

# Beyond EM: Bayesian Techniques for HLT

Hal Daumé III

me@hal3.name

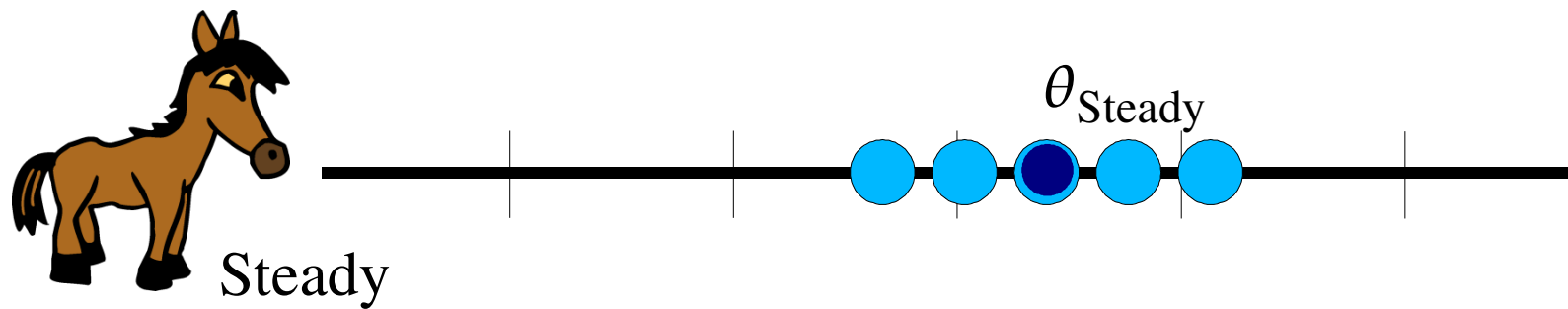
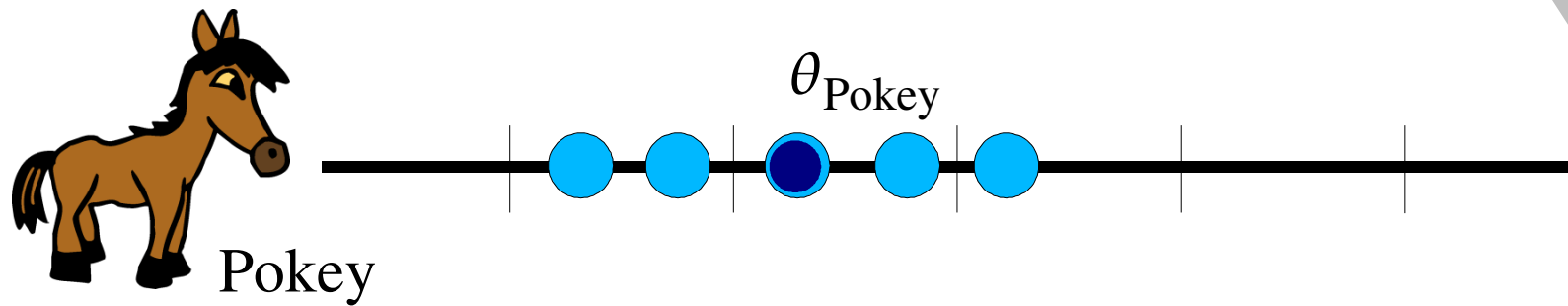
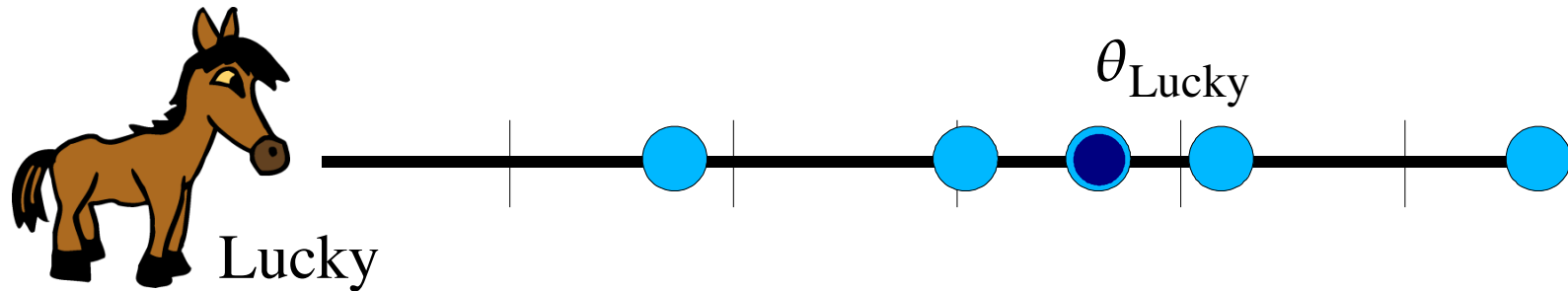


<http://bayes.hal3.name/>

Acknowledgments: David Blei, Yee-Whye Teh, Aaron D'Souza



# Horse Racing



# Who Should Be Here?

**“My EM converges to garbage!”**

**“I want to integrate domain knowledge.”**

**“My independence assumptions  
don't factor nicely!”**

**“Bayesian techniques are  
worthless...  
too hard...  
too slow...”**

# Tutorial Goals

**Understand when to be Bayesian**

**Know the natural prior distributions**

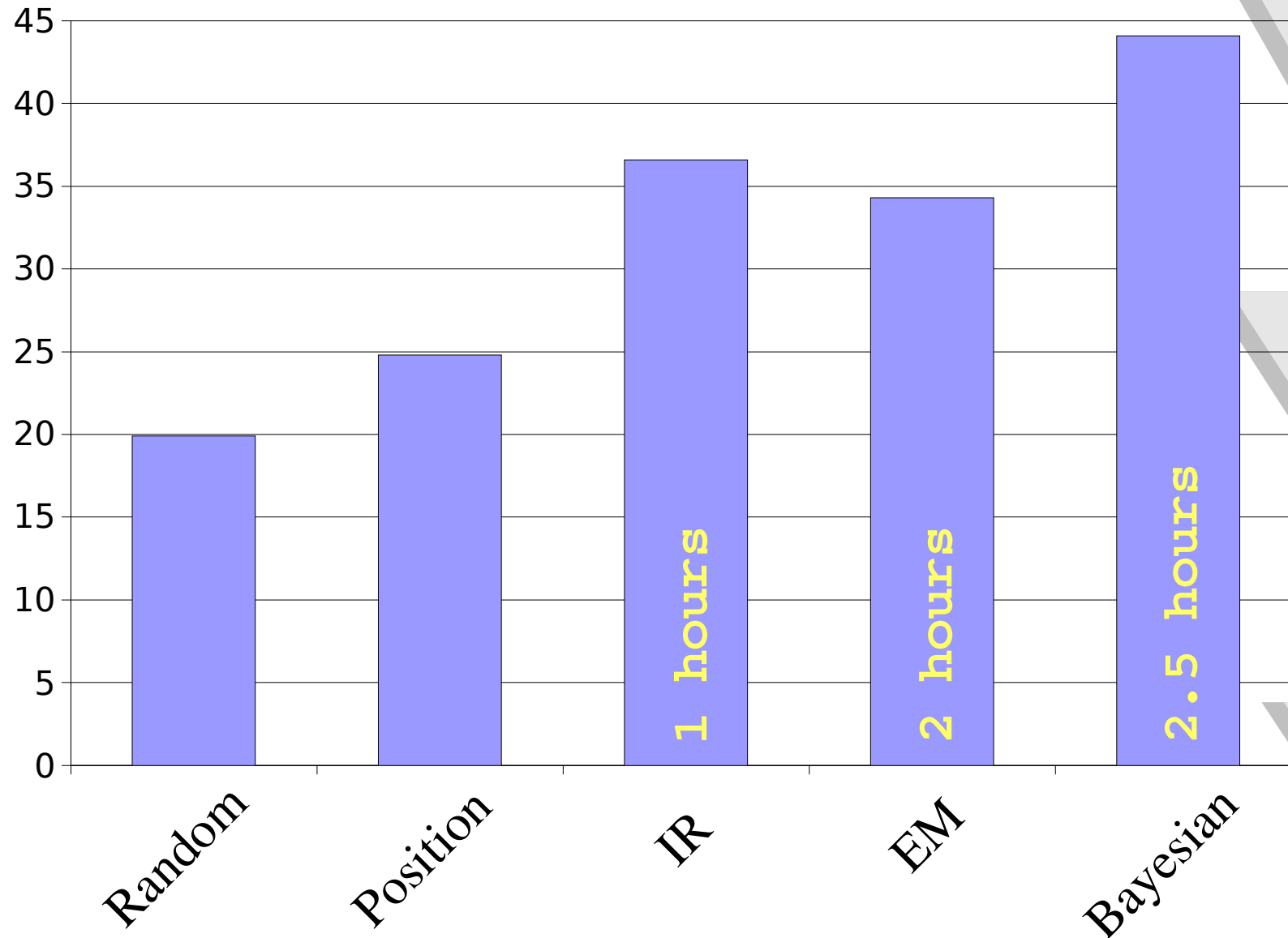
**Draw complex graphical models**

**Implement a Gibbs sampler for LDA**

**Read NIPS/UAI/etc. papers**

# Empirical Motivation

## Mean Average Precision



# Model for Q-F Summarization

- Suppose a document D is relevant to two queries, Q1 and Q2
  - Mark each sentence with the degree to which it is about:
    - Q1
    - Q2
    - D, but not Q1 nor Q2
    - General English
  - Now, mark each word in that sentence with an absolute judgment about where it came from
    - Sentences which are more like Q1 are more likely to have words from Q1
    - General English words are likely to be consistent across the whole corpus
    - Document-specific words are likely to be consistent across the whole document
    - Query-specific words are likely to be consistent across all documents relevant to a given query

Iraq's National Assembly approved a list of Cabinet members for a transitional government Thursday, three months after national elections.	(0.5, 0.2, 0.2, 0.1)
Three ministries – Defense, Oil and Electricity – were filled with temporary appointments because of a last minute failure to reach a compromise.	(0.1, 0.6, 0.1, 0.1)
Prime minister Ibrahim al-Jaafari assumed his post with the creation of his government	(0.2, 0.2, 0.3, 0.3)
The approval of the Cabinet represents the end of a major political impasse in the country.	(0.4, 0.4, 0.1, 0.1)
On Wednesday, al-Jaafari told a news conference that he had submitted his proposal Cabinet to President Jalal Talabani, who had to approve the names before the transitional National Assembly voted on them.	(0.1, 0.2, 0.5, 0.1)
Al-Jaafari's announcement came a short time after gunment shot and killed an assembly member on her doorstep in Baghdad...	(0.2, 0.4, 0.1, 0.3)

# Tutorial Outline

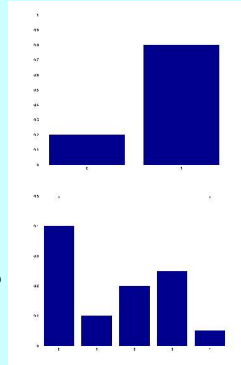
- Introduction to the Bayesian Paradigm
- Background Material
  - Graphical Models
  - Maximum Likelihood
  - Expectation Maximization
- Priors, priors, priors (subjective, conjugate, reference, etc.)
- Inference Problem and Solutions
  - Summing
  - Monte Carlo
  - Markov Chain Monte Carlo
  - Laplace Approximation
  - Variational Approximation
  - Message Passing...
- Survey of Popular Models
- Pointers to Literature
- Conclusions

# A Brief Refresher

## Distributions

Binomial

Binary



$$\text{Bin}(x | N, \theta) \propto \theta^n (1 - \theta)^{N-n}$$

Multinomial

K classes

$$\text{Mult}(\bar{x} | \bar{\theta}) \propto \prod \theta_k^{x_k}$$

## Expectations:

$$E_{x \sim p}[f(x)] = \begin{cases} \sum_{x \in X} p(x) f(x) & X \text{ is discrete} \\ \int_X dx p(x) f(x) & X \text{ is continuous} \end{cases}$$

## Probability Calculus:

$$p(x_{1:N}) = \prod_n p(x_n | x_{1:n-1}) \quad p(a | b) = \frac{p(a) p(b | a)}{p(b)}$$



# The Bayesian Paradigm

- Every statistical problem has *data* and *parameters*
- Find a probability *distribution* of the *parameters* given the data using Bayes' Rule:

Diagram illustrating Bayes' Rule with callouts for the components of the equation:

- Prior**:  $P(\text{params})$
- Likelihood**:  $P(\text{data} | \text{params})$
- Posterior**:  $P(\text{params} | \text{data})$
- Marginal**:  $P(\text{data})$

$$P(\text{params} | \text{data}) = \frac{P(\text{params}) P(\text{data} | \text{params})}{P(\text{data})}$$

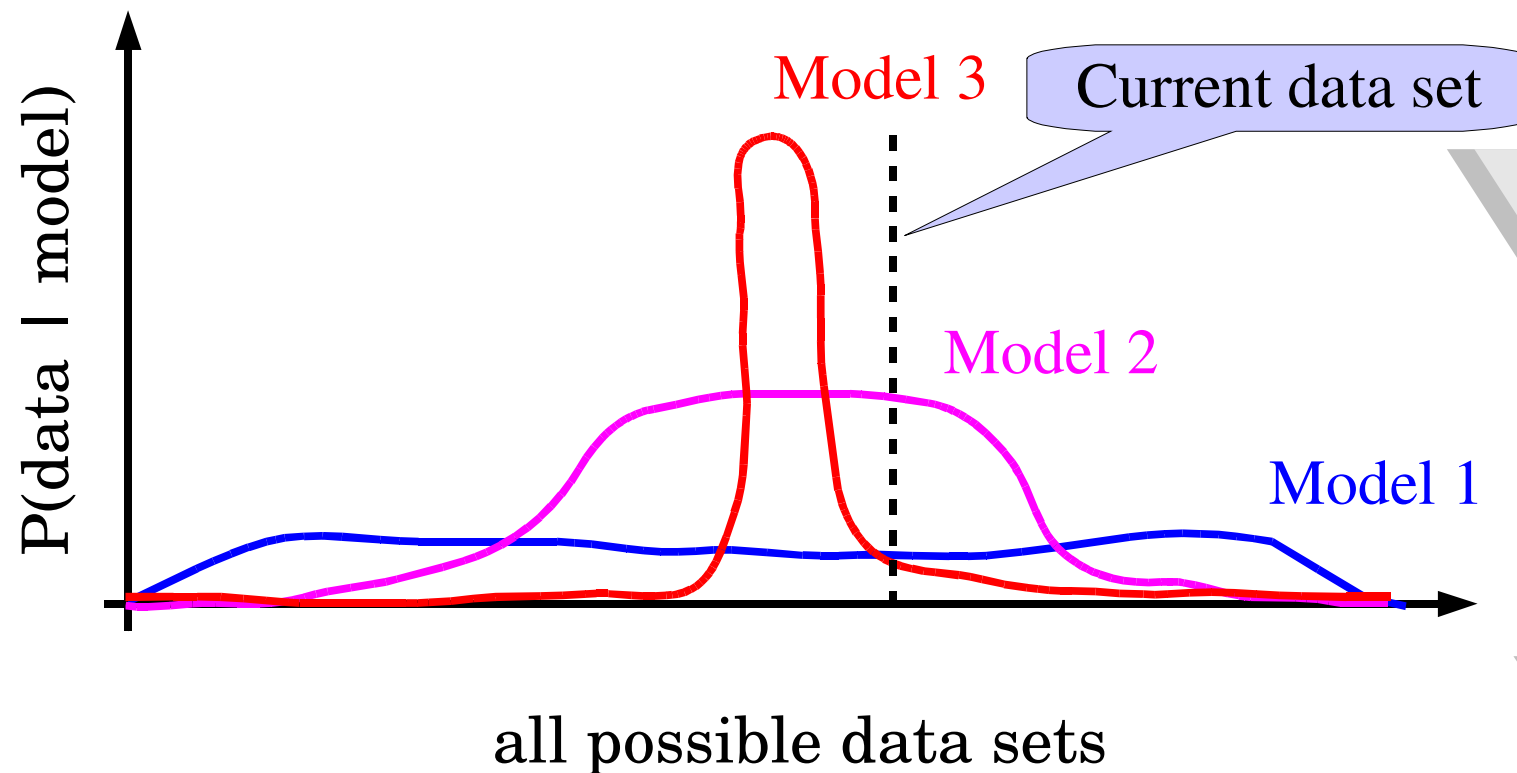
- Use the posterior to:
  - Predict unseen data (machine learning)
  - Reach scientific conclusions (statistics)
  - Make optimal decisions (Bayesian decision theory)

# Models, Parameters and Data

- Model = Our explanation of the world (data)
    - Examples: maximum entropy models, IBM model 1, trigram LM
  - Parameters = All unknown aspects of the model
    - Examples: “lambda” parameters, T-table,  $p(\text{ate} \mid \text{the man})$
  - Data = All observed variables
- 
- Inference problems:
    - Estimate parameters (or their distribution)
    - Estimate missing data (prediction)
    - Find a good model

# What is a *Good Model*?

- We can consider models by looking at the probability that they generate our data set (the marginal likelihood of the data):

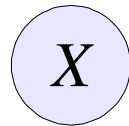


# Tutorial Outline

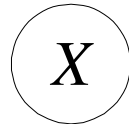
- Introduction to the Bayesian Paradigm
- **Background Material**
  - **Graphical Models**
  - Maximum Likelihood
  - Expectation Maximization
- Priors, priors, priors (subjective, conjugate, reference, etc.)
- Inference Problem and Solutions
  - Summing
  - Monte Carlo
  - Markov Chain Monte Carlo
  - Laplace Approximation
  - Variational Approximation
  - Message Passing...
- Survey of Popular Models
- Pointers to Literature
- Conclusions

# Graphical Models

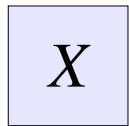
- Convenient notation for representing probability distributions and conditional independence assumptions



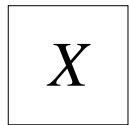
A observed random variable



A unobserved/hidden random variable



A observed/known parameter



A unobserved/unknown parameter

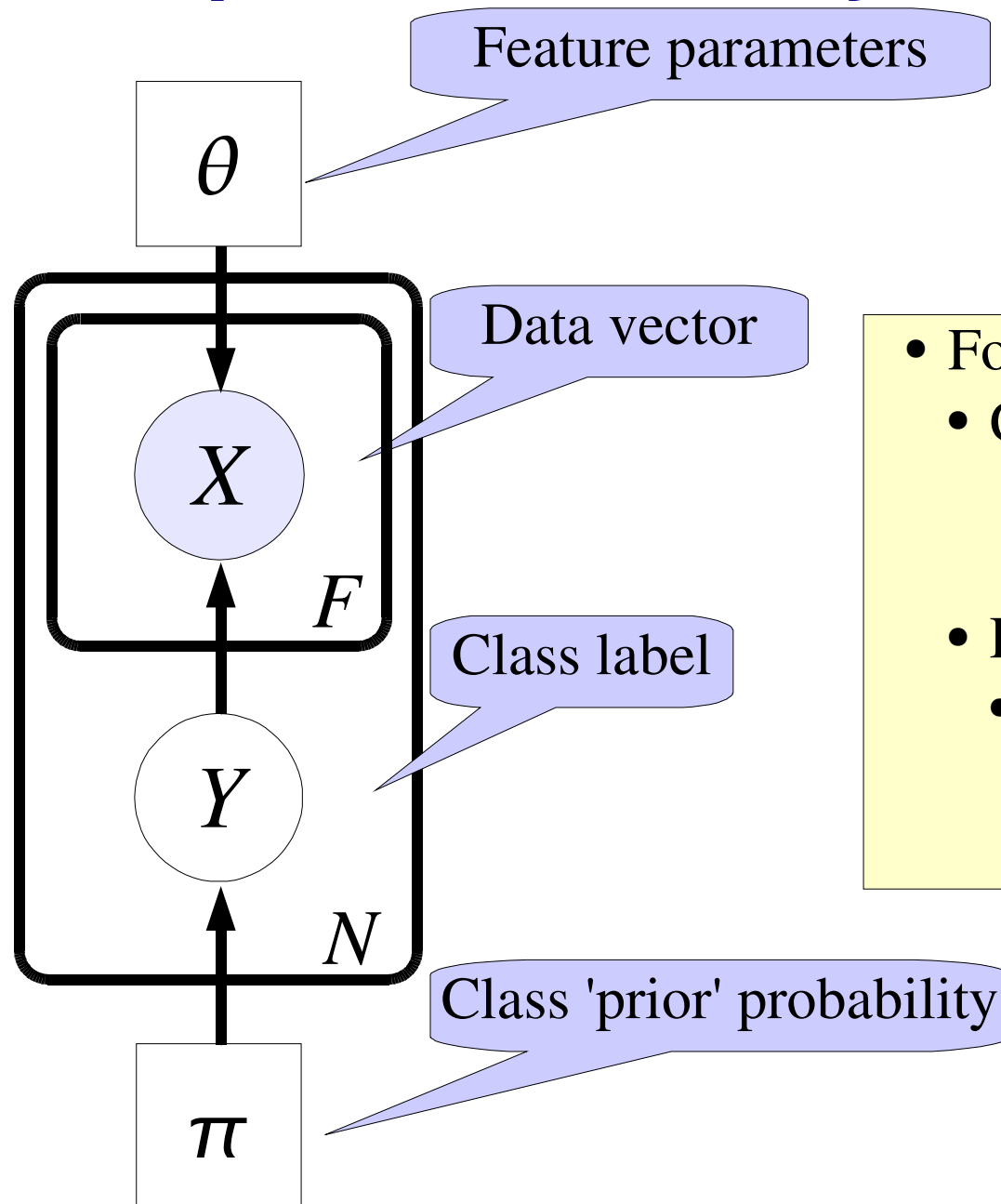


A submodel replicated  $N$  times



An indication of conditional dependence

# Example 1: Naïve Bayes



$$X \mid \theta, Y \sim \text{Binomial}(X \mid \theta^Y)$$

$$Y \mid \pi \sim \text{Multinomial}(\pi)$$

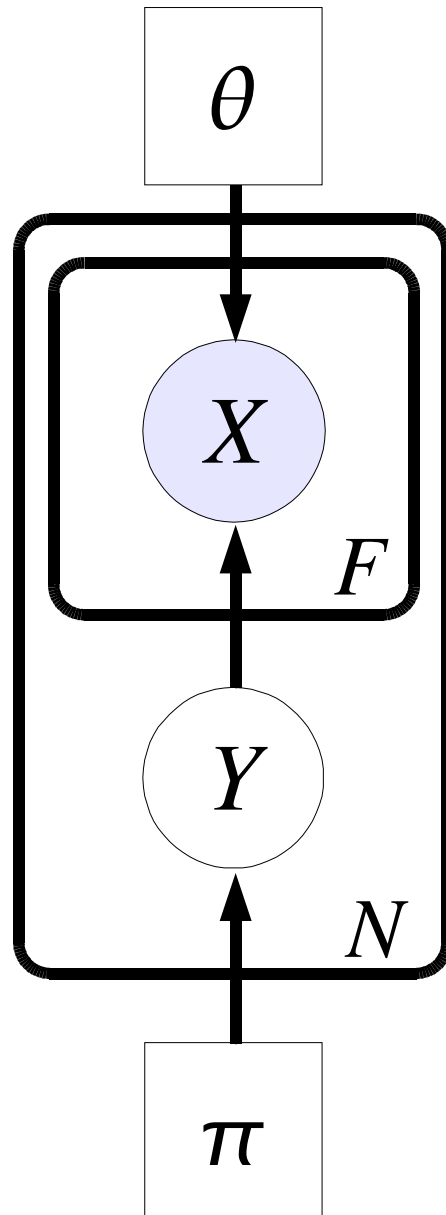
- For each example  $n$ :
- Choose a class  $Y$  by:

$$p(Y=y \mid \pi) \propto \pi_y$$

- For each feature  $f$ :
- Choose  $X$  by:

$$p(X_f \mid \theta^Y) \propto \theta_f^Y$$

# Example 1: Naïve Bayes



$$p(D | \theta, \pi)$$

$$= \prod_n p(y_n | \pi) \prod_f p(x_{nf} | y_n, \theta)$$

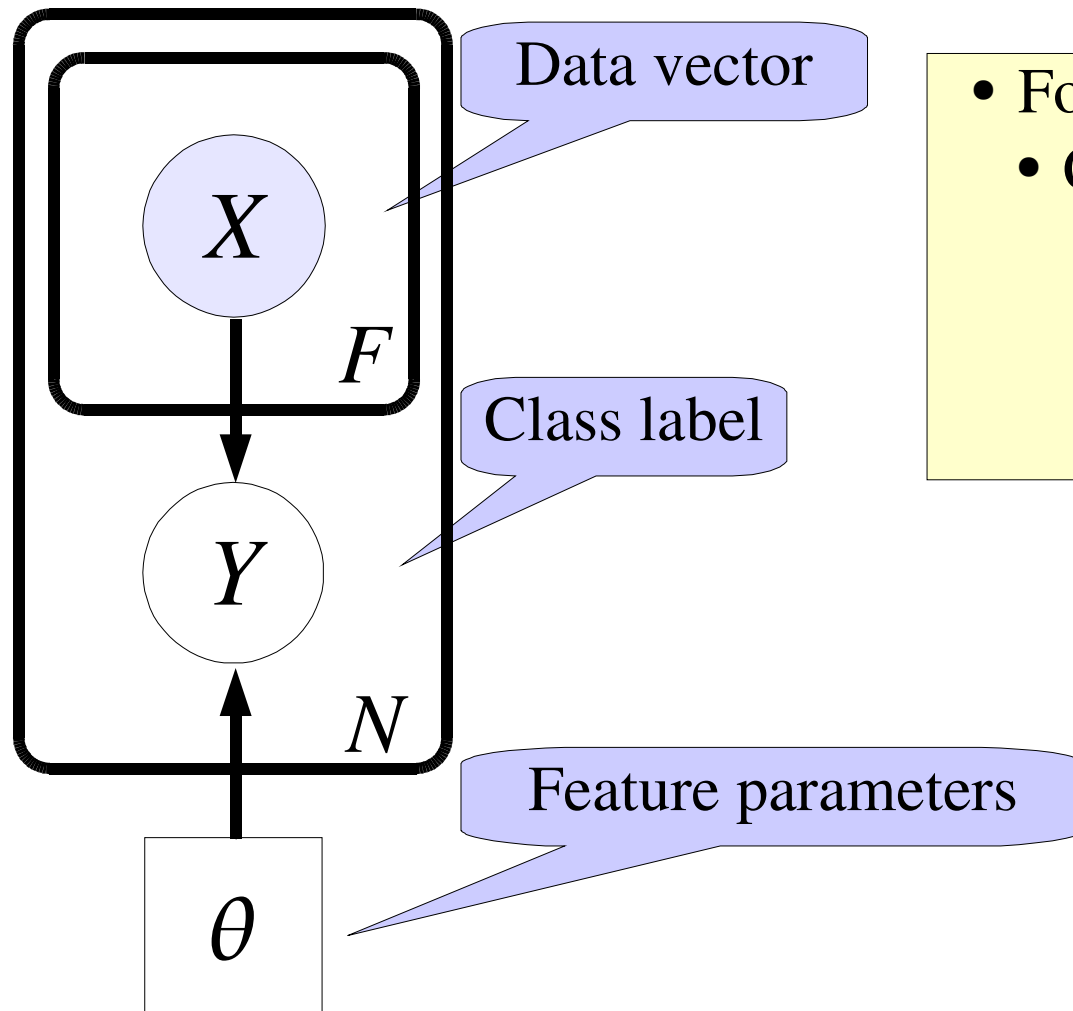
$$= \prod_n \underbrace{\pi^{y_n} (1 - \pi)^{1 - y_n}}_{\text{if } y_n = 1} \underbrace{\prod_f \prod_v \theta_{y_n f v}^{x_{nfv}}}_{\text{if } x_{nfv} = 1}$$

$$\begin{array}{ll} \pi & \text{if } y_n = 1 \\ 1 - \pi & \text{if } y_n = 0 \end{array}$$

$$\begin{array}{ll} \theta_{yfv} & \\ & \text{if } x_{nfv} = 1 \end{array}$$

$\theta_{yfv}$  = probability that feature  $f$  takes value  $v$  if the class is  $y$

