

Latent Dirichlet Allocation



Topics

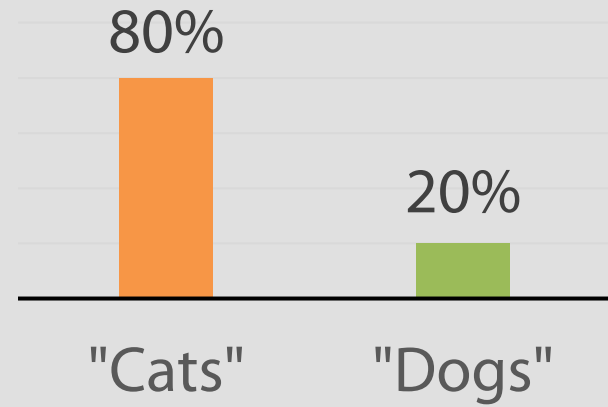
Document is a distribution over topics



Topics

Document is a distribution over topics

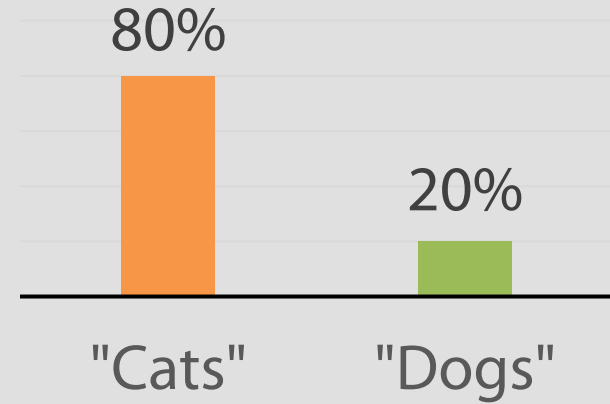
Document



Topics

Document is a distribution over topics

Document



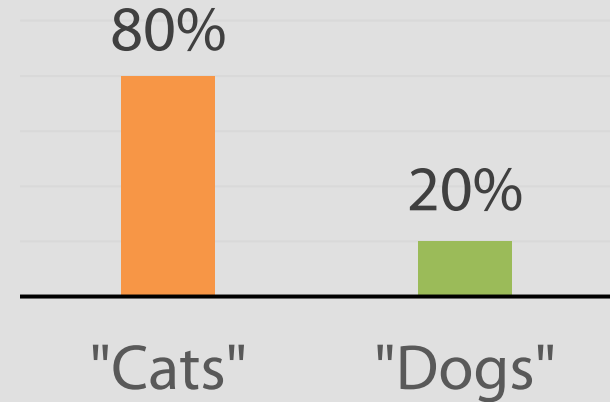
