

CMPUT463/563
Probabilistic Graphical Models

Supervised Learning: Bayesian Networks

Lili Mou

Dept. Computing Science, University of Alberta

lmou@ualberta.ca

Outline

- Supervised learning: all variables are observed in training
- MLE of BN decomposes to every conditional probability
- For tabular parametrization
 - MLE: counting
 - MAP: add- $(\alpha - 1)$ smoothing with Dirichlet prior
 - Bayesian: add- α smoothing with Dirichlet prior
- Hierarchical Bayes
- Application: pLSI, LDA

Bayesian Networks

BN over variables X_1, \dots, X_N is a DAG with joint probability

$$p(X_1, \dots, X_N) = \prod_{i=1}^N p(X_i | \text{Par}(X_i))$$

Supervised training: all variables are observed in each training sample. $\mathcal{D} = \{(x_1^{(m)}, \dots, x_N^{(m)})\}_{m=1}^M$ for M samples

Maximum likelihood estimation: Assuming parameters θ are unknown constants; choosing such parameters that maximize the probability of data

$$\mathcal{L}(\theta; \mathcal{D}) = P(\mathcal{D}; \theta)$$

Likelihood (of parameters) is just a fancy name of probability of data. It's also common to say the likelihood of data.

MLE for BN

maximize $P(\mathcal{D}; \boldsymbol{\theta}) \Leftrightarrow \text{maximize } \log P(\mathcal{D}; \boldsymbol{\theta})$

$$\Leftrightarrow \text{maximize } \log \prod_{m=1}^M \prod_{n=1}^N P(x_n^{(m)} | \text{Par}(x_n^{(m)}))$$

$$\Leftrightarrow \text{maximize } \sum_{n=1}^M \left[\sum_{m=1}^M \log P(x_i^{(m)} | \text{Par}(x_i^{(m)}; \boldsymbol{\theta}_i)) \right]$$

The optimization decomposes for the parameters $\boldsymbol{\theta}_i$ of each variable X_i

In other words, the **training of a parameter** in BN only concerns the observations **involving the parameter**

Such decomposition works for any BN parametrization (discrete/continuous) and for weight sharing.

Categorical (multinomial) distribution:

Tabular BN

$$\hat{\pi}_k^{(\text{MLE})} = \frac{N_k}{\sum_{k'} N_{k'}} \quad N_k \text{ is the count of } k\text{th category}$$

- **Moment**

- n th-order raw moment $\mathbb{E}[X^n]$
 - E.g., mean = 1st-order raw moment
- n th-order central moment $\mathbb{E}[(X - \mathbb{E}[X])^n]$
 - E.g., variance = 2nd-order central moment

- **Moment matching**

- Tune model parameters (esp. moment parameters) s.t. model moment matches empirical moment
- E.g., π_k is the moment parameter of $\mathbb{E}[X_k]$ for $X_k = 1 \{X = k\}$

$$\frac{N_k}{\sum_{k'} N_{k'}} \text{ is the empirical moment}$$

- For categorical distribution (and a wider range, known as **exponential family**), **moment matching is equivalent to MLE**

Max a posteriori estimation

Bayesian interpretation of probability: everything unknown is a random variable.

- Parameters are treated as random variables with **Prior distribution** $P(\theta)$, the belief of θ before seeing data
- Under some parameter, the **likelihood** is still the probability of data $P(\mathcal{D} | \theta)$. Here, $P(\mathcal{D} | \theta)$ is computationally the same as $P(\mathcal{D}; \theta)$. When using “|”, we treat θ as a random variable; when using “;”, we treat θ as an unknown constants.
- **Posterior distribution:** $P(\theta | \mathcal{D})$, the belief of θ after seeing data
- **Max a posteriori (MAP) estimation** $\operatorname{argmax}_{\theta} P(\theta | \mathcal{D})$

