf{x}; \text{vec}(\mathbf{x} \mathbf{x}^T)]$$

$$A(\eta) = \frac{1}{2} \mu^T \Sigma^{-1} \mu + \log|\Sigma| = -\frac{1}{2} \text{tr}(\eta_2 \eta_1 \eta_1^T) - \frac{1}{2} \log(-2\eta_2)$$

$$h(\mathbf{x}) = (2\pi)^{-k/2}$$

- Note: a  $k$ -dimensional Gaussian is a  $(\mathcal{A}, \mathcal{A}^{\mathcal{L}})$ -parameter distribution with a  $(\mathcal{A}, \mathcal{A}^{\mathcal{L}})$ -element vector of sufficient statistics (but because of symmetry and positivity, parameters are constrained and have lower degree of freedom)

# Multinomial distribution



- For a binary vector random variable  $\mathbf{X} \sim \text{multi}(\mathbf{X} | \boldsymbol{\pi})$ ,

$$p(\mathbf{x}|\boldsymbol{\pi}) = \pi_1^{x_1} \pi_2^{x_2} \dots \pi_K^{x_K} = \exp\left\{\sum_K x^k \ln \pi_k\right\}$$

$$= \exp\left\{\sum_{k=1}^{K-1} x^k \ln \pi_k + \left(1 - \sum_{k=1}^{K-1} x^k\right) \ln\left(1 - \sum_{k=1}^{K-1} \pi_k\right)\right\}$$

$$= \exp\left\{\sum_{k=1}^{K-1} x^k \ln\left(\frac{\pi_k}{1 - \sum_{k=1}^{K-1} \pi_k}\right) + \ln\left(1 - \sum_{k=1}^{K-1} \pi_k\right)\right\}$$

- Exponential family representation

$$\eta = \left[\ln\left(\frac{\pi_k}{\pi_K}\right); \mathbf{0}\right]$$

$$T(\mathbf{x}) = [\mathbf{x}]$$

$$A(\eta) = -\ln\left(1 - \sum_{k=1}^{K-1} \pi_k\right) = \ln\left(\sum_{k=1}^K e^{\eta_k}\right)$$

$$h(\mathbf{x}) = 1$$



## Why exponential family?

- Moment generating property

$$\begin{aligned}\frac{dA}{d\eta} &= \frac{d}{d\eta} \log Z(\eta) = \frac{1}{Z(\eta)} \frac{d}{d\eta} Z(\eta) \\ &= \frac{1}{Z(\eta)} \frac{d}{d\eta} \int h(x) \exp\{\eta^T T(x)\} dx \\ &= \int T(x) \frac{h(x) \exp\{\eta^T T(x)\}}{Z(\eta)} dx \\ &= E[T(x)]\end{aligned}$$

$$\begin{aligned}\frac{d^2 A}{d\eta^2} &= \int T^2(x) \frac{h(x) \exp\{\eta^T T(x)\}}{Z(\eta)} dx - \left[ \int T(x) \frac{h(x) \exp\{\eta^T T(x)\}}{Z(\eta)} dx \right]^2 \\ &= E[T^2(x)] - E^2[T(x)] \\ &= \text{Var}[T(x)]\end{aligned}$$



## Moment estimation

- We can easily compute moments of any exponential family distribution by taking the derivatives of the log normalizer  $A(\eta)$ .
- The  $q^{\text{th}}$  derivative gives the  $q^{\text{th}}$  centered moment.

$$\frac{dA(\eta)}{d\eta} = \text{mean}$$

$$\frac{d^2 A(\eta)}{d\eta^2} = \text{variance}$$

...

- When the sufficient statistic is a stacked vector, partial derivatives need to be considered.

