# Document models

Ayush Tewari

November 24th, 2025

Adapted from Carl Edward Rasmussen

# Key concepts

- a simple document model

# Key concepts

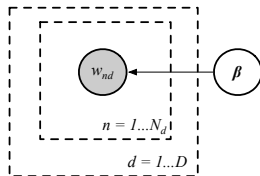- a simple document model
- a mixture model for document

# Key concepts

- a simple document model
- a mixture model for document
- fitting the mixture model with EM

# A really simple document model

Consider a collection of D documents from a vocabulary of M words.

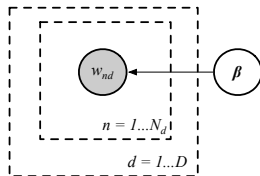- $N_d$: number of words in document d.



---

[1]It's a categorical distribution if we observe the sequence of words in the document, it's a multinomial if we only observe the counts.

# A really simple document model

Consider a collection of D documents from a vocabulary of M words.

- $N_d$: number of words in document d.
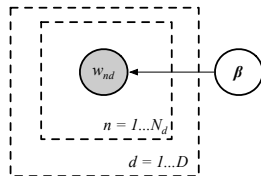- $w_{nd}$: n-th word in document d ($w_{nd} \in \{1 \ldots M\}$).



---

[1]It's a categorical distribution if we observe the sequence of words in the document, it's a multinomial if we only observe the counts.

# A really simple document model

Consider a collection of D documents from a vocabulary of M words.

- $N_d$: number of words in document d.
- $w_{nd}$: n-th word in document d ($w_{nd} \in \{1 \dots M\}$).
- $w_{nd} \sim \mathrm{Cat}(\boldsymbol{\beta})$: each word is drawn from a discrete categorical distribution with parameters $\boldsymbol{\beta}$
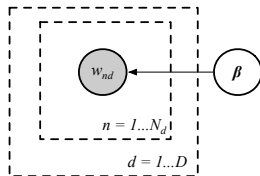


---

[1]It's a categorical distribution if we observe the sequence of words in the document, it's a multinomial if we only observe the counts.

# A really simple document model

Consider a collection of D documents from a vocabulary of M words.

- $N_d$: number of words in document d.
- $w_{nd}$: n-th word in document d ($w_{nd} \in \{1 \ldots M\}$).
- $w_{nd} \sim \mathrm{Cat}(\beta)$: each word is drawn from a discrete categorical distribution with parameters $\beta$
- $\beta = [\beta_1, \ldots, \beta_M]^\top$: parameters of a categorical / multinomial distribution[1] over the M vocabulary words.
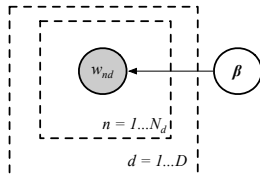


---

[1] It's a categorical distribution if we observe the sequence of words in the document, it's a multinomial if we only observe the counts.

# A really simple document model

Modelling D documents from a vocabulary of M unique words.

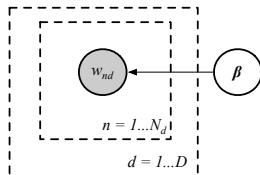- $N_d$: number of words in document d.
- $w_{nd}$: n-th word in document d ($w_{nd} \in \{1 \dots M\}$).
- $w_{nd} \sim \mathrm{Cat}(\beta)$: each word is drawn from a discrete categorical distribution with parameters $\beta$

# A really simple document model

Modelling D documents from a vocabulary of M unique words.

- $N_d$: number of words in document d.
- $w_{nd}$: n-th word in document d ($w_{nd} \in \{1 \dots M\}$).
- $w_{nd} \sim \mathrm{Cat}(\beta)$: each word is drawn from a discrete categorical distribution with parameters $\beta$



We can fit $\beta$ by maximising the likelihood:

# A really simple document model

Modelling D documents from a vocabulary of M unique words.

