

# Introduction to topic models: Building up towards LDA

Vivi Nastase

Summer semester 2012  
ICL, University of Heidelberg

# Motivation

Assumption: the data was generated according to some model that we try to retrieve.

# Motivation

Assumption: the data was generated according to some model that we try to retrieve.

The **observation**  $\mathcal{O}$ : the texts.

The **model**  $\mu$ : parameters and probability distributions that generated the data – the one that fits the data best.

$$p(\mathcal{O}, \mu) = \underbrace{p(\mathcal{O}|\mu)}_{\text{likelihood}} \quad p(\mu) = \underbrace{p(\mu|\mathcal{O})}_{\text{posterior}} \quad \underbrace{p(\mathcal{O})}_{\text{evidence}}$$

## Motivation

Assumption: the data was generated according to some model that we try to retrieve.

The **observation**  $\mathcal{O}$ : the texts.

The **model**  $\mu$ : parameters and probability distributions that generated the data – the one that fits the data best.

$$p(\mathcal{O}, \mu) = \underbrace{p(\mathcal{O}|\mu)}_{\text{likelihood}} \underbrace{p(\mu)}_{\text{prior}} = \underbrace{p(\mu|\mathcal{O})}_{\text{posterior}} \underbrace{p(\mathcal{O})}_{\text{evidence}}$$

**Best fit:**

- maximize likelihood
- maximize the posterior
- compute the probability distribution of the posterior

# Maximum likelihood estimation (ML)

Likelihood:

$$L(\mu|\mathcal{O}) = p(\mathcal{O}|\mu) = \prod_{x \in \mathcal{O}} p(x|\mu)$$

$$\mathcal{L}(\mu|\mathcal{O}) = \log L(\mu|\mathcal{O}) = \sum_{x \in \mathcal{O}} \log p(x|\mu)$$

## Maximum likelihood estimation (ML)

Likelihood:

$$L(\mu|\mathcal{O}) = p(\mathcal{O}|\mu) = \prod_{x \in \mathcal{O}} p(x|\mu)$$

$$\mathcal{L}(\mu|\mathcal{O}) = \log L(\mu|\mathcal{O}) = \sum_{x \in \mathcal{O}} \log p(x|\mu)$$

The model that maximizes the likelihood:

$$\mu_{ML} = \operatorname{argmax}_{\mu} \mathcal{L}(\mu|\mathcal{O}) = \operatorname{argmax}_{\mu} \sum_{x \in \mathcal{O}} \log p(x|\mu)$$

# Maximum likelihood estimation (ML)

Likelihood:

$$L(\mu|\mathcal{O}) = p(\mathcal{O}|\mu) = \prod_{x \in \mathcal{O}} p(x|\mu)$$

$$\mathcal{L}(\mu|\mathcal{O}) = \log L(\mu|\mathcal{O}) = \sum_{x \in \mathcal{O}} \log p(x|\mu)$$

The model that maximizes the likelihood:

$$\mu_{ML} = \operatorname{argmax}_{\mu} \mathcal{L}(\mu|\mathcal{O}) = \operatorname{argmax}_{\mu} \sum_{x \in \mathcal{O}} \log p(x|\mu)$$

Maximum:

$$\frac{\partial \mathcal{L}(\mu|\mathcal{O})}{\partial \mu_k} = 0; \forall \mu_k \in \mu$$