# Example 2: Maximum Entropy



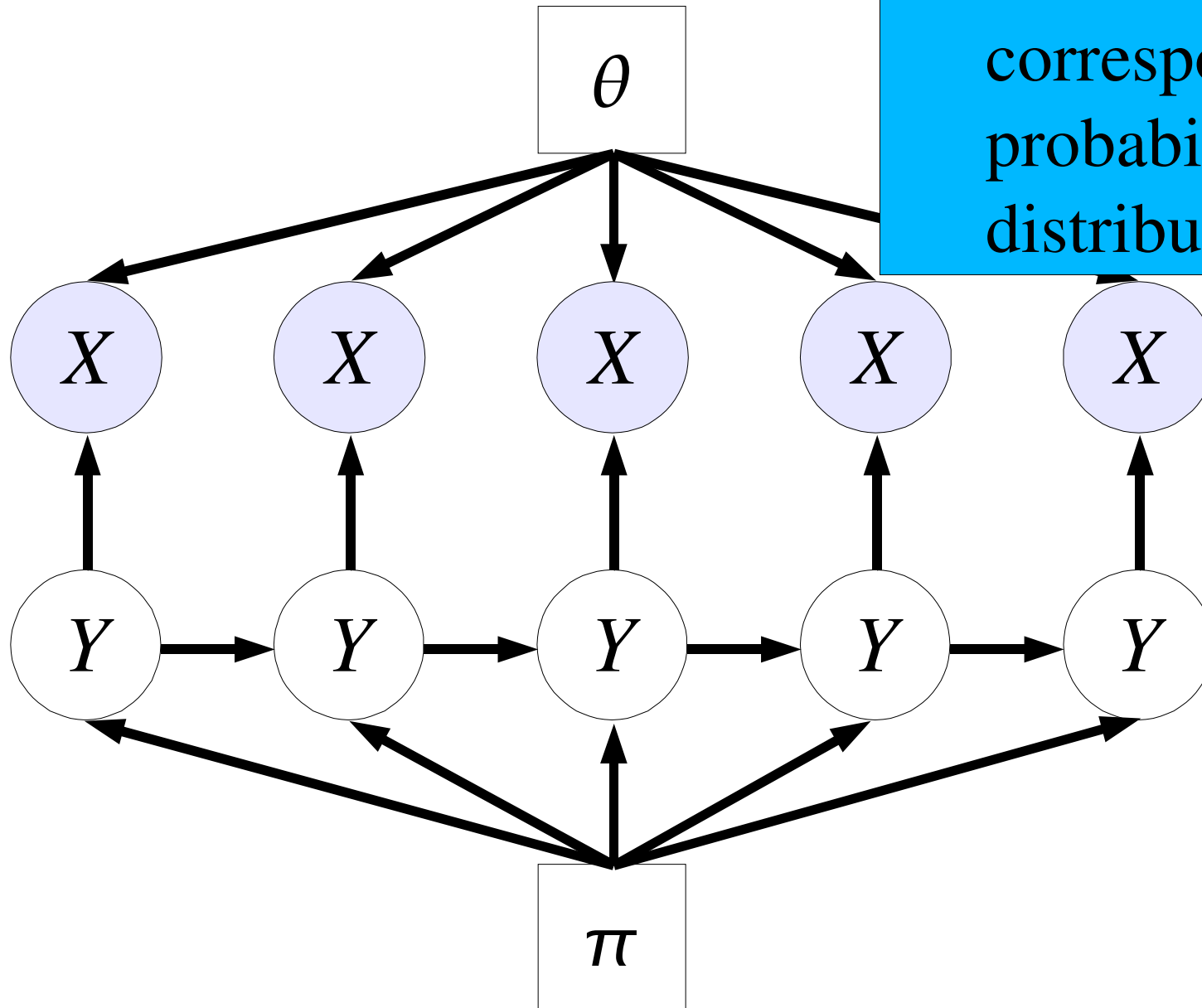
- For each example  $n$ :
- Choose a class  $Y$  by:

$$p(Y=y \mid X, \theta) \propto \exp\left[\sum_f X_f \theta_f\right]$$



# Example 3: Hidden Markov Models

Task: Write out corresponding probability distribution.



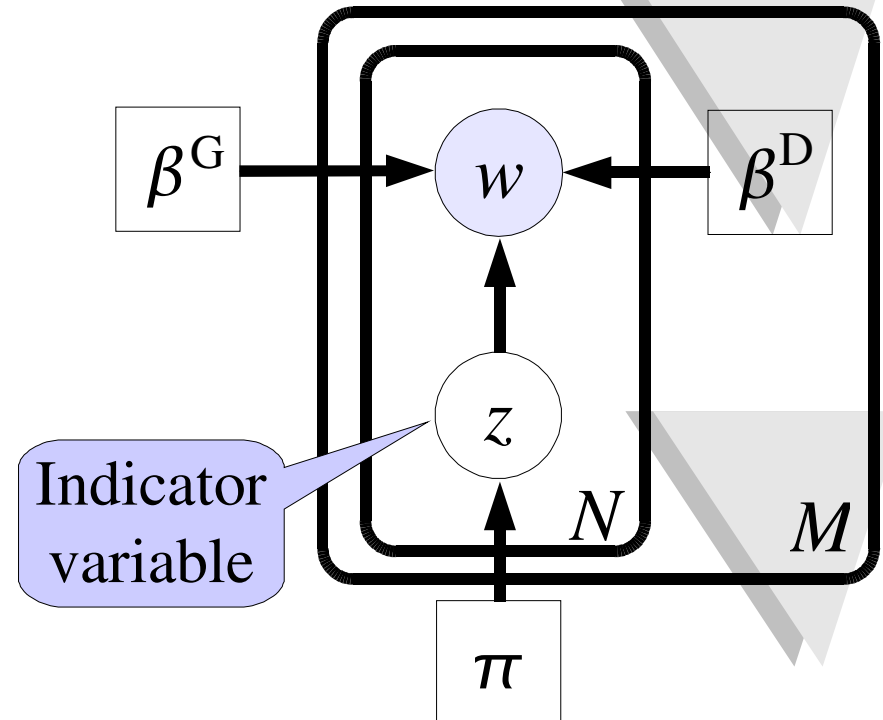
# Example for Summarization

- Consider a stupid summarization model:
  - Each word in a document is drawn independently
  - Each word is drawn either from a *general English* model, or a document specific model
  - We don't know which words are drawn from which

$$p(w \mid \pi, \beta^G, \beta^D) = \prod_m \prod_n \sum_{z_{mn}} p(z_{mn} \mid \pi)$$

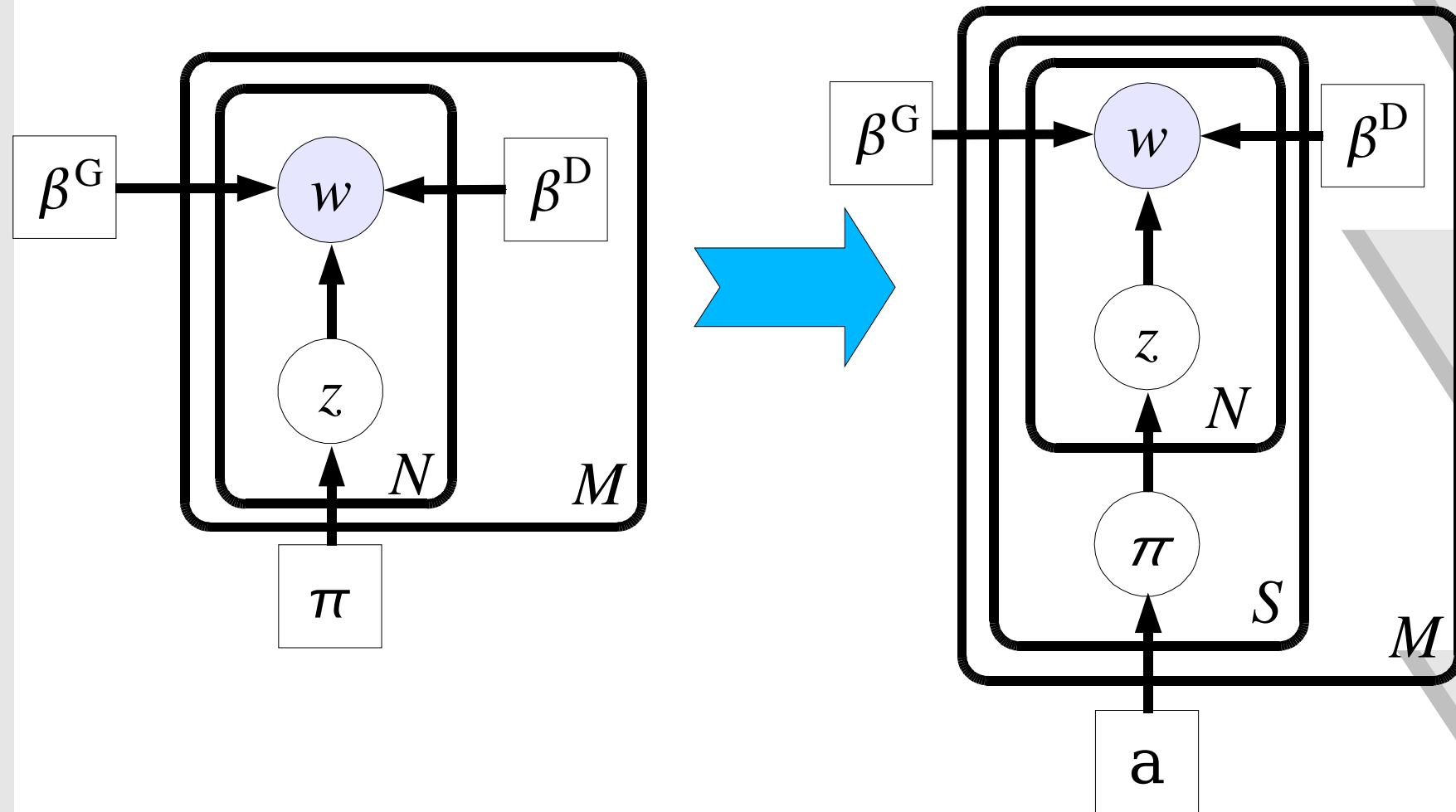
$$p(w \mid \beta^G)^{z_{mn}}$$

$$p(w \mid \beta_m^D)^{1 - z_{mn}}$$



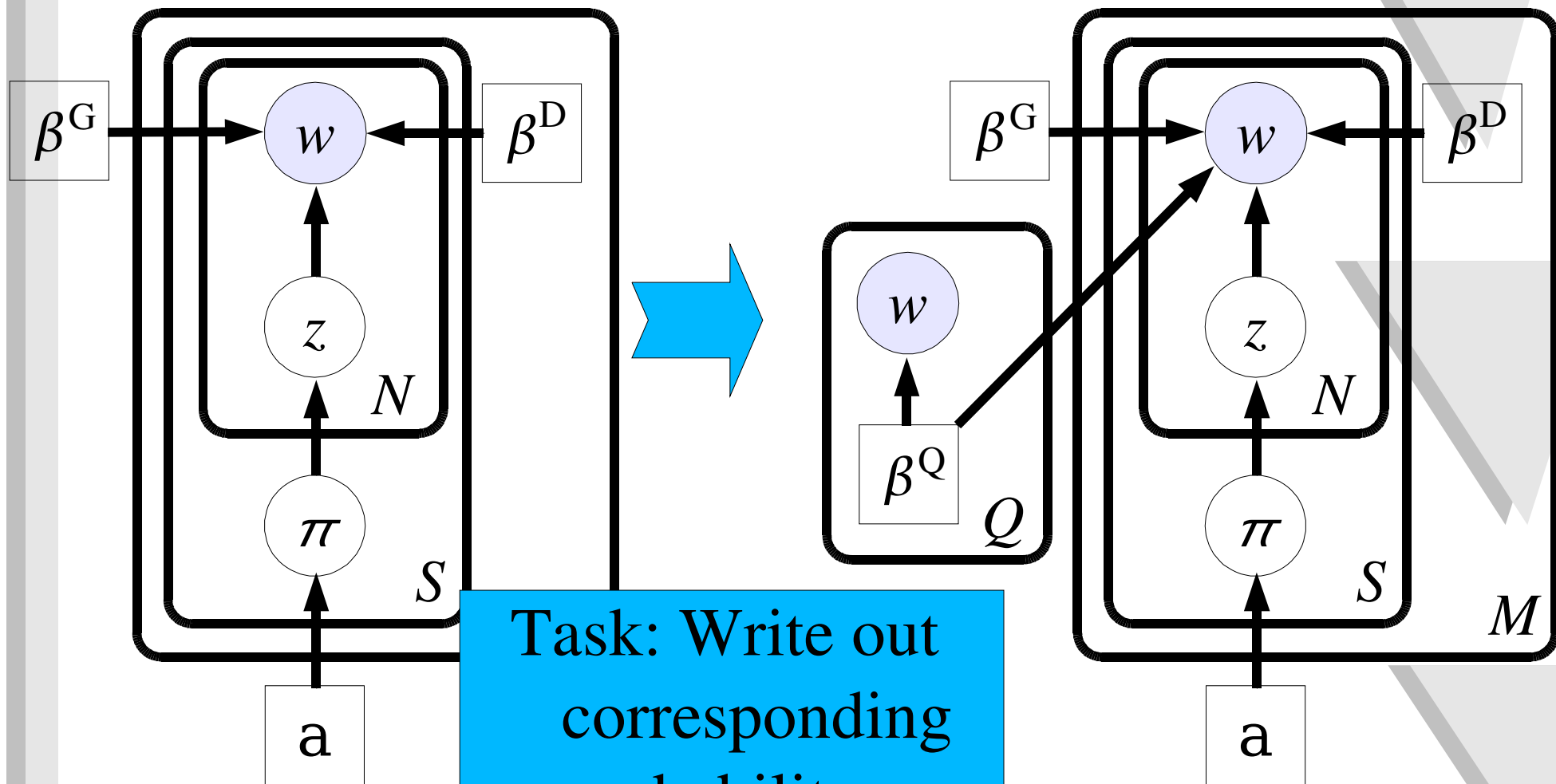
# Fun with Graphical Models

- Easy to propose extensions to the model: add sentences!



# Fun with Graphical Models

- Add queries!



Task: Write out corresponding probability distribution.

# Tutorial Outline

- Introduction to the Bayesian Paradigm
- **Background Material**
  - Graphical Models
  - **Maximum Likelihood**
  - Expectation Maximization
- Priors, priors, priors (subjective, conjugate, reference, etc.)
- Inference Problem and Solutions
  - Summing
  - Monte Carlo
  - Markov Chain Monte Carlo
  - Laplace Approximation
  - Variational Approximation
  - Message Passing...
- Survey of Popular Models
- Pointers to Literature
- Conclusions

# Maximum Likelihood Estimators (MLE)

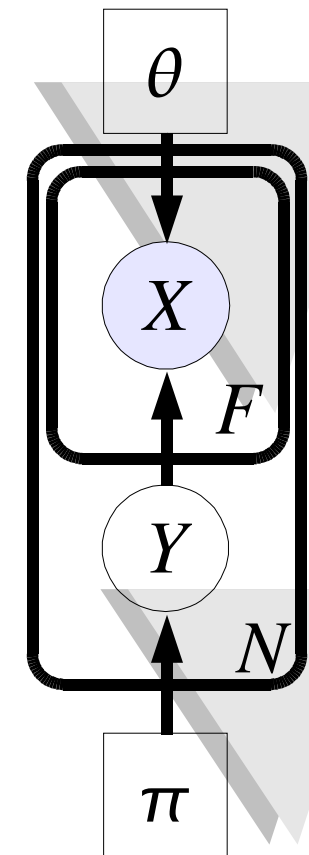
- Take a parameterized model and some data
- Find the parameters that maximize the likelihood of that data (i.e., the 'probability' of the parameters given the data):

$$L(\theta, \pi \mid X_{1:N}, Y_{1:N}) = \prod_{n=1}^N \left( \prod_{k=1}^K \pi_k^{Y_{nk}} (1 - \pi_k)^{1 - Y_{nk}} \right) \left( \prod_{f=1}^F (\theta_f^{Y_n})^{X_{nf}} (1 - \theta_f^{Y_n})^{1 - X_{nf}} \right)$$

$$l(\theta, \pi) = \sum_n \sum_k (Y_{nk} \log \pi_k + (1 - Y_{nk}) \log (1 - \pi_k)) + \sum_n \sum_f (X_{nf} \log \theta_f^{Y_n} + (1 - X_{nf}) \log (1 - \theta_f^{Y_n}))$$

$$\frac{\partial l}{\partial \pi} = \sum_n \sum_k \left[ \frac{Y_{nk}}{\pi_k} - \frac{1 - Y_{nk}}{1 - \pi_k} \right]$$

$$\frac{\partial l}{\partial \theta^k} = \sum_{n: Y_n=k} \sum_f \left[ \frac{X_{nf}}{\theta_f^k} - \frac{1 - X_{nf}}{1 - \theta_f^k} \right]$$



# Tutorial Outline

- Introduction to the Bayesian Paradigm
- **Background Material**
  - Graphical Models
  - Maximum Likelihood
  - **Expectation Maximization**
- Priors, priors, priors (subjective, conjugate, reference, etc.)
- Inference Problem and Solutions
  - Summing
  - Monte Carlo
  - Markov Chain Monte Carlo
  - Laplace Approximation
  - Variational Approximation
  - Message Passing...
- Survey of Popular Models
- Pointers to Literature
- Conclusions

# MLE with hidden variables

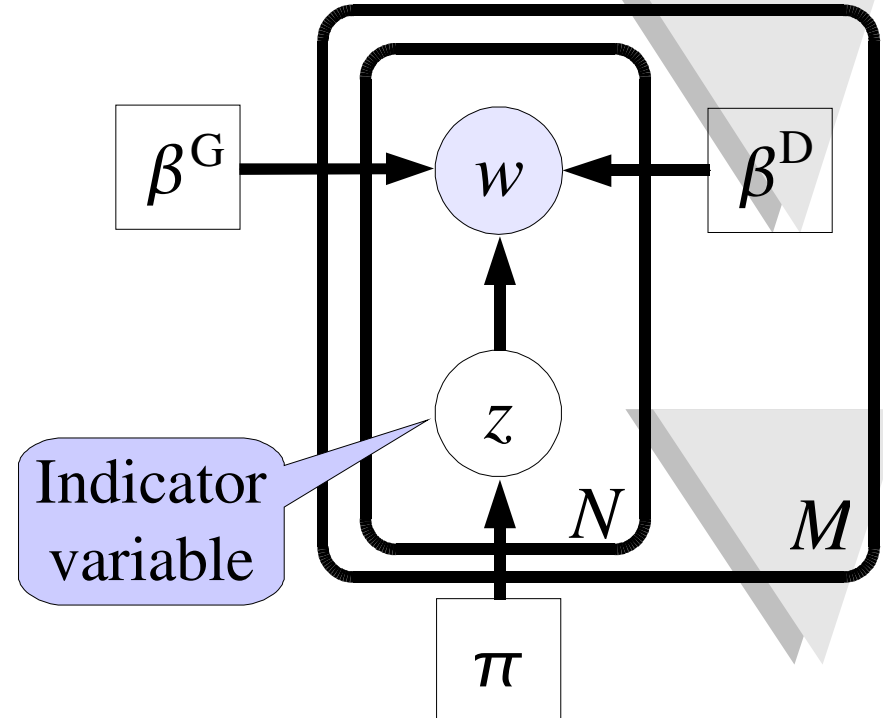
- Consider a stupid summarization model:
  - Each word in a document is drawn independently
  - Each word is draw either from a *general English* model, or a document specific model
  - We don't know which words are drawn from which

$$p(w \mid \pi, \beta^G, \beta^D) = \prod_m \prod_n \sum_{z_{mn}} p(z_{mn} \mid \pi) p(w \mid \beta^G)^{z_{mn}} p(w \mid \beta_m^D)^{1-z_{mn}}$$

- Compute log likelihood:

$$l(\pi, \beta \mid w) = \sum_m \sum_n \log \sum_{z_{mn}} \dots$$

- Uh oh! Logs can't go inside sums!





# Expectation Maximization

- We would like to move the log inside the sum, but can we?
- Jensen's Inequality to the rescue:

$$\begin{aligned}
 \log p(x|\theta) &= \log \int_z dz p(x, z|\theta) \\
 &= \log \int_z dz q(z) \frac{p(x, z|\theta)}{q(z)} \\
 &\geq \int_z dz q(z) \log \frac{p(x, z|\theta)}{q(z)} \\
 &= \int_z q(z) \log p(x, z|\theta) - \int_z q(z) \log q(z) \\
 &= \mathbf{E}_{z \sim q} \{ \log p(x, z|\theta) \} - \mathbf{E}_{z \sim q} \{ \log q(z) \}
 \end{aligned}$$

- For *any* distribution  $Q$  (with the same support)
- How should we choose  $Q$ ?

