

Latent Dirichlet Allocation LDA

Тематическая модель

$$p(\theta | \alpha) = \text{Dir}(\theta | \alpha) = \frac{\Gamma(\sum_j \alpha_j)}{\prod_j \Gamma(\alpha_j)} \prod_{j=1}^m \theta_j^{\alpha_j - 1}$$

$$\theta \in \mathbb{R}^m, \theta_j \geq 0, \sum_{j=1}^m \theta_j = 1$$

$$= B^{-1}(\alpha_1, \dots, \alpha_m)$$

$\ln \theta_j$ - гостомонические статистики

$$\mathbb{E} \ln \theta_j = - \frac{\partial}{\partial \alpha_j} \ln B^{-1}(\alpha_1, \dots, \alpha_m) = - \frac{\partial}{\partial \alpha_j} \ln \Gamma(\sum_j \alpha_j) +$$

$$+ \frac{\partial}{\partial \alpha_j} \ln \Gamma(\alpha_j) = \psi(\alpha_j) - \psi(\sum_j \alpha_j), \psi(x) = \frac{d}{dx} \ln \Gamma(x)$$

D -гоукменгов из N_d -слов

$\{1 \dots W\} \ni w_{d,n}$ - номер n -го ^{слова} ~~гоукменгов~~ в гок. d

$\{1 \dots T\} \ni z_{d,n}$ - номер ~~менн~~ гок n -го слова в гок. d .

$S^T \ni \theta_d$ - прогрок гоукменгов

$\varphi_{w,z}$ - вер-ть слова w в ~~гоукменгов~~ ^{D} ~~менн~~ z

$$\varphi \in \mathbb{R}^{W \times T} \quad \begin{matrix} \uparrow \\ T \\ \downarrow \\ W \end{matrix} \quad p(w, z, \theta | \varphi, \alpha) = \prod_{d=1}^D \left(p(\theta_d | \alpha) \cdot \prod_{n=1}^{N_d} p(z_{dn} | \theta_d) \cdot p(w_{dn} | z_{dn}, \varphi) \right) =$$

$$= \prod_{d=1}^D \left(\text{Dir}(\theta_d | \alpha) \prod_{n=1}^{N_d} \prod_{t=1}^T \theta_{dt}^{[z_{dn}=t]} \varphi_{w_{dn} z_{dn}} \{ \varphi_{w_{dn}}^{[z_{dn}=t]} \} \right) =$$

$$= \prod_{d=1}^D \left(\text{Dir}(\theta_d | \alpha) \prod_{n=1}^{N_d} \prod_{t=1}^T (\theta_{dt} \varphi_{w_{dn}})^{[z_{dn}=t]} \right)$$

$$p(w | \varphi, \alpha) \rightarrow \max_{\varphi} \quad , \quad \ln p(w | \varphi, \alpha) \rightarrow \max_{\varphi}$$

$$E\text{-step} \quad p(w, z, \theta | \Phi, \alpha) \approx q(z) q(\theta)$$

$$M\text{-step} \quad \mathbb{E}_{z, \theta} \ln p(w, z, \theta | \Phi, \alpha)$$

$$\ln q(\theta) = \mathbb{E}_{q(z)} \ln p(w, z, \theta | \Phi, \alpha) + C =$$

$$= \mathbb{E}_{q(z)} \left(\sum_{d=1}^D \sum_{t=1}^T (\alpha-1) \ln \theta_{dt} + \sum_{n=1}^{N_d} \sum_{t=1}^T \sum_{d=1}^D [z_{dn}=t] \ln \theta_{dt} \right) + C =$$

$$= \sum_{d=1}^D \sum_{t=1}^T \ln \theta_{dt} [(\alpha-1) + \sum_{n=1}^{N_d} \mathbb{E}_{z_{dn}} [z_{dn}=t]] + C =$$

$$= \sum_{d=1}^D \sum_{t=1}^T \ln \theta_{dt} [(\alpha-1) + \sum_{n=1}^{N_d} \delta_{dnt}] + C$$

$$q(\theta) = \prod_{d=1}^D \prod_{t=1}^T q(\theta_{dt}), \quad q(\theta_{dt}) = \text{Dir}(\theta_{dt} | \alpha + \sum_{n=1}^{N_d} \delta_{dnt})$$

$$q(\theta) = \prod_{d=1}^D q(\theta_d) = \prod_{d=1}^D \text{Dir}(\theta_d | \alpha_d)$$

$$\ln q(z) = \mathbb{E}_{q(\theta)} \ln p(w, z, \theta | \Phi, \alpha) + C =$$

$$= \mathbb{E}_{q(\theta)} \left(\sum_{d=1}^D \sum_{t=1}^T (\alpha-1) \ln \theta_{dt} + \sum_{d=1}^D \sum_{n=1}^{N_d} \sum_{t=1}^T [z_{dn} \neq t] (\ln \theta_{dt} + \ln \Phi_{w_{dn}, t}) \right) + C$$

$$= \sum_{d=1}^D \sum_{n=1}^{N_d} \sum_{t=1}^T [z_{dn}=t] (\mathbb{E}_{q(\theta_{dt})} \ln \theta_{dt} + \ln \Phi_{w_{dn}, t}) + C$$

$$q(z) = \prod_{d=1}^D \prod_{n=1}^{N_d} q(z_{dn})$$

$$\ln q(z_{dn}=t) = \mathbb{E}_{q(\theta_{dt})} \ln \theta_{dt} + \ln \Phi_{w_{dn}, t} + C$$

$$q(z_{dn}=t) = \frac{\Phi_{w_{dn}, t} \exp(\mathbb{E}_{q(\theta_{dt})} \ln \theta_{dt})}{\sum_{s=1}^T \Phi_{w_{dn}, s} \exp(\mathbb{E}_{q(\theta_{ds})} \ln \theta_{ds})} = \delta_{dnt}$$

наиболее вероятная тема t в словаре d

M-step

$$\mathbb{E}_{z, \theta} \ln p(w, z, \theta | \varphi, \alpha) = \mathbb{E}_{z, \theta} \sum_{d=1}^D [\ln p(\theta_d | \alpha) +$$

$$+ \sum_{n=1}^{N_d} \sum_{t=1}^T [z_{dn}=t] (\ln \theta_{dt} + \ln \varphi_{w_{dn}, t})] = C +$$

$$+ \mathbb{E}_{z, \theta} \sum_{d=1}^D \sum_{n=1}^N \sum_{t=1}^T [z_{dn}=t] \ln \varphi_{w_{dn}, t} = C +$$

$$+ \sum_{d, n, t} \ln \varphi_{w_{dn}, t} \delta_{dnt} \rightarrow \max_{\varphi}$$

$$\varphi: \sum_{w=1}^V \varphi_{wt} = 1$$

$$\sum_{d, n, t} \delta_{dnt} \ln \varphi_{w_{dn}, t} \rightarrow \sum_{t=1}^T \lambda_t (\sum_{w=1}^V \varphi_{wt} - 1) \rightarrow \text{ext}_2$$

$$\frac{\partial \mathcal{L}}{\partial \varphi_{wt}} = \sum_{(d, n): w_{dn}=w} \sum_{t=1}^T \delta_{dnt} \frac{1}{\varphi_{wt}} - \lambda_t = 0$$

$$\varphi_{wt} = \frac{1}{\lambda_t} \sum_{(d, n): w_{dn}=w} \sum_{t=1}^T \delta_{dnt}$$

$$\lambda_t \varphi_{wt} = \sum_{(d, n): w_{dn}=w} \delta_{dnt} \quad | \quad \sum_{w=1}^V$$

$$\lambda_t = \sum_{w=1}^V \sum_{(d, n): w_{dn}=w} \delta_{dnt} = \sum_{d=1}^D \sum_{n=1}^N \delta_{dnt}$$

$$\varphi_{wt} = \frac{\sum_{d=1}^D \sum_{n=1}^N \delta_{dnt}}{\sum_{d=1}^D \sum_{n=1}^N \delta_{dnt}}$$

гора броуниана
слова в меню

мематический процесс по гуггеншта

$$q(d_*) q(z) \approx p(\theta_{d_*}, z | w_{d_*}, \Phi_d)$$

↘ Dir

$$p(w, z, \theta | \Phi, \alpha) = \prod_d (p(\theta_d | \alpha) \prod_{n=1}^{N_d} p(z_{dn} | \theta_d) \cdot p(w_{dn} | z_{dn}, \Phi))$$

$$p(\theta, z, \Phi | w, \alpha, \beta) \approx q(\theta) q(z) q(\Phi)$$

$$p(\theta_d | \alpha) = DP, \mu c$$

$$p(\theta_d | \theta_{d-1}, \alpha)$$

18.11.16 sumo sem

смесь Сторгенма

$$p(x, z, T | \mu, \Sigma, \gamma, \pi) = \prod_{n,k} [\pi_k N(x_n | \mu_k, \frac{1}{z_n} \Sigma_k) \cdot G(z_n | \frac{\gamma_k}{2}, \frac{\gamma_k}{2})]^{t_{nk}}$$

$q(z, T) \approx q(z) q(T)$
 $\approx q(z | T) q(T)$

$p(z, T), p(x | z, T)$, сравнение

$$\downarrow$$

$$p(z, T | x) = q(z, T) \text{ аналитический буг}$$

$$p(x | z, T) = \prod_{n,k} [N(x_n | \mu_k, \frac{1}{z_n} \Sigma_k)]^{t_{nk}} =$$

$$= \prod_{n,k} \left[\frac{z_n^{d/2}}{(2\pi)^{d/2} \sqrt{\det \Sigma_k}} \exp\left(-\frac{z_n}{2} (x_n - \mu_k)^T \Sigma_k^{-1} (x_n - \mu_k)\right) \right]^{t_{nk}} \neq$$

$$\neq \prod_{n,k} \left[\frac{z_n^c}{c} \exp(-z_n c) \right]^{t_{nk}}$$

$$p(z, T) = \prod_{n,k} [\pi_k G(z_n | \frac{\gamma_k}{2}, \frac{\gamma_k}{2})]^{t_{nk}} = \prod_{n,k} \left[\frac{z_n^c}{c} \exp(-z_n c) \right]^{t_{nk}}$$

$$p(w, z, \theta | \alpha, \beta) = \prod_{d=1}^D \text{Dir}(\theta_d | \alpha) \prod_{n=1}^N \prod_{t=1}^T (\theta_{dt} \phi_{tw_{dn}})^{[z_{dn}=t]}$$

$$p(w | z, \theta), \quad p(z, \theta) ?$$

$$p(w | z, \theta) = \prod_{d,n,t} \phi_{tw_{dn}}^{[z_{dn}=t]} = \prod_{d,n,t} c^{[z_{dn}=t]}$$

$$p(z, \theta) = \prod_d \text{Dir}(\theta_d | \alpha) \prod_{n,t} \theta_{dt}^{[z_{dn}=t]} =$$

$$= \prod_d \frac{\Gamma(\sum_t \alpha_t)}{\prod_t \Gamma(\alpha_t)} \prod_t \theta_{dt}^{\alpha_t - 1} \prod_{n,t} \theta_{dt}^{[z_{dn}=t]} =$$

$$= \frac{1}{c} \prod_{d,t} \theta_{dt}^c \prod_{n,t} \theta_{dt}^{[z_{dn}=t]} = \frac{1}{c} \prod_t \left[\prod_d \theta_{dt}^c \prod_n \theta_{dt}^{[z_{dn}=t]} \right] =$$

$$= \frac{1}{c} \prod_{d,n,t} \theta_{dt}^{c + [z_{dn}=t]} = \frac{1}{c} \prod_{d,n,t} \theta_{dt}^{c + [z_{dn}=t]}$$

$$\# \quad p(w, z, \theta, \phi | \alpha, \beta) = \prod_d \text{Dir}(\theta_d | \alpha) \prod_{n,t} [\theta_{dt} \phi_{tw_{dn}}]^{[z_{dn}=t]}$$

$$\prod_t \text{Dir}(\phi_t | \beta)$$

$$p(w | z, \theta, \phi), \quad p(z, \theta, \phi) ?$$

$$p(w | z, \theta, \phi) = \frac{1}{c} \prod_d \prod_t \theta_{dt}^c \prod_{n,t} \theta_{dt}^{[z_{dn}=t]} \prod_t \frac{1}{c} \prod_w \phi_{tw}^c$$

$$p(w | z, \theta, \phi) = \prod_{d,n,t} \phi_{tw_{dn}}^{[z_{dn}=t]}$$

$$q(z) q(\theta, \phi) = q(z) q(\theta) q(\phi)$$

$$E: \quad q(z, \theta, \phi) = q(z) \cdot q(\theta) q(\phi)$$

$$\ln q(\phi) = \mathbb{E}_{q(z)q(\theta)} \ln q p(w, z, \theta, \phi | \alpha, \beta) + C \Leftrightarrow$$

$$= \mathbb{E}_{q(z)q(\theta)} \sum_{d=1}^D \sum_{n=1}^N \sum_{t=1}^T \ln \left(\prod_{w=1}^V \phi_{tw}^{[z_{dn}=t][w_{dn}=w]} \right)$$

$$\ominus \sum_t \sum_w (\beta-1) \ln \Phi_{tw} + \sum_{d,n,t,w} \delta_{dn} \mathbb{I}[w_{dn}=w] \ln \Phi_{tw}^{\beta} =$$

$$= \sum_{t,w} [\ln \Phi_{tw} (\sum_{d,n} \delta_{dn} \mathbb{I}[w_{dn}=w] + \beta - 1)] +$$

$$q(\Phi) = \prod_{t,w} \text{Dir}(\Phi_t | \beta + \sum_{d,n} \delta_{dn} \mathbb{I}[w_{dn}=w])$$

$$p(w, z, \theta | \Phi, \alpha) = \prod_d \argmax_{\theta_d} \text{Dir}(\theta_d | \alpha)$$

$$\prod_n \prod_t [\hat{\theta}_{dt} \Phi_{tw_{dn}}]^{z_{dn}=t}$$

Иерархические процессы Дирихле.

$$G \sim DP(G_0, \alpha), \quad G = \sum_{k=1}^{\infty} \pi_k \delta_{\theta_k}$$

$$\theta_k \sim G_0, \quad \pi_k = v_k \prod_{j=1}^{k-1} (1 - v_j), \quad v_k \sim \text{Beta}(1, \alpha)$$

$$p(x, z, \pi, \theta) = \left[\prod_n \prod_k p(x_n | \theta_k)^{\mathbb{I}[z_n=k]} \pi_k^{\mathbb{I}[z_n=k]} \right] \underbrace{p(\pi) p(\theta)}_{[z_n=k]} =$$

$$= \left[\prod_k p_{G_0}(\theta_k) \text{Beta}(v_k | 1, \alpha) \right] \underbrace{\prod_n \prod_k [p(x_n | \theta_k) v_k \prod_{j=1}^{k-1} (1 - v_j)]}_{\sim p(\theta)} \underbrace{\pi_k}_{\sim p(\pi)}$$

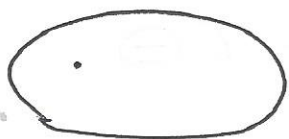
$$p(w, z, \theta, \Phi | \alpha, \beta) = \prod_t p(\Phi_t | \beta) \prod_d p(\theta_d | \alpha)$$

$$\prod_n \underbrace{p(z_{dn} | \theta_d)}_{\sim p(z_n | \pi)} \underbrace{p(w_{dn} | z_{dn}, \Phi)}_{\sim p(x_n | z_n, \theta_{z_n})}$$

слово - монета

$$G \sim DP(\alpha, H), \quad H = \text{Dir}(\eta)$$

$$\Phi_t \quad \overline{}_w$$



$$\Phi_t \sim \text{Dir}(\Phi_t | \eta)$$

$$v_t \sim \text{Beta}(1, \alpha)$$

$$\beta_t = v_t \prod_{j=1}^{t-1} (1 - v_j)$$

$$G = \sum_{t=1}^{\infty} \beta_t \delta_{\Phi_t}$$

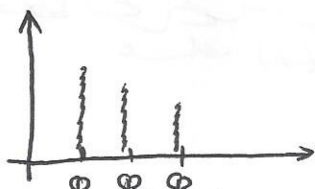
$$G_d \sim DP(\gamma, G)$$

\Rightarrow

$$\Phi_n \sim G$$

$$\pi_{dn} \sim \text{Beta}(1, \gamma)$$

$$\theta_{dn} = \pi_{dn} \prod_{j=1}^{n-1} (1 - \pi_{dj})$$



$$G_d = \sum_k \delta_{\Phi_k} \theta_{dk}$$

$$\pi_{dt} \sim \text{Beta}(\beta_t, \beta(1 - \sum_{k=1}^t \beta_k))$$

$$\Theta_{dt} = \pi_{dt} \prod_{j=1}^{t-1} (1 - \pi_{dj})$$

$$\begin{aligned} p(w, z, \theta, \Phi | \alpha, \beta) &= \prod_t \text{Dir}(\Phi_t | \eta) \text{Beta}(\nu_t | 1, \alpha) \\ &\prod_d \prod_t \text{Beta}(\pi_{dt} | \beta_t, \beta(1 - \sum_k \beta_k)) \cdot \prod_n \prod_t (\Theta_{dt})^{[z_{dn}=t]} \\ &\cdot \prod_t \prod_w \Phi_{tw}^{[z_{dn}=t][w_{dn}=w]} \end{aligned}$$