Prior Distribution of Categorical Distribution

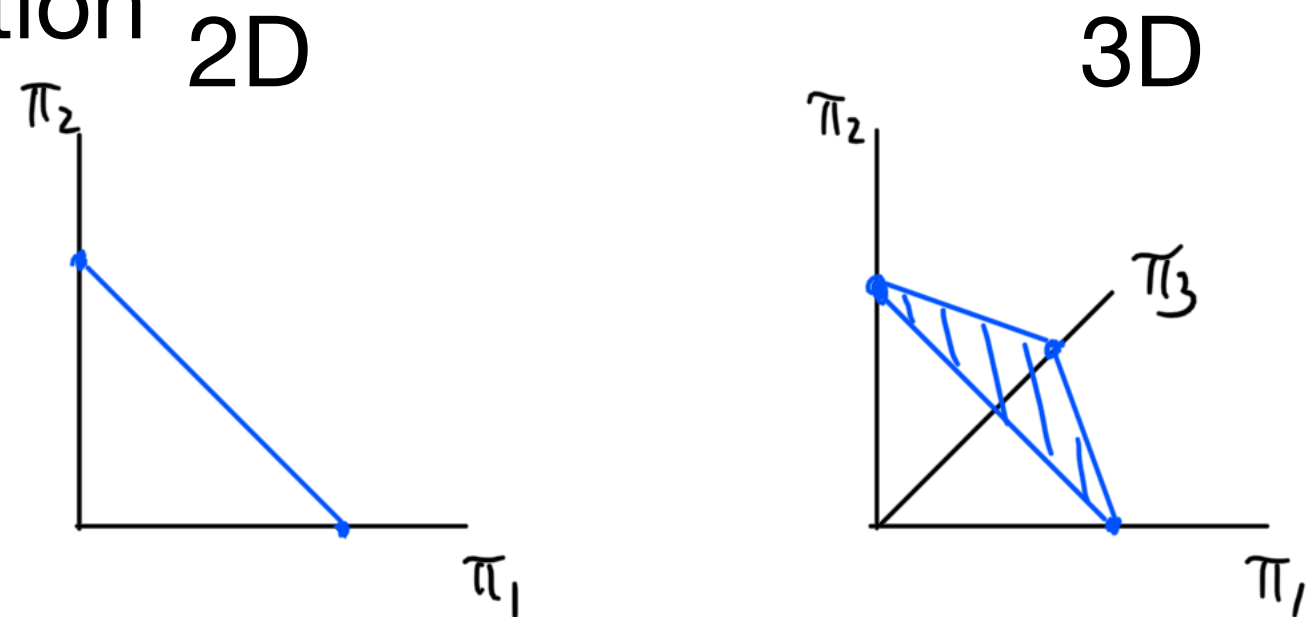
- K -way categorical distribution is fully characterized by $\pi_k = \Pr[X = k]$ for $k = 1, \dots, K$, such that

$$\pi_k \geq 0 \quad \text{and} \quad \sum_k \pi_k = 1$$

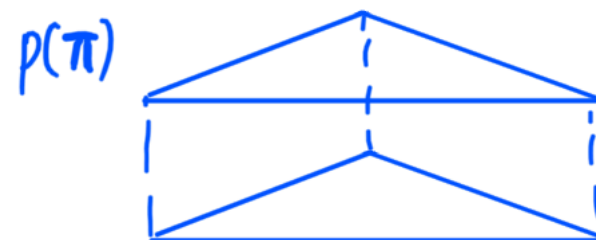
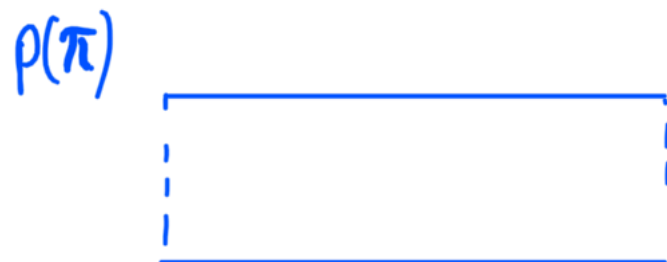
- Consider one-hot representation

$$X = (X_1, \dots, X_K)$$

$$P(X) = \prod_{k=1}^K \pi_k^{X_k}$$



- A prior distribution is a distribution over the blue area (known as **simplex**). E.g., uniform over the simple