# Expectation Maximization

- If we set  $q(z) = p(z | x, \theta)$  then the lower bound becomes an equality:

$$\begin{aligned}
 \int_{\mathcal{Z}} dz q(z) \log \frac{p(x, z | \theta)}{q(z)} &= \int_{\mathcal{Z}} dz p(x | z, \theta) \log \frac{p(x, z | \theta)}{p(x | z, \theta)} \\
 &= \int_{\mathcal{Z}} dz p(x | z, \theta) \log \frac{p(z | x, \theta) p(x | \theta)}{p(x | z, \theta)} \\
 &= \int_{\mathcal{Z}} dz p(x | z, \theta) \log p(x | \theta) \\
 &= \log p(x | \theta) \int_{\mathcal{Z}} dz p(x | z, \theta) \\
 &= \log p(x | \theta)
 \end{aligned}$$

- So, when computing  $\mathbf{E}_{z \sim q} \{\log p(x, z | \theta)\}$ , the expectation should be taken with respect to the true posterior

# EM in Practice

- Recall, we wanted to estimate parameters for:

$$\begin{aligned} p(w \mid \pi, \beta^G, \beta^D) &= \prod_m \prod_n \sum_{z_{mn}} p(z_{mn} \mid \pi) p(w \mid \beta^G)^{z_{mn}} p(w \mid \beta_m^D)^{1-z_{mn}} \\ &= \prod_m \prod_n \mathbf{E}_{z_{mn} \sim \pi} \{ p(w \mid \beta^G)^{z_{mn}} p(w \mid \beta_m^D)^{1-z_{mn}} \} \end{aligned}$$

- So we replace the hidden variables with their expectations:

$$l(\beta \mid w) \geq \sum_m \sum_n \mathbf{E}\{z_{mn}\} \log p(w \mid \beta^G) + (1 - \mathbf{E}\{z_{mn}\}) \log p(w \mid \beta_m^D)$$

- All we need to do is calculate the expectations:

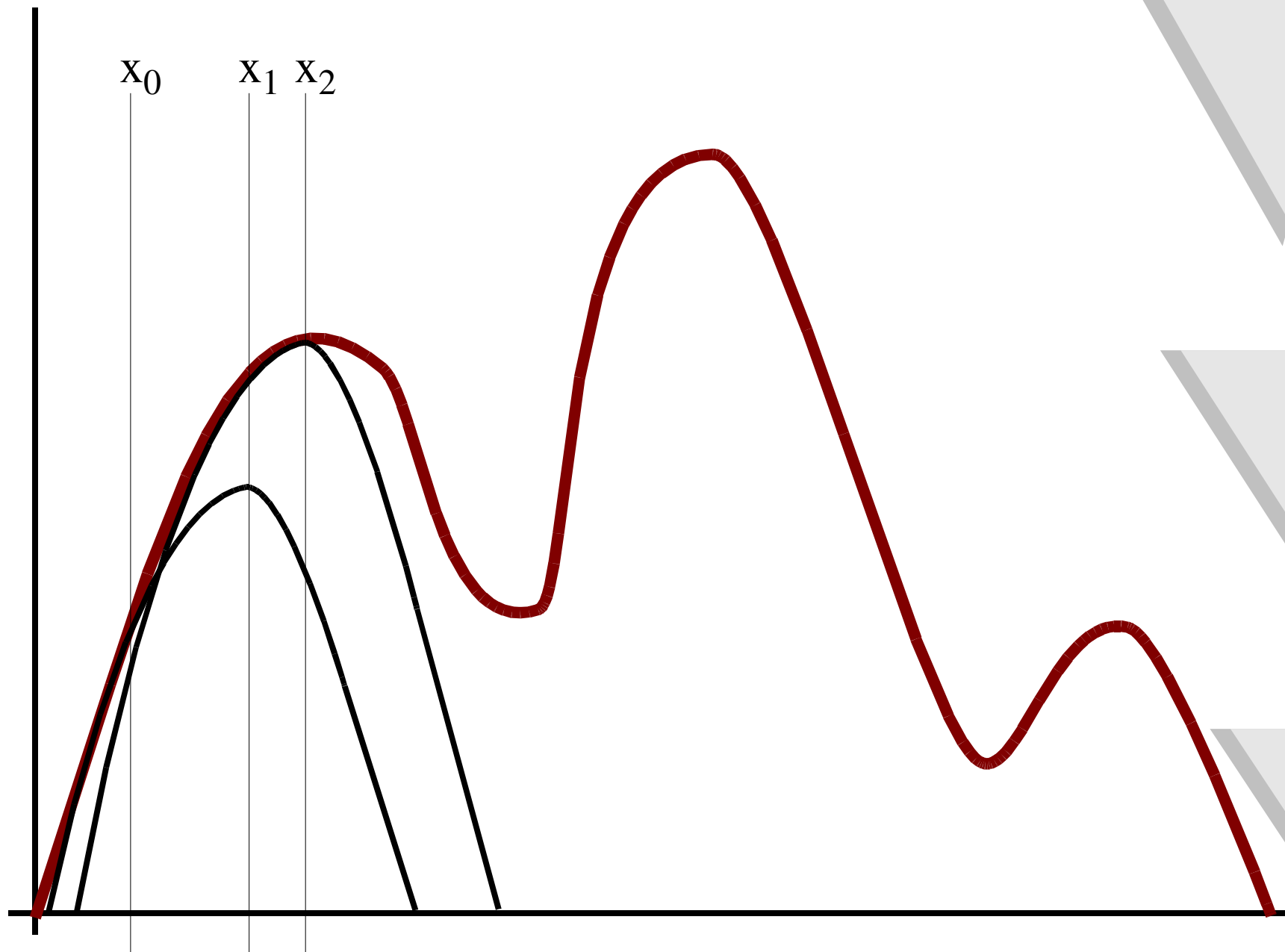
$$\mathbf{E}\{z_{mn}\} \propto p(z_{mn}=1 \mid \pi) p(w \mid \beta^G)$$

- And now the computation proceeds as in the no-hidden-variable setting

# EM Summed Up

- Initialize parameters however you desire
- Repeat:
  - *E-STEP*:  
Compute expectations of hidden variables under the current parameter settings
  - *M-STEP*:  
Optimize parameters given those expectation
  
- This procedure is guaranteed to:
  - Converge to a (local) maximum
  - Monotonically increase the incomplete log-likelihood

# EM Graphically



# EM on our simple model

- Suppose we have three words: {A, B, C}
- Document 1 = [A B], Document 2 = [A C]
- Initialized uniformly

- E-step:  $\mathbf{E}\{z_{mn}\} \propto p(z_{mn}=1 \mid \pi) p(w \mid \beta^G)$

$$\mathbf{E}\{z_{11}\} = \frac{\pi \beta_A^G}{\pi \beta_A^G + (1-\pi) \beta_{1A}^D} = \frac{0.5 * 1/3}{0.5 * 1/3 + 0.5 * 1/3} = 0.5$$

$$\mathbf{E}\{z_{12}\} = \mathbf{E}\{z_{21}\} = \mathbf{E}\{z_{22}\} = 0.5$$

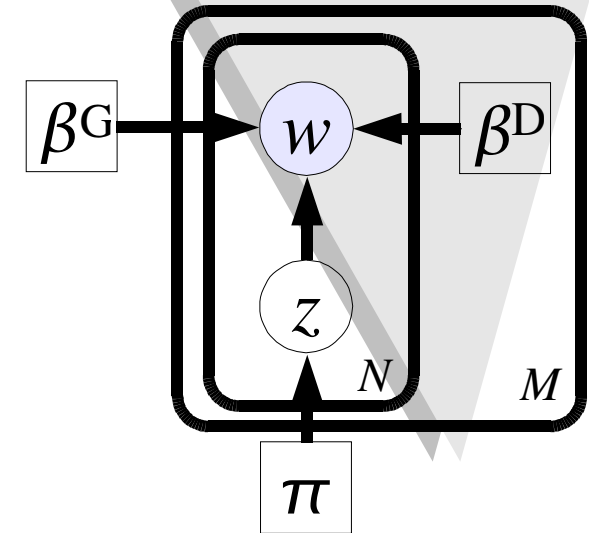
- M-step:

$$\beta_A^G = \frac{1}{Z} [\mathbf{E}\{z_{11}\} + \mathbf{E}\{z_{21}\}] = \frac{1}{2} \quad \beta_B^G = \frac{1}{Z} [\mathbf{E}\{z_{12}\}] = \frac{1}{4} \quad \beta_C^G = \frac{1}{Z} [\mathbf{E}\{z_{22}\}] = \frac{1}{4}$$

$$\beta_{1A}^D = \frac{1}{Z} [1 - \mathbf{E}\{z_{11}\}] = \frac{1}{2} \quad \beta_{1B}^D = \frac{1}{Z} [1 - \mathbf{E}\{z_{12}\}] = \frac{1}{2} \quad \beta_{1C}^D = 0$$

$$\beta_{2A}^D = \frac{1}{Z} [1 - \mathbf{E}\{z_{21}\}] = \frac{1}{2} \quad \beta_{2B}^D = 0 \quad \beta_{2C}^D = \frac{1}{Z} [1 - \mathbf{E}\{z_{22}\}] = \frac{1}{2}$$

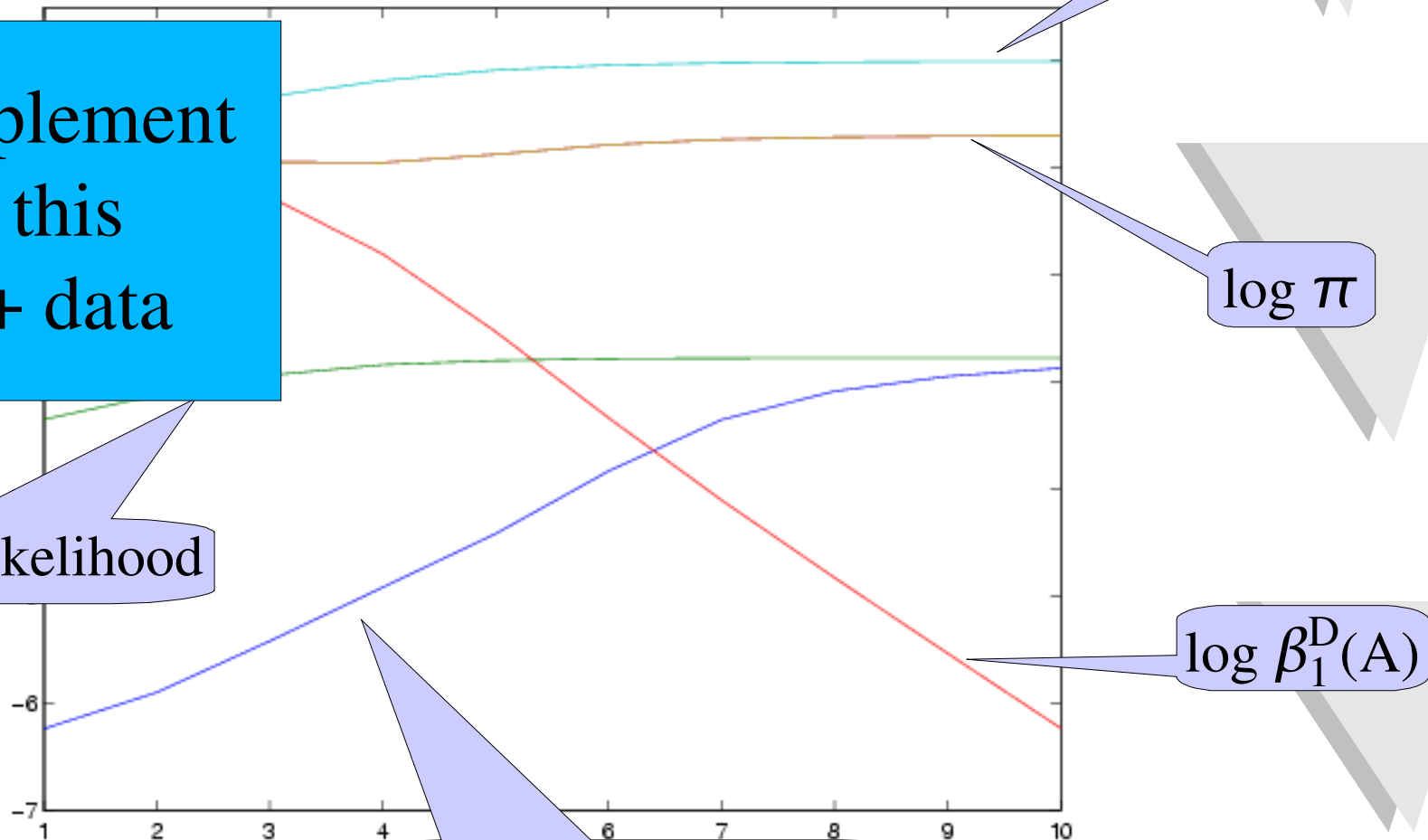
$$\pi = \frac{\mathbf{E}\{z_{11}\} + \mathbf{E}\{z_{21}\}}{\mathbf{E}\{z_{11}\} + \mathbf{E}\{z_{21}\} + \mathbf{E}\{z_{12}\} + \mathbf{E}\{z_{22}\}} = \frac{1}{2}$$



# EM on our simple model

- Suppose we have three words: {A, B, C}
- Document 1 = [A B], Document 2 = [A C]
- Initialized uniformly

Task: Implement EM for this model + data



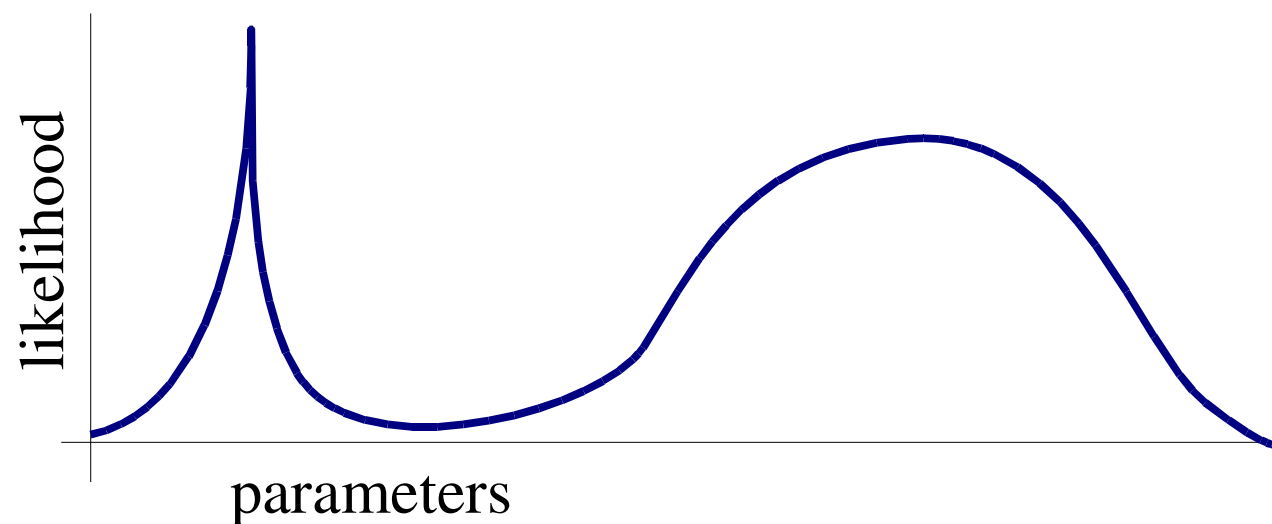
# Problems with Maximum Likelihood

**Powerful model  $\Rightarrow$  Worthless results**  
(due to overfitting...)

**Theoretically unjustified**  
(some would argue...)

**Computationally Expensive**  
(all that cross-validation...)

**Background knowledge is 0/1**





# Tutorial Outline

- Introduction to the Bayesian Paradigm
- Background Material
  - Graphical Models
  - Maximum Likelihood
  - Expectation Maximization
- Priors, priors, priors (subjective, conjugate, reference, etc.)
- Inference Problem and Solutions
  - Summing
  - Monte Carlo
  - Markov Chain Monte Carlo
  - Laplace Approximation
  - Variational Approximation
  - Message Passing...
- Survey of Popular Models
- Pointers to Literature
- Conclusions

# What is a Prior?

- Recall Bayes' Rule:

The diagram shows the Bayes' Rule equation with four callout boxes: 'Posterior' pointing to the left side of the equation, 'Prior' pointing to the numerator's first term, 'Likelihood' pointing to the numerator's second term, and 'Marginal' pointing to the denominator's integral term.

$$P(\theta | D) = \frac{P(\theta) P(D | \theta)}{\int_{\Theta} d\theta P(\theta) P(D | \theta)}$$

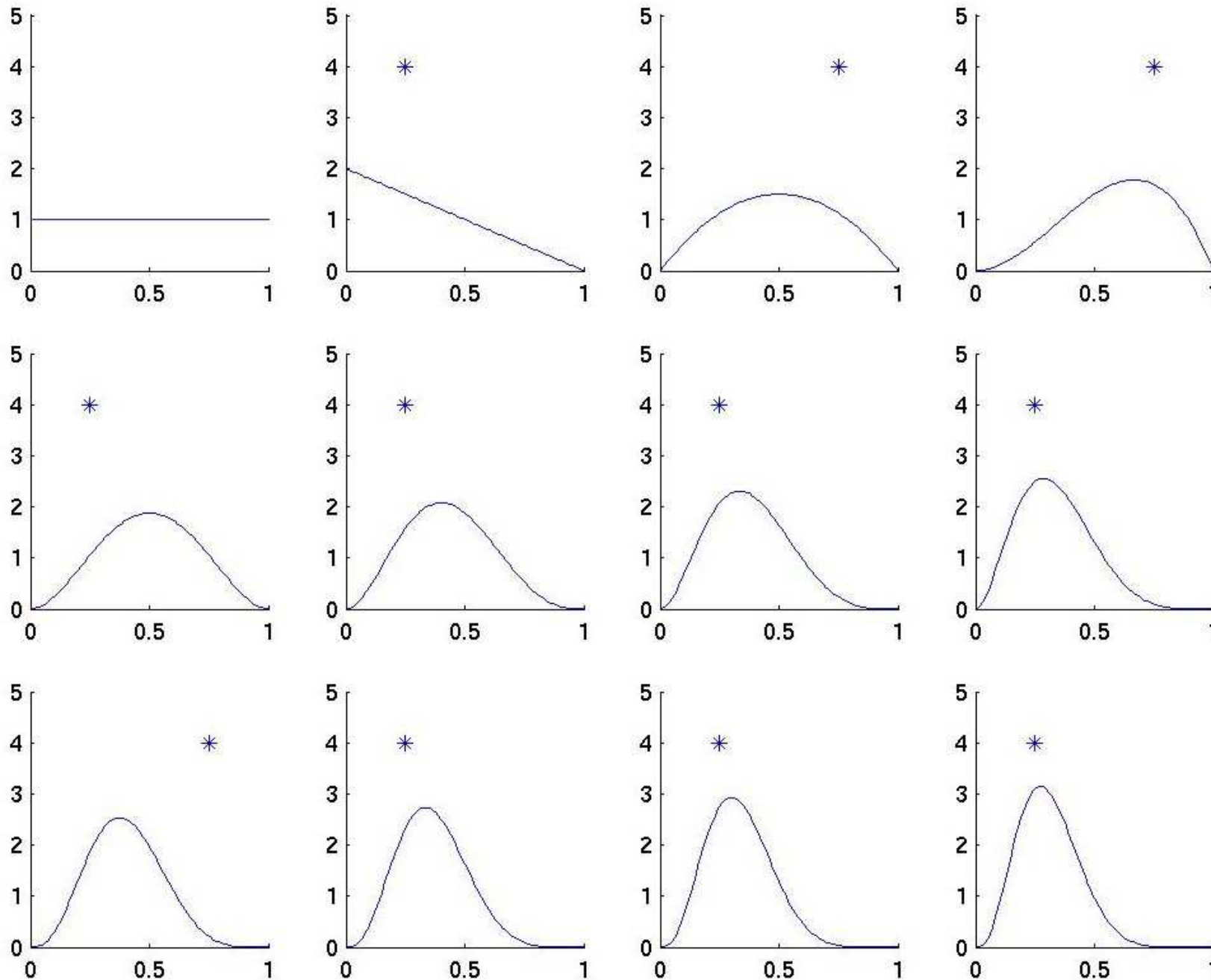
- A *prior* is a specification of our *beliefs* about the values parameters can take, *before seeing any data*

# How Does the Posterior Behave?

Take sequence of data  $x_1, \dots, x_N \dots$

$$\begin{aligned}
 p(\theta) &= \text{just the prior} \\
 p(\theta | x_1) &= \frac{p(\theta) p(x_1 | \theta)}{\int d\theta p(\theta) p(x_1 | \theta)} \\
 p(\theta | x_1, x_2) &= \frac{p(\theta | x_1) p(x_2 | \theta)}{\int d\theta p(\theta | x_1) p(x_2 | \theta)} \\
 &\vdots \\
 p(\theta | x_{1:N}) &= \frac{p(\theta | x_{1:N-1}) p(x_N | \theta)}{\int d\theta p(\theta | x_{1:N-1}) p(x_N | \theta)} \\
 &= \frac{p(\theta) \prod_n p(x_n | \theta)}{\int d\theta p(\theta) \prod_n p(x_n | \theta)}
 \end{aligned}$$

# Binomial Example



# Specifying Priors

- A prior is a map  $\pi$  that:
  - Assigns to every setting of parameters a real value
  - Integrates to 1 over the parameter space
- Such a beast can be difficult to describe! Tools:
  - When the parameters are discrete, we can set them by hand
  - Otherwise, we will often choose a *parametric* prior  $\pi(\theta) = \pi(\theta | \alpha)$  and deal with the *hyper-parameters*
  - Or choose a set of priors and integrate over them (robust Bayes)
  - ...

# Empirical Bayes

- Specify a class of priors (typically a functional form):

$$\Gamma = \{ \pi : \pi(\theta) = g(\theta | \alpha) \}$$

- Estimate the prior by maximizing the marginal likelihood:

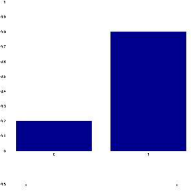
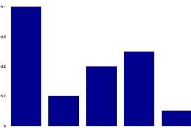
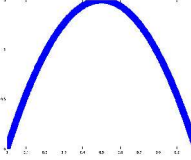
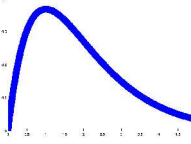
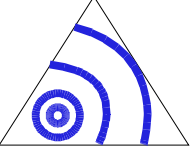
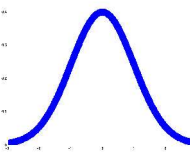
$$\begin{aligned} \hat{\pi} &= \max_{\pi \in \Gamma} p(X | \pi) \\ &= \max_{\alpha \in A} \int_{\Theta} d\theta \pi(\theta | \alpha) p(X | \theta) \end{aligned}$$

# Conjugate (convenient) Priors

- Recall:  $p(\theta | x_{1:N}) = \frac{p(\theta) \prod_n p(x_n | \theta)}{\int d\theta p(\theta) \prod_n p(x_n | \theta)}$
- Given a distribution  $p(\mathbf{x} | \theta)$
- And a prior  $\pi(\theta | \alpha)$
- The prior is *conjugate* if:

$$p(\theta | \alpha, \mathbf{x}) = \frac{\pi(\theta | \alpha) p(\mathbf{x} | \theta)}{\int_{\Theta} F^{\pi(\alpha)}(\theta) p(\mathbf{x} | \theta)} = \pi(\theta | \hat{\alpha})$$

# Summary of Distributions

<u>Distribution</u>	<u>Domain</u>	<u>Picture</u>	<u>Parametric Form</u>
Binomial	Binary		$Bin(x   N, \theta) \propto \theta^n (1 - \theta)^{N-n}$
Multinomial	K classes		$Mult(\bar{x}   \bar{\theta}) \propto \prod \theta_k^{x_k}$
Beta	[0,1]		$Beta(\theta   \alpha, \beta) \propto \theta^{\alpha-1} (1 - \theta)^{\beta-1}$
Gamma	[0, ∞)		$Gam(x   a, b) \propto x^{-a-1} \exp(-bx)$
Dirichlet	Simplex		$Dir(\bar{\theta}   \bar{\alpha}) \propto \prod \theta_k^{\alpha_k-1}$
Gaussian	Reals		$Nor(x   \mu, \sigma^2) \propto \exp\left(-\frac{(x-\mu)^2}{2\sigma^2}\right)$



# Binomial and Beta Distributions

- Binomial distribution models flips of coins (domain= $\{0,1\}$ ):

- Probability that a coin, bias  $\theta$ , flipped  $N$  times will come up  $x$  heads

- Parameters:  $N \in \mathbb{N}^+$ ,  $\theta \in [0,1]$

- Distribution:  $Bin(x | N, \theta) = \binom{N}{x} \theta^x (1-\theta)^{N-x}$

- Moments:  $\mu = N\theta$ ,  $var = N\theta(1-\theta)$

- Beta distribution models nothing (we care about) (domain= $[0,1]$ ):

- Parameters:  $\alpha \in \mathbb{R}^+$ ,  $\beta \in \mathbb{R}^+$

- Distribution:  $Beta(\theta | \alpha, \beta) = \frac{\Gamma(\alpha+\beta)}{\Gamma(\alpha)\Gamma(\beta)} \theta^{\alpha-1} (1-\theta)^{\beta-1}$

- Moments:  $\mu = \frac{\alpha}{\alpha+\beta}$ ,  $var = \frac{\alpha\beta}{(\alpha+\beta)^2(\alpha+\beta+1)}$

- Beta is conjugate to binomial:

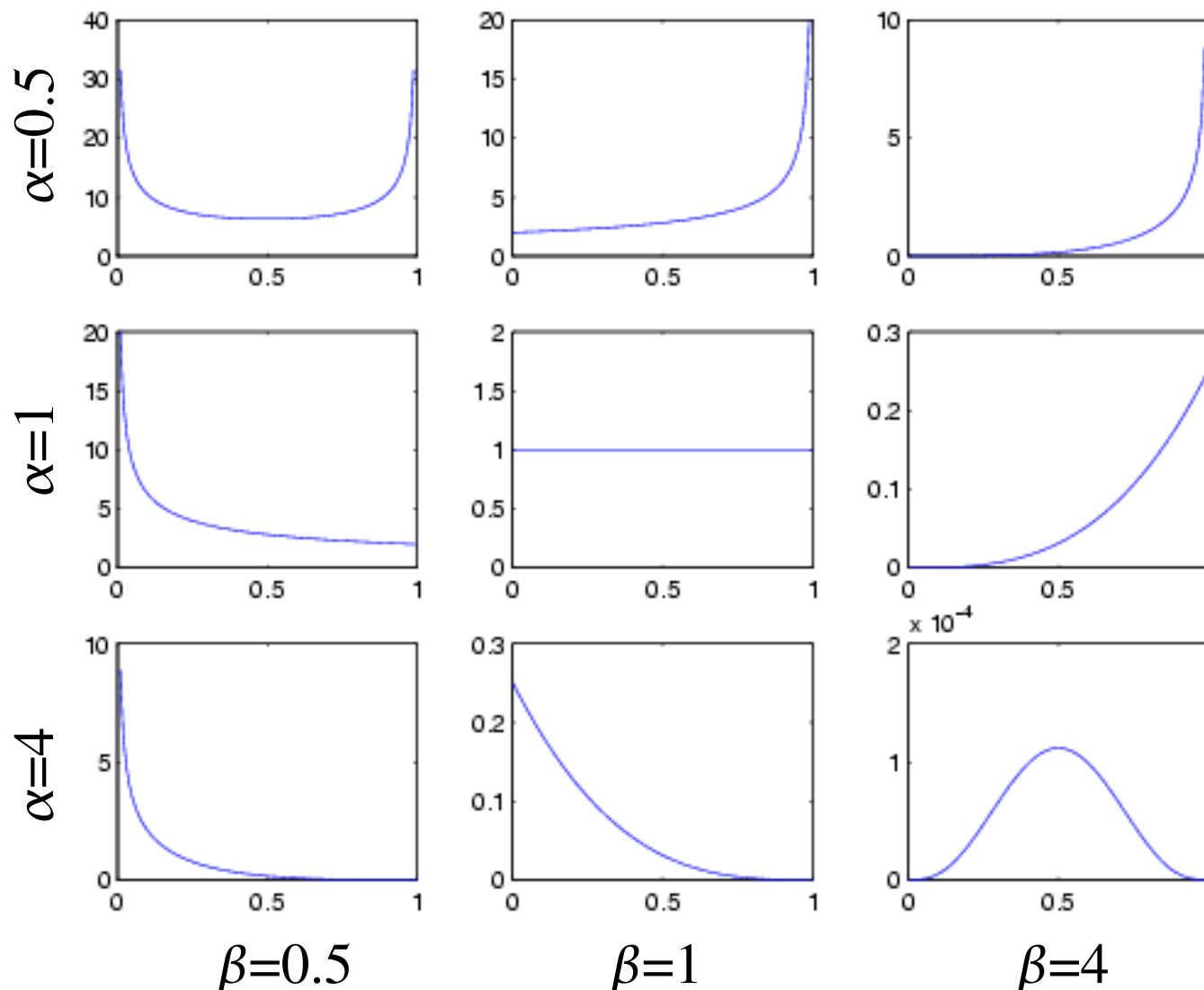
- Posterior parameters:  $\hat{\alpha} = \alpha + x$ ,  $\hat{\beta} = \beta + N - x$

- Marginal distribution:

$$p(x | \alpha, \beta) = \frac{\Gamma(\alpha+\beta)}{\Gamma(\alpha)\Gamma(\beta)} \binom{N}{x} \frac{\Gamma(\alpha+x)\Gamma(\beta+N-x)}{\Gamma(\alpha+\beta+N)}$$

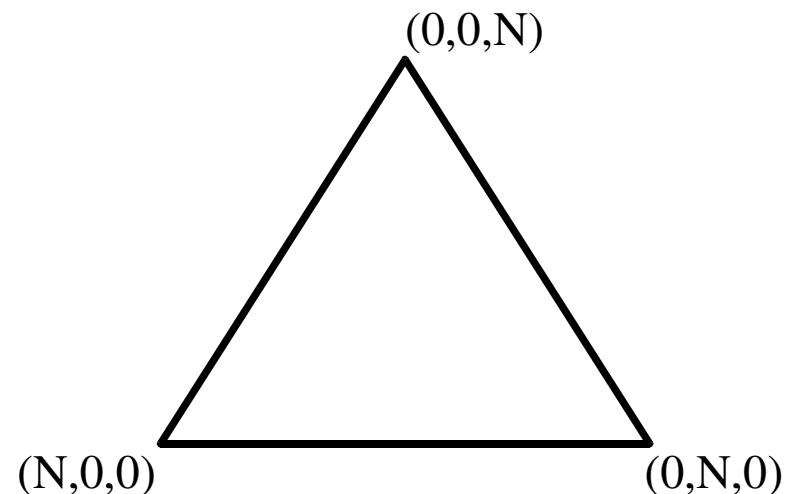
# Beta Distribution Examples

$$\text{Beta}(\theta | \alpha, \beta) = \frac{\Gamma(\alpha + \beta)}{\Gamma(\alpha)\Gamma(\beta)} \theta^{\alpha-1} (1-\theta)^{\beta-1}$$



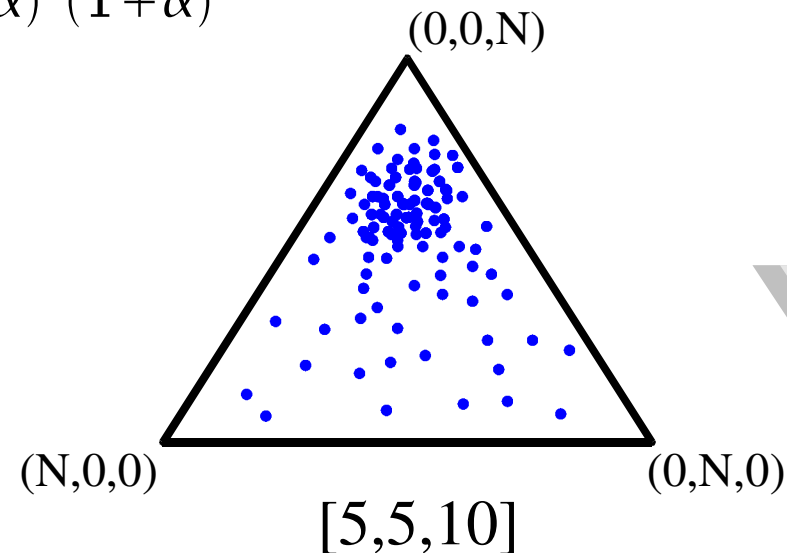
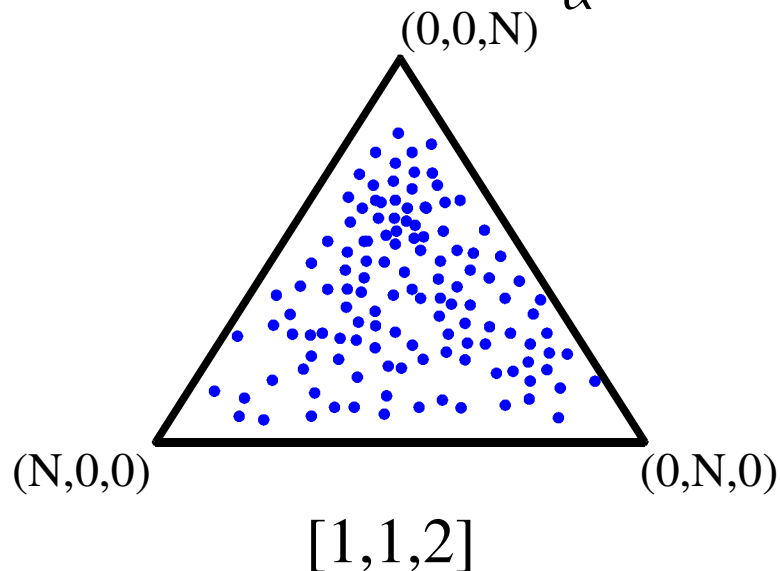
# Multinomial Distribution

- A distribution over counts of  $K > 1$  discrete events (words)
- Domain:  $\langle \mathbf{x}_1, \dots, \mathbf{x}_K \rangle \in \mathbb{N}^K$
- Parameters:  $\langle \theta_1, \dots, \theta_K \rangle \in \Delta_K = \{ \theta_{1:K} : \theta_k \geq 0, \sum_k \theta_k = 1 \}$
- Distribution:  $Mult(\bar{\mathbf{x}} | \bar{\theta}) = \frac{\Gamma(\sum_k x_k + 1)}{\prod_k \Gamma(x_k + 1)} \prod_k \theta_k^{x_k}$
- Moments:  $\langle \theta_1, \dots, \theta_K \rangle \in \Delta_K = \{ \theta_{1:K} : \theta_k \geq 0, \sum_k \theta_k = 1 \}$



# Dirichlet Distribution

- A distribution over a probability simplex
- Domain:  $\langle \theta_1, \dots, \theta_K \rangle \in \delta^K$
- Parameters:  $\langle \alpha_1, \dots, \alpha_K \rangle \in (\mathbb{R}^+)^K$ ,  $\hat{\alpha} = \sum_k \alpha_k$
- Distribution:  $Dir(\bar{\theta} | \bar{\alpha}) = \frac{\Gamma(\hat{\alpha})}{\prod_k \Gamma(\alpha_k)} \prod_k \theta_k^{\alpha_k - 1}$
- Moments:  $\mu_k = \frac{\alpha_k}{\hat{\alpha}}$ ,  $var_k = \frac{\alpha_k(\hat{\alpha} - \alpha_k)}{(\hat{\alpha})^2(1 + \hat{\alpha})}$



# Multinomial/Dirichlet Pair

➤ Multinomial distribution:  $Mult(\bar{\mathbf{x}} | \bar{\theta}) = \frac{\Gamma(\sum \mathbf{x}_k + 1)}{\prod \Gamma(\mathbf{x}_k + 1)} \prod \theta_k^{\mathbf{x}_k}$

➤ Dirichlet distribution:  $Dir(\bar{\theta} | \bar{\alpha}) = \frac{\Gamma(\hat{\alpha})}{\prod_k \Gamma(\alpha_k)} \prod_k \theta_k^{\alpha_k - 1}$

➤ Posterior hyper-parameters:

$$\langle \hat{\alpha}_1, \dots, \hat{\alpha}_K \rangle = \langle \alpha_1 + \mathbf{x}_1, \dots, \alpha_K + \mathbf{x}_K \rangle$$

➤ Marginal Distribution:

$$p(\bar{\mathbf{x}} | \bar{\alpha}) = \frac{\Gamma(\sum \mathbf{x}_k + 1)}{\prod \Gamma(\mathbf{x}_k + 1)} \frac{\Gamma(\hat{\alpha})}{\prod \Gamma(\alpha_k)} \frac{\prod \Gamma(\alpha_k + \mathbf{x}_k)}{\Gamma(\hat{\alpha} + \sum \mathbf{x}_k)}$$

# Gaussian/Gaussian-Gamma

- Gaussian distribution:  $Nor(x | \mu, \sigma^2) = (2\pi\sigma^2)^{-1/2} \exp\left(-\frac{(x-\mu)^2}{2\sigma^2}\right)$
- Gaussian prior:  $Nor(\mu | m, s^2)$
- Gamma prior:  $Gam(\sigma | a, b) = \frac{1}{b^a \Gamma(a)} \sigma^{-2a-1} \exp\left(-\frac{1}{b\sigma^2}\right)$

$$a > 0, b > 0, \text{ domain} = \mathbb{R}^+$$

- Posterior hyper-parameters:

$$\begin{aligned} \mathring{s} &= \left( \frac{1}{s^2} + \frac{1}{\sigma^2} \right)^{-1/2} & \mathring{m} &= \frac{m/s^2 + \sum_i x_i/\sigma^2}{1/s^2 + N/\sigma^2} \\ \mathring{a} &= a + 1/2 & \mathring{b} &= \left( b^{-1} + \frac{1}{2} \sum_i (x_i - \bar{x})^2 \right)^{-1} \end{aligned}$$

- Marginal distribution:

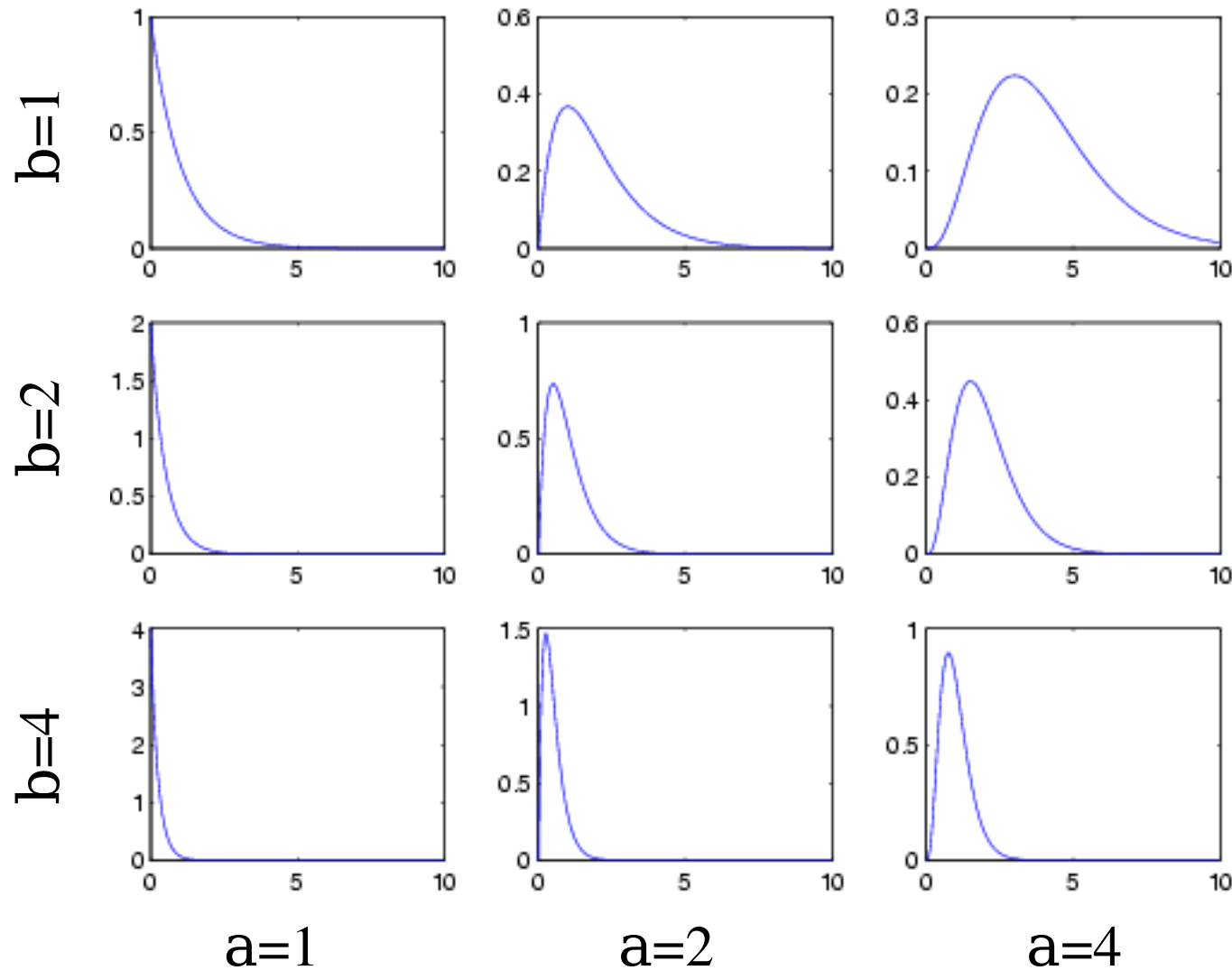
$$p(x | m, s^2, a, b) = \text{StuT}(m, a, b)$$

Non-standard  
Student's T distribution

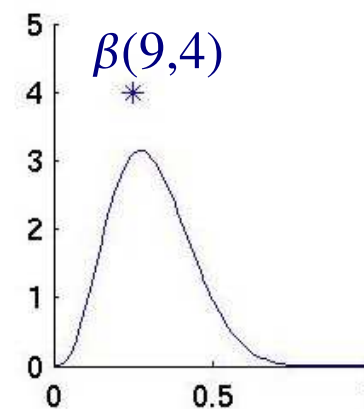
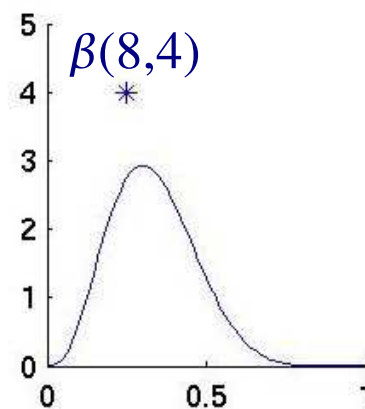
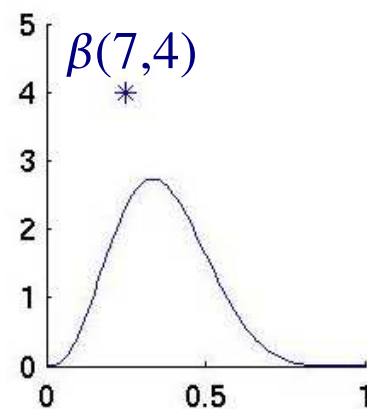
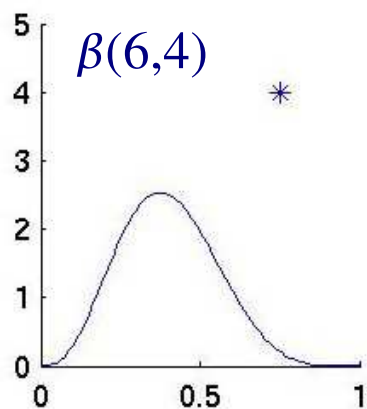
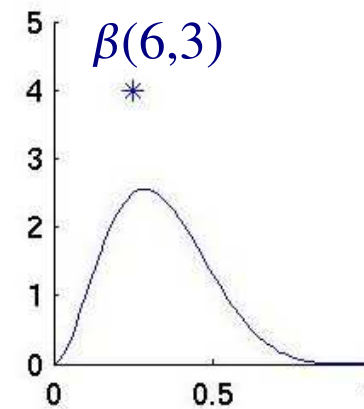
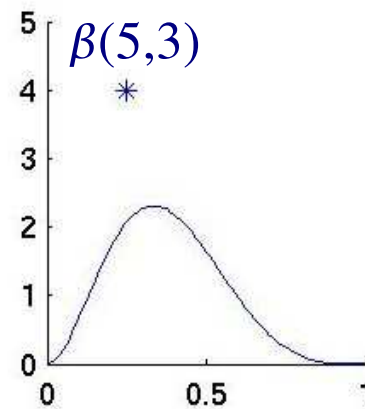
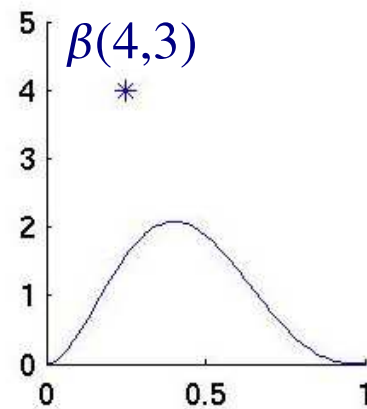
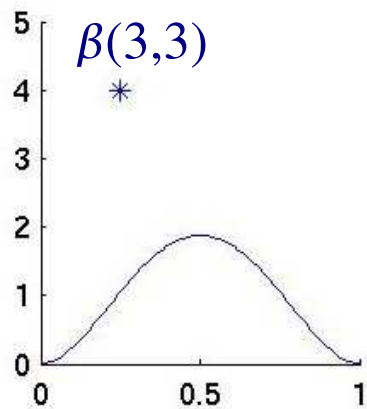
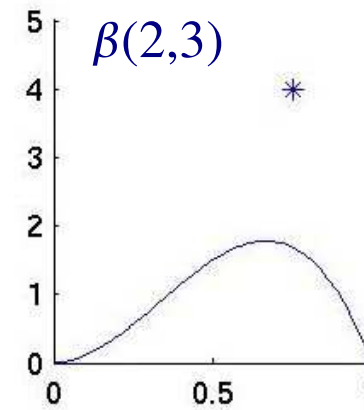
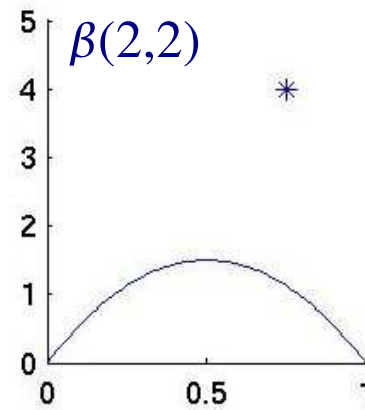
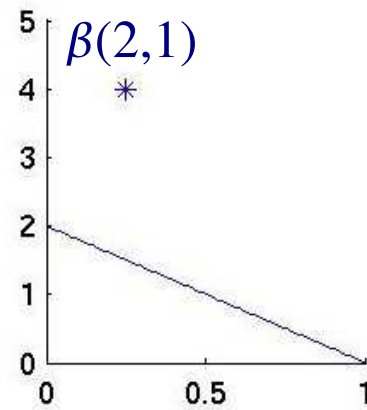
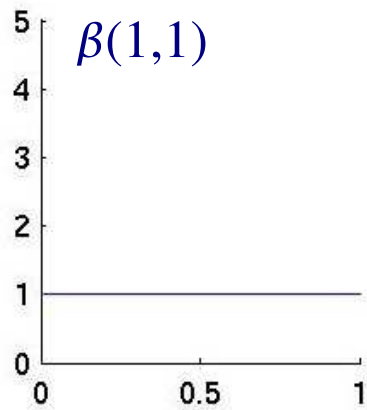
# Gamma Distribution

$$\text{Gam}(x | a, b) = \frac{b^a}{\Gamma(a)} x^{-a-1} \exp(-bx)$$

$$\begin{aligned} \mu &= a/b \\ \text{var} &= a/b^2 \end{aligned}$$



# Conjugate Priors in Action

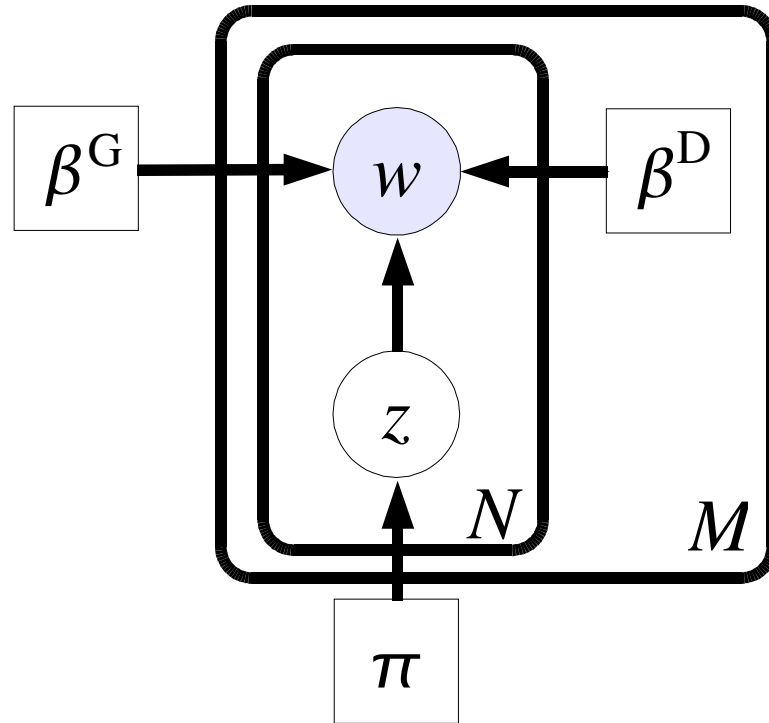




# Tutorial Outline

- Introduction to the Bayesian Paradigm
- Background Material
  - Graphical Models
  - Maximum Likelihood
  - Expectation Maximization
- Priors, priors, priors (subjective, conjugate, reference, etc.)
- **Inference Problem and Solutions**
  - Summing
  - Monte Carlo
  - Markov Chain Monte Carlo
  - Laplace Approximation
  - Variational Approximation
  - Message Passing...
- Survey of Popular Models
- Pointers to Literature
- Conclusions

# Recall our summarization model

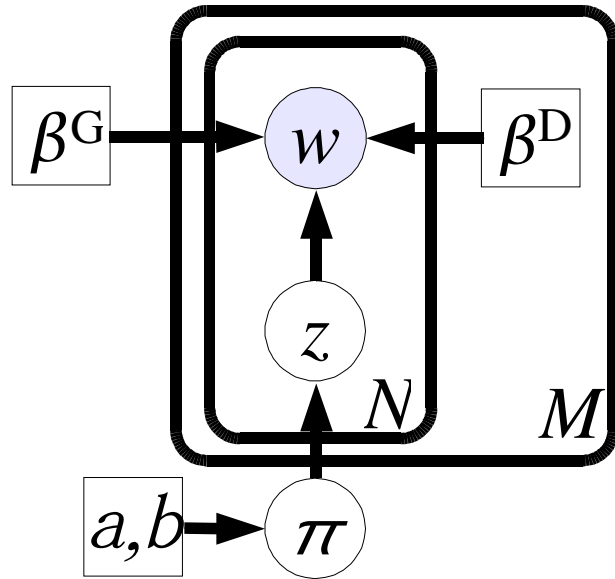


$$z \mid \pi \sim \text{Bin}(\pi)$$

$$w \mid z, \beta \sim \text{Mult}(\beta^G)^z \text{Mult}(\beta^D)^{1-z}$$

- The problem was that we don't believe that it's okay for  $\pi$  to go to 0 or 1
- Solution?  
Put a prior on  $\pi$ !
- What's a good prior?