## Maximum likelihood estimation (ML) – example

$\mathcal{O}$ :  $N$  coin tossing experiments,  $N = n_h + n_t$



## Maximum likelihood estimation (ML) – example

$\mathcal{O}$ :  $N$  coin tossing experiments,  $N = n_h + n_t$

$$\mu: p_h \quad p(X = x_i | p_h) = p_h^{x_i} (1 - p_h)^{1 - x_i}$$

$x_i = 1$  for heads,  $x_i = 0$  for tails

## Maximum likelihood estimation (ML) – example

$\mathcal{O}$ :  $N$  coin tossing experiments,  $N = n_h + n_t$

$\mu$ :  $p_h$       $p(X = x_i | p_h) = p_h^{x_i} (1 - p_h)^{1-x_i}$

$x_i = 1$  for heads,  $x_i = 0$  for tails

$$\mathcal{L}(p_h | N) = \sum_{i=1, N} \log p(X = x_i | p_h)$$

## Maximum likelihood estimation (ML) – example

$\mathcal{O}$ :  $N$  coin tossing experiments,  $N = n_h + n_t$

$\mu$ :  $p_h$       $p(X = x_i | p_h) = p_h^{x_i} (1 - p_h)^{1 - x_i}$

$x_i = 1$  for heads,  $x_i = 0$  for tails

$$\begin{aligned}\mathcal{L}(p_h | N) &= \sum_{i=1, N} \log p(X = x_i | p_h) \\ &= \sum_{i=1, N} \log(p_h^{x_i} (1 - p_h)^{1 - x_i}) \\ &= \sum_{i=1, N} x_i \log p_h + (1 - x_i) \log(1 - p_h) \\ &= n_h \log p_h + n_t \log(1 - p_h)\end{aligned}$$

$$\begin{aligned}p_{h \text{ ML}} &= \operatorname{argmax}_p \mathcal{L}(p_h | N) \\ \rightarrow \frac{\partial \mathcal{L}(p_h | N)}{\partial p_h} &= \frac{n_h}{p_{h \text{ ML}}} - \frac{n_t}{1 - p_{h \text{ ML}}} = 0\end{aligned}$$

## Maximum a posteriori estimation (MAP)

Similar to ML estimation, but incorporates some prior knowledge about the parameters.

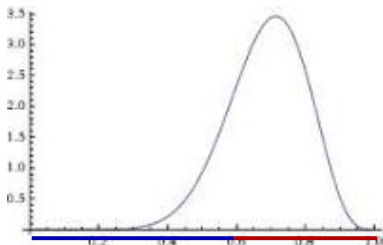
$$p(\mathcal{O}, \mu) = \underbrace{p(\mathcal{O}|\mu)}_{\text{likelihood}} \underbrace{p(\mu)}_{\text{prior}} = \underbrace{p(\mu|\mathcal{O})}_{\text{posterior}} \underbrace{p(\mathcal{O})}_{\text{evidence}}$$
$$p(\mu|\mathcal{O}) = \frac{p(\mathcal{O}|\mu) p(\mu)}{p(\mathcal{O})}$$

## Maximum a posteriori estimation (MAP)

Similar to ML estimation, but incorporates some prior knowledge about the parameters.

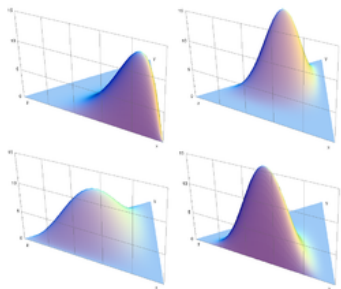
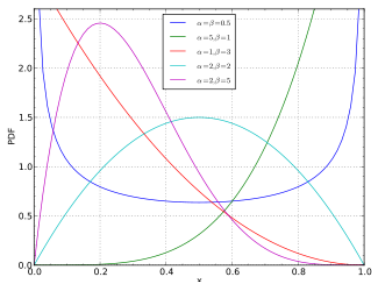
$$\begin{aligned}\mu_{MAP} &= \operatorname{argmax}_{\mu} p(\mu|\mathcal{O}) \\ &= \operatorname{argmax}_{\mu} \frac{p(\mathcal{O}|\mu)p(\mu)}{p(\mathcal{O})} \\ &= \operatorname{argmax}_{\mu} \underbrace{p(\mathcal{O}|\mu)}_{\text{likelihood}} \underbrace{p(\mu)}_{\text{prior}} \\ &= \operatorname{argmax}_{\mu} \mathcal{L}(\mathcal{O}|\mu) + \log p(\mu)\end{aligned}$$