- $N_d$: number of words in document d.
- $w_{nd}$: n-th word in document d ($w_{nd} \in \{1 \dots M\}$).
- $w_{nd} \sim \text{Cat}(\beta)$: each word is drawn from a discrete categorical distribution with parameters $\beta$
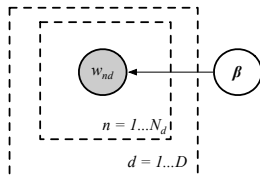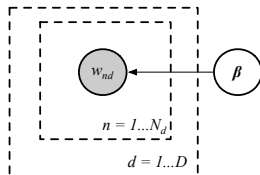


We can fit $\beta$ by maximising the likelihood:

$$\hat{\beta} = \text{argmax}_\beta \prod_{d=1}^{D} \prod_{n}^{N_d} \text{Cat}(w_{nd}|\beta)$$

# A really simple document model

Modelling D documents from a vocabulary of M unique words.

- $N_d$: number of words in document d.
- $w_{nd}$: n-th word in document d ($w_{nd} \in \{1 \ldots M\}$).
- $w_{nd} \sim \mathrm{Cat}(\boldsymbol{\beta})$: each word is drawn from a discrete categorical distribution with parameters $\boldsymbol{\beta}$



We can fit $\boldsymbol{\beta}$ by maximising the likelihood:

$$\hat{\boldsymbol{\beta}} = \mathrm{argmax}_{\boldsymbol{\beta}} \prod_{d=1}^{D} \prod_{n}^{N_d} \mathrm{Cat}(w_{nd}|\boldsymbol{\beta})$$

$$= \mathrm{argmax}_{\boldsymbol{\beta}} \, \mathrm{Mult}(c_1, \ldots, c_M | \boldsymbol{\beta}, N)$$

# A really simple document model

Modelling D documents from a vocabulary of M unique words.

- $N_d$: number of words in document d.
- $w_{nd}$: n-th word in document d ($w_{nd} \in \{1 \dots M\}$).
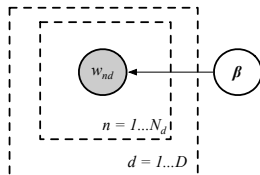- $w_{nd} \sim \mathrm{Cat}(\boldsymbol{\beta})$: each word is drawn from a discrete categorical distribution with parameters $\boldsymbol{\beta}$



We can fit $\boldsymbol{\beta}$ by maximising the likelihood:

$$\hat{\boldsymbol{\beta}} = \mathrm{argmax}_{\boldsymbol{\beta}} \prod_{d=1}^{D} \prod_{n}^{N_d} \mathrm{Cat}(w_{nd}|\boldsymbol{\beta})$$

$$= \mathrm{argmax}_{\boldsymbol{\beta}} \, \mathrm{Mult}(c_1, \dots, c_M|\boldsymbol{\beta}, N)$$

$$\boxed{\hat{\beta}_m = \frac{c_m}{N} = \frac{c_m}{\sum_{\ell=1}^{M} c_\ell}}$$

- $N = \sum_{d=1}^{D} N_d$: total number of words in the collection.
- $c_m = \sum_{d=1}^{D} \sum_{n=1}^{N_d} \mathbb{I}(w_{nd} = m)$: total count of vocabulary word m.

# Maximum Likelihood and Lagrange multipliers

In maximum likelihood learning, we want to maximize the (log) likelihood

- Likelihood: $p(w|\beta) = \prod_{d=1}^{D} \prod_{n=1}^{N_d} \beta_{w_{nd}} = \prod_{m=1}^{M} \beta_m^{c_m}$

# Maximum Likelihood and Lagrange multipliers

In maximum likelihood learning, we want to maximize the (log) likelihood

- Likelihood: $p(w|\boldsymbol{\beta}) = \prod_{d=1}^{D} \prod_{n=1}^{N_d} \beta_{w_{nd}} = \prod_{m=1}^{M} \beta_m^{c_m}$
- Log-Likelihood: $\log p(w|\boldsymbol{\beta}) = \sum_{m=1}^{M} c_m \log \beta_m$

# Maximum Likelihood and Lagrange multipliers

In maximum likelihood learning, we want to maximize the (log) likelihood

- Likelihood: $p(w|\beta) = \prod_{d=1}^{D} \prod_{n=1}^{N_d} \beta_{w_{nd}} = \prod_{m=1}^{M} \beta_m^{c_m}$
- Log-Likelihood: $\log p(w|\beta) = \sum_{m=1}^{M} c_m \log \beta_m$
- **Constraint:** $\sum_{m=1}^{M} \beta_m = 1$.