# Bayesianified summarization model



$$\pi \mid a, b \sim \text{Beta}(a, b)$$

$$z \mid \pi \sim \text{Bin}(\pi)$$

$$w \mid z, \beta \sim \text{Mult}(\beta^G)^z \text{Mult}(\beta^D)^{1-z}$$

$$p(D \mid \beta, a, b) = \int_U d\pi \frac{\Gamma(a+b)}{\Gamma(a)\Gamma(b)} \pi^{a-1} (1-\pi)^{b-1}$$

Conjugacy does not help because of the hidden variables

$$\prod_m \prod_n \sum_{z_{mn} \in \{0,1\}} \pi^{z_{mn}} (1-\pi)^{1-z_{mn}}$$

$$\prod_v (\beta_v^G)^{z_{mn} w_{mnv}} (\beta_{dv}^D)^{(1-z_{mn}) w_{mnv}}$$

# Interesting Inference Questions

- Predict values of unobserved data:

$$P(U | D) \propto \int_{\Theta} d\pi(\theta) P(D | \theta) P(U | \theta)$$

- Compute data likelihood:

$$P(D) \propto \int_{\Theta} d\pi(\theta) P(D | \theta)$$

- Maximize marginal likelihood:

$$P(\alpha | D) \propto \int_{\Theta} d\pi(\theta | \alpha) P(D | \theta)$$

- Estimate posterior:

$$P(\theta | D) = \frac{\pi(\theta) P(D | \theta)}{P(D)}$$

- **GENERAL FORM:**

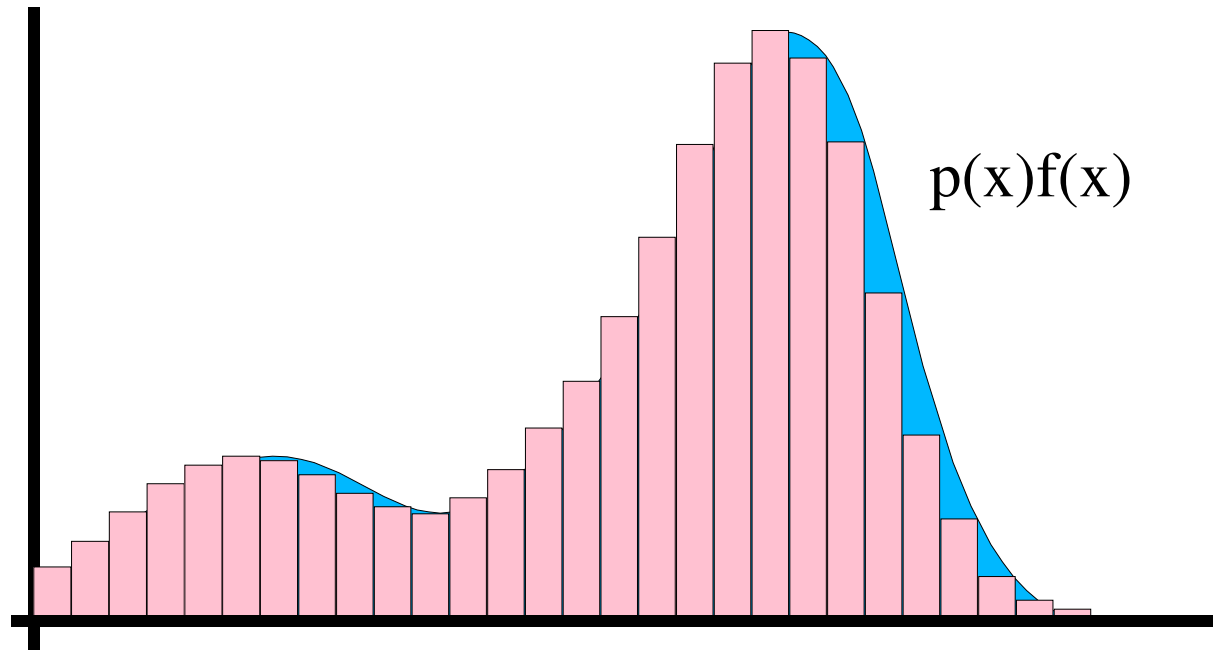
$$F = \int_{\mathcal{X}} dx p(\mathbf{x}) f(\mathbf{x}) = \mathbf{E}_{\mathbf{x} \sim p} \{ f(\mathbf{x}) \}$$

# Tutorial Outline

- Introduction to the Bayesian Paradigm
- Background Material
  - Graphical Models
  - Maximum Likelihood
  - Expectation Maximization
- Priors, priors, priors (subjective, conjugate, reference, etc.)
- Inference Problem and Solutions
  - Summing
  - Monte Carlo
  - Markov Chain Monte Carlo
  - Laplace Approximation
  - Variational Approximation
  - Message Passing...
- Survey of Popular Models
- Pointers to Literature
- Conclusions

# Integration by Summation

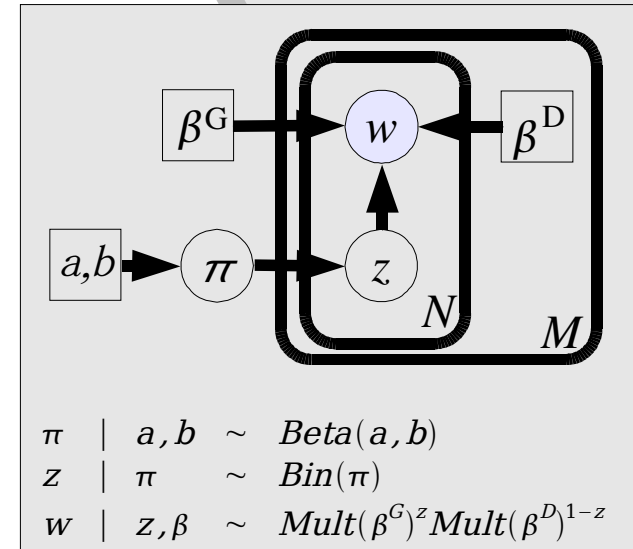
- Remember your 9<sup>th</sup> grade math:



$$F = \int_x dx p(x) f(x) \approx \frac{1}{R} \sum_{x \in R} p(x) f(x)$$

# Summing in our Model

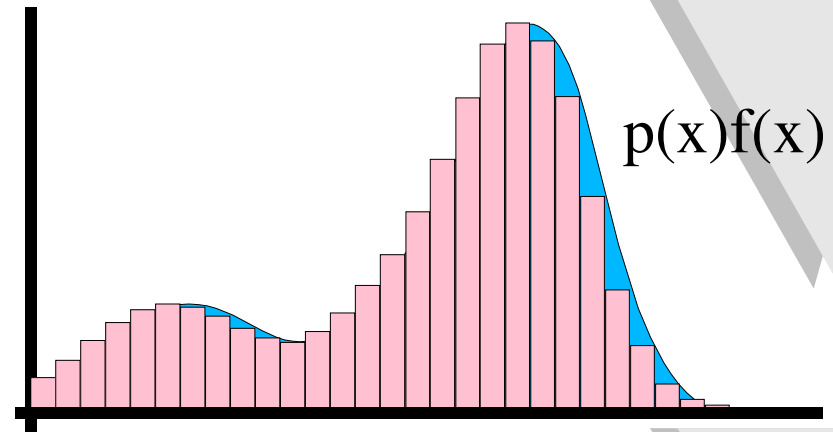
- Simply rewrite the integral as a sum:



$$\begin{aligned}
 p(D | \beta, a, b) &= \int_U d\pi \frac{\Gamma(a+b)}{\Gamma(a)\Gamma(b)} \pi^{a-1} (1-\pi)^{b-1} \\
 &\quad \prod_m \prod_n \sum_{z_{mn} \in \{0,1\}} \pi^{z_{mn}} (1-\pi)^{1-z_{mn}} \prod_v (\beta_v^G)^{z_{mn} w_{mnv}} (\beta_{dv}^D)^{(1-z_{mn}) w_{mnv}} \\
 &\approx \frac{\sum_{p=1}^{100} \Gamma(a+b)}{\Gamma(a)\Gamma(b)} (p/100)^{a-1} (1-p/100)^{b-1} \\
 &\quad \prod_m \prod_n \sum_{z_{mn} \in \{0,1\}} (p/100)^{z_{mn}} (1-p/100)^{1-z_{mn}} \\
 &\quad \prod_v (\beta_v^G)^{z_{mn} w_{mnv}} (\beta_{dv}^D)^{(1-z_{mn}) w_{mnv}}
 \end{aligned}$$

# Integration by Summation

- Pros:
  - Easy to implement
  - Arbitrarily accurate
- Cons:
  - Only works for doubly-bounded regions
  - Intractable for  $>1$  or  $>2$  dimensions
  - Difficult to choose granularity
- Idea: let's choose  $R$  differently



$$F = \int_{\mathcal{X}} dx p(\mathbf{x}) f(\mathbf{x}) \approx \frac{1}{R} \sum_{\mathbf{x} \in R} p(\mathbf{x}) f(\mathbf{x})$$

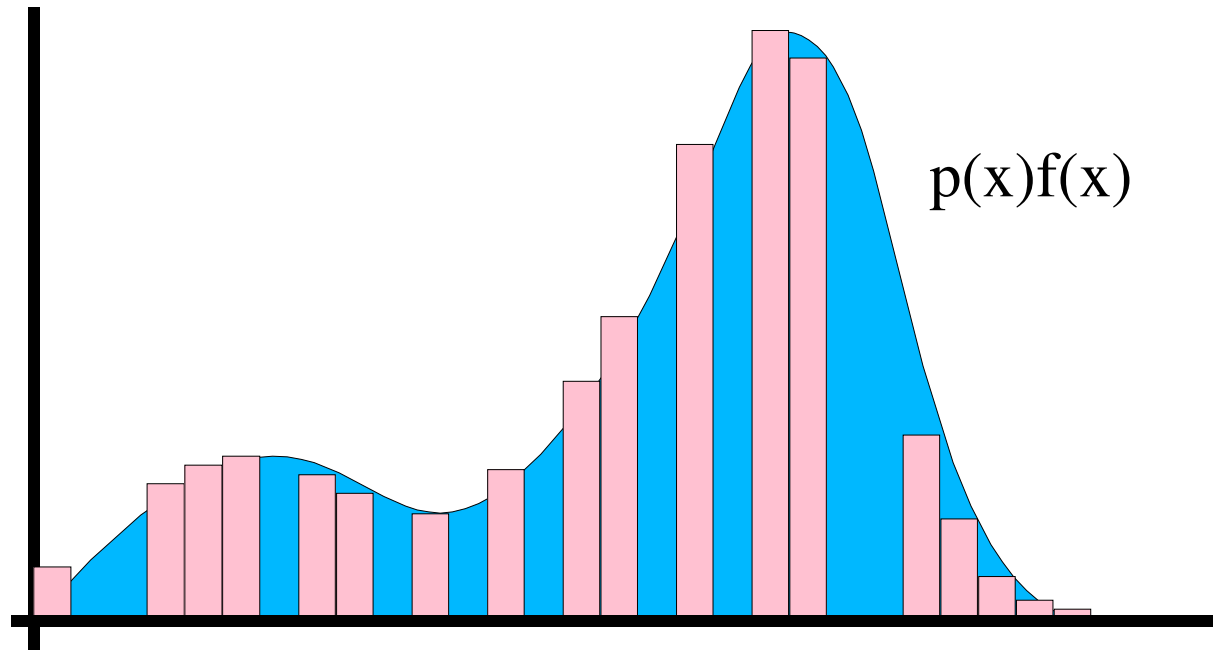


# Tutorial Outline

- Introduction to the Bayesian Paradigm
- Background Material
  - Graphical Models
  - Maximum Likelihood
  - Expectation Maximization
- Priors, priors, priors (subjective, conjugate, reference, etc.)
- Inference Problem and Solutions
  - Summing
  - Monte Carlo
  - Markov Chain Monte Carlo
  - Laplace Approximation
  - Variational Approximation
  - Message Passing...
- Survey of Popular Models
- Pointers to Literature
- Conclusions

# Monte Carlo Integration

- Uniform sampling:
  - Let  $R$  be a (multi)set of points drawn uniformly at random

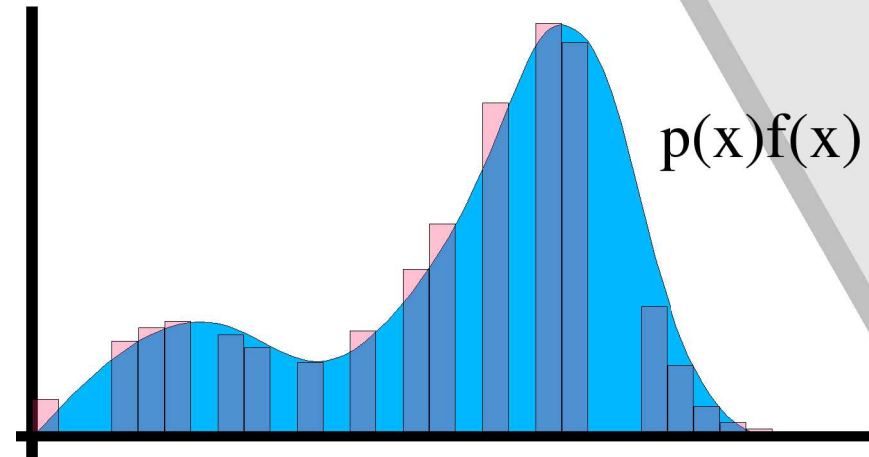


$$F = \int_X dx p(x) f(x) \approx \frac{1}{R} \sum_{x \in R} p(x) f(x)$$

# Uniform Sampling

## ➤ Pros:

- Can now work in arbitrarily high dimensions (in theory)
- Choice is now size of  $R$ , not the width of windows



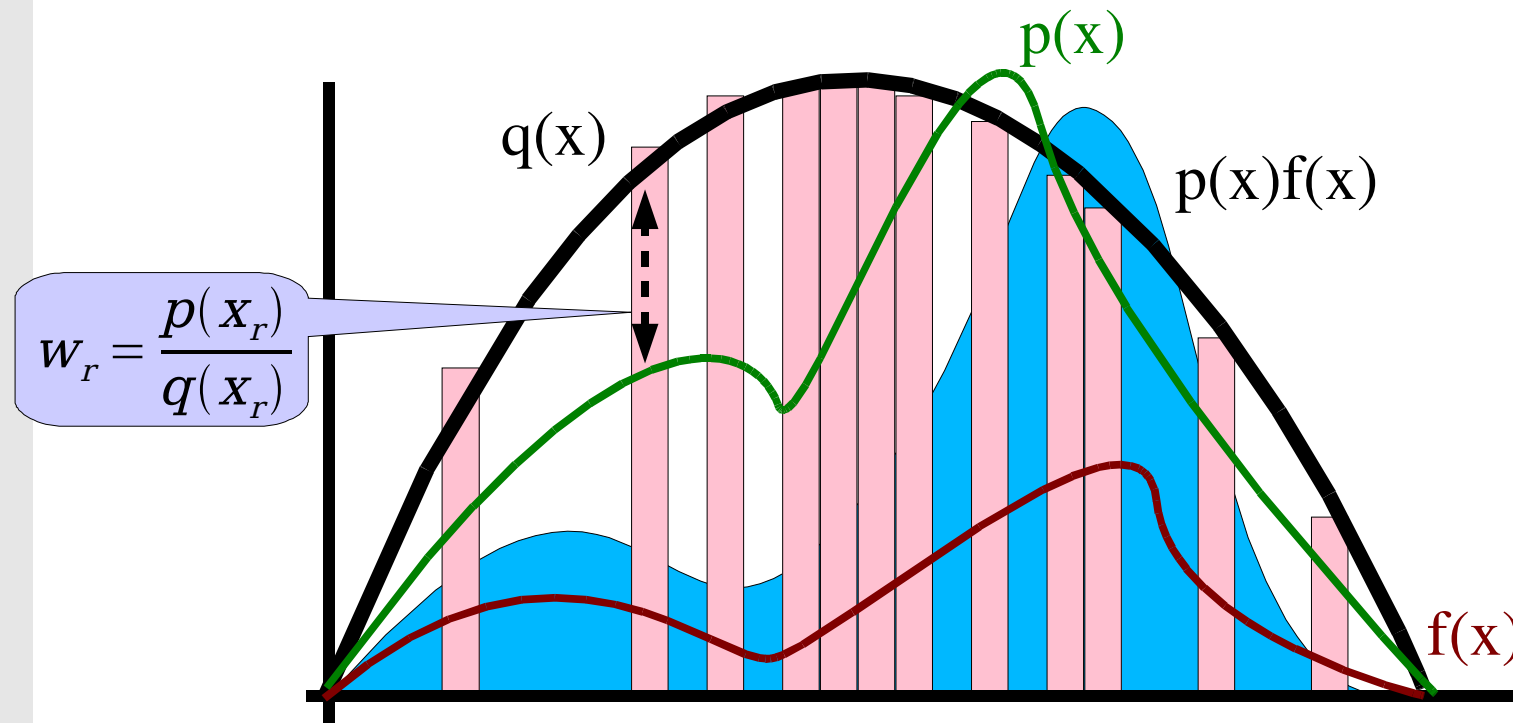
## ➤ Cons:

- Number of samples required to get near the mode of a spiky distribution is huge:  $R \sim 2^{D/2}$
- True distribution is rarely uniform

$$F = \int_x dx p(x) f(x) \approx \frac{1}{R} \sum_{x \in R} p(x) f(x)$$

# Importance Sampling

- Let  $R$  be a set of points drawn from a proposal distribution  $q$

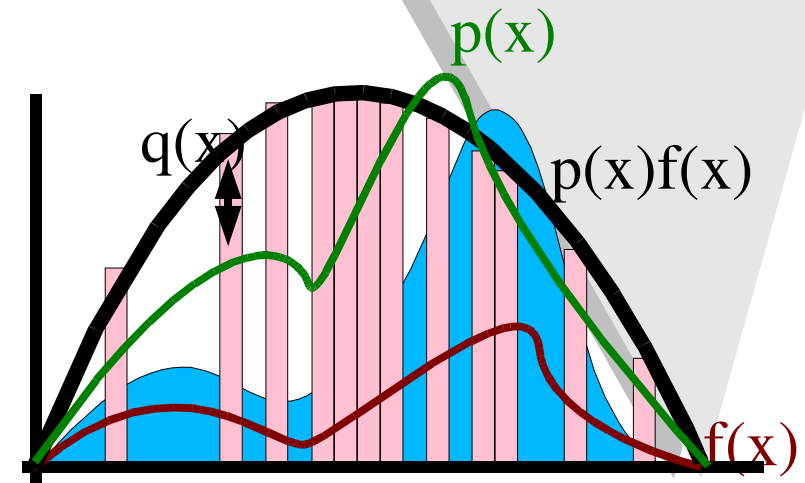


$$F = \int_X dx p(x) f(x) \approx \frac{\sum_r w_r f(x_r)}{\sum_r w_r}, \quad w_r = \frac{p(x_r)}{q(x_r)}$$

# Importance Sampling

- Pros:
  - If  $q$  can be constructed similar to  $p$ , then good samples can be had
  - Can scale better than uniform sampling (not saying much)
- Cons:
  - Very sensitive to choice of  $q$
  - Hard to evaluate whether it has converged
  - Still a lot of samples required:

$$\text{IS: } R \sim \exp\sqrt{2D} \qquad \text{US: } R \sim 2^{D/2}$$



$$F \approx \frac{\sum w_r f(x_r)}{\sum w_r}, \quad w_r = \frac{p(x_r)}{q(x_r)}$$

# Tutorial Outline

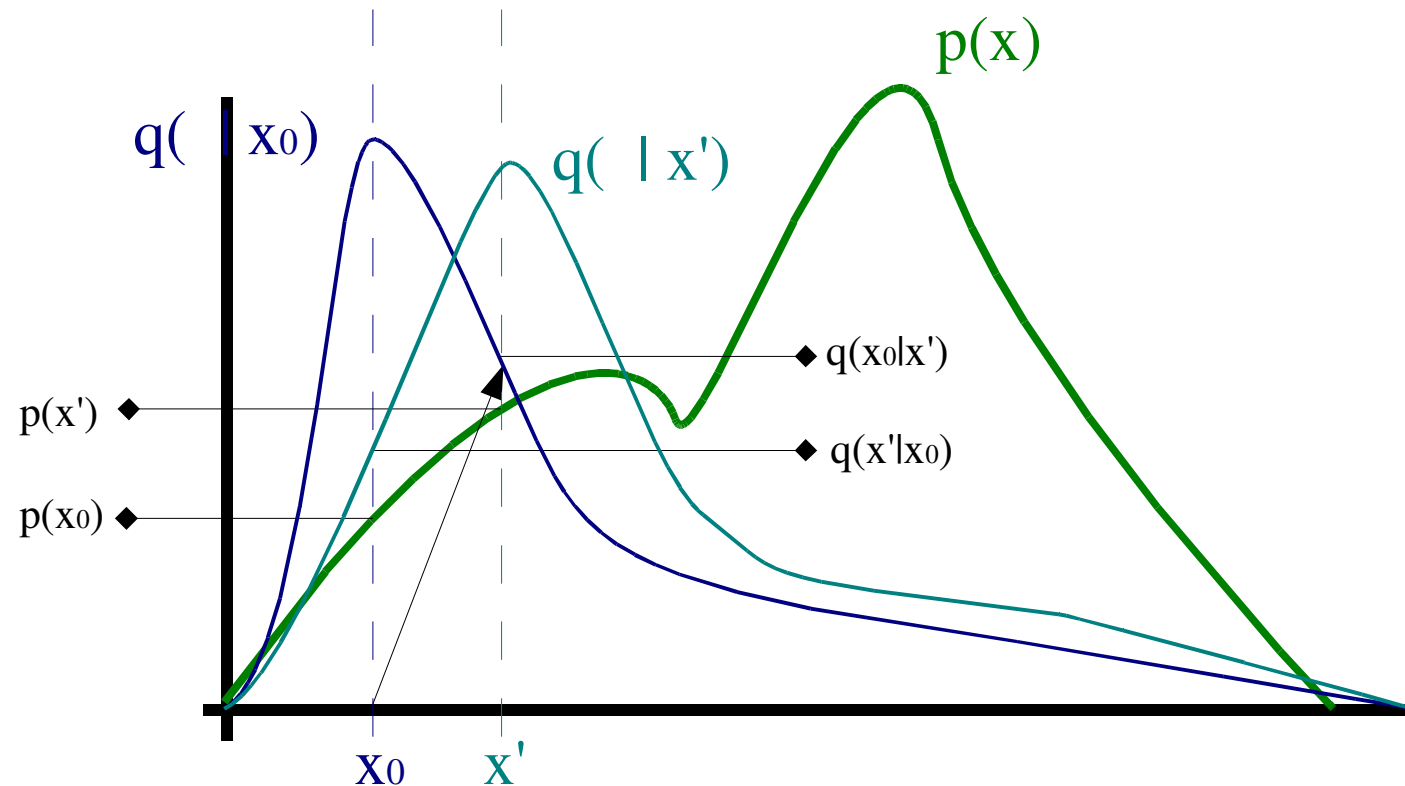
- Introduction to the Bayesian Paradigm
- Background Material
  - Graphical Models
  - Maximum Likelihood
  - Expectation Maximization
- Priors, priors, priors (subjective, conjugate, reference, etc.)
- Inference Problem and Solutions
  - Summing
  - Monte Carlo
  - Markov Chain Monte Carlo
  - Laplace Approximation
  - Variational Approximation
  - Message Passing...
- Survey of Popular Models
- Pointers to Literature
- Conclusions

# Markov Chain Monte Carlo

- Monte Carlo methods suffer because the proposal density needs to be similar to the true density everywhere
- MCMC methods get around this problem by changing the proposal density after each sample
  
- General framework:
  - Choose a proposal density  $q(\cdot | x)$  parameterized by location  $x$
  - Initialize state  $x$  arbitrarily
  - Repeatedly sample by:
    - Propose a new state  $x'$  from  $q(x' | x)$
    - Either accept or reject this new state
      - If accepted, set  $x = x'$
  
- *New problem:* samples are no longer independent!