## About priors – Beta/Dirichlet



$$p(p_h | \alpha_h, \alpha_t) = \text{Beta}(p_h | \alpha_h, \alpha_t) = \frac{1}{B(\alpha_h, \alpha_t)} p_h^{\alpha_h - 1} (1 - p_h)^{\alpha_t - 1}$$

## About priors – Beta/Dirichlet



$$p(p_h | \alpha_h, \alpha_t) = \text{Beta}(p_h | \alpha_h, \alpha_t) = \frac{1}{B(\alpha_h, \alpha_t)} p_h^{\alpha_h - 1} (1 - p_h)^{\alpha_t - 1}$$

## Maximum a posteriori estimation (MAP) – example

$N$  coin tossing experiments, head/tails;  $N = n_h + n_t$

$p_h$  – probability to get a head

Prior belief about the coin:  $p(p_h) = \text{Beta}(p_h|\alpha_h, \alpha_t)$

$$\begin{aligned}\mu_{MAP} &= \operatorname{argmax}_{\mu} \mathcal{L}(p_h|N) + \log p(p_h) \\ &\rightarrow \frac{n_h}{p_h} - \frac{n_t}{1-p_h} + \frac{\alpha_h - 1}{p_h} - \frac{\alpha_t - 1}{1-p_h} = 0 \\ \mu_{MAP} &= \frac{n_h + \alpha_h - 1}{N + \alpha_h - 1 + \alpha_t - 1}\end{aligned}$$



## Bayesian estimation

Similar to MAP estimation, but instead of maximizing to obtain the model, we induce a distribution of the model:

$$p(\mu|\mathcal{O}) = \frac{p(\mathcal{O}|\mu)p(\mu)}{p(\mathcal{O})}$$

The probability of the observation ( $\mathcal{O}$ ) is the expected value, according to all possible model variations, given its prior:

$$p(\mathcal{O}) = \int_{\mu \in \theta} p(\mathcal{O}|\mu)p(\mu)d\mu$$

## More about priors – conjugate distributions

$$p(\mu|\mathcal{O}) = \frac{p(\mathcal{O}|\mu)p(\mu)}{\int_{\mu \in \theta} p(\mathcal{O}|\mu)p(\mu)d\mu}$$

## More about priors – conjugate distributions

$$p(\mu|\mathcal{O}) = \frac{p(\mathcal{O}|\mu)p(\mu)}{\int_{\mu \in \theta} p(\mathcal{O}|\mu)p(\mu)d\mu}$$

A conjugate prior  $p(\mu)$  of a likelihood  $p(\mathcal{O}|\mu)$  is a distribution that results in a posterior distribution  $p(\mu|\mathcal{O})$  with the same functional form as the prior, and a parametrisation that incorporates the observations  $\mathcal{O}$ . (they are conjugate if they have the same functional form)

Discrete	Continuous
Bernoulli	Beta
$p(x) = q^x(1 - q)^{1-x}$	$p(q) = \frac{1}{B(\alpha, \beta)} q^{\alpha-1}(1 - q)^{\beta-1}$
$B(\alpha, \beta) = \frac{\Gamma(\alpha)\Gamma(\beta)}{\Gamma(\alpha + \beta)} = \int_0^1 q^{\alpha-1}(1 - q)^{\beta-1}$	

## More about priors – conjugate distributions

$$p(\mu|\mathcal{O}) = \frac{p(\mathcal{O}|\mu)p(\mu)}{\int_{\mu \in \theta} p(\mathcal{O}|\mu)p(\mu)d\mu}$$

Discrete	Continuous
Bernoulli	Beta
$p(x) = q^x(1-q)^{1-x}$	$p(q) = \frac{1}{B(\alpha,\beta)} q^{\alpha-1}(1-q)^{\beta-1}$
Multinomial	Dirichlet
...	...
Poisson	Gamma
$p(x) = \frac{\lambda^x}{x!} e^{-\lambda}$	$p(\lambda) = \frac{1}{\theta^k} \frac{\lambda^{k-1}}{\Gamma(k)} e^{-\frac{\lambda}{\theta}}$