# Maximum Likelihood and Lagrange multipliers

In maximum likelihood learning, we want to maximize the (log) likelihood

- Likelihood: $p(w|\beta) = \prod_{d=1}^{D} \prod_{n=1}^{N_d} \beta_{w_{nd}} = \prod_{m=1}^{M} \beta_m^{c_m}$
- Log-Likelihood: $\log p(w|\beta) = \sum_{m=1}^{M} c_m \log \beta_m$
- **Constraint:** $\sum_{m=1}^{M} \beta_m = 1$.

An easy way to do this optimization is to add the Lagrange multiplier to the cost:

$$F = \sum_{m=1}^{M} c_m \log \beta_m + \lambda(1 - \sum_{m=1}^{M} \beta_m),$$

# Maximum Likelihood and Lagrange multipliers

In maximum likelihood learning, we want to maximize the (log) likelihood

- Likelihood: $p(w|\boldsymbol{\beta}) = \prod_{d=1}^{D} \prod_{n=1}^{N_d} \beta_{w_{nd}} = \prod_{m=1}^{M} \beta_m^{c_m}$
- Log-Likelihood: $\log p(w|\boldsymbol{\beta}) = \sum_{m=1}^{M} c_m \log \beta_m$
- **Constraint:** $\sum_{m=1}^{M} \beta_m = 1$.

An easy way to do this optimization is to add the Lagrange multiplier to the cost:

$$F = \sum_{m=1}^{M} c_m \log \beta_m + \lambda(1 - \sum_{m=1}^{M} \beta_m),$$

taking derivatives and setting to zero, we obtain

$$\frac{\partial F}{\partial \beta_m} = \frac{c_m}{\beta_m} - \lambda = 0 \Rightarrow \beta_m = \frac{c_m}{\lambda} \text{ and } \frac{\partial F}{\partial \lambda} = 0 \Rightarrow \sum_{m=1}^{M} \beta_m = 1,$$

# Maximum Likelihood and Lagrange multipliers

In maximum likelihood learning, we want to maximize the (log) likelihood

- Likelihood: $p(w|\boldsymbol{\beta}) = \prod_{d=1}^{D} \prod_{n=1}^{N_d} \beta_{w_{nd}} = \prod_{m=1}^{M} \beta_m^{c_m}$
- Log-Likelihood: $\log p(w|\boldsymbol{\beta}) = \sum_{m=1}^{M} c_m \log \beta_m$
- **Constraint:** $\sum_{m=1}^{M} \beta_m = 1$.

An easy way to do this optimization is to add the Lagrange multiplier to the cost:

$$F = \sum_{m=1}^{M} c_m \log \beta_m + \lambda(1 - \sum_{m=1}^{M} \beta_m),$$

taking derivatives and setting to zero, we obtain

$$\frac{\partial F}{\partial \beta_m} = \frac{c_m}{\beta_m} - \lambda = 0 \Rightarrow \beta_m = \frac{c_m}{\lambda} \text{ and } \frac{\partial F}{\partial \lambda} = 0 \Rightarrow \sum_{m=1}^{M} \beta_m = 1,$$

which we combine to $\beta_m = c_m/N$, where $N$ is the total number of words.

# Limitations of the really simple document model

- Document $d$ is the result of sampling $N_d$ words from the categorical distribution with parameters $\beta$.

# Limitations of the really simple document model

- Document d is the result of sampling $N_d$ words from the categorical distribution with parameters $\beta$.
- $\beta$ estimated by maximum likelihood reflects the aggregation of all documents.

# Limitations of the really simple document model

- Document d is the result of sampling $N_d$ words from the categorical distribution with parameters $\beta$.
- $\beta$ estimated by maximum likelihood reflects the aggregation of all documents.
- All documents are therefore modelled by the global word frequency distribution.