# Metropolis-Hastings Sampling

- Accept new states with probability:  $\min \left\{ 1, \frac{p(x')}{p(x)} \frac{q(x|x')}{q(x'|x)} \right\}$
- Only put every  $N^{\text{th}}$  sample into  $R$



$$F = \int_{\mathcal{X}} dx p(x) f(x) \approx \frac{1}{R} \sum_{x \in R} f(x)$$



# MH in our Model

- Invent a proposal distribution  $q$

$$\log a' \mid a \sim \text{Nor}(\log(a), 1)$$

$$\log b' \mid b \sim \text{Nor}(\log(b), 1)$$

$$\sigma(\pi') \mid \pi \sim \text{Nor}(\sigma(\pi), 1)$$

$$z_{mn}' \mid z \sim \text{Bin}(0.5)$$

- Or, condition on all variables:

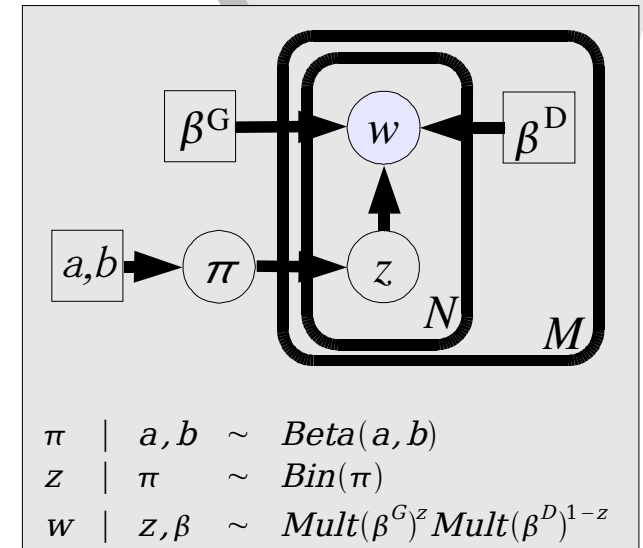
$$\log a' \sim \text{Nor}(\log(a), 1)$$

$$\log b' \sim \text{Nor}(\log(b), 1)$$

$$\sigma(\pi') \sim \text{Beta}(\pi' \mid a, b) \prod_{m,n} \text{Bin}(z_{mn} \mid \pi')$$

$$z_{mn}' \sim \text{Bin}(z_{mn}' \mid \pi) p(w_{mn} \mid z_{mn}', \beta)$$

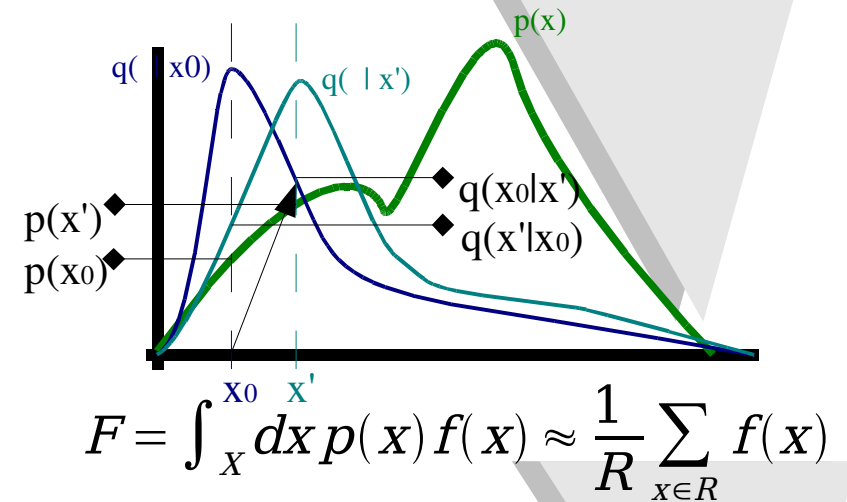
- Now we can compute expectations of  $z$  easily and use these for the M-step of EM



# Metropolis-Hastings Sampling

## ➤ Pros:

- No longer need to specify a universally good proposal distribution; only locally good
- Simple proposal distributions can go far



## ➤ Cons:

- Hard to tell how far to space samples:
  - Suppose we use spherical proposals and, then we need at least

$$N \geq (\sigma_{max} / \sigma_{min})^2$$

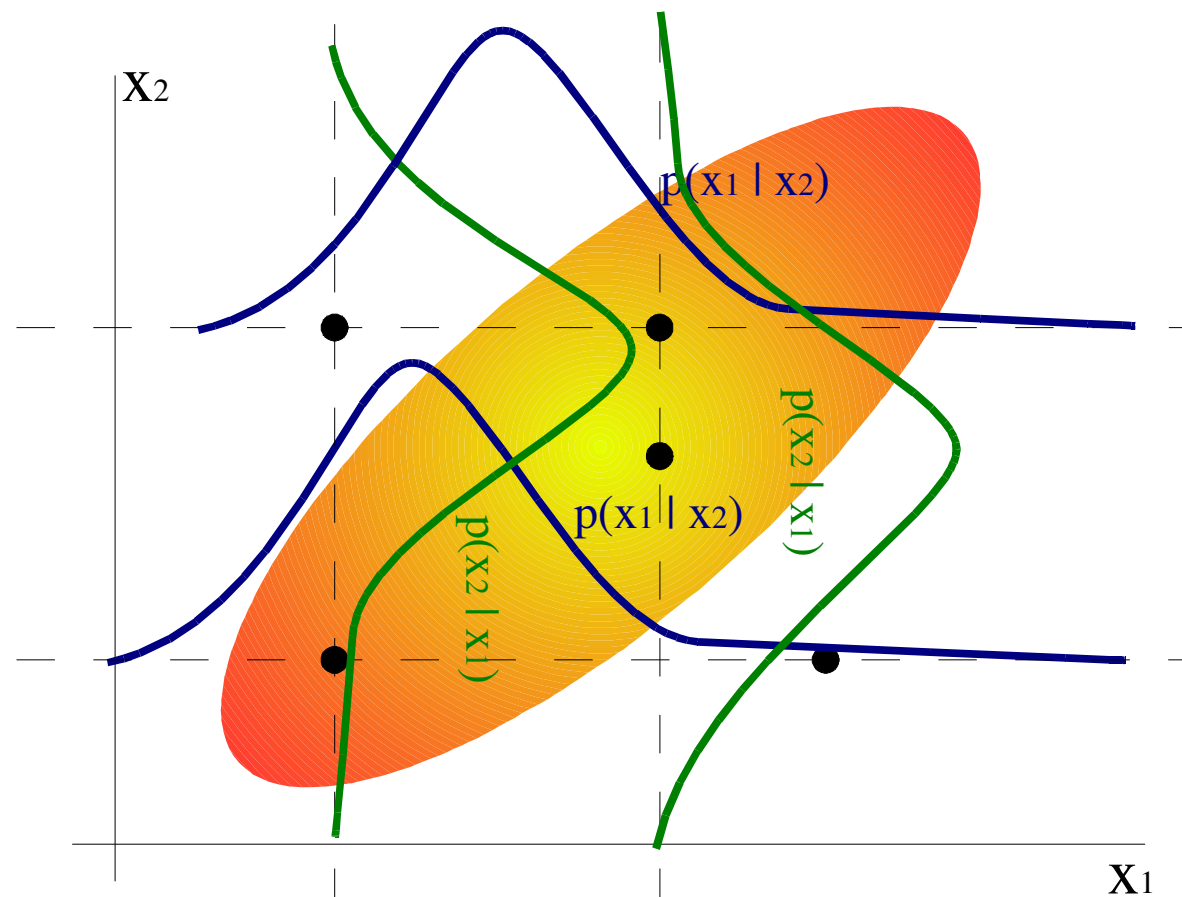
where *sigmas* are lengths of the major density in  $p$

- Auto-correlation to track this:

$$r_k = \frac{\sum_{i=1}^{N-k} (\mathbf{x}_i - \bar{\mathbf{x}})(\mathbf{x}_{i+k} - \bar{\mathbf{x}})}{\sum_{i=1}^N (\mathbf{x}_i - \bar{\mathbf{x}})^2}$$

# Gibbs Sampling

- Defined only for multidimensional problems
- Useful when you can take out one variable and explicitly sample the rest



$$F = \int_{\mathcal{X}} dx p(x) f(x) \approx \frac{1}{R} \sum_{x \in R} f(x)$$

# Gibbs Sampling

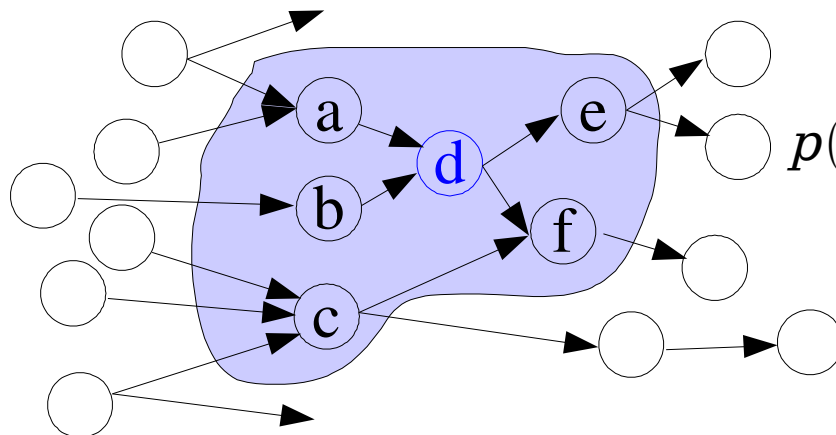
- Typically our params are:  $\bar{\theta} = \langle \theta_1, \dots, \theta_D \rangle$
- If, for each  $i$ , we can draw a sample from:

$$p(\theta_i | \theta_{-i}) = p(\theta_i | \theta_1, \dots, \theta_{i-1}, \theta_{i+1}, \dots, \theta_D)$$

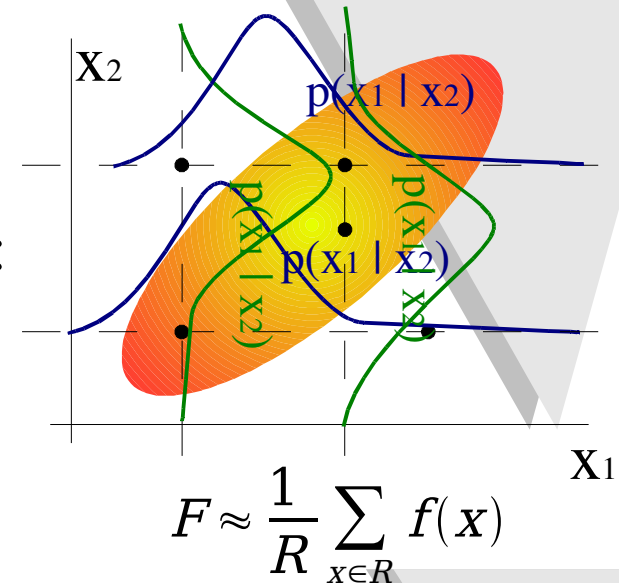
then we *can* use Gibbs sampling

- In graphical models, only depends on the *Markov blanket*:

$$p(\theta_i | \theta_{-i}) = p(\theta_i | \text{par}(\theta_i)) \prod_{j: \theta_j \in \text{par}(\theta_i)} p(\theta_j | \text{par}(\theta_j))$$



$$p(d | \theta_{-d}) = p(d | a, b) p(e | d) p(f | d, c)$$



$$F \approx \frac{1}{R} \sum_{x \in R} f(x)$$

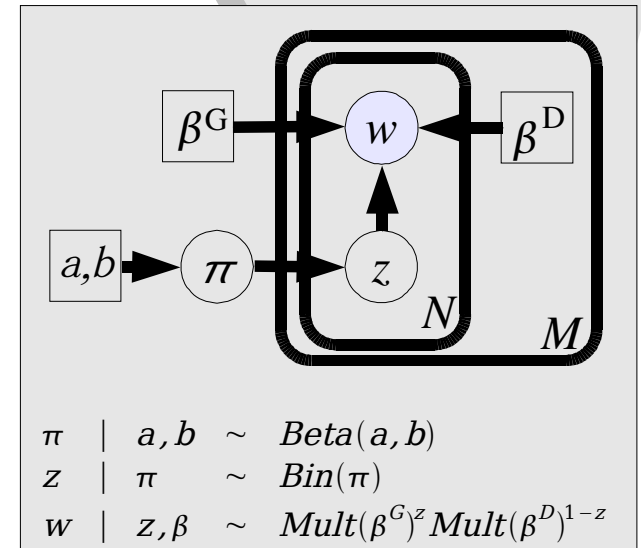
# Gibbs in our Model

- Compute conditional probabilities

$$a, b \mid \neg a, b \sim \text{Beta}(\pi \mid a, b)$$

$$\pi \mid \neg \pi \sim \text{Beta}(\pi \mid a, b) \prod_{m, n} \text{Bin}(z_{mn} \mid \pi)$$

$$Z_{mn} \mid \neg Z_{mn} \sim \text{Bin}(z_{mn} \mid \pi) p(w_{mn} \mid z_{mn}, \beta)$$



- Now we can compute expectations of  $z$  easily and use these for the M-step of EM
- Alternatively, we could propose values for LMs in the sampling

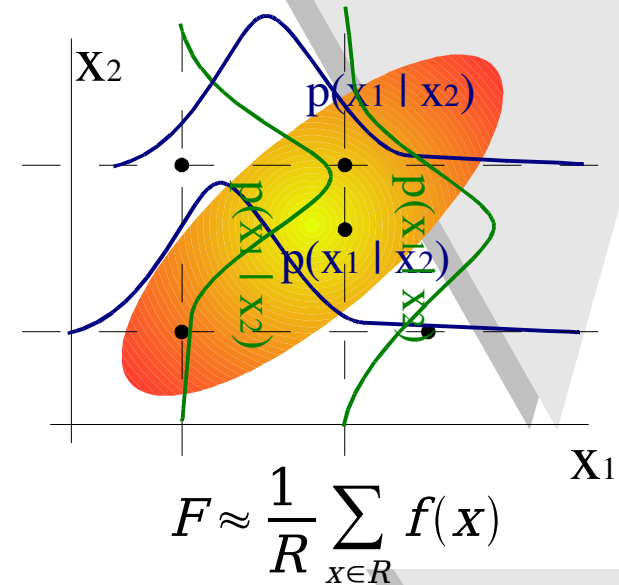
# Gibbs Sampling

## ➤ Pros:

- Designed to work in high dimensional spaces
- Terribly simple to implement
- Automatable

## ➤ Cons:

- Hard to judge convergence, can require many many samples to get an independent one (often worse than MH)
- Only applicable when conditional distributions are 'nice'
  - (Though there are ways around this)

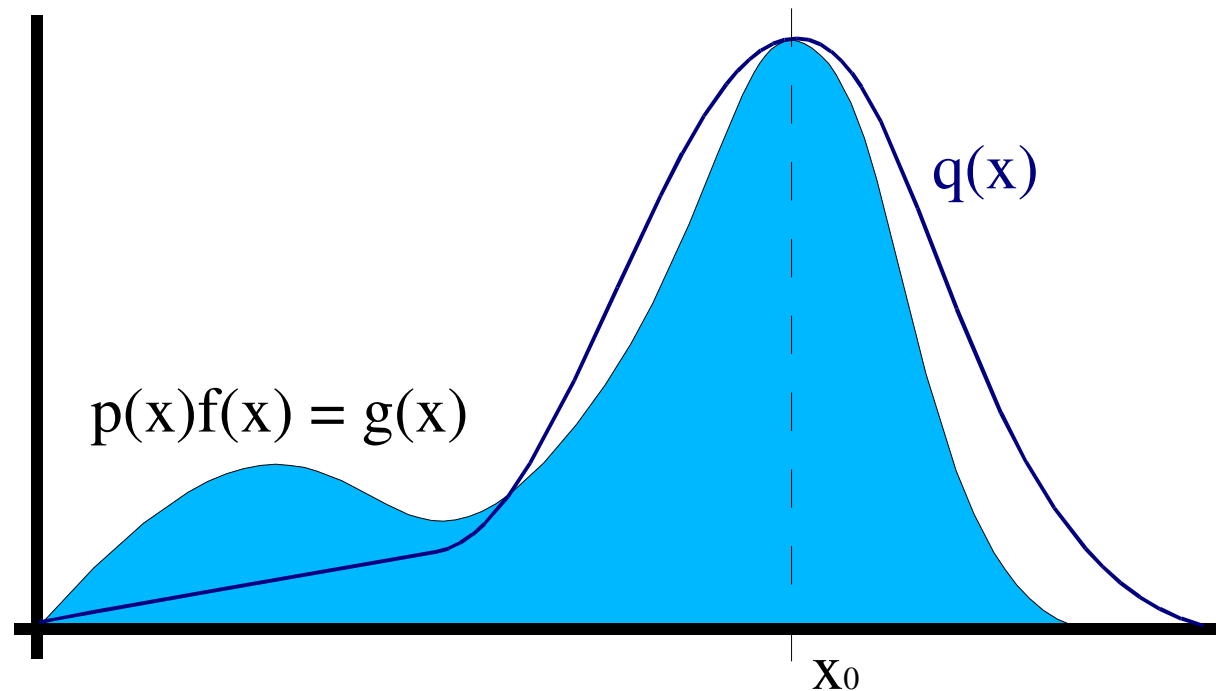


# Tutorial Outline

- Introduction to the Bayesian Paradigm
- Background Material
  - Graphical Models
  - Maximum Likelihood
  - Expectation Maximization
- Priors, priors, priors (subjective, conjugate, reference, etc.)
- Inference Problem and Solutions
  - Summing
  - Monte Carlo
  - Markov Chain Monte Carlo
  - Laplace Approximation
  - Variational Approximation
  - Message Passing...
- Survey of Popular Models
- Pointers to Literature
- Conclusions

# Laplace (Saddlepoint) Approximation

- Idea: approximate the expectation by a quadratic (Taylor expansion) and use the normalizing constant from the resulting Gaussian distribution



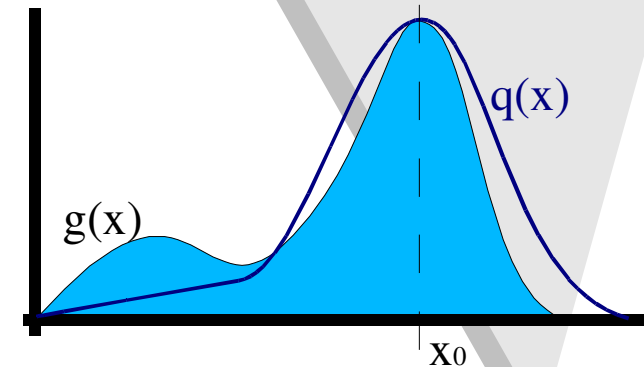
$$F = \int_{\mathbf{x}} d\mathbf{x} p(\mathbf{x}) f(\mathbf{x}) \approx g(\mathbf{x}_0) \sqrt{\frac{2\pi}{c}} \quad , \quad c = - \left[ \frac{\partial^2}{\partial \mathbf{x}^2} \ln g(\mathbf{x}) \right]_{\mathbf{x}=\mathbf{x}_0}$$



# Laplace Approximation

- Find a mode  $\mathbf{x}_0$  of the high-dimensional distribution  $g$
- Approximate  $\ln g(\mathbf{x})$  by a Taylor expansion around this mode:

$$\ln g(\bar{\mathbf{x}}) \approx \ln g(\bar{\mathbf{x}}_0) - \frac{1}{2} (\bar{\mathbf{x}} - \bar{\mathbf{x}}_0)^T \mathbf{A} (\bar{\mathbf{x}} - \bar{\mathbf{x}}_0)$$



$$F \approx g(\mathbf{x}_0) \sqrt{2\pi/c}$$

$$c = -\left[ \partial^2 \ln g(\mathbf{x}) / \partial \mathbf{x}^2 \right]_{\mathbf{x}=\mathbf{x}_0}$$

- Compute the matrix  $\mathbf{A}$  of second derivatives

$$A_{ij} = -\left[ \frac{\partial^2}{\partial x_i \partial x_j} \ln g(\bar{\mathbf{x}}) \right]_{\bar{\mathbf{x}}=\bar{\mathbf{x}}_0}$$

- The exponential form is a Gaussian distribution; use the Gaussian normalizing constant:

$$F = \int_{\mathbb{R}^D} d\mathbf{x} g(\mathbf{x}) \approx g(\bar{\mathbf{x}}_0) \sqrt{\frac{(2\pi)^D}{\det \mathbf{A}}}$$

# Laplace in our Model

➤ Compute second derivatives:

$$\int_U d\pi Z_{ab} \pi^{a-1} (1-\pi)^{b-1} \prod_m \prod_n \sum_{z_{mn}} \pi^{z_{mn}} (1-\pi)^{1-z_{mn}} p(w_{mn} | z_{mn}, \beta)$$

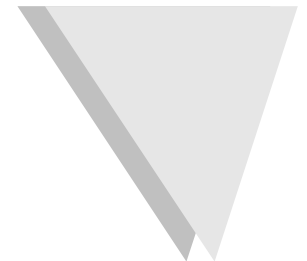
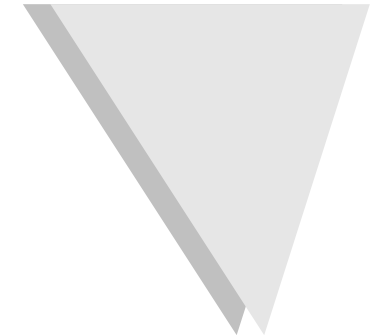
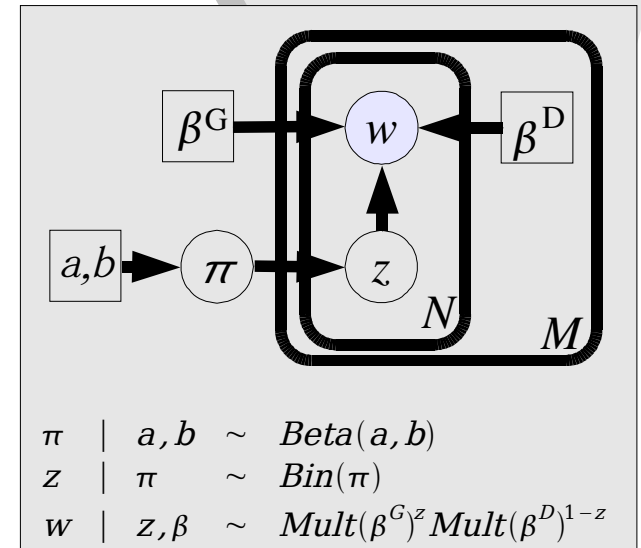
$g(\pi)$

$$\frac{\partial \log g}{\partial \pi} = \frac{(a-1)}{\pi} - \frac{(b-1)}{(1-\pi)} + \sum_{m,n} \left[ \frac{z_{mn}}{\pi} - \frac{(1-z_{mn})}{(1-\pi)} \right]$$

$$\frac{\partial^2 \log g}{\partial \pi^2} = -\frac{(a-1)}{\pi^2} - \frac{(b-1)}{(1-\pi)^2} - \sum_{m,n} \left[ \frac{z_{mn}}{\pi^2} + \frac{(1-z_{mn})}{(1-\pi)^2} \right]$$

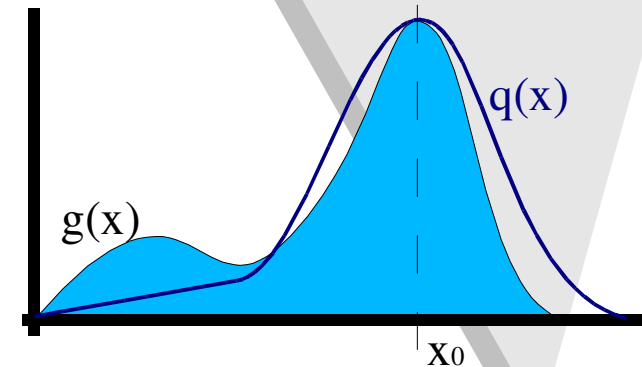
$$\frac{\pi_0}{(1-\pi_0)} = \frac{a-1 + \sum_{m,n} z_{mn}}{b-1 + \sum_{m,n} (1-z_{mn})}$$

$$F = \int_x dx p(x) f(x) \approx g(x_0) \sqrt{\frac{2\pi}{c}} \quad , \quad c = - \left[ \frac{\partial^2}{\partial x^2} \ln g(x) \right]_{x=x_0}$$



# Laplace Approximation

- Pros:
  - Deterministic
  - Efficient if  $\mathbf{A}$  is of a suitable form (i.e., diagonal or block-diagonal)
  - Can apply transformations to make quadratic approximation more reasonable
  
- Cons:
  - Poor fit for multimodal distributions
  - Often,  $\det \mathbf{A}$  cannot be found efficiently



$$F \approx g(x_0) \sqrt{2\pi/c}$$

$$c = -\left[\partial^2 \ln g(x) / \partial x^2\right]_{x=x_0}$$

# Tutorial Outline

- Introduction to the Bayesian Paradigm
- Background Material
  - Graphical Models
  - Maximum Likelihood
  - Expectation Maximization
- Priors, priors, priors (subjective, conjugate, reference, etc.)
- **Inference Problem and Solutions**
  - Summing
  - Monte Carlo
  - Markov Chain Monte Carlo
  - Laplace Approximation
  - **Variational Approximation**
  - Message Passing...
- Survey of Popular Models
- Pointers to Literature
- Conclusions

# Variational Approximation

- Basic idea: replace intractable  $p$  with tractable  $q$
- Old Problem:
  - We cannot come up with a good, single,  $q$  to approximate  $p$
- Key Idea:
  - Consider a *family* of distributions  $Q = \{q(\cdot | \phi) : \phi \in \Phi\}$  with 'variational parameters'  $\phi$
  - Choose a member  $q$  from  $Q$  that is closest to  $p$
- New problems:
  - How do we choose  $Q$ ?
  - How do we measure 'closeness' between  $q$  and  $p$ ?