## Probability of the observation – coin tossing

$$p(\mathcal{O}|\alpha_h, \alpha_t) = \int_0^1 p(\mathcal{O}|p_h)p(p_h|\alpha_h, \alpha_t)dp_h$$

## Probability of the observation – coin tossing

$$\begin{aligned} p(\mathcal{O}|\alpha_h, \alpha_t) &= \int_0^1 p(\mathcal{O}|p_h)p(p_h|\alpha_h, \alpha_t)dp_h \\ &= \int_0^1 p_h^{n_h}(1-p_h)^{n_t} \frac{1}{B(\alpha_h, \alpha_t)} p_h^{\alpha_h-1}(1-p_h)^{\alpha_t-1} dp_h \\ &= \frac{1}{B(\alpha_h, \alpha_t)} \int_0^1 p_h^{n_h+\alpha_h-1}(1-p_h)^{n_t+\alpha_t-1} dp_h \\ &= \frac{1}{B(\alpha_h, \alpha_t)} B(n_h + \alpha_h, n_t + \alpha_t) \end{aligned}$$

## Probability of the observation – coin tossing

$$\begin{aligned} p(\mathcal{O}|\alpha_h, \alpha_t) &= \int_0^1 p(\mathcal{O}|p_h)p(p_h|\alpha_h, \alpha_t)dp_h \\ &= \int_0^1 p_h^{n_h}(1-p_h)^{1-n_t} \frac{1}{B(\alpha_h, \alpha_t)} p_h^{\alpha_h-1}(1-p_h)^{\alpha_t-1} dp_h \\ &= \frac{1}{B(\alpha_h, \alpha_t)} \int_0^1 p_h^{n_h+\alpha_h-1}(1-p_h)^{n_t+\alpha_t-1} dp_h \\ &= \frac{1}{B(\alpha_h, \alpha_t)} B(n_h + \alpha_h, n_t + \alpha_t) \end{aligned}$$

$$B(\alpha_h, \alpha_t) = \frac{\Gamma(\alpha_h)\Gamma(\alpha_t)}{\Gamma(\alpha_h + \alpha_t)}$$

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

## Predicting a new observation

$$p(x = 1 | \mathcal{O}, \alpha_h, \alpha_t) = \frac{p(x = 1, \mathcal{O} | \alpha_h, \alpha_t)}{p(\mathcal{O} | \alpha_h, \alpha_t)}$$

$$p(\mathcal{O} | \alpha_h, \alpha_t) = \frac{1}{B(\alpha_h, \alpha_t)} B(n_h + \alpha_h, n_t + \alpha_t)$$



## Predicting a new observation

$$\begin{aligned} p(x = 1 | \mathcal{O}, \alpha_h, \alpha_t) &= \frac{p(x = 1, \mathcal{O} | \alpha_h, \alpha_t)}{p(\mathcal{O} | \alpha_h, \alpha_t)} \\ &= \frac{\frac{1}{B(\alpha_h, \alpha_t)} B(n_h + 1 + \alpha_h, n_t + \alpha_t)}{\frac{1}{B(\alpha_h, \alpha_t)} B(n_h + \alpha_h, n_t + \alpha_t)} \end{aligned}$$