Topic is a distribution over words



Topics

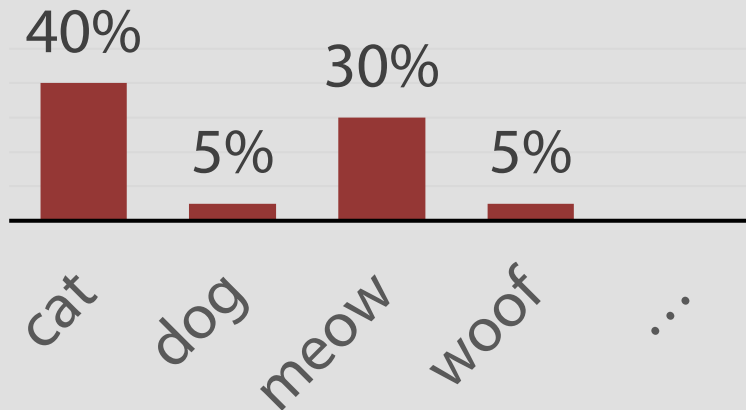
Document is a distribution over topics

Document

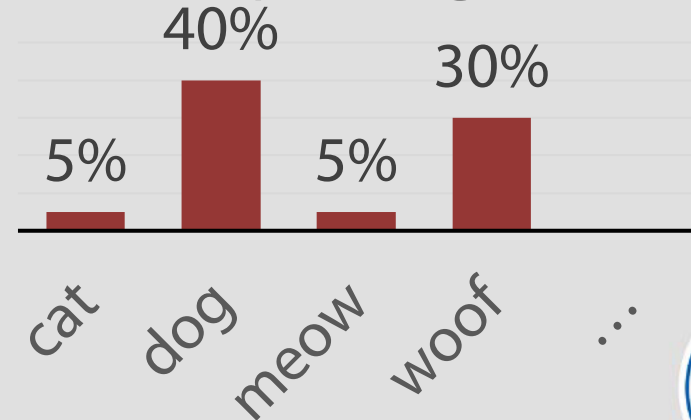


Topic is a distribution over words

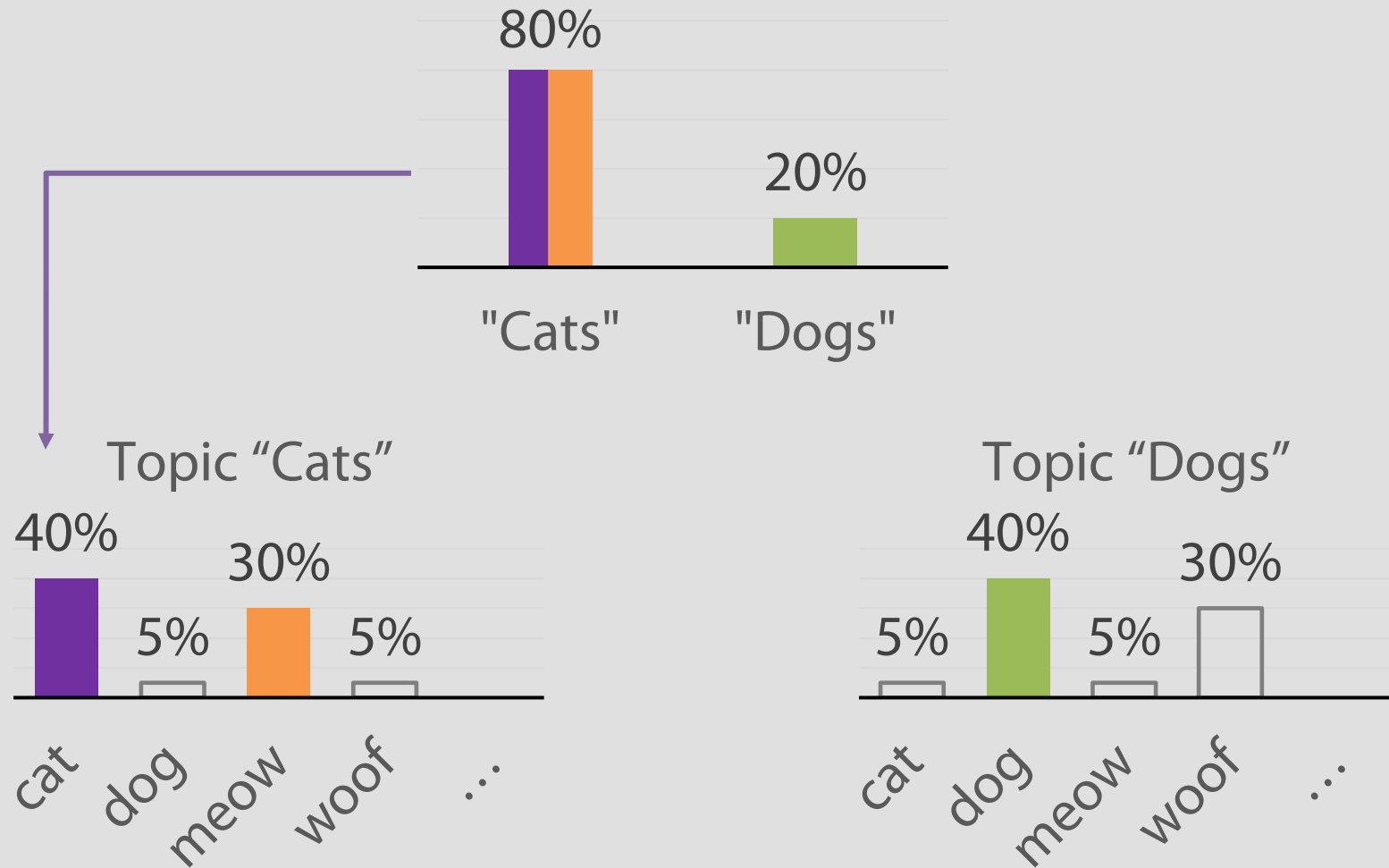
Topic "Cats"



Topic "Dogs"