## Moment vs canonical parameters

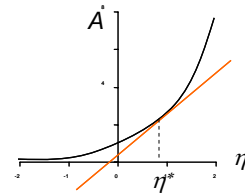


- The moment parameter  $\mu$  can be derived from the natural (canonical) parameter

$$\frac{dA(\eta)}{d\eta} = E[T(\mathcal{X})] \stackrel{\text{def}}{=} \mu$$

- $A(\eta)$  is convex since

$$\frac{d^2 A(\eta)}{d\eta^2} = \text{Var}[T(\mathcal{X})] > 0$$



- Hence we can invert the relationship and infer the canonical parameter from the moment parameter (1-to-1):

$$\eta = \psi(\mu) \stackrel{\text{def}}{=}$$

- A distribution in the exponential family can be parameterized not only by  $\eta$  – the canonical parameterization, but also by  $\mu$  – the moment parameterization.

## MLE for Exponential Family



- For *iid* data, the log-likelihood is

$$\begin{aligned} \ell(\eta; D) &= \log \prod_n h(x_n) \exp\{\eta^T T(x_n) - A(\eta)\} \\ &= \sum_n \log h(x_n) + \left( \eta^T \sum_n T(x_n) \right) - NA(\eta) \end{aligned}$$

- Take derivatives and set to zero:

$$\begin{aligned} \frac{\partial \ell}{\partial \eta} &= \sum_n T(x_n) - N \frac{\partial A(\eta)}{\partial \eta} = 0 \\ \frac{\partial A(\eta)}{\partial \eta} &= \frac{1}{N} \sum_n T(x_n) \\ \Rightarrow \hat{\mu}_{MLE} &= \frac{1}{N} \sum_n T(x_n) \end{aligned}$$

- This amounts to **moment matching**.
- We can infer the canonical parameters using  $\hat{\eta}_{MLE} = \psi(\hat{\mu}_{MLE})$



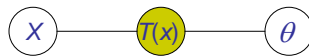
## Sufficiency

- For  $p(x|\theta)$ ,  $T(x)$  is **sufficient** for  $\theta$  if there is no information in  $X$  regarding  $\theta$  beyond that in  $T(x)$ .
  - We can throw away  $X$  for the purpose of inference w.r.t.  $\theta$ .

- Bayesian view  $X \rightarrow T(x) \rightarrow \theta$   $p(\theta|T(x), x) = p(\theta|T(x))$

- Frequentist view  $\theta \rightarrow T(x) \rightarrow X$   $p(x|T(x), \theta) = p(x|T(x))$

- The Neyman factorization theorem



- $T(x)$  is **sufficient** for  $\theta$  if

$$p(x, T(x), \theta) = \psi_1(T(x), \theta) \psi_2(x, T(x))$$

$$\Rightarrow p(x|\theta) = g(T(x), \theta) h(x, T(x))$$



## Examples

- Gaussian:

$$\eta = \left[ \Sigma^{-1} \mu; -\frac{1}{2} \text{vec}(\Sigma^{-1}) \right]$$

$$T(x) = \left[ x; \text{vec}(xx^T) \right]$$

$$A(\eta) = \frac{1}{2} \mu^T \Sigma^{-1} \mu + \frac{1}{2} \log |\Sigma|$$

$$h(x) = (2\pi)^{-k/2}$$

$$\Rightarrow \mu_{MLE} = \frac{1}{N} \sum_n T_1(x_n) = \frac{1}{N} \sum_n x_n$$

- Multinomial:

$$\eta = \left[ \ln \left( \frac{\pi_k}{\pi_K} \right); 0 \right]$$

$$T(x) = [x]$$

$$A(\eta) = -\ln \left( 1 - \sum_{k=1}^{K-1} \pi_k \right) = \ln \left( \sum_{k=1}^K e^{\eta_k} \right)$$

$$h(x) = 1$$

$$\Rightarrow \mu_{MLE} = \frac{1}{N} \sum_n x_n$$

- Poisson:

$$\eta = \log \lambda$$

$$T(x) = x$$

$$A(\eta) = \lambda = e^\eta$$

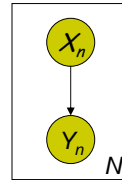
$$h(x) = \frac{1}{x!}$$

$$\Rightarrow \mu_{MLE} = \frac{1}{N} \sum_n x_n$$

# Generalized Linear Models (GLIMs)

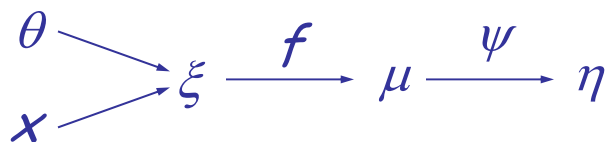


- The graphical model
  - Linear regression
  - Discriminative linear classification
  - Commonality:
    - model  $E(Y) = \mu = f(\theta^T X)$ 
      - What is  $p(\cdot)$ , the cond. dist. Of  $Y$ ?
      - What is  $f(\cdot)$ , the response function?



- GLIM
  - The observed input  $x$  is assumed to enter into the model via a linear combination of its elements  $\xi = \theta^T x$
  - The conditional mean  $\mu$  is represented as a function  $f(\xi)$  of  $\xi$ , where  $f$  is known as the response function
  - The observed output  $y$  is assumed to be characterized by an exponential family distribution with conditional mean  $\mu$ .

## GLIM, cont.



$$p(y | \eta) = h(x) \exp\{\eta^T y - A(\eta)\}$$

$$\Rightarrow p(y | \eta) = h(x) \exp\left\{\frac{1}{\phi} (\eta^T y - A(\eta))\right\}$$

- The choice of exp family is constrained by the nature of the data  $Y$ 
  - Example:  $y$  is a continuous vector  $\rightarrow$  multivariate Gaussian
  - $y$  is a class label  $\rightarrow$  Bernoulli or multinomial
- The choice of the response function
  - Following some mild constrains, e.g.,  $[0,1]$ . Positivity ...
  - Canonical response function:  $f = \psi^{-1}(\cdot)$ 
    - In this case  $\theta^T x$  directly corresponds to canonical parameter  $\eta$ .

## MLE for GLIMs with natural response



- Log-likelihood

$$\ell = \sum_n \log h(y_n) + \sum_n (\theta^T x_n y_n - A(\eta_n))$$

- Derivative of Log-likelihood

$$\begin{aligned} \frac{d\ell}{d\theta} &= \sum_n \left( x_n y_n - \frac{dA(\eta_n)}{d\eta_n} \frac{d\eta_n}{d\theta} \right) \\ &= \sum_n (y_n - \mu_n) x_n \\ &= X^T (y - \mu) \end{aligned}$$

This is a fixed point function because  $\mu$  is a function of  $\theta$

- Online learning for canonical GLIMs

- Stochastic gradient ascent = least mean squares (LMS) algorithm:

$$\begin{aligned} \theta^{t+1} &= \theta^t + \rho (y_n - \mu_n^t) x_n \\ \text{where } \mu_n^t &= (\theta^t)^T x_n \text{ and } \rho \text{ is a step size} \end{aligned}$$

## Batch learning for canonical GLIMs



- The Hessian matrix

$$\begin{aligned} H &= \frac{d^2 \ell}{d\theta d\theta^T} = \frac{d}{d\theta^T} \sum_n (y_n - \mu_n) x_n = \sum_n x_n \frac{d\mu_n}{d\theta^T} \\ &= - \sum_n x_n \frac{d\mu_n}{d\eta_n} \frac{d\eta_n}{d\theta^T} \\ &= - \sum_n x_n \frac{d\mu_n}{d\eta_n} x_n^T \text{ since } \eta_n = \theta^T x_n \\ &= -X^T W X \end{aligned}$$

where  $X = [x_n^T]$  is the design matrix and

$$W = \text{diag} \left( \frac{d\mu_1}{d\eta_1}, \dots, \frac{d\mu_N}{d\eta_N} \right)$$

which can be computed by calculating the 2<sup>nd</sup> derivative of  $A(\eta_n)$