- Dirichlet distribution

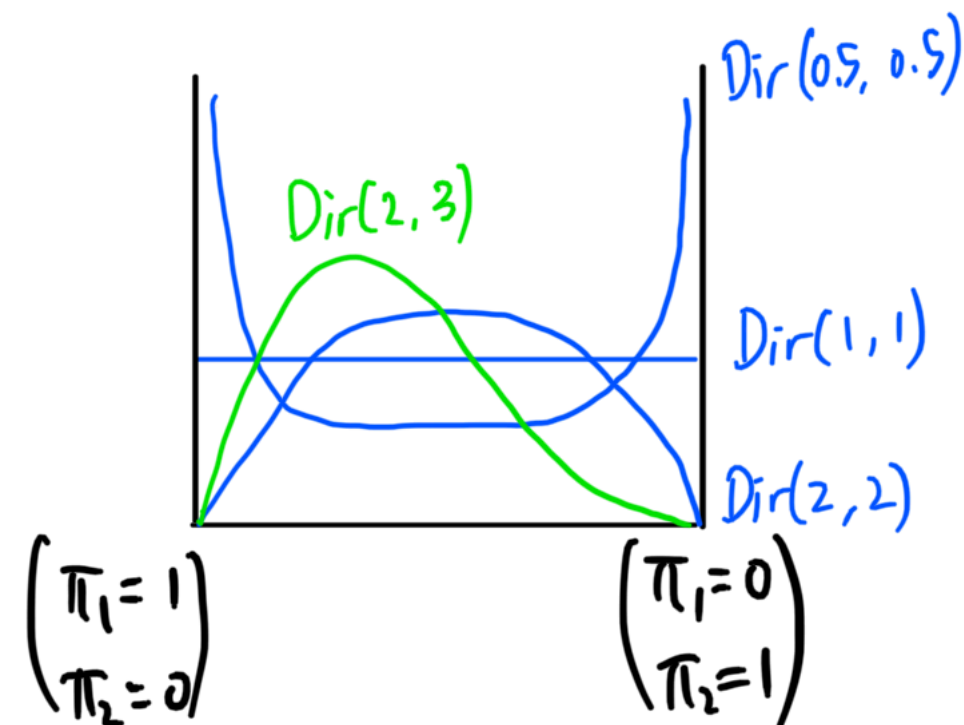
Dirichlet Distribution

$$P(\boldsymbol{\pi}; \boldsymbol{\alpha}) = \frac{\Gamma(\sum_{k=1}^K \alpha_k)}{\prod_{k=1}^K \Gamma(\alpha_k)} \prod_{k=1}^K \pi_k^{\alpha_k - 1}$$

