# Document models

Carl Edward Rasmussen
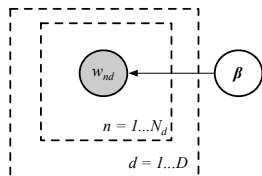
November 15th, 2017

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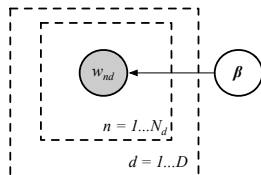
  Summing beta_1 to beta_M gives unity.



small box is one document

---

[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 ($w_{nd} \in \{1 \dots M\}$).
- $w_{nd} \sim \text{Cat}(\beta)$: each word is drawn from a discrete categorical distribution with parameters $\beta$



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

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

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

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

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

# Maximum Likelihood and Lagrange multipliers

In maximum likelihood learning, we want to maximize the (log) likelihood
(Notice the absence of multinomial coefficient, hence modelling sequence but not count.)

$$p(\mathbf{w}|\boldsymbol{\beta}) = \prod_{n=1}^{D} \prod_{n=1}^{N_d} \beta_{w_{nd}} = \prod_{m=1}^{M} \beta_m^{c_m}, \text{ or } \log p(\mathbf{w}|\boldsymbol{\beta}) = \sum_{m=1}^{M} c_m \log \beta_m,$$

subject to the normalizing constraint that $\sum_{m=1}^{M} \beta_m = 1$.
An easy way to do this optimization is to add the Lagrange multiplier to the cost

$$F = \sum_{m=1}^{M} c_m \log \beta_m + \lambda(1 - \sum_{m=1}^{M} \beta_m),$$
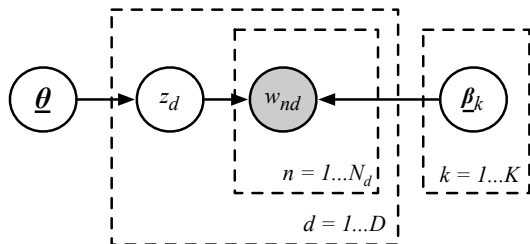
taking derivatives and setting to zero, we obtain

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

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 $\beta$.
- $\beta$ estimated by maximum likelihood reflects the aggregation of all documents.
- All documents are therefore modelled by the global word frequency distribution.  bad!
- This generative model does not specialise.
- We would like a model where different documents might be about different *topics*.

# A mixture of categoricals model



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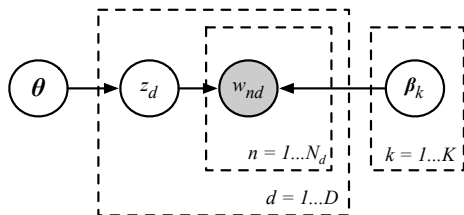$$w_{nd}|z_d \sim \text{Cat}(\underline{\beta}_{z_d})$$

We want to allow for a mixture of K categoricals parametrised by $\underline{\beta}_1, \ldots, \underline{\beta}_K$.
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 $\theta = [\theta_1, \ldots, \theta_K]$ 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 θ and β?

# A mixture of categoricals model: the likelihood



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

prob that
the observed
words appear
this way, given
the model parameters.

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

$$= \prod_{d=1}^{D} \sum_{k=1}^{K} p(\mathbf{w}_d, z_d = k|\boldsymbol{\theta}, \boldsymbol{\beta}) \quad \text{sum rule of prob.}$$

$$= \prod_{d=1}^{D} \sum_{k=1}^{K} p(z_d = k|\boldsymbol{\theta}) \underset{\wedge}{p(\mathbf{w}_d|z_d = k, \boldsymbol{\beta}_k)} \quad \begin{array}{l} \text{beta\_k and theta} \\ \text{are independent.} \end{array}$$

beta_k · · · theta

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

Assume all words in document d are independent.

# EM and Mixtures of Categoricals 📝

In the mixture model, the likelihood is:

sum over latent variable z.

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

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

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

prior         likelihood

M-step: Maximize

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

Why has d become
a summation?

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

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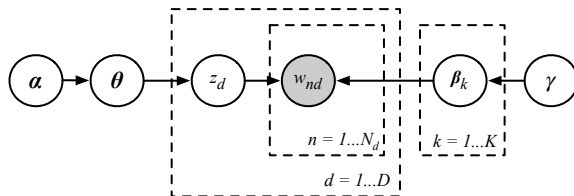
📝

# 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.

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

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

$$\hat{\beta}_{km} \leftarrow \operatorname{argmax}_{\beta_{km}} F(R, \theta, \beta) + \sum_{k'=1}^{K} \lambda_{k'}(1 - \sum_{m'=1}^{M} \beta_{k'm'})$$

$$= \frac{\sum_{d=1}^{D} r_{kd} c_{md}}{\sum_{m'=1}^{M} \sum_{d=1}^{D} r_{kd} c_{m'd}}$$

# A Bayesian mixture of categoricals model



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

With the EM algorithm we have essentially estimated $\theta$ and $\beta$ by maximum likelihood. An alternative, Bayesian treatment infers these parameters starting from priors, e.g.:

- $\theta \sim \text{Dir}(\alpha)$ 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.