## Predicting a new observation

$$\begin{aligned} p(x = 1 | \mathcal{O}, \alpha_h, \alpha_t) &= \frac{p(x = 1, \mathcal{O} | \alpha_h, \alpha_t)}{p(\mathcal{O} | \alpha_h, \alpha_t)} \\ &= \frac{\frac{1}{B(\alpha_h, \alpha_t)} B(n_h + 1 + \alpha_h, n_t + \alpha_t)}{\frac{1}{B(\alpha_h, \alpha_t)} B(n_h + \alpha_h, n_t + \alpha_t)} \\ &= \frac{B(n_h + 1 + \alpha_h, n_t + \alpha_t)}{B(n_h + \alpha_h, n_t + \alpha_t)} \\ &= \frac{\Gamma(n_h + 1 + \alpha_h) \Gamma(n_t + \alpha_t)}{\Gamma(N + 1 + \alpha_h + \alpha_t)} \\ &= \frac{\Gamma(n_h + \alpha_h) \Gamma(n_t + \alpha_t)}{\Gamma(N + \alpha_h + \alpha_t)} \end{aligned}$$

## Predicting a new observation

$$\begin{aligned} p(x = 1 | \mathcal{O}, \alpha_h, \alpha_t) &= \frac{p(x = 1, \mathcal{O} | \alpha_h, \alpha_t)}{p(\mathcal{O} | \alpha_h, \alpha_t)} \\ &= \frac{\frac{1}{B(\alpha_h, \alpha_t)} B(n_h + 1 + \alpha_h, n_t + \alpha_t)}{\frac{1}{B(\alpha_h, \alpha_t)} B(n_h + \alpha_h, n_t + \alpha_t)} \\ &= \frac{B(n_h + 1 + \alpha_h, n_t + \alpha_t)}{B(n_h + \alpha_h, n_t + \alpha_t)} \\ &= \frac{\Gamma(n_h + 1 + \alpha_h) \Gamma(n_t + \alpha_t)}{\Gamma(N + 1 + \alpha_h + \alpha_t)} \\ &= \frac{\Gamma(n_h + \alpha_h) \Gamma(n_t + \alpha_t)}{\Gamma(N + \alpha_h + \alpha_t)} \\ &= \frac{n_h + \alpha_h}{N + \alpha_h + \alpha_t} \end{aligned}$$

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

## Bayesian estimation – example

$\mathcal{O}$ : A sequence of  $N$  coin tosses,  $N = n_h + n_t$

$\mu$ :  $p_h | \alpha_h, \alpha_t$

$$\begin{aligned} p(\mu | \mathcal{O}) &= \frac{p(\mathcal{O} | \mu) p(\mu)}{p(\mathcal{O})} \\ &= \frac{\prod_{i=1, N} p(X = x_i | p_h) p(p_h | \alpha_h, \alpha_t)}{\int_0^1 \prod_{i=1, N} p(X = x_i | p_h) p(p_h | \alpha_h, \alpha_t) dp_h} \end{aligned}$$

## Bayesian estimation – example

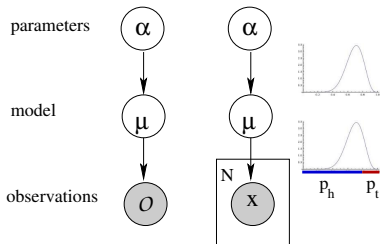
$\mathcal{O}$ : A sequence of  $N$  coin tosses,  $N = n_h + n_t$

$\mu$ :  $p_h | \alpha_h, \alpha_t$

$$\begin{aligned} p(\mu | \mathcal{O}) &= \frac{p(\mathcal{O} | \mu) p(\mu)}{p(\mathcal{O})} \\ &= \frac{\prod_{i=1, N} p(X = x_i | p_h) p(p_h | \alpha_h, \alpha_t)}{\int_0^1 \prod_{i=1, N} p(X = x_i | p_h) p(p_h | \alpha_h, \alpha_t) dp_h} \\ &= \frac{p_h^{n_h} (1 - p_h)^{n_t} \frac{1}{B(\alpha_h, \alpha_t)} p_h^{\alpha_h - 1} (1 - p_h)^{\alpha_t - 1}}{\frac{1}{B(\alpha_h, \alpha_t)} B(n_h + \alpha_h, n_t + \alpha_t)} \\ &= \frac{1}{B(n_h + \alpha_h, n_t + \alpha_t)} p_h^{n_h + \alpha_h - 1} (1 - p_h)^{n_t + \alpha_t - 1} \\ &= \text{Beta}(p_h | n_h + \alpha_h, n_t + \alpha_t) \end{aligned}$$