# Iteratively Reweighted Least Squares (IRLS)



- Recall Newton-Raphson methods with cost function  $\mathcal{J}$

$$\theta^{t+1} = \theta^t - H^{-1} \nabla_{\theta} \mathcal{J}$$

- We now have

$$\nabla_{\theta} \mathcal{J} = X^T (y - \mu)$$

$$H = -X^T W X$$

- Now:

$$\theta^{t+1} = \theta^t + H^{-1} \nabla_{\theta} \mathcal{J}$$

$$= (X^T W^t X)^{-1} [X^T W^t X \theta^t + X^T (y - \mu^t)]$$

- 

$$= (X^T W^t X)^{-1} X^T W^t z^t$$

where the adjusted response is  $z^t = X \theta^t + (W^t)^{-1} (y - \mu^t)$

- This can be understood as solving the following "Iteratively reweighted least squares" problem

$$\theta^{t+1} = \arg \min_{\theta} (z - X\theta)^T W (z - X\theta)$$

# Example 1: logistic regression (sigmoid classifier)

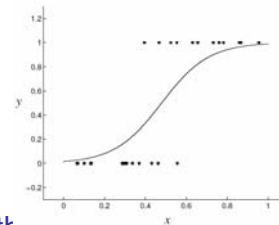


- The condition distribution: a Bernoulli

$$p(y | x) = \mu(x)^y (1 - \mu(x))^{1-y}$$

where  $\mu$  is a logistic function

$$\mu(x) = \frac{1}{1 + e^{-\eta(x)}}$$



- $p(y|x)$  is an exponential family function, with

- mean:  $E[y | x] = \mu = \frac{1}{1 + e^{-\eta(x)}}$

- and canonical response function  $\eta = \xi = \theta^T x$

- IRLS

$$\frac{d\mu}{d\eta} = \mu(1 - \mu)$$

$$W = \begin{pmatrix} \mu_1(1 - \mu_1) & & \\ & \ddots & \\ & & \mu_N(1 - \mu_N) \end{pmatrix}$$



# Logistic regression: practical issues



- It is very common to use *regularized* maximum likelihood.

$$p(y = \pm 1 | x, \theta) = \frac{1}{1 + e^{-y\theta^T x}} = \sigma(y\theta^T x)$$

$$p(\theta) \sim \text{Normal}(0, \lambda^{-1} I)$$

$$l(\theta) = \sum_n \log(\sigma(y_n \theta^T x_n)) - \frac{\lambda}{2} \theta^T \theta$$

- IRLS takes  $\mathcal{O}(Nd)$  per iteration, where  $N$  = number of training cases and  $d$  = dimension of input  $x$ .
- Quasi-Newton methods, that approximate the Hessian, work faster.
- Conjugate gradient takes  $\mathcal{O}(Nd)$  per iteration, and usually works best in practice.
- Stochastic gradient descent can also be used if  $N$  is large c.f. perceptron rule:

$$\nabla_{\theta} \ell = (1 - \sigma(y_n \theta^T x_n)) y_n x_n - \lambda \theta$$

# Example 2: linear regression



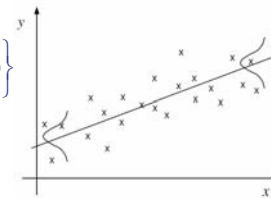
- The condition distribution: a Gaussian

$$p(y | x, \theta, \Sigma) = \frac{1}{(2\pi)^{k/2} |\Sigma|^{1/2}} \exp\left\{-\frac{1}{2} (y - \mu(x))^T \Sigma^{-1} (y - \mu(x))\right\}$$

**Recall**  $\Rightarrow h(x) \exp\left\{-\frac{1}{2} \Sigma^{-1} (\eta^T(x) y - \mathcal{A}(\eta))\right\}$

where  $\mu$  is a linear function

$$\mu(x) = \theta^T x = \eta(x)$$



- $p(y | x)$  is an exponential family function, with

- mean:  $E[y | x] = \mu = \theta^T x$

- and canonical response function  $\eta_1 = \xi = \theta^T x$

- IRLS  $\frac{d\mu}{d\eta} = 1 \Rightarrow \begin{aligned} \theta^{t+1} &= (X^T W^t X)^{-1} X^T W^t z^t \\ &= (X^T X)^{-1} X^T (X \theta^t + (y - \mu^t)) \\ &= \theta^t + (X^T X)^{-1} X^T (y - \mu^t) \end{aligned} \xrightarrow{t \rightarrow \infty} \theta = (X^T X)^{-1} X^T y$

