

## روش تخصیص دیریکله پنهان ترکیبی (compound latent Dirichlet allocation)

در مدل LDA قبلی، تنها داده‌های یک نوع سرطان به مدل داده می‌شد که برای تجزیه و تحلیل مقایسه‌ای بین انواع دیگر سرطان‌ها دارای نقص است. مجموعه داده از یک نوع سرطان با دیگر نوع‌های سرطان دارای امضاهای جهش مشترک است، اما نسبت امضاهای جهش بین نوع‌های سرطان به صورت نسبی متفاوت است. در مدل cLDA فرض بر این است که برای ژنوم‌های هر نوع سرطان ترکیبی از امضاهای جهش پنهان مشترک با دیگر انواع سرطان وجود دارد. با تعیین این امضاهای جهش می‌توان گفت هر امضای جهش به چه میزان در هر ژنوم فردی تاثیر گذاشته است. نسبت امضاهای جهش برای مجموعه ژنوم‌های یک نوع سرطان یکسان است.

تعاریف:

واژگان (vocabulary)  $V$ :  $v$  نوع جهش در دیکشنری جهش