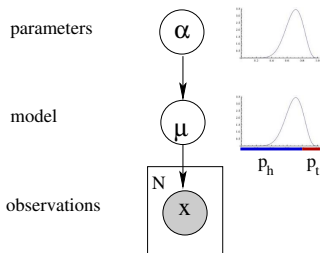
# Bayesian network

Coin tossing

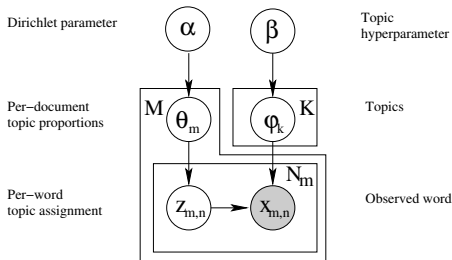


# Latent Dirichlet Allocation

## Coin tossing

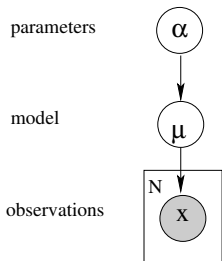


## Latent Dirichlet Allocation

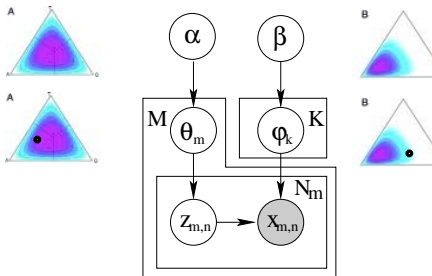


# Latent Dirichlet Allocation

Coin tossing

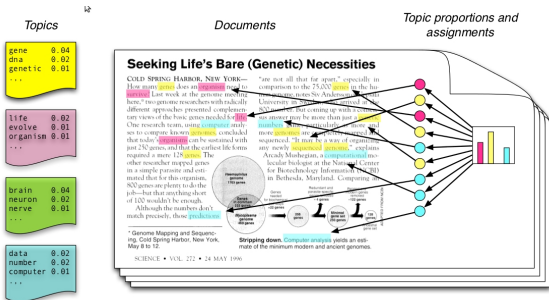


Latent Dirichlet Allocation





# Topic models



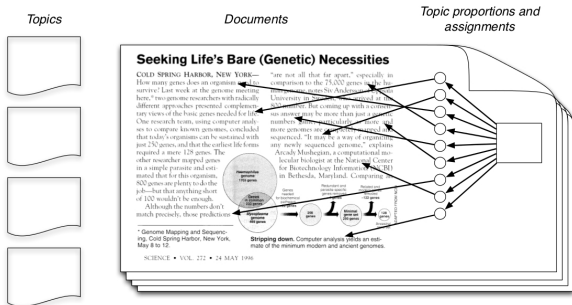
Each document is a mixture of topics:

$$\sum_k p(z_m = k) = \sum_k \theta_{m,k} = 1$$

Each word is drawn from one of its document's topics:

$$p(w_{m,n}) = \sum_k p(w_{m,n} | z_{m,n} = k) p(z_{m,n} = k) = \sum_k \phi_k(w_{m,n}) \theta_{m,k}$$

# Topic models



The **observations** are the documents:  $\mathbf{w}_m$ ,  $m \in 1, M$

We need to infer the **model**, i.e the underlying topic structure, i.e. the topic assignments  $z_{m,n}$ , the topic  $\theta_m$ ,  $m \in 1, M$  and word distributions  $\phi_k$ ,  $k \in 1, K$