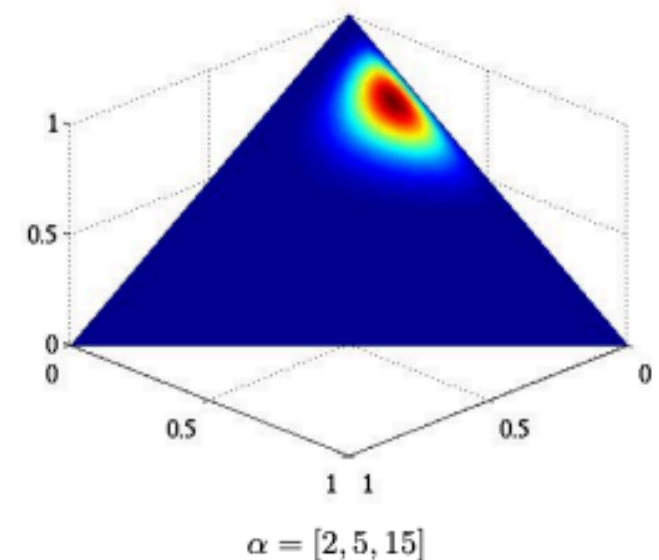
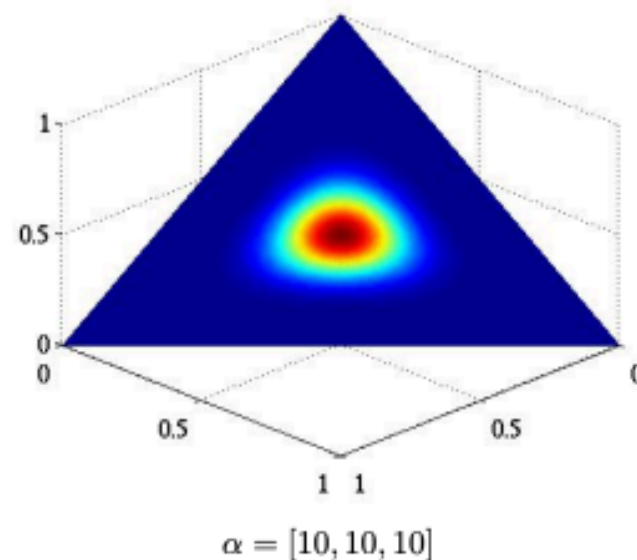
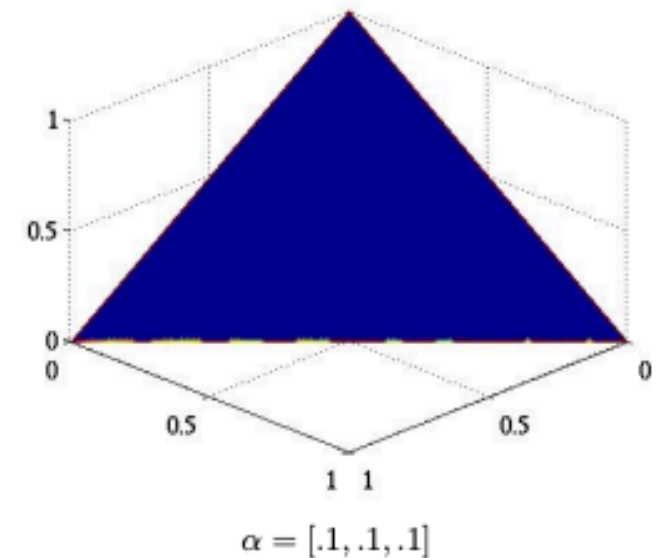
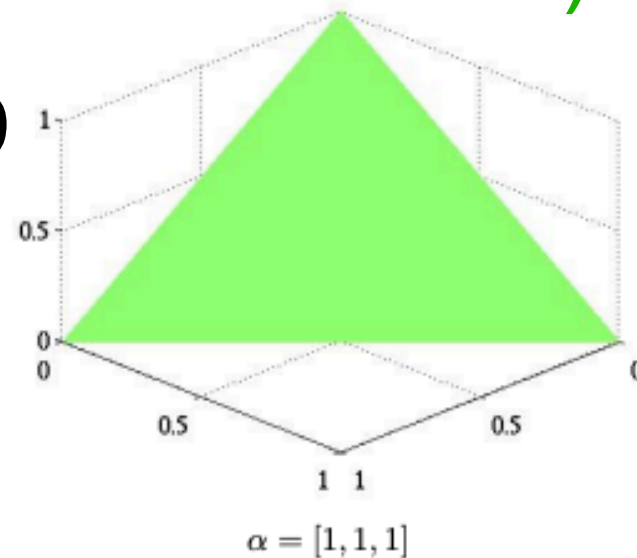
Very analogous to
Categorical distribution

Normalizing factor
(Least important, need not be memorized)