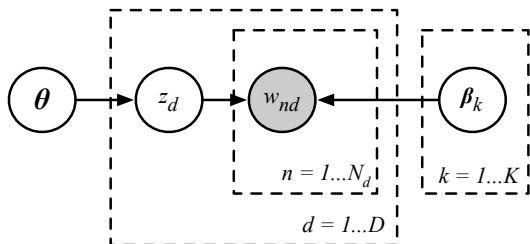# Limitations of the really simple document model

- Document d is the result of sampling $N_d$ words from the categorical distribution with parameters $\beta$.
- $\beta$ estimated by maximum likelihood reflects the aggregation of all documents.
- All documents are therefore modelled by the global word frequency distribution.
- **This generative model does not specialise.**

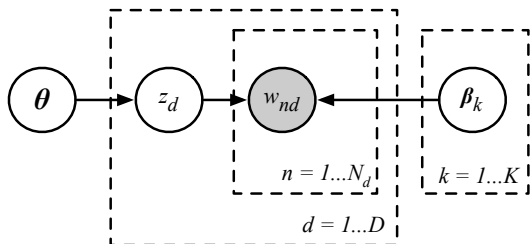# Limitations of the really simple document model

- Document d is the result of sampling $N_d$ words from the categorical distribution with parameters $\beta$.
- $\beta$ estimated by maximum likelihood reflects the aggregation of all documents.
- All documents are therefore modelled by the global word frequency distribution.
- **This generative model does not specialise.**
- We would like a model where different documents might be about different *topics*.
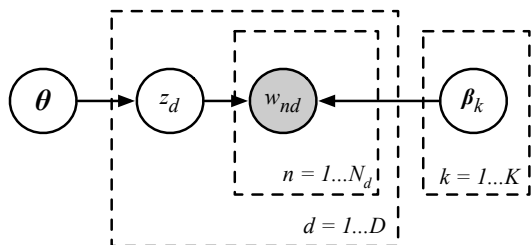
# A mixture of categoricals model

# A mixture of categoricals model



$$z_{\mathtt{d}} \quad \sim \quad \mathrm{Cat}(\boldsymbol{\theta})$$

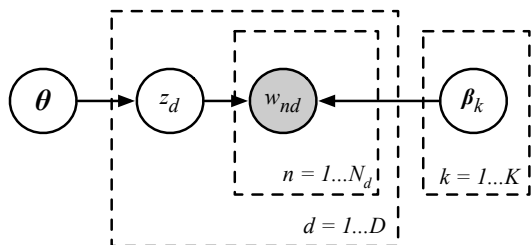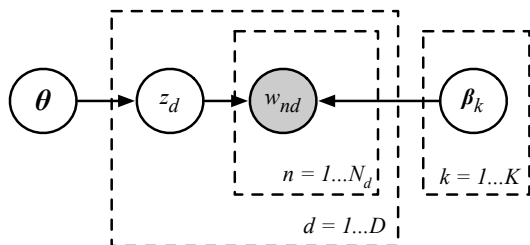$$w_{\mathtt{nd}}|z_{\mathtt{d}} \quad \sim \quad \mathrm{Cat}(\boldsymbol{\beta}_{z_{\mathtt{d}}})$$

# A mixture of categoricals model



$$z_d \sim \text{Cat}(\theta)$$
$$w_{nd}|z_d \sim \text{Cat}(\beta_{z_d})$$

We want to allow for a mixture of K categoricals parametrised by $\beta_1, \ldots, \beta_K$.

# A mixture of categoricals model



$$z_d \sim \text{Cat}(\theta)$$
$$w_{nd}|z_d \sim \text{Cat}(\beta_{z_d})$$

We want to allow for a mixture of K categoricals parametrised by $\beta_1, \ldots, \beta_K$.
Each of those categorical distributions corresponds to a *document category*.
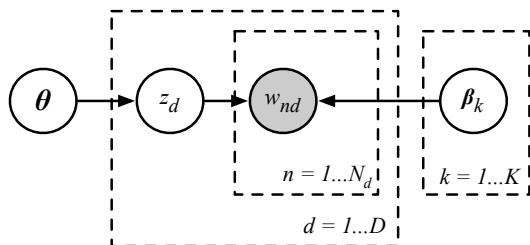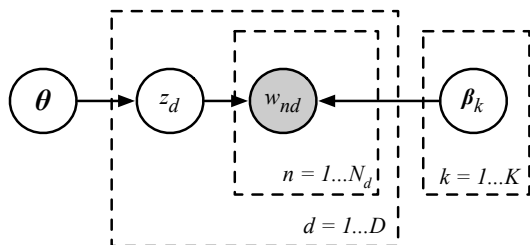
# A mixture of categoricals model