Steepest descent

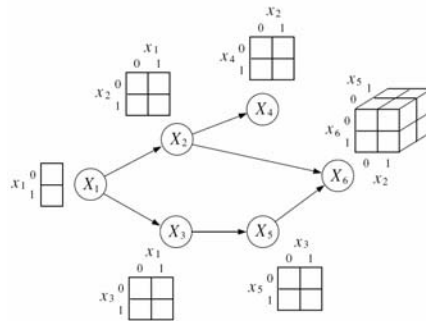
Normal equation

## MLE for general BNs



- If we assume the parameters for each CPD are globally independent, and all nodes are fully observed, then the log-likelihood function decomposes into a sum of local terms, one per node:

$$\ell(\theta; D) = \log p(D | \theta) = \log \prod_n \left( \prod_i p(x_{n,i} | \mathbf{x}_{\pi_i}, \theta_i) \right) = \sum_i \left( \sum_n \log p(x_{n,i} | \mathbf{x}_{\pi_i}, \theta_i) \right)$$



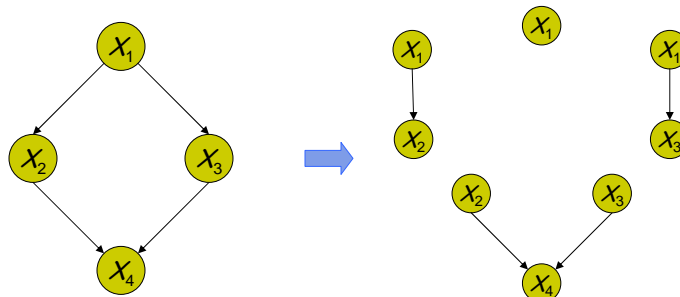
## Example: A directed model



- Consider the distribution defined by the directed acyclic GM:

$$p(x | \theta) = p(x_1 | \theta_1) p(x_2 | x_1, \theta_2) p(x_3 | x_1, \theta_3) p(x_4 | x_2, x_3, \theta_4)$$

- This is exactly like learning four separate small BNs, each of which consists of a node and its parents.



## MLE for BNs with tabular CPDs



- Assume each CPD is represented as a table (multinomial) where

$$\theta_{ijk} \stackrel{\text{def}}{=} p(X_i = j \mid X_{\pi_i} = k)$$

- Note that in case of multiple parents,  $X_{\pi_i}$  will have a composite state, and the CPD will be a high-dimensional table
- The sufficient statistics are counts of family configurations



$$n_{ijk} \stackrel{\text{def}}{=} \sum_n x_{n,i}^j x_{n,\pi_i}^k$$

- The log-likelihood is

$$\ell(\theta; D) = \log \prod_{i,j,k} \theta_{ijk}^{n_{ijk}} = \sum_{i,j,k} n_{ijk} \log \theta_{ijk}$$

- Using a Lagrange multiplier to enforce so  $\sum_j \theta_{ijk} = 1$  we get

$$\theta_{ijk}^{ML} = \frac{n_{ijk}}{\sum_{i,j,k} n_{i,j,k}}$$

## MLE and Kulback-Leibler divergence



- KL divergence

$$D(q(x) \parallel p(x)) = \sum_x q(x) \log \frac{q(x)}{p(x)}$$

- Empirical distribution

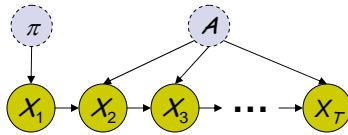
$$\tilde{p}(x) \stackrel{\text{def}}{=} \frac{1}{N} \sum_{n=1}^N \delta(x, x_n)$$