# Recall EM and Jensen's Inequality

- Jensen gives us:

$$\begin{aligned}
 \log p(\mathbf{x} | \theta) &= \log \int_{\mathbf{z}} d\mathbf{z} p(\mathbf{x}, \mathbf{z} | \theta) \\
 &= \log \int_{\mathbf{z}} d\mathbf{z} q(\mathbf{z}) \frac{p(\mathbf{x}, \mathbf{z} | \theta)}{q(\mathbf{z})} \\
 &\geq \int_{\mathbf{z}} d\mathbf{z} q(\mathbf{z}) \log \frac{p(\mathbf{x}, \mathbf{z} | \theta)}{q(\mathbf{z})} \\
 &= \int_{\mathbf{z}} q(\mathbf{z}) \log p(\mathbf{x}, \mathbf{z} | \theta) - \int_{\mathbf{z}} q(\mathbf{z}) \log q(\mathbf{z}) \\
 &= \underbrace{\mathbf{E}_{\mathbf{z} \sim q} \{ \log p(\mathbf{x}, \mathbf{z} | \theta) \} - \mathbf{E}_{\mathbf{z} \sim q} \{ \log q(\mathbf{z}) \}}_{\mathcal{L}(\mathbf{x} | \theta)}
 \end{aligned}$$

- Where we chose  $q(\mathbf{z}) = p(\mathbf{z} | \mathbf{x}, \theta)$  to turn the inequality into an equality. But we can also compute:

$$\log p(\mathbf{x} | \theta) = \mathcal{L} + KL(q(\mathbf{z}) || p(\mathbf{z} | \mathbf{x}, \theta))$$

for *any* choice of  $q$

# Variational EM

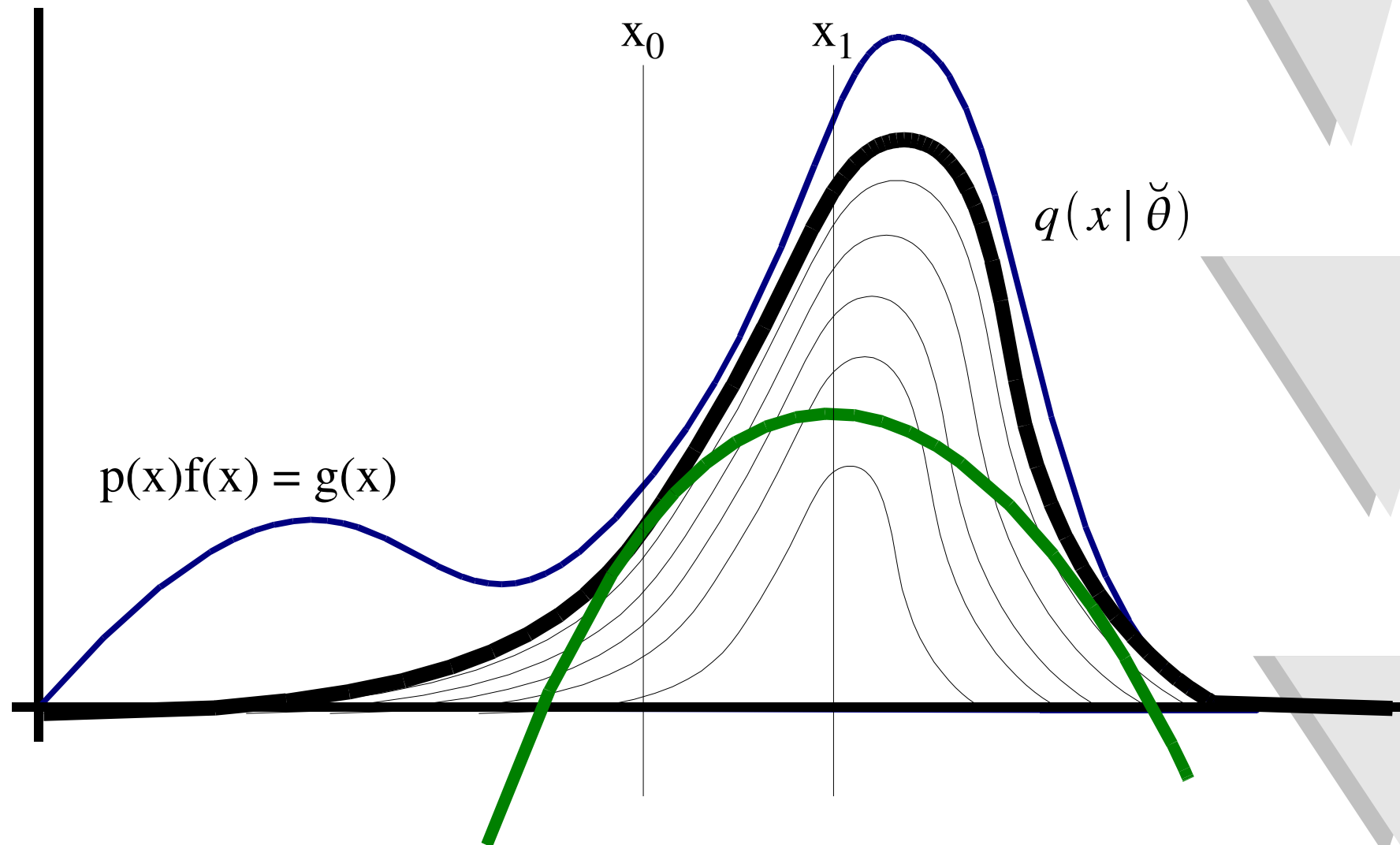
- Parameterize  $q$  and directly optimize:

$$\log p(x | \theta) = \mathbf{E}_{z \sim q} \{ \log p(x, z | \theta) \} - \mathbf{E}_{z \sim q} \{ \log q(z) \} + KL(q(z | \check{\theta}) || p(z | x, \theta))$$

- Iterate:

- **V-Step:** Compute variational parameters  $\check{\theta}$  to minimize KL
  - **E-Step:** Compute expectations of hidden variables wrt  $q(\check{\theta})$
  - **M-Step:** Maximize  $\mathcal{L}$  wrt true parameters  $\theta$
- 
- Art: inventing  $q$  so that this is all tractable

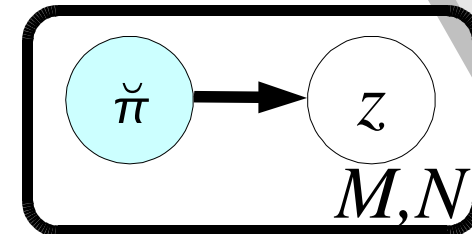
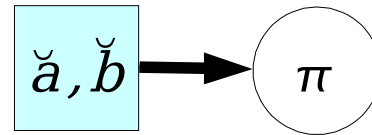
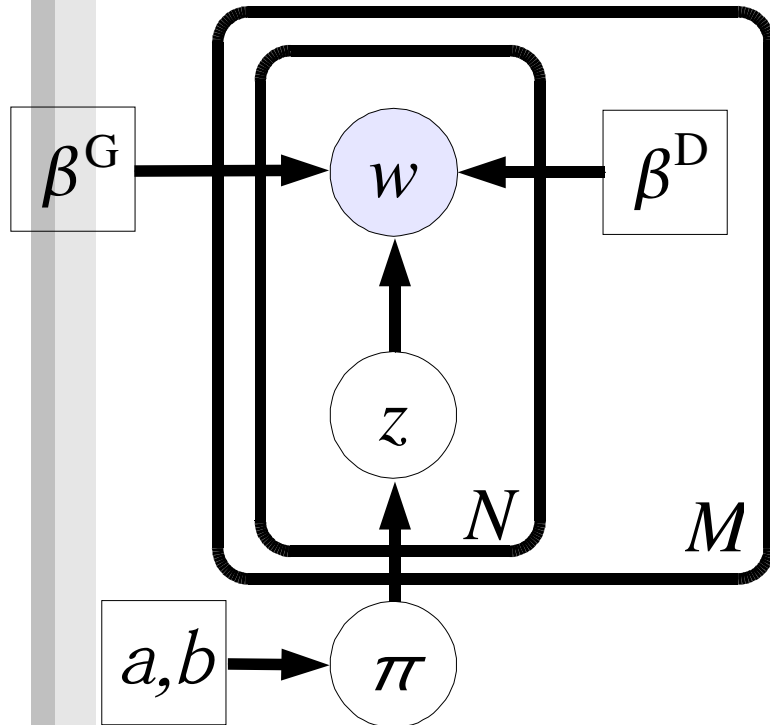
# Variational EM in Pictures





# Variational: Choosing Q

➤ Mixture model:



$$\begin{aligned} \pi &| \check{a}, \check{b} \sim \text{Beta}(\check{a}, \check{b}) \\ z &| \check{\pi} \sim \text{Bin}(\check{\pi}) \end{aligned}$$

$$q(\pi, z | \check{a}, \check{b}, \check{\pi}) = \frac{\Gamma(\check{a} + \check{b})}{\Gamma(\check{a})\Gamma(\check{b})} \prod_{m,n} \prod_i \pi_{mni}^{\check{a}_i - 1} \check{\pi}_{mn}^{z_{nmi}}$$

$$\begin{aligned} p(w, \pi, z | \rho, a, b) = & \frac{\Gamma(a+b)}{\Gamma(a)\Gamma(b)} \pi^{a-1} (1-\pi)^{b-1} \\ & \prod_{m,n} \pi^{z_{nm}} (1-\pi)^{1-z_{nm}} \prod_i \prod_v \left[ \beta_v^i \right]^{z_{mni}} \end{aligned}$$

Key:  $\pi$  and  $z$  are now not tied in the  $q$  distribution!

# VEM in our Model

- Iterate:
  - Optimize variational parameters:

$$\check{\pi}_{mni} \propto \exp[\check{E}_i + \omega_{mni}]$$

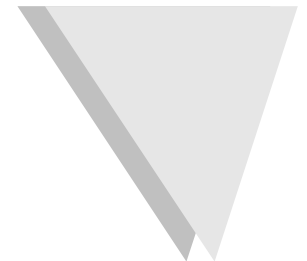
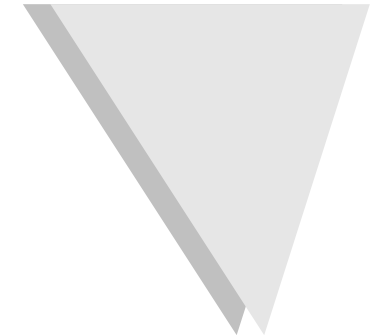
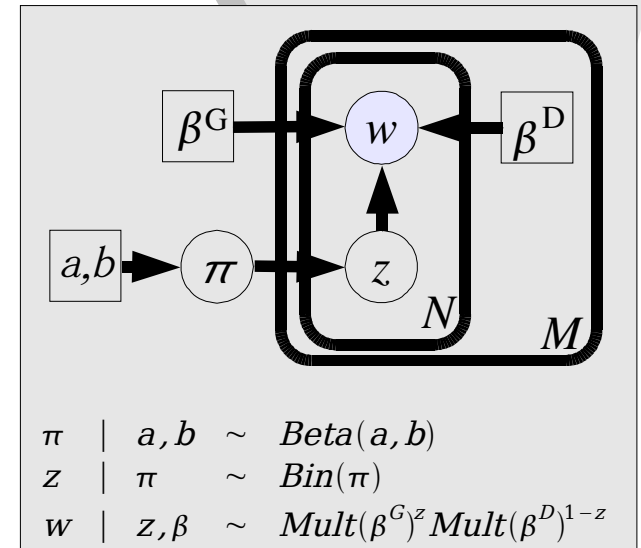
$$\check{a}_i = a_i + \sum_{m,n} \check{\pi}_{mni}$$

$$\check{E}_i = \Psi(\check{a}_i) - \Psi\left(\sum_i \check{a}_i\right) \quad \omega_{mni} = \sum_j w_{mnj} \log \beta_j^i$$

- Optimize model parameters:

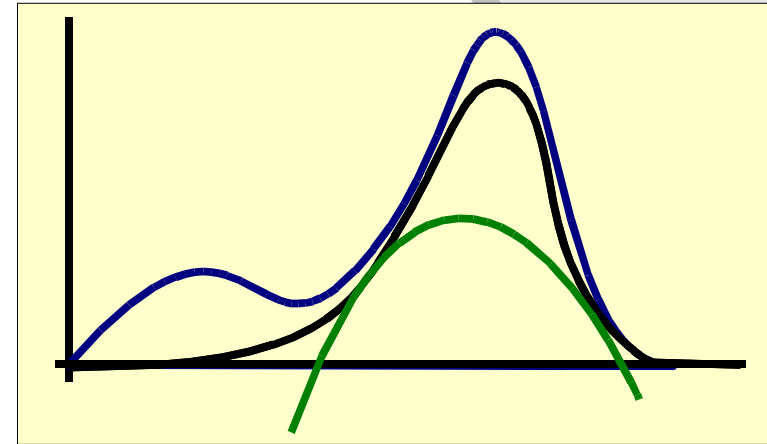
$$\beta_v^i \propto \sum_{m,n} \check{\pi}_{mni} w_{mnv}$$

$a, b \sim$  generic optimization techniques



# Variational EM Summed Up

- Steps:
  - Write down conditional likelihood and choose an approximating distribution (eg, by factoring everything) with variational parameters
  - Iterate between optimizing the VPs and model parameters
  
- Pros:
  - Efficient, deterministic, often quite accurate
- Cons:
  - At it's heart, still a mode-based technique
  - Often underestimates the spread of a distribution
  - Approximation is *local*

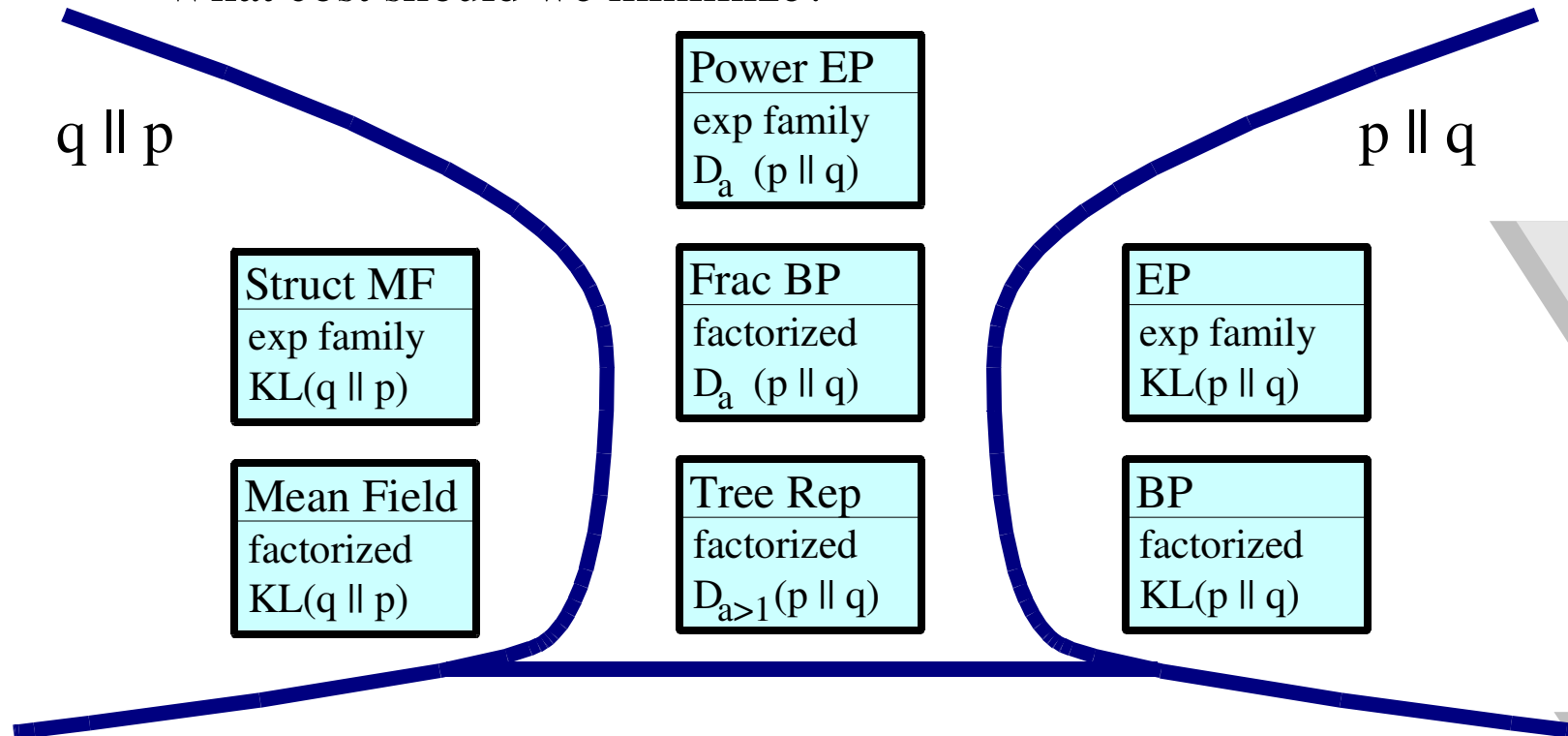


# Tutorial Outline

- Introduction to the Bayesian Paradigm
- Background Material
  - Graphical Models
  - Maximum Likelihood
  - Expectation Maximization
- Priors, priors, priors (subjective, conjugate, reference, etc.)
- Inference Problem and Solutions
  - Summing
  - Monte Carlo
  - Markov Chain Monte Carlo
  - Laplace Approximation
  - Variational Approximation
  - Message Passing...
- Survey of Popular Models
- Pointers to Literature
- Conclusions

# Message Passing Algorithms

- Two major choices:
  - What approximating distribution should we use?
  - What cost should we minimize?



$$D_a(p \parallel q) = \frac{1}{\beta(1-\beta)} \int dx \beta p^{x+(1-\beta)q} - p^\beta q^\beta$$

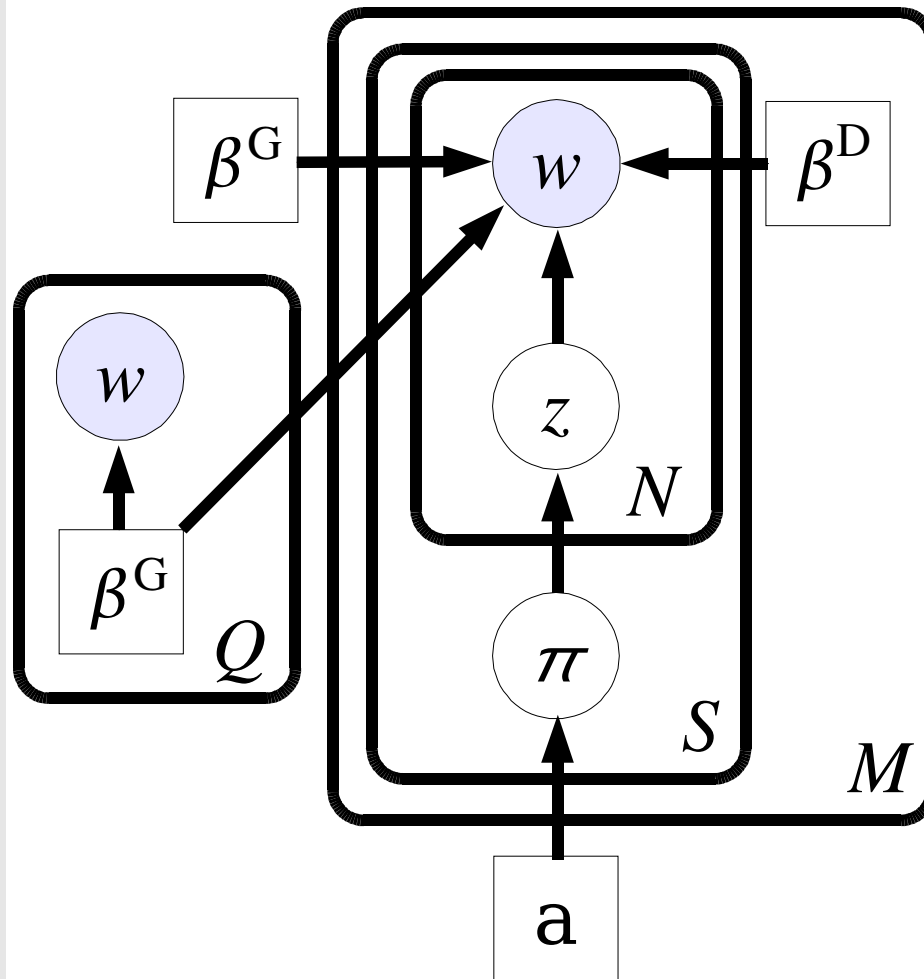
$$\beta = \frac{1}{2}(1+a)$$

$$KL(p \parallel q) = D_1(p \parallel q)$$

$$KL(q \parallel p) = D_{-1}(p \parallel q)$$

# Empirical Evaluation of Methods

- Query-focused summarization model:



$$\begin{aligned}
 w_{qn}^Q &\sim \text{Mult}(\beta_q^Q) \\
 \pi_{ms} &\sim \text{Dir}(a) \\
 z_{msn} &\sim \text{Mult}(\pi_{ms}) \\
 w_{msn} &\sim \text{Mult}(\beta^G)^{z_{msn1}}
 \end{aligned}$$

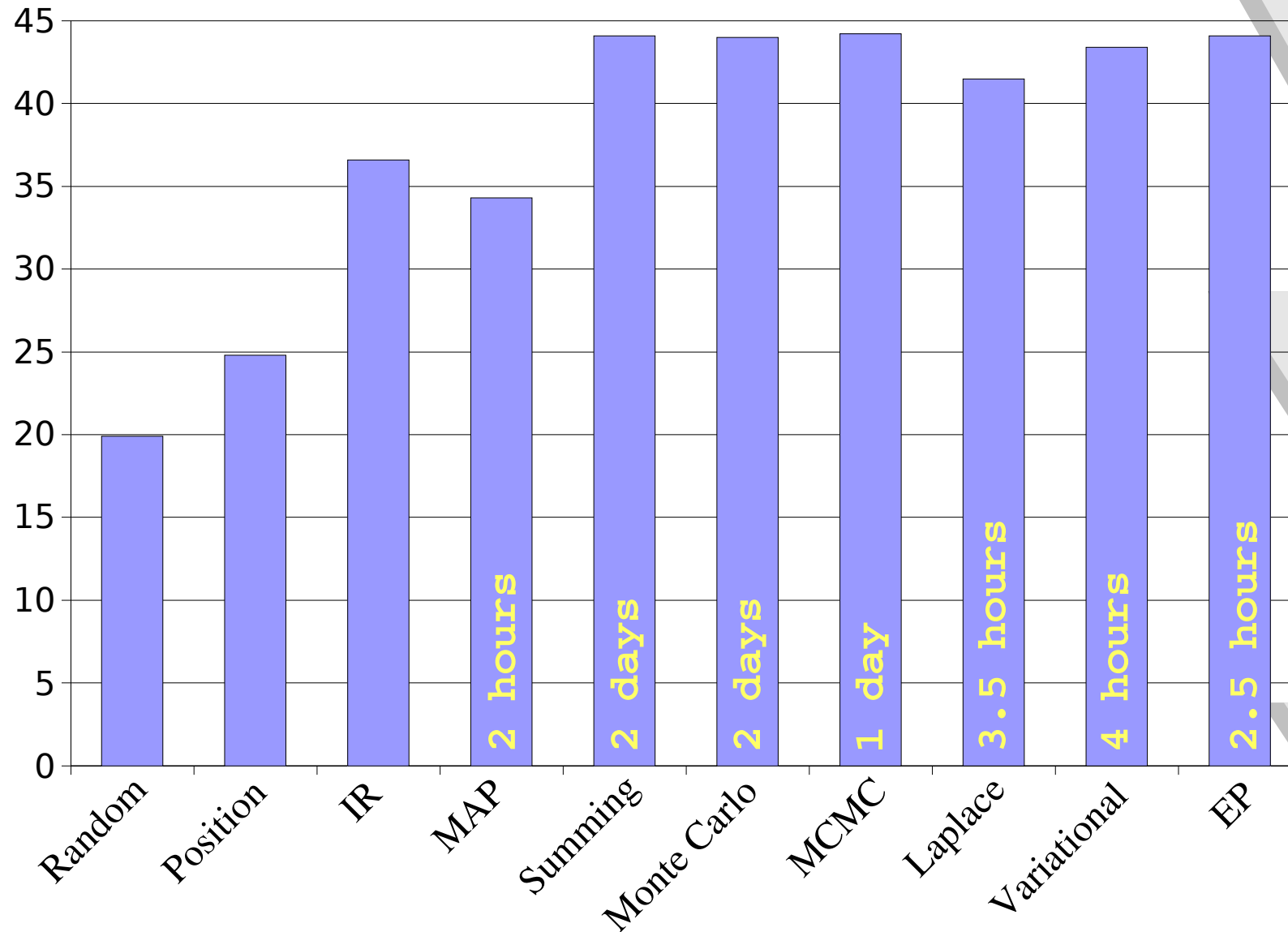
$$\prod_m \text{Mult}(\beta_m^D)^{z_{msn}(m+1)} \\
 \prod_q \text{Mult}(\beta_q^Q)^{z_{msn}(q+M+1)}$$

# Evaluation Data

- All TREC data
  - Queries 51-350 and 401-450 (35k words)
  - All relevant documents (43k docs, 2.1m sents, 65.8m words)
  - Asked 7 annotators to select up to 4 sentences for an extract
    - Each annotated 25 queries (166 total)
  - Systems produce ranked lists of sentences
    - Compared on mean average precision, mean reciprocal rank and precision at 2
  
- Computation Time:
  - MAP-EM (2 hours)
  - Summing (2 days)
  - Monte Carlo (2 days)
  - MCMC (1 day)
  - Laplace (5 hours)
  - Variational (4 hours)
  - EP (2.5 hours)

# Evaluation Results

## Mean Average Precision





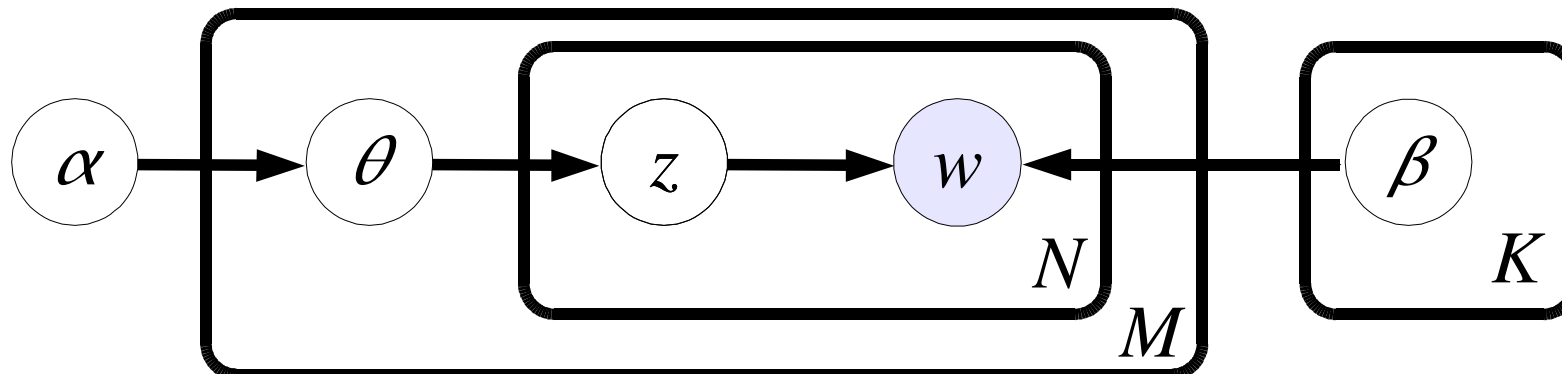
# Tutorial Outline

- Introduction to the Bayesian Paradigm
- Background Material
  - Graphical Models
  - Maximum Likelihood
  - Expectation Maximization
- Priors, priors, priors (subjective, conjugate, reference, etc.)
- Inference Problem and Solutions
  - Summing
  - Monte Carlo
  - Markov Chain Monte Carlo
  - Laplace Approximation
  - Variational Approximation
  - Message Passing...
- **Survey of Popular Models**
- Pointers to Literature
- Conclusions