$$z_d \sim \mathrm{Cat}(\theta)$$
$$w_{nd}|z_d \sim \mathrm{Cat}(\beta_{z_d})$$

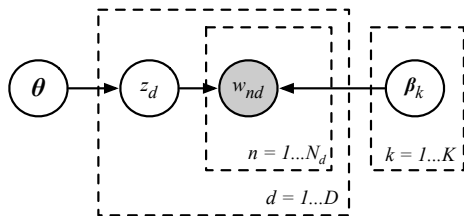We want to allow for a mixture of K categoricals parametrised by $\beta_1, \ldots, \beta_K$.
Each of those categorical distributions corresponds to a *document category*.

- $z_d \in \{1, \ldots, K\}$ assigns document d to one of the K categories.

# A mixture of categoricals model



$$z_d \sim \text{Cat}(\boldsymbol{\theta})$$
$$w_{nd}|z_d \sim \text{Cat}(\boldsymbol{\beta}_{z_d})$$

We want to allow for a mixture of K categoricals parametrised by $\boldsymbol{\beta}_1, \ldots, \boldsymbol{\beta}_K$.
Each of those categorical distributions corresponds to a *document category*.

- $z_d \in \{1, \ldots, K\}$ assigns document d to one of the K categories.
- $\theta_k = p(z_d = k)$ is the probability any document d is assigned to category k.
- so $\boldsymbol{\theta} = [\theta_1, \ldots, \theta_K]$ is the parameter of a categorical distribution over K categories.

# A mixture of categoricals model



$$z_d \sim \text{Cat}(\boldsymbol{\theta})$$
$$w_{nd}|z_d \sim \text{Cat}(\boldsymbol{\beta}_{z_d})$$

We want to allow for a mixture of K categoricals parametrised by $\boldsymbol{\beta}_1, \ldots, \boldsymbol{\beta}_K$. Each of those categorical distributions corresponds to a *document category*.

- $z_d \in \{1, \ldots, K\}$ assigns document d to one of the K categories.
- $\theta_k = p(z_d = k)$ is the probability any document d is assigned to category k.
- so $\boldsymbol{\theta} = [\theta_1, \ldots, \theta_K]$ is the parameter of a categorical distribution over K categories.

We have introduced a new set of *hidden* variables $z_d$.

- How do we fit those variables? What do we do with them?
- Are these variables interesting? Or are we only interested in $\boldsymbol{\theta}$ and $\boldsymbol{\beta}$?

# A mixture of categoricals model: the likelihood



$$z_d \sim \text{Cat}(\grave{})$$
$$w_{nd}|z_d \sim \text{Cat}(\text{fi}_{z_d})$$

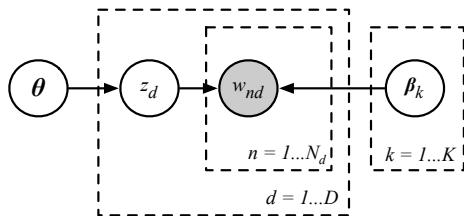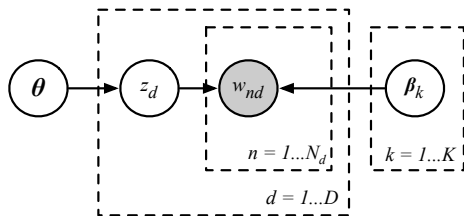$$p(w|\theta, \beta) = \prod_{d=1}^{D} p(w_d|\theta, \beta)$$

# A mixture of categoricals model: the likelihood



$$z_d \sim \text{Cat}(\grave{})$$
$$w_{nd}|z_d \sim \text{Cat}(fi_{z_d})$$

$$
\begin{aligned}
p(w|\theta, \beta) &= \prod_{d=1}^{D} p(w_d|\theta, \beta) \\
&= \prod_{d=1}^{D} \sum_{k=1}^{K} p(w_d, z_d = k|\theta, \beta)
\end{aligned}
$$