2D



3D

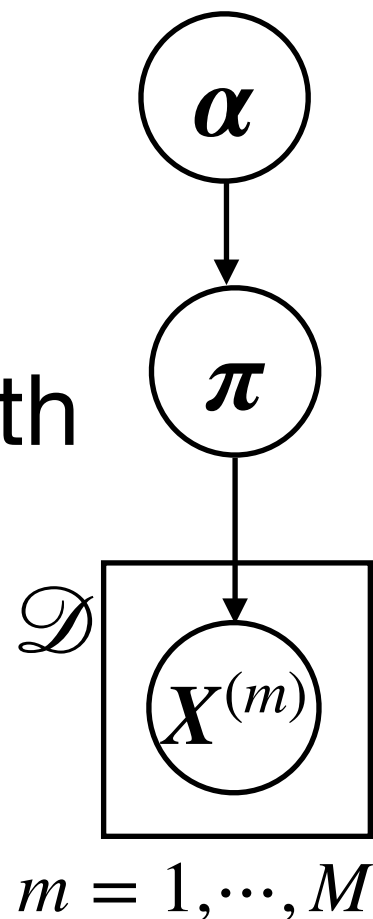


Max a posterior inference

- Why Dirichlet? Equipped with key controls
 - Preference to certain categories; preference strength
- Computationally convenient

$$P(\boldsymbol{\pi} | \mathcal{D}) \propto P(\boldsymbol{\pi})P(\mathcal{D} | \boldsymbol{\pi}) \propto \pi_k^{\alpha_k-1} \pi_k^{N_k} = \pi_k^{\alpha_k+N_k-1}$$

N_k is the count of $X = k$ in the training data

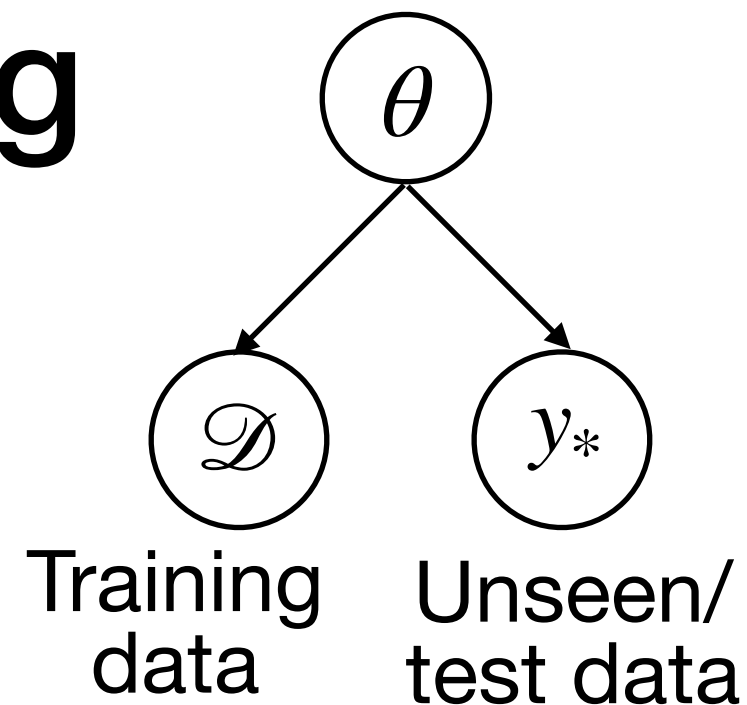


- In other words, $\boldsymbol{\pi} | \mathcal{D} \sim \text{Dir}(\alpha_1 + N_1, \dots, \alpha_K + N_K)$
 - The posterior has the same form as the prior, known as **conjugate prior**.
 - It also has the same form as MLE. Thus, MAP w/ Dirichlet prior is **equivalent to add- $(\alpha - 1)$ smoothing** for MLE

$$\hat{\pi}_k^{\text{MAP}} = \operatorname{argmax} P(\boldsymbol{\pi} | \mathcal{D}) = \frac{N_k + \alpha_k - 1}{\sum_{k'} (N_{k'} + \alpha_{k'} - 1)}$$