- Where  $\delta(x, x_n)$  is a Kronecker delta function

- $\text{Max}_{\theta}(\text{MLE}) \equiv \text{Min}_{\theta}(\text{KL})$

$$\begin{aligned} D(\tilde{p}(x) \parallel p(x|\theta)) &= \sum_x \tilde{p}(x) \log \frac{\tilde{p}(x)}{p(x|\theta)} \\ &= \sum_x \tilde{p}(x) \log \tilde{p}(x) - \sum_x \tilde{p}(x) \log p(x|\theta) \\ &= \sum_x \tilde{p}(x) \log \tilde{p}(x) - \frac{1}{N} \sum_n \log p(x_n|\theta) \\ &= C + \frac{1}{N} \ell(\theta; D) \end{aligned}$$

## Parameter sharing



- Consider a time-invariant (stationary) 1<sup>st</sup>-order Markov model
  - Initial state probability vector:  $\pi_k = p(X_1^k = 1)$
  - State transition probability matrix:  $A_{ij} = p(X_t^j = 1 | X_{t-1}^i = 1)$
- The joint:  $p(X_{1:T} | \theta) = p(x_1 | \pi) \prod_{t=2}^T p(x_t | x_{t-1})$
- The log-likelihood:  $\ell(\theta; \mathcal{D}) = \sum_n \log p(x_{n,1} | \pi) + \sum_n \sum_{t=2}^T \log p(x_{n,t} | x_{n,t-1}, A)$
- Again, we optimize each parameter separately
  - $\pi$  is a multinomial frequency vector, and we've seen it before
  - What about  $A$ ?

## Learning a Markov chain transition matrix



- $A$  is a stochastic matrix:  $\sum_j A_{ij} = 1$
- Each row of  $A$  is multinomial distribution.
- So **MLE** of  $A_{ij}$  is the fraction of transitions from  $i$  to  $j$

$$A_{ij}^{ML} = \frac{\#(i \rightarrow j)}{\#(i \rightarrow \bullet)} = \frac{\sum_n \sum_{t=2}^T x_{n,t-1}^i x_{n,t}^j}{\sum_n \sum_{t=2}^T x_{n,t-1}^i}$$

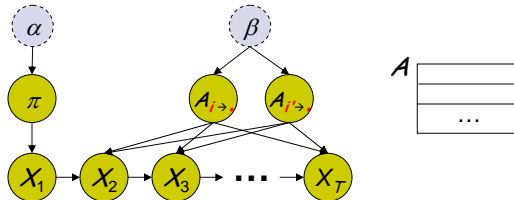
- Application:
  - if the states  $X_t$  represent words, this is called a *bigram language model*
- Sparse data problem:
  - If  $i \rightarrow j$  did not occur in data, we will have  $A_{ij} = 0$ , then any further sequence with word pair  $i \rightarrow j$  will have zero probability.
  - A standard hack: *backoff smoothing* or *deleted interpolation*

$$\tilde{A}_{i \rightarrow \bullet} = \lambda \eta_i + (1 - \lambda) A_{i \rightarrow \bullet}^{ML}$$

# Bayesian language model



- Global and local parameter independence



- The posterior of  $A_{i \rightarrow \cdot}$  and  $A_{i' \rightarrow \cdot}$  is factorized despite v-structure on  $X_p$  because  $X_{p-1}$  acts like a **multiplexer**
- Assign a Dirichlet prior  $\beta_i$  to each row of the transition matrix:

$$A_{ij}^{Bayes} \stackrel{\text{def}}{=} p(j | i, D, \beta_i) = \frac{\#(i \rightarrow j) + \beta_{i,k}}{\#(i \rightarrow \bullet) + |\beta_i|} = \lambda_i \beta_{i,k} + (1 - \lambda_i) A_{ij}^{ML}, \text{ where } \lambda_i = \frac{|\beta_i|}{|\beta_i| + \#(i \rightarrow \bullet)}$$

- We could consider more realistic priors, e.g., mixtures of Dirichlets to account for types of words (adjectives, verbs, etc.)