# A mixture of categoricals model: the likelihood



$$z_d \sim \mathrm{Cat}(\grave{})$$
$$w_{nd}|z_d \sim \mathrm{Cat}(\mathrm{fi}_{z_d})$$

$$
\begin{aligned}
p(w|\theta, \beta) &= \prod_{d=1}^{D} p(w_d|\theta, \beta) \\
&= \prod_{d=1}^{D} \sum_{k=1}^{K} p(w_d, z_d = k|\theta, \beta) \\
&= \prod_{d=1}^{D} \sum_{k=1}^{K} p(z_d = k|\theta) p(w_d|z_d = k, \beta_k)
\end{aligned}
$$

# A mixture of categoricals model: the likelihood



$$
\begin{aligned}
z_d &\sim \text{Cat}(\theta) \\
w_{nd}|z_d &\sim \text{Cat}(\beta_{z_d})
\end{aligned}
$$

$$
\begin{aligned}
p(w|\theta, \beta) &= \prod_{d=1}^{D} p(w_d|\theta, \beta) \\
&= \prod_{d=1}^{D} \sum_{k=1}^{K} p(w_d, z_d = k|\theta, \beta) \\
&= \prod_{d=1}^{D} \sum_{k=1}^{K} p(z_d = k|\theta) p(w_d|z_d = k, \beta_k) \\
&= \prod_{d=1}^{D} \sum_{k=1}^{K} p(z_d = k|\theta) \prod_{n=1}^{N_d} p(w_{nd}|z_d = k, \beta_k)
\end{aligned}
$$

# EM and Mixtures of Categoricals

In the mixture model, the likelihood is:

$$p(w|\theta, \beta) = \prod_{d=1}^{D} \sum_{k=1}^{K} p(z_d = k|\theta) \prod_{n=1}^{N_d} p(w_{nd}|z_d = k, \beta_k)$$

# EM and Mixtures of Categoricals

In the mixture model, the likelihood is:

$$p(w|\theta, \beta) = \prod_{d=1}^{D} \sum_{k=1}^{K} p(z_d = k|\theta) \prod_{n=1}^{N_d} p(w_{nd}|z_d = k, \beta_k)$$

E-step: for each $d$, set $q$ to the posterior (where $c_{md} = \sum_{n=1}^{N_d} \mathbb{I}(w_{nd} = m)$):

$$q(z_d = k) \propto p(z_d = k|\theta) \prod_{n=1}^{N_d} p(w_{nd}|\beta_{k,w_n}) = \theta_k \, \mathrm{Mult}(c_{1d}, \ldots, c_{Md}|\beta_k, N_d) \stackrel{\mathrm{def}}{=} r_{kd}$$

# EM and Mixtures of Categoricals

In the mixture model, the likelihood is:

$$p(w|\theta, \beta) = \prod_{d=1}^{D} \sum_{k=1}^{K} p(z_d = k|\theta) \prod_{n=1}^{N_d} p(w_{nd}|z_d = k, \beta_k)$$

E-step: for each $d$, set $q$ to the posterior (where $c_{md} = \sum_{n=1}^{N_d} \mathbb{I}(w_{nd} = m)$):

$$q(z_d = k) \propto p(z_d = k|\theta) \prod_{n=1}^{N_d} p(w_{nd}|\beta_{k,w_n}) = \theta_k \, \mathrm{Mult}(c_{1d}, \ldots, c_{Md}|\beta_k, N_d) \overset{\mathrm{def}}{=} r_{kd}$$

M-step: Maximize

$$\sum_{d=1}^{D} \sum_{k=1}^{K} q(z_d = k) \log p(w, z_d) = \sum_{k,d} r_{kd} \log \left[ p(z_d = k|\theta) \prod_{n=1}^{N_d} p(w_{nd}|\beta_{k,w_{nd}}) \right]$$

# EM and Mixtures of Categoricals

In the mixture model, the likelihood is:

$$p(w|\theta, \beta) = \prod_{d=1}^{D} \sum_{k=1}^{K} p(z_d = k|\theta) \prod_{n=1}^{N_d} p(w_{nd}|z_d = k, \beta_k)$$

E-step: for each $d$, set $q$ to the posterior (where $c_{md} = \sum_{n=1}^{N_d} \mathbb{I}(w_{nd} = m)$):

$$q(z_d = k) \propto p(z_d = k|\theta) \prod_{n=1}^{N_d} p(w_{nd}|\beta_{k,w_n}) = \theta_k \, \text{Mult}(c_{1d}, \ldots, c_{Md}|\beta_k, N_d) \stackrel{\text{def}}{=} r_{kd}$$