**Priors:**

$\theta \sim$  distribution with hyperparameter  $\alpha$

$\phi \sim$  distribution with hyperparameter  $\beta$

## Topic models – Latent Dirichlet Allocation

$$p(\theta|\alpha) = \frac{1}{B(\alpha)} \prod_k \theta_k^{\alpha_k - 1}$$

$$\sum_k \theta_{m,k} = 1$$

$\alpha$  controls the mean shape and sparsity of  $\theta$

The topic proportions ( $\theta_m$ ) are a K-dimensional Dirichlet

$z_{m,n}$  are multinomial distributions from  $\theta_m$

$$p(z_{m,n}|\theta_m) = \frac{N!}{\prod_{k=1}^K n_k!} \prod_{k=1}^K \theta_{m,k}^{n_k}$$

## Topic models – Latent Dirichlet Allocation

$$p(\phi|\beta) = \frac{1}{B(\beta)} \prod_{\nu} \phi_{\nu}^{\beta_{\nu}-1}$$

$$\sum_{\nu} \phi_{k,\nu} = 1$$

$\beta$  controls the mean shape and sparsity of  $\phi$

The topics ( $\phi_k$ ) are a V-dimensional Dirichlet

$w_{m,n}$  are multinomial distributions from  $\phi_{z_{m,n}}$

$$p(w_{m,n}|\phi_k) = \frac{V!}{\prod_{\nu=1}^V n_{\nu}!} \prod_{\nu=1}^V \phi_{k,\nu}^{n_{\nu}}$$

## Topic models

Remember:  $p(x = 1|\mathcal{O}, \alpha_h, \alpha_t) = \frac{p(x=1, \mathcal{O}|\alpha_h, \alpha_t)}{p(\mathcal{O}|\alpha_h, \alpha_t)} = \frac{n_h + \alpha_h}{N + \alpha_h + \alpha_t}$

$$p(z_i = k, w_i | \mathbf{z}_{-i}) = p(z_i = k | \mathbf{z}_{-i}) p(w_i | z_i)$$

$$\begin{aligned} p(z_i = k | \mathbf{z}_{-i}, \alpha) &= p(z_i = k | \mathbf{z}_{-i}, \alpha) \\ &= \frac{nz_{m,-i}^k + \alpha_k}{\sum_{j=1}^K (nz_{m,-i}^j + \alpha_j)} \\ &= \hat{\theta}_{m,k} \quad \text{an approximation of } \theta_{m,k} \end{aligned}$$

$$\begin{aligned} p(w_i | z_i, \beta) &= \frac{nw_{w_i,-i}^{z_i} + \beta_{w_i}}{\sum_{l=1}^V (nw_l^{z_i} + \beta_l)} \\ &= \hat{\phi}_k(w_i) \quad \text{an approximation of } \phi_k(w_i) \end{aligned}$$

Sample  $z_i$  and assign it at position  $i$ :  $z_i \propto \hat{\theta}_{m,k} \hat{\phi}_k(w_i)$

# Inference via Gibbs sampling

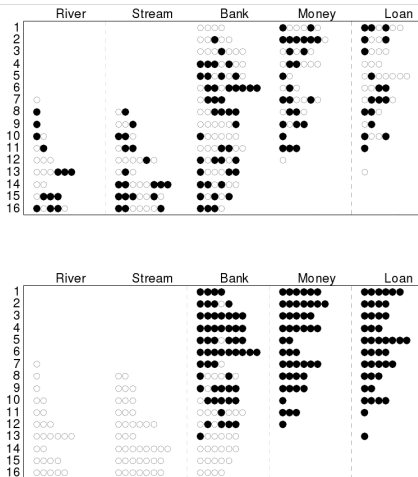


Figure 7. An example of the Gibbs sampling procedure.

Sample  $z_i$  and assign it at position  $i$ :  $z_i \propto \hat{\theta}_{m,k} \hat{\phi}_k(w_i)$

## References

- *Probabilistic topic models*, Mark Steyvers, Tom Griffiths
- *Parameter estimation for text analysis*, Gregor Heinrich
- *Topic Models*, David Blei (tutorial, [videolectures.net](http://videolectures.net))