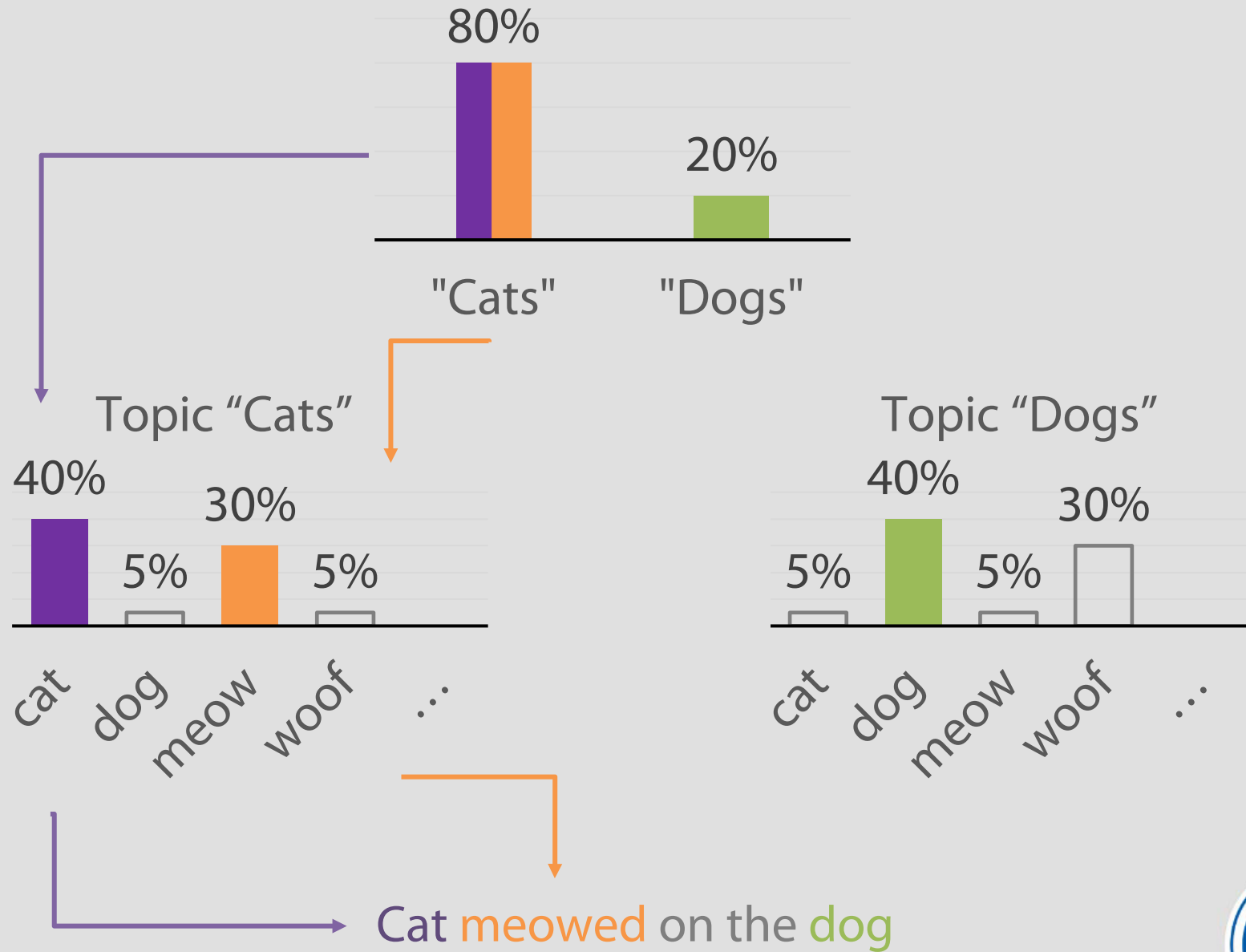
Text generation



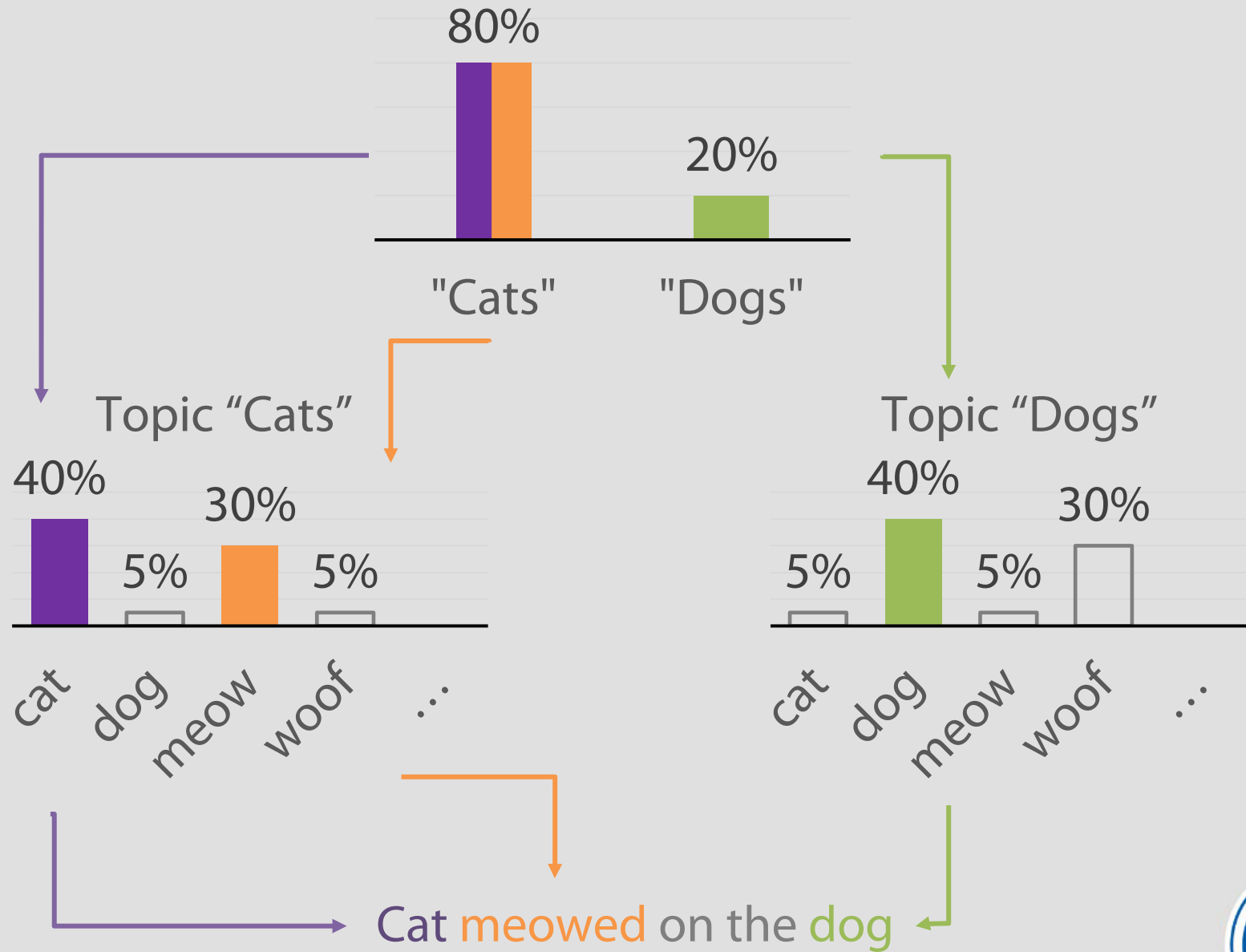
Cat meowed on the dog



Text generation

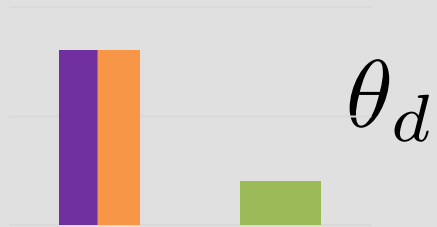


Text generation



Model

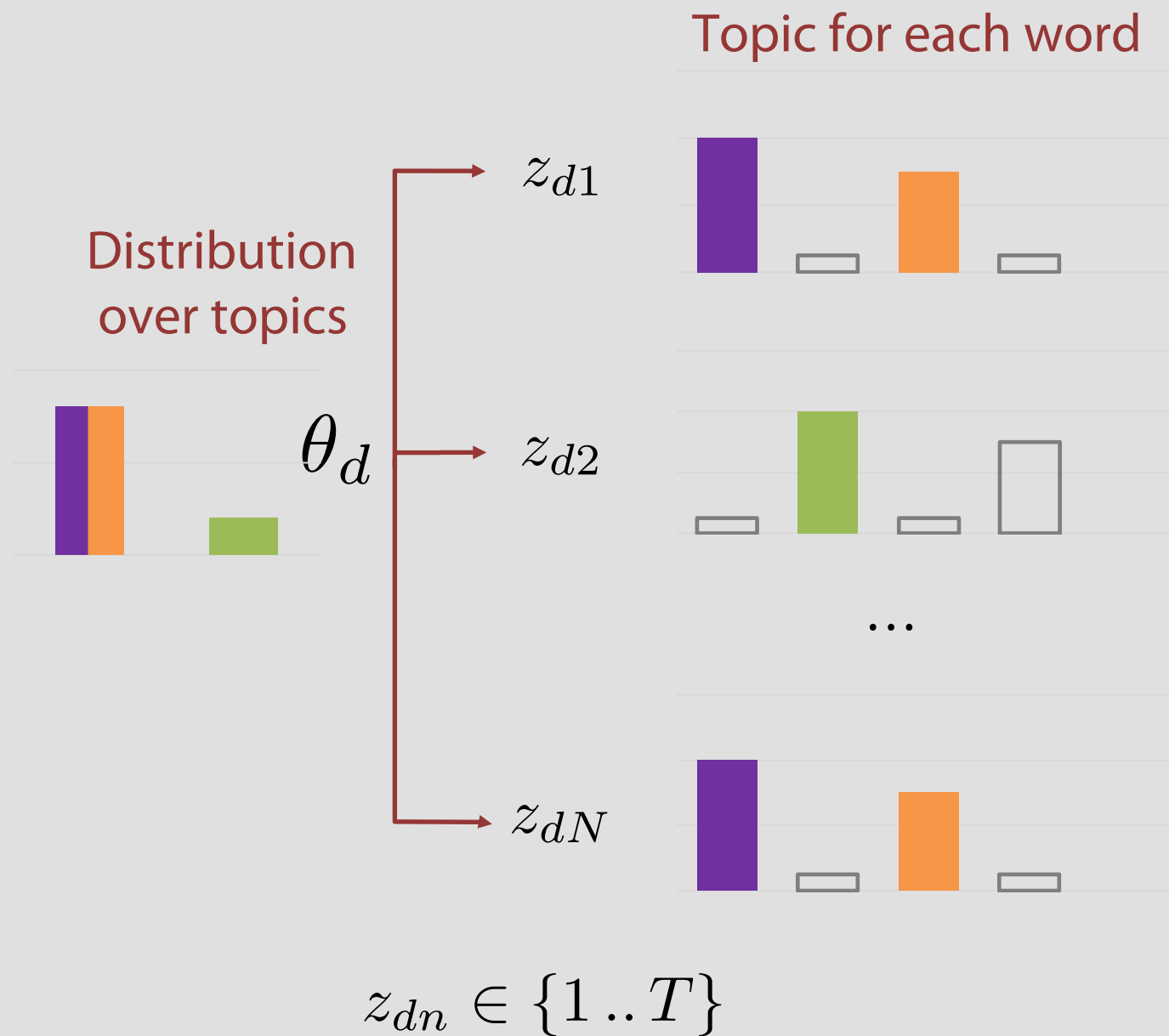
Distribution
over topics



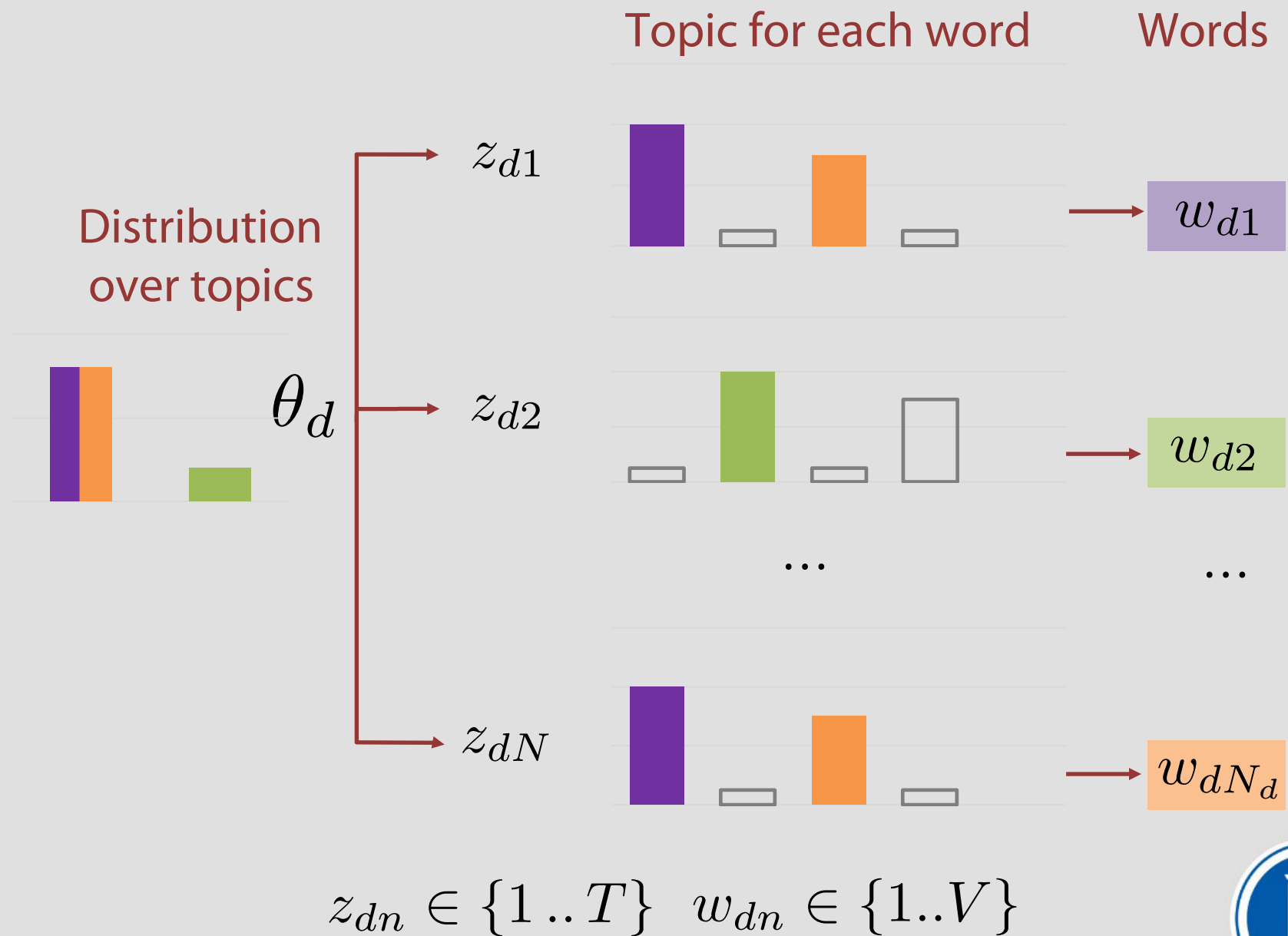
θ_d



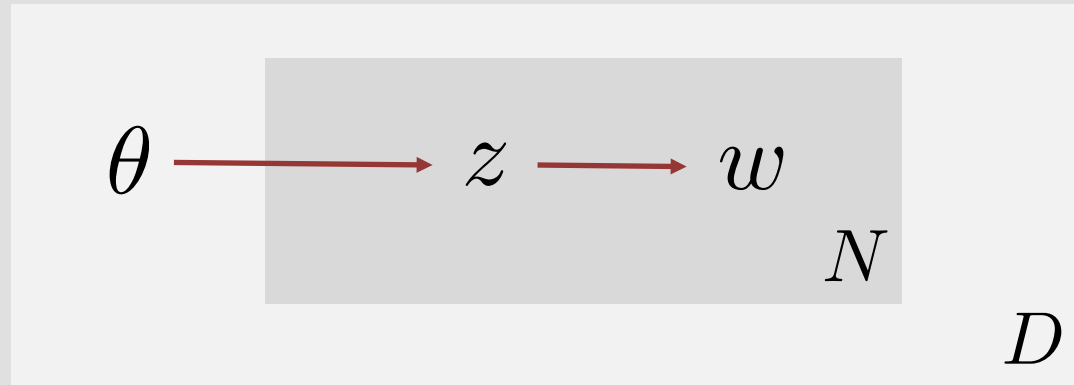
Model



Model



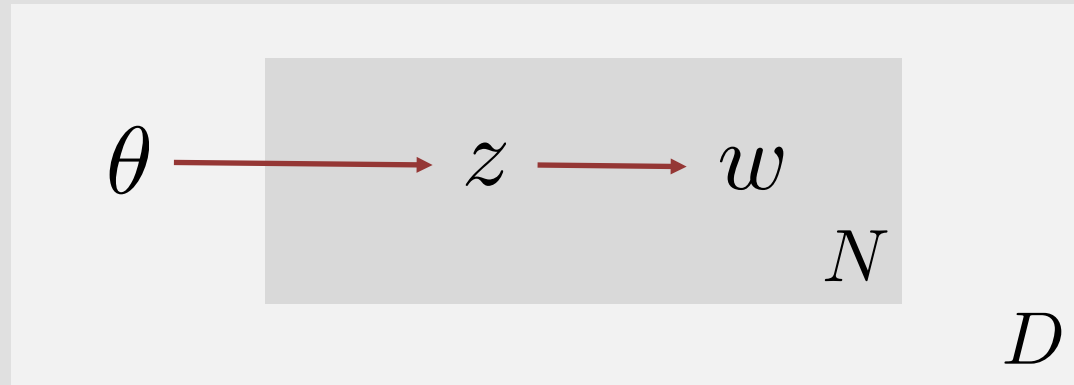
LDA Model



$$p(W, Z, \Theta) = \prod_{d=1}^D p(\theta_d) \prod_{n=1}^{N_d} p(z_{dn} | \theta_d) p(w_{dn} | z_{dn})$$



LDA Model

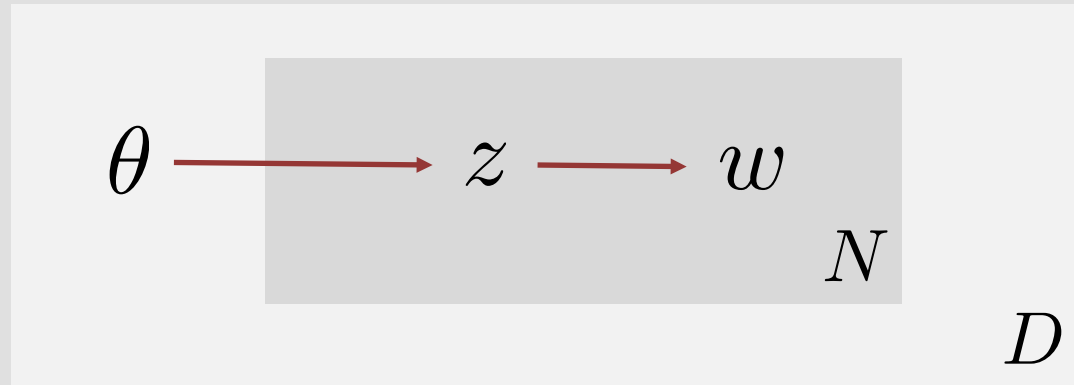


$$p(W, Z, \Theta) = \prod_{d=1}^D p(\theta_d) \prod_{n=1}^{N_d} p(z_{dn} | \theta_d) p(w_{dn} | z_{dn})$$

for each document \uparrow



LDA Model

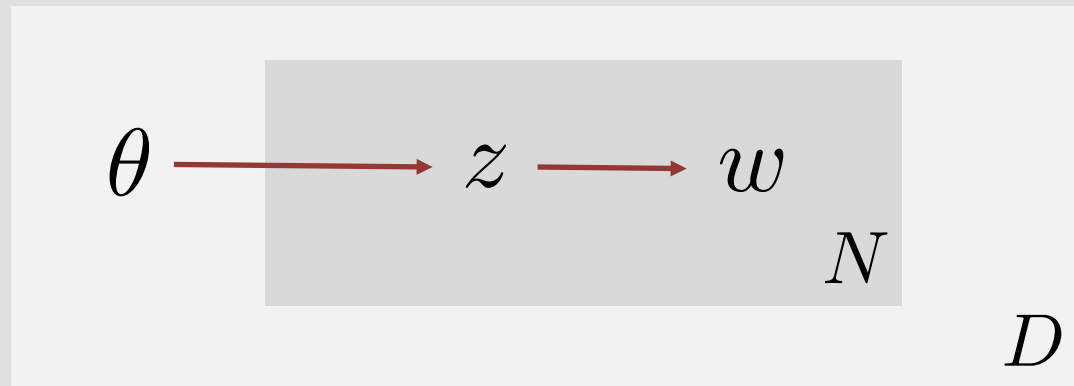


$$p(W, Z, \Theta) = \prod_{d=1}^D p(\theta_d) \prod_{n=1}^{N_d} p(z_{dn} | \theta_d) p(w_{dn} | z_{dn})$$

for each document \uparrow \uparrow generate topic probabilities



LDA Model



for each word

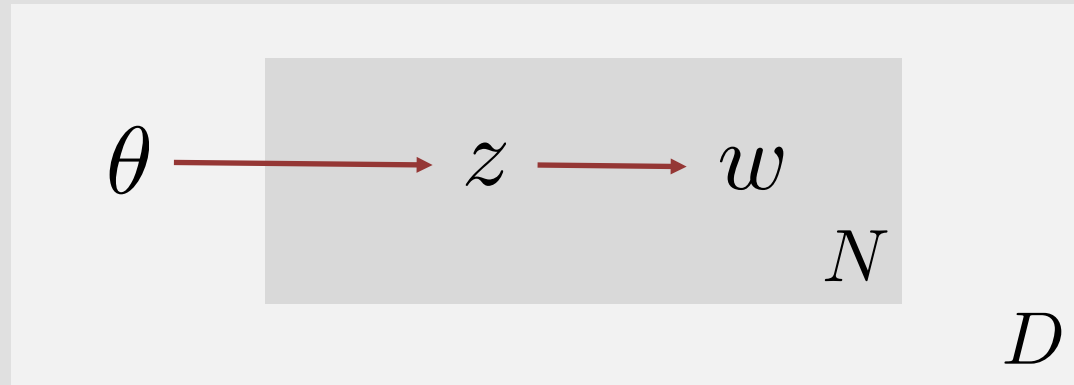
$$p(W, Z, \Theta) = \prod_{d=1}^D p(\theta_d) \prod_{n=1}^{N_d} p(z_{dn} | \theta_d) p(w_{dn} | z_{dn})$$

for each document

generate topic probabilities



LDA Model



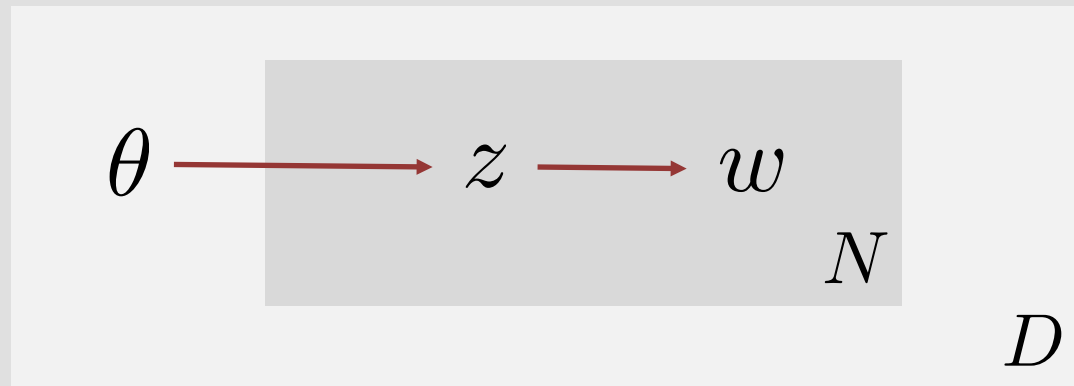
for each word select topic

$$p(W, Z, \Theta) = \prod_{d=1}^D p(\theta_d) \prod_{n=1}^{N_d} p(z_{dn} | \theta_d) p(w_{dn} | z_{dn})$$

for each document generate topic probabilities



LDA Model



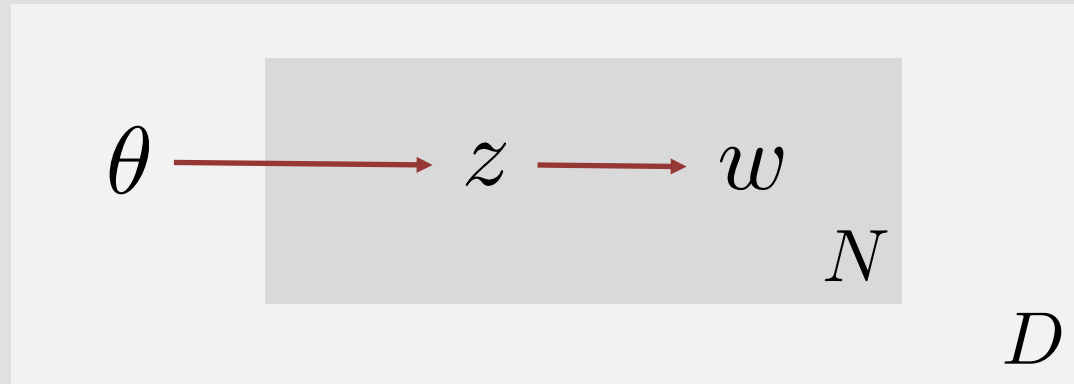
for each word select topic select word from topic

$$p(W, Z, \Theta) = \prod_{d=1}^D p(\theta_d) \prod_{n=1}^{N_d} p(z_{dn} | \theta_d) p(w_{dn} | z_{dn})$$

for each document generate topic probabilities



LDA Model



for each word select topic select word from topic

$$p(W, Z, \Theta) = \prod_{d=1}^D p(\theta_d) \prod_{n=1}^{N_d} p(z_{dn} | \theta_d) p(w_{dn} | z_{dn})$$

for each document generate topic probabilities



LDA Model

$$p(\textcolor{red}{W}, \textcolor{green}{Z}, \textcolor{blue}{\Theta}) = \prod_{d=1}^D \textcolor{blue}{p}(\theta_d) \prod_{n=1}^{N_d} \textcolor{green}{p}(z_{dn} | \theta_d) \textcolor{red}{p}(w_{dn} | z_{dn})$$



LDA Model

$$p(\textcolor{red}{W}, \textcolor{green}{Z}, \textcolor{blue}{\Theta}) = \prod_{d=1}^D \textcolor{blue}{p(\theta_d)} \prod_{n=1}^{N_d} \textcolor{green}{p(z_{dn}|\theta_d)} \textcolor{red}{p(w_{dn}|z_{dn})}$$

$$\textcolor{blue}{p(\theta_d)} \sim \text{Dir}(\alpha)$$



LDA Model

$$p(\mathbf{W}, \mathbf{Z}, \Theta) = \prod_{d=1}^D p(\theta_d) \prod_{n=1}^{N_d} p(z_{dn} | \theta_d) p(w_{dn} | z_{dn})$$

$$p(\theta_d) \sim \text{Dir}(\alpha)$$

$$p(z_{dn} | \theta_d) = \theta_{dz_{dn}}$$



LDA Model

$$p(\mathbf{W}, \mathbf{Z}, \Theta) = \prod_{d=1}^D p(\theta_d) \prod_{n=1}^{N_d} p(z_{dn} | \theta_d) p(w_{dn} | z_{dn})$$

$$p(\theta_d) \sim \text{Dir}(\alpha)$$

$$p(z_{dn} | \theta_d) = \theta_{dz_{dn}}$$

$$p(w_{dn} | z_{dn}) = \Phi_{z_{dn} w_{dn}}$$



LDA Model

$$p(\mathbf{W}, \mathbf{Z}, \Theta) = \prod_{d=1}^D p(\theta_d) \prod_{n=1}^{N_d} p(z_{dn} | \theta_d) p(w_{dn} | z_{dn})$$

$$p(\theta_d) \sim \text{Dir}(\alpha)$$

$$p(z_{dn} | \theta_d) = \theta_{dz_{dn}}$$

$$p(w_{dn} | z_{dn}) = \Phi_{z_{dn} w_{dn}} \longleftarrow \sum_w \Phi_{tw} = 1$$

Constraints:

$$\Phi_{tw} \geq 0$$



LDA Model

Known: W data

Unknown: Φ parameters, distribution over words for each topic

Unknown: Z latent variables, topic of each word

Unknown: Θ latent variables, distribution over topics for each document



ТЕХНИЧЕСКИЙ СЛАЙД (15 мин на доску)

- ВЫВОД ФОРМУЛ VAR. EM НА ДОСКЕ