Bayesian Learning

- Prior $P(\theta)$
- Likelihood $P(\mathcal{D} | \theta)$
- Posteriori $P(\theta | \mathcal{D}) \propto P(\theta)P(\mathcal{D} | \theta)$



- Recall: **Max a posteriori (MAP) inference**

$$\hat{\theta}^{(\text{MAP})} = \operatorname{argmax}_{\theta} P(\theta | \mathcal{D}) \quad \hat{y}_*^{(\text{MAP})} = P(y | x_*, \hat{\theta}^{(\text{MAP})})$$

- **Bayesian learning:**
 - No specific θ is chosen
 - Everything unrelated to the final prediction should be marginalized out

$$\hat{y}_*^{(\text{Bayesian})} = \int P(y | x_*, \theta) P(\theta | \mathcal{D}) d\theta \propto \int P(y | x_*, \theta) P(\theta) P(\mathcal{D} | \theta) d\theta$$

Bayesian Learning for Dirichlet-Categorical

$$P(y^* | \mathcal{D}) = \int_{\Delta} P(y^* | \pi) P(\pi | \mathcal{D}) d\pi$$

$$= \int_{\Delta} \pi_{y^*} \frac{\Gamma(\sum_{k=1}^K \alpha'_k)}{\prod_{k=1}^K \Gamma(\alpha'_k)} \prod_{k=1}^K \pi_k^{\alpha'_k - 1} d\pi$$

$$= \frac{\Gamma(\sum_{k=1}^K \alpha'_k)}{\prod_{k=1}^K \Gamma(\alpha'_k)} \int_{\Delta} \prod_{k=1}^K \pi_k^{\mathbb{1}\{y^*=k\} + \alpha'_k - 1} d\pi$$

$\propto \text{Dir}(\mathbb{1}\{y^*=k\} + \alpha'_k - 1)$

$$= \frac{\Gamma(\sum_{k=1}^K \alpha'_k)}{\prod_{k=1}^K \Gamma(\alpha'_k)} \cdot \frac{\prod_{k=1}^K \Gamma(\mathbb{1}\{y^*=k\} + \alpha'_k)}{\Gamma(\sum_{k=1}^K (\mathbb{1}\{y^*=k\} + \alpha'_k))}$$

$= 1$ for exactly one term

$$= \frac{\Gamma(\sum_{k=1}^K \alpha'_k)}{\prod_{k=1}^K \Gamma(\alpha'_k)} \cdot \frac{\prod_{k=1}^K \Gamma(\mathbb{1}\{y^*=k\} + \alpha'_k)}{\Gamma(1 + \sum_{k=1}^K \alpha'_k)}$$

$$= \frac{\Gamma(\sum_{k=1}^K \alpha'_k)}{\prod_{k=1}^K \Gamma(\alpha'_k)} \cdot \frac{\alpha'_{y^*} \cdot \prod_{k=1}^K \Gamma(\alpha'_k)}{\sum_{k=1}^K \alpha'_k \cdot \Gamma(\sum_{k=1}^K \alpha'_k)}$$

$$= \frac{\alpha'_{y^*}}{\sum_{k=1}^K \alpha'_k} = \frac{\alpha_{y^*} + N_{y^*}}{\sum_{k=1}^K (\alpha_k + N_k)}$$

$$[\alpha'_k = \alpha_k + N_k, \text{ prior} + \text{count}]$$

Note:

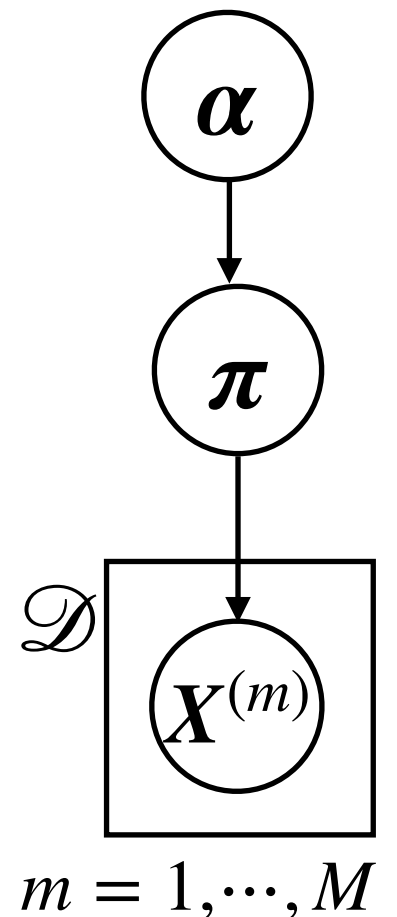
$$\left[\frac{\Gamma(\sum_{k=1}^K \mathbb{1}\{y^*=k\} + \alpha'_k)}{\prod_{k=1}^K \Gamma(\mathbb{1}\{y^*=k\} + \alpha'_k)} \prod_{k=1}^K \pi_k^{\mathbb{1}\{y^*=k\} + \alpha'_k - 1} d\pi = 1 \right]$$

$\text{Dir}(\mathbb{1}\{y^*=k\} + \alpha'_k - 1)$

$$\Gamma(x+1) = x \Gamma(x)$$

Hierarchical Bayes

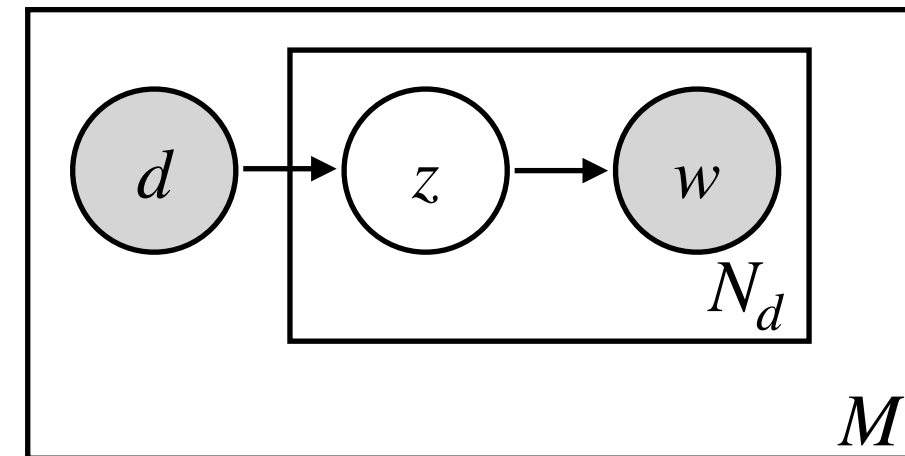
- Sometimes the prior distribution has hyper-parameters Prior $P(\theta | \alpha)$
- How can we handle them?
 - Maximum likelihood estimation (Type-II MLE, empirical Bayes)
 - Max a posteriori estimation (Type-II MAP)
 - Full Bayesian



Probabilistic Latent Semantic Indexing (pLSA)

[Hofmann, SIGIR'99]

- Documents: $d \in \{d_1, \dots, d_M\}$
- Latent topics: $z = \{z_1, \dots, z_K\}$
 - Think of sport, politics, etc.
- Words $w = \{w_1, \dots, w_{|V|}\}$



- Select a document d
- Pick a latent class z with probability $P(z | d)$
- Generate a word $P(w | z)$

Latent Dirichlet Allocation

[Blei, Ng, Jordan, JMLR'03]

- For each document of length N_d
 - Choose $\theta \sim \text{Dir}(\alpha)$
 - For each word w_n
 - $z_n \sim \text{cat}(\theta)$
 - $w_n \sim P(w_n | z_n, \beta)$

where $\beta_{ij} = \Pr[w = j | z = i]$

- A more complicate variant