# Latent Dirichlet Allocation

[Blei, Ng + Jordan, JMLR 03]

- Unigram model of documents
- Each document is a *mixture* over topics
- Each topic is a *mixture* over words
- **Generative model for each document (M total):**
  - Choose a single topic mixture:  $\theta \sim \text{Dir}(\alpha)$
  - For each word (N total):
    - Choose a topic for this word:  $z \sim \text{Mult}(\theta)$
    - Choose the word itself:  $w \sim \text{Mult}(\beta^z)$

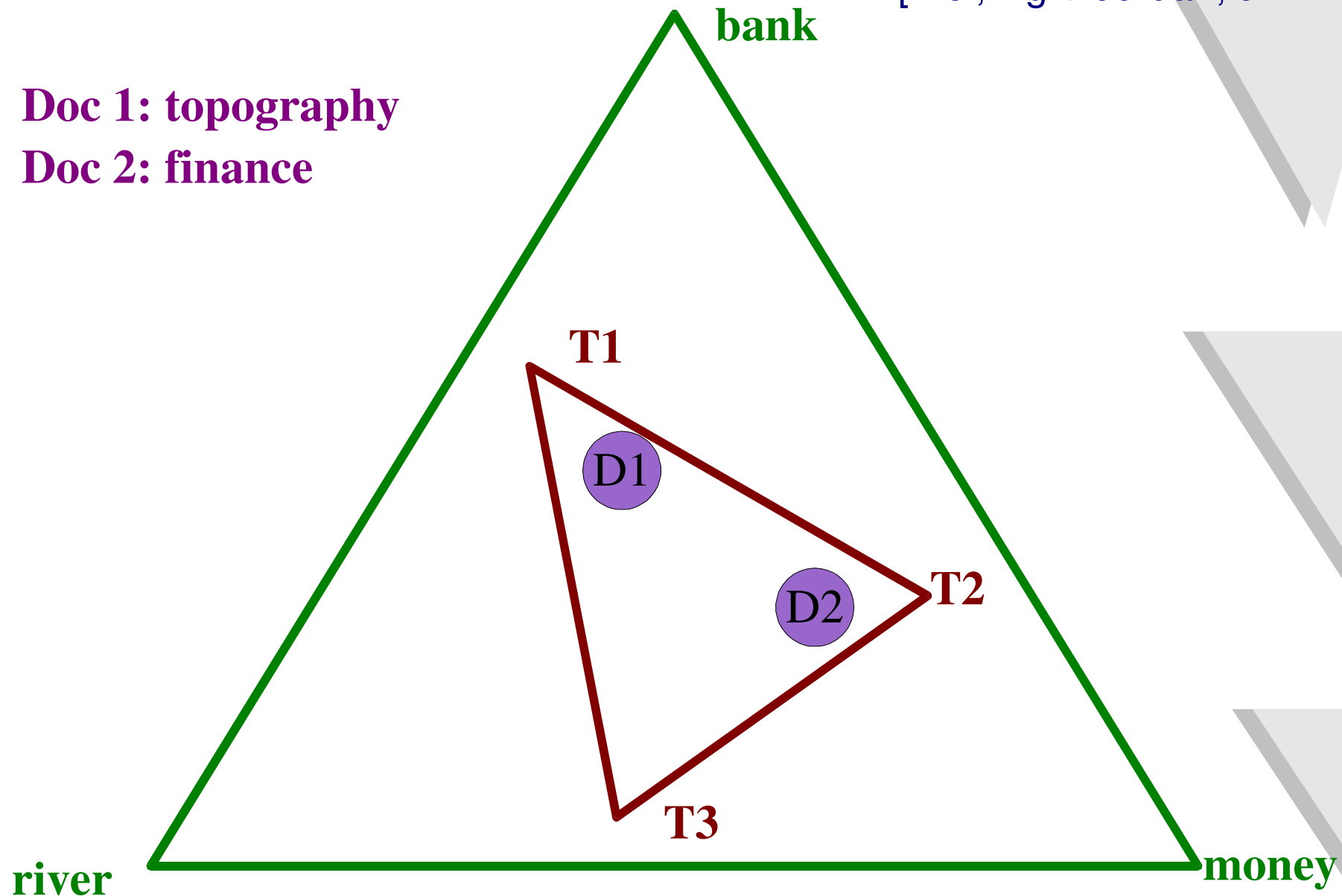


# LDA: Geometric Interpretation

[Blei, Ng + Jordan, JMLR 03]

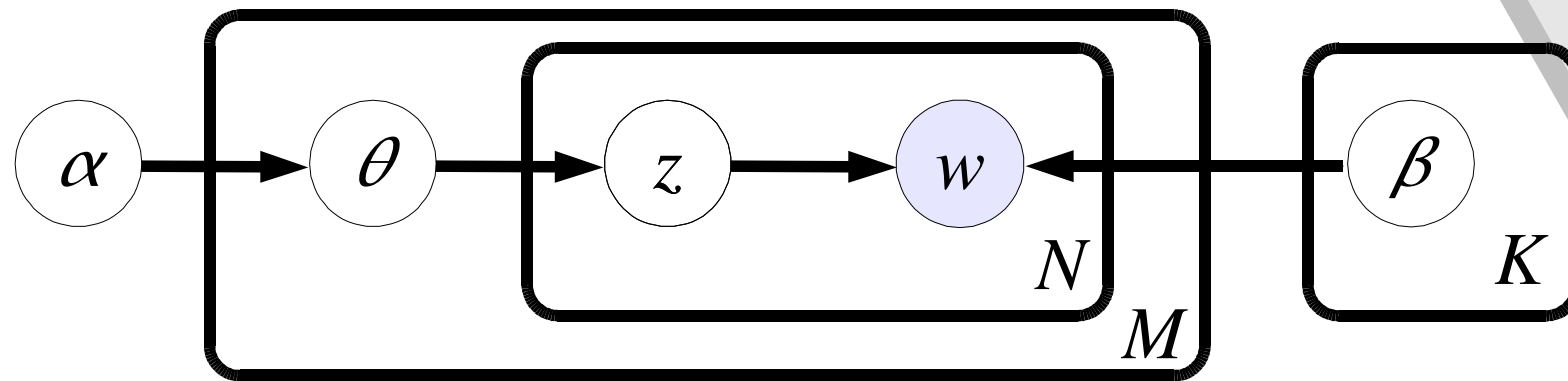
Doc 1: topography

Doc 2: finance



# LDA: Inference

[Blei, Ng + Jordan, JMLR 03]



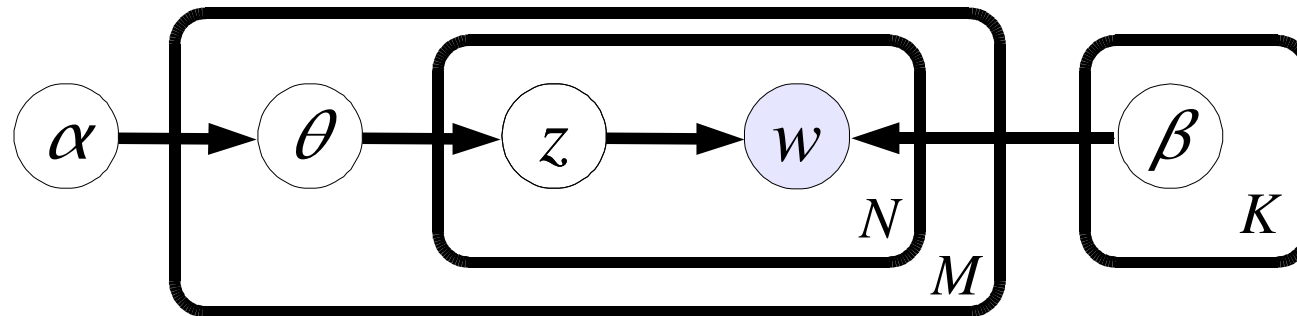
$$P(D) = \int_{\Delta_V} dP(\beta) \int_{\mathbb{R}^+} dP(\alpha) \prod_{m=1}^M \int_{\Delta_K} d\theta \left[ \frac{\Gamma(K\alpha)}{\Gamma(\alpha)^K} \prod_{j=1}^K \theta_j^{\alpha-1} \right]$$

$$\prod_{n=1}^N \sum_{z_{mn}=1}^K \prod_{i=1}^{|V|} \prod_{j=1}^K \beta_{ji} \mathbf{1}[w_{mn}=i] \mathbf{1}[z_{mn}=j]$$

**Desired: either  $\beta$ s or  $z$ s**

# LDA: Naïve Gibbs Sampler

[Griffiths + Tenenbaum, CogSci 03]



$$\alpha \sim P(\alpha) \prod_m \text{Dir}(\theta_m | \alpha)$$

$$\beta_j \sim P(\beta_j) \prod_{mn} \text{Mult}(w_{mn} | \beta_j) \mathbf{1}[z_{mn}=j]$$

$$\theta_m \sim \text{Dir}(\theta_m | \alpha) \prod_n \text{Mult}(z_{mn} | \theta_m)$$

$$z_{mn} \sim \text{Mult}(z_{mn} | \theta_m) \text{Mult}(w_{mn} | \beta_{z_{mn}})$$

Can collapse this step!

# LDA Results

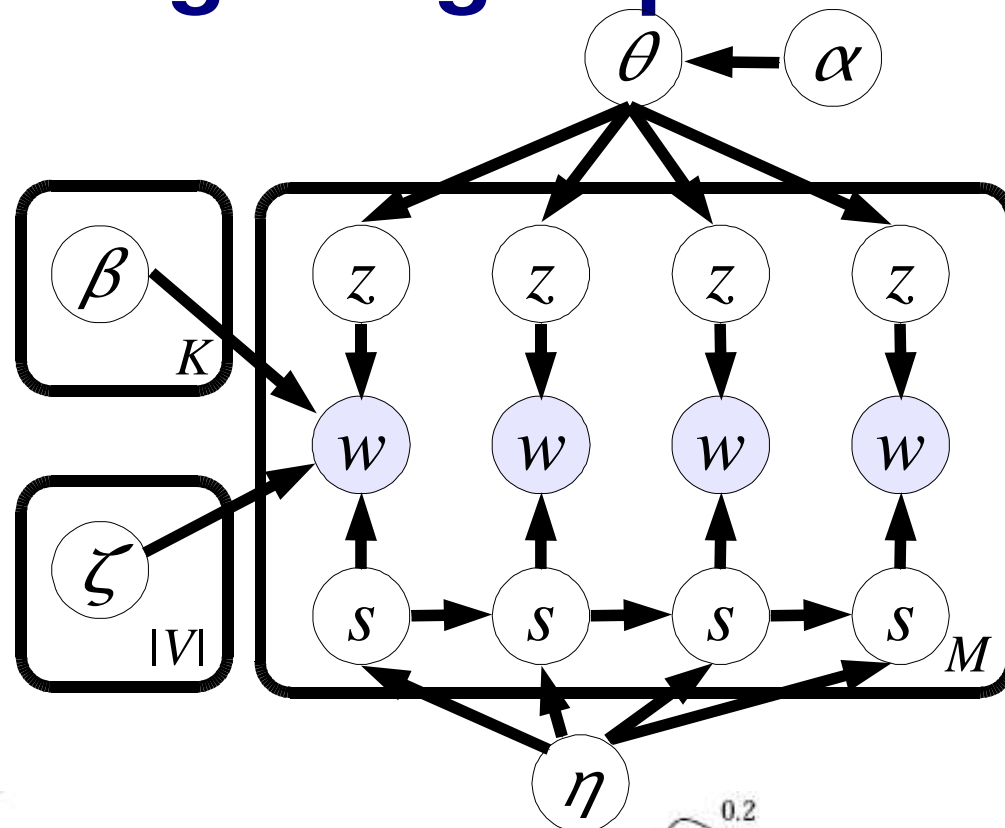
[Blei, Ng + Jordan, JMLR 03]

"Arts"	"Budgets"	"Children"	"Education"
NEW	MILLION	CHILDREN	SCHOOL
FILM	TAX	WOMEN	STUDENTS
SHOW	PROGRAM	PEOPLE	SCHOOLS
MUSIC	BUDGET	CHILD	EDUCATION
MOVIE	BILLION	YEARS	TEACHERS
PLAY	FEDERAL	FAMILIES	HIGH
MUSICAL	YEAR	WORK	PUBLIC
BEST	SPENDING	PARENTS	TEACHER
ACTOR	NEW	SAYS	BENNETT
FIRST	STATE	FAMILY	MANIGAT
YORK	PLAN	WELFARE	NAMPHY
OPERA	MONEY	MEN	STATE
THEATER	PROGRAMS	PERCENT	PRESIDENT
ACTRESS	GOVERNMENT	CARE	ELEMENTARY
LOVE	CONGRESS	LIFE	HAITI

The William Randolph Hearst Foundation will give \$1.25 million to Lincoln Center, Metropolitan Opera Co., New York Philharmonic and Juilliard School. "Our board felt that we had a real opportunity to make a mark on the future of the performing arts with these grants an act every bit as important as our traditional areas of support in health, medical research, education and the social services," Hearst Foundation President Randolph A. Hearst said Monday in announcing the grants. Lincoln Center's share will be \$200,000 for its new building, which will house young artists and provide new public facilities. The Metropolitan Opera Co. and New York Philharmonic will receive \$400,000 each. The Juilliard School, where music and the performing arts are taught, will get \$250,000. The Hearst Foundation, a leading supporter of the Lincoln Center Consolidated Corporate Fund, will make its usual annual \$100,000 donation, too.

# Integrating Topics and Syntax

[Griffiths, Steyvers, Blei +  
Tenenbaum, NIPS 2004]



For each document  $M$ :  
Choose a topic mixture

For each word  $N$ :

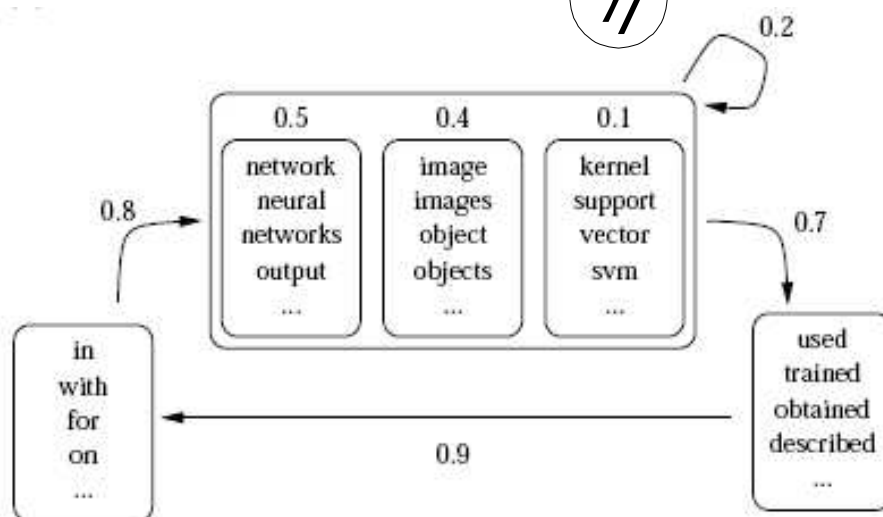
Choose topic  $z$

Choose class  $s$

Choose  $w$  from:

$\beta_z$  if  $s=0$

$\zeta_s$  otherwise



network used for images

image obtained with kernel

output described with objects

neural network trained with svm images

# LDA versus Topics+Syntax

LDA

the	the	the	the	the	a	the	the	the
blood	,	,	of	a	the	,	,	,
,	and	and	,	of	of	of	a	a
of	of	of	to	,	,	a	of	in
body	a	in	in	in	in	and	and	game
heart	in	land	and	to	water	in	drink	ball
and	trees	to	classes	picture	is	story	alcohol	and
in	tree	farmers	government	film	and	is	to	team
to	with	for	a	image	matter	to	bottle	to
is	on	farm	state	lens	are	as	in	play

Topics

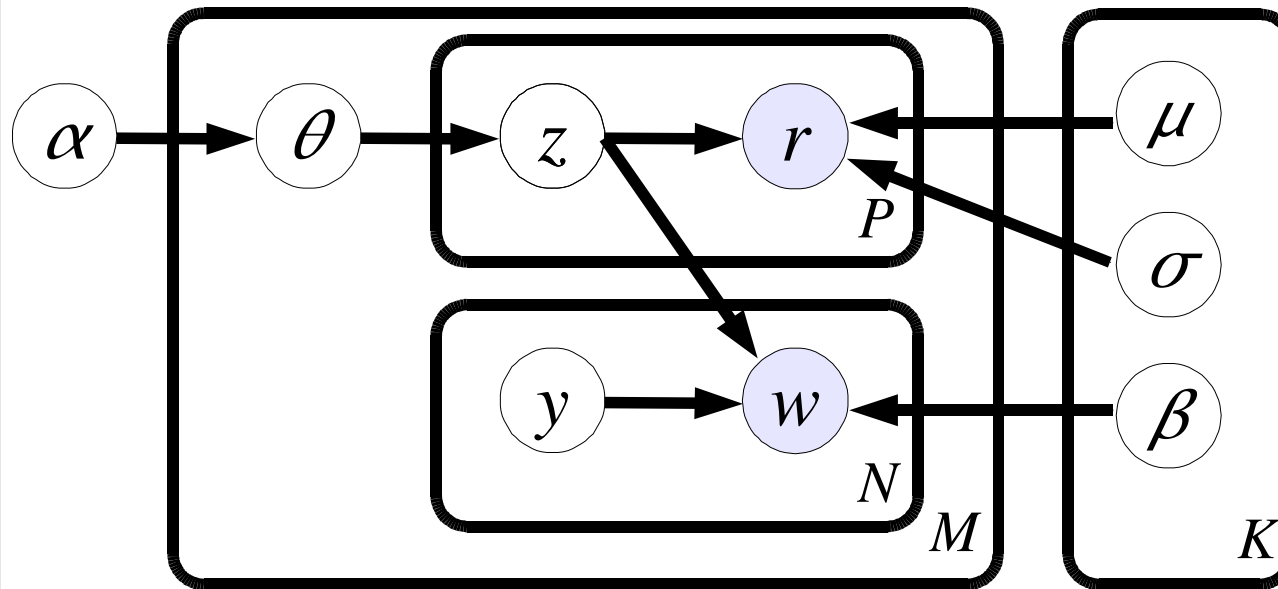
blood	forest	farmers	government	light	water	story	drugs	ball
heart	trees	land	state	eye	matter	stories	drug	game
pressure	forests	crops	federal	lens	molecules	poem	alcohol	team
body	land	farm	public	image	liquid	characters	people	*
lungs	soil	food	local	mirror	particles	poetry	drinking	baseball
oxygen	areas	people	act	eyes	gas	character	person	players
vessels	park	farming	states	glass	solid	author	effects	football
arteries	wildlife	wheat	national	object	substance	poems	marijuana	player
*	area	farms	laws	objects	temperature	life	body	field
breathing	rain	corn	department	lenses	changes	poet	use	basketball

Syntax

the	in	he	*	be	said	can	time	,
a	for	it	new	have	made	would	way	;
his	to	you	other	see	used	will	years	(
this	on	they	first	make	came	could	day	:
their	with	i	same	do	went	may	part	)
these	at	she	great	know	found	had	number	
your	by	we	good	get	called	must	kind	
her	from	there	small	go		do	place	
my	as	this	little	take		have		
some	into	who	old	find		did		



# Matching Words and Pictures



[Barnard, Duygulu, de Freitas, Forsyth, Blei + Jordan, JMLR 2003]

1. People, tree
2. Sky, jet
3. Sky, clouds
4. Sky, mountain
5. Plane, jet
6. Plane, jet

For each image/caption pair  $M$

Draw a topic mixture  $\theta \sim \text{Dir}(\alpha)$

For each image region  $P$

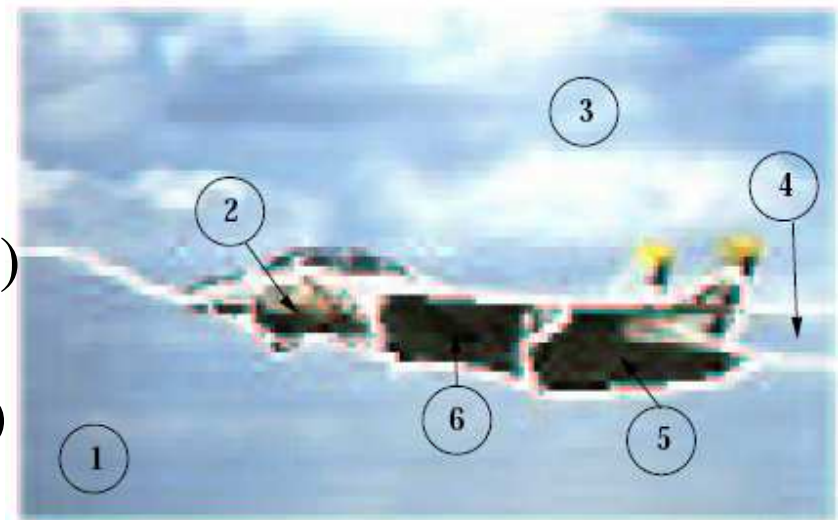
Draw a topic  $z \sim \text{Mult}(\theta)$

Draw the region  $r \sim \text{Gaussian}(\mu, \sigma^2)$

For each word  $N$

Draw a image region  $y \sim \text{Unif}(1..P)$

Draw the word  $w \sim \text{Mult}(\beta_{zy})$



# Matching Words and Pictures

[Barnard. Duygulu, de Freitas,  
Forsyth, Blei + Jordan, JMLR 2003]



**True caption**

market people

**Corr-LDA**

people market pattern textile display



**True caption**

scotland water

**Corr-LDA**

scotland water flowers hills tree



**True caption**

sky tree water

**Corr-LDA**

tree water sky people buildings



**True caption**

birds tree

**Corr-LDA**

birds nest leaves branch tree



**True caption**

fish reefs water

**Corr-LDA**

fish water ocean tree coral



**True caption**

clouds jet plane

**Corr-LDA**

sky plane jet mountain clouds

# Conclusions

- Bayesian methods provide efficient, effective models
- Graphical models are an easy language
- Plug and play of Multinomial/Dirichlet/Beta/Gamma leads to models that admit efficient Gibbs sampling methods
- For faster inference, the variational approximation is effective
- Bayesian models of text problems is largely unexplored
- Many topics not discussed:
  - Alternative inference techniques (belief/expectation propagation)
  - Classifiers/discriminative models (Gaussian Processes  $\approx$  SVMs)
  - Infinite models (Dirichlet Processes, Chinese Restaurant Processes)

# Bayes in Action (NLP/IR/Text)

Blei, Ng + Jordan, *Latent Dirichlet allocation*, JMLR03.

Barnard, Duygulu, de Freitas, Forsyth, Blei + Jordan. *Matching words and pictures*. JMLR03.

Daumé III + Marcu, *Bayesian Query-Focused Summarization*, ACL06.

Griffiths, Steyvers, Blei, Tenenbaum, *Integrating topics and syntax*. NIPS04.

McCallum, Corrada-Emmanuel + Wang, *Topic and Role Discovery in Social Networks*. IJCAI05.

Zhang, Callan + Minka, *Novelty and Redundancy Detection in Adaptive Filtering*. SIGIR02.

# For Further Information (Books)

James O. Berger, *Statistical Decision Theory and Bayesian Analysis*.  
Springer, 1985.

David MacKay, *Information Theory, Inference and Learning Algorithms*. Cambridge University Press, 2003.

Larry Wasserman, *All of Statistics: A Concise Course in Statistical Inference*. Springer, 2003.

Christopher Bishop, *Pattern Recognition and Machine Learning*.  
Springer, 2006.

# For Further Information (Tutorials)

Andrieu, de Freitas, Doucet + Jordan, *An Introduction to MCMC for Machine Learning*. ML 2003

Wainwright + Jordan, *Graphical models, exponential families and variational inference*. UCB Stat TR#649, 2003.

Murphy, *A Brief Introduction to Graphical Models and Bayesian Networks*. [www.cs.ubc.ca/~murphyk/Bayes/bayes.html](http://www.cs.ubc.ca/~murphyk/Bayes/bayes.html)

Minka, *Using lower bounds to approximate integrals*. 2003.

[www.research.microsoft.com/~minka/papers/rem.html](http://www.research.microsoft.com/~minka/papers/rem.html).

# Other References

Lawrence, *Fast sparse Gaussian process methods: the informative vector machine*. NIPS 2003.

Minka, *Expectation Propagation for Approximate Bayesian Inference*. UAI 2001.

Minka, *Divergence Measures and Message Passing*. AI-Stats 2005.

Neal, *Markov chain sampling methods for Dirichlet process mixture models*, TR. 9815, Dept. of Statistics, University of Toronto.

<http://bayes.hal3.name/>

<http://nlpers.blogspot.com>