M-step: Maximize

$$\sum_{d=1}^{D} \sum_{k=1}^{K} q(z_d = k) \log p(w, z_d) = \sum_{k,d} r_{kd} \log \left[ p(z_d = k|\theta) \prod_{n=1}^{N_d} p(w_{nd}|\beta_{k,w_{nd}}) \right]$$

$$= \sum_{k,d} r_{kd} \left( \log \prod_{m=1}^{M} \beta_{km}^{c_{md}} + \log \theta_k \right)$$

# EM and Mixtures of Categoricals

In the mixture model, the likelihood is:

$$p(w|\theta, \beta) = \prod_{d=1}^{D} \sum_{k=1}^{K} p(z_d = k|\theta) \prod_{n=1}^{N_d} p(w_{nd}|z_d = k, \beta_k)$$

E-step: for each $d$, set $q$ to the posterior (where $c_{md} = \sum_{n=1}^{N_d} \mathbb{I}(w_{nd} = m)$):

$$q(z_d = k) \propto p(z_d = k|\theta) \prod_{n=1}^{N_d} p(w_{nd}|\beta_{k,w_n}) = \theta_k \, \mathrm{Mult}(c_{1d}, \ldots, c_{Md}|\beta_k, N_d) \overset{\mathrm{def}}{=} r_{kd}$$

M-step: Maximize

$$\sum_{d=1}^{D} \sum_{k=1}^{K} q(z_d = k) \log p(w, z_d) = \sum_{k,d} r_{kd} \log \left[ p(z_d = k|\theta) \prod_{n=1}^{N_d} p(w_{nd}|\beta_{k,w_{nd}}) \right]$$

$$= \sum_{k,d} r_{kd} \left( \log \prod_{m=1}^{M} \beta_{km}^{c_{md}} + \log \theta_k \right)$$

$$= \sum_{k,d} r_{kd} \left( \sum_{m=1}^{M} c_{md} \log \beta_{km} + \log \theta_k \right) \overset{\mathrm{def}}{=} F(R, \theta, \beta)$$

# EM: M step for mixture model

$$F(R, \theta, \beta) = \sum_{k,d} r_{kd} \left( \sum_{m=1}^{M} c_{md} \log \beta_{km} + \log \theta_k \right)$$

Need Lagrange multipliers to constrain the maximization of $F$ and ensure proper distributions.

# EM: M step for mixture model

$$F(R, \theta, \beta) = \sum_{k,d} r_{kd} \left( \sum_{m=1}^{M} c_{md} \log \beta_{km} + \log \theta_k \right)$$

Need Lagrange multipliers to constrain the maximization of F and ensure proper distributions.

**Update for $\theta$ (Topic Proportions):**

$$\hat{\theta}_k \leftarrow \operatorname{argmax}_{\theta_k} F(R, \theta, \beta) + \lambda \left( 1 - \sum_{k'=1}^{K} \theta_{k'} \right)$$

$$= \frac{\sum_{d=1}^{D} r_{kd}}{\sum_{k'=1}^{K} \sum_{d=1}^{D} r_{k'd}} = \boxed{\frac{\sum_{d=1}^{D} r_{kd}}{D}}$$

# EM: M step for mixture model

$$F(R, \theta, \beta) = \sum_{k,d} r_{kd} \left( \sum_{m=1}^{M} c_{md} \log \beta_{km} + \log \theta_k \right)$$

Need Lagrange multipliers to constrain the maximization of $F$ and ensure proper distributions.
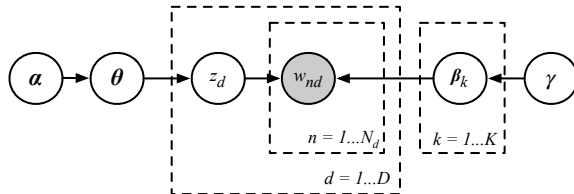
**Update for $\theta$ (Topic Proportions):**

$$\hat{\theta}_k \leftarrow \mathrm{argmax}_{\theta_k} \ F(R, \theta, \beta) + \lambda(1 - \sum_{k'=1}^{K} \theta_{k'})$$

$$= \frac{\sum_{d=1}^{D} r_{kd}}{\sum_{k'=1}^{K} \sum_{d=1}^{D} r_{k'd}} = \boxed{\frac{\sum_{d=1}^{D} r_{kd}}{D}}$$

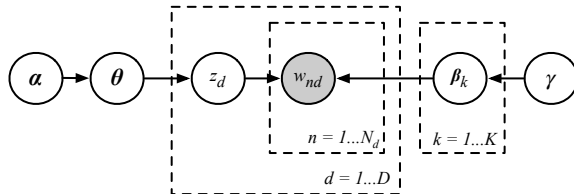**Update for $\beta$ (Word Probabilities):**

$$\hat{\beta}_{km} \leftarrow \mathrm{argmax}_{\beta_{km}} \ F(R, \theta, \beta) + \sum_{k'=1}^{K} \lambda_{k'}(1 - \sum_{m'=1}^{M} \beta_{k'm'})$$
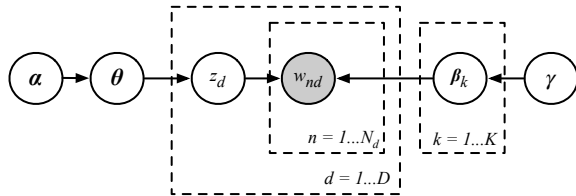
# A Bayesian mixture of categoricals model

# A Bayesian mixture of categoricals model



$$\theta \sim \mathrm{Dir}(\alpha)$$
$$\beta_k \sim \mathrm{Dir}(\gamma)$$
$$z_d | \theta \sim \mathrm{Cat}(`)$$
$$w_{nd} | z_d, \beta \sim \mathrm{Cat}(\mathrm{fi}_{z_d})$$

# A Bayesian mixture of categoricals model



$$\theta \sim \mathrm{Dir}(\alpha)$$
$$\beta_k \sim \mathrm{Dir}(\gamma)$$
$$z_d | \theta \sim \mathrm{Cat}(`)$$
$$w_{nd} | z_d, \beta \sim \mathrm{Cat}(\mathrm{fi}_{z_d})$$

- With the EM algorithm we have essentially estimated $\theta$ and $\beta$ by maximum likelihood.

# A Bayesian mixture of categoricals model



$$\theta \sim \mathrm{Dir}(\alpha)$$
$$\beta_k \sim \mathrm{Dir}(\gamma)$$
$$z_d | \theta \sim \mathrm{Cat}(`)$$
$$w_{nd} | z_d, \beta \sim \mathrm{Cat}(fi_{z_d})$$

- With the EM algorithm we have essentially estimated $\theta$ and $\beta$ by maximum likelihood.

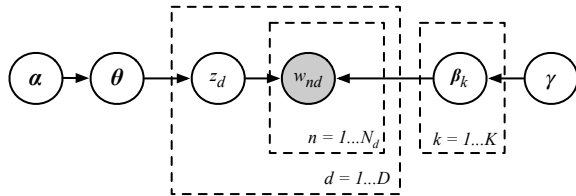- An alternative, Bayesian treatment infers these parameters starting from priors:
    - $\theta \sim \mathrm{Dir}(\alpha)$ is a symmetric Dirichlet over category probabilities.
    - $\beta_k \sim \mathrm{Dir}(\gamma)$ are symmetric Dirichlets over vocabulary probabilities.

# A Bayesian mixture of categoricals model



$$\theta \sim \text{Dir}(\alpha)$$
$$\beta_k \sim \text{Dir}(\gamma)$$
$$z_d|\theta \sim \text{Cat}(`)$$
$$w_{nd}|z_d, \beta \sim \text{Cat}(fi_{z_d})$$

- With the EM algorithm we have essentially estimated $\theta$ and $\beta$ by maximum likelihood.
- An alternative, Bayesian treatment infers these parameters starting from priors:
    - $\theta \sim \text{Dir}(\alpha)$ is a symmetric Dirichlet over category probabilities.
    - $\beta_k \sim \text{Dir}(\gamma)$ are symmetric Dirichlets over vocabulary probabilities.

What is different?

- We no longer want to compute a point estimate of $\theta$ or $\beta$.
- We are now interested in computing the *posterior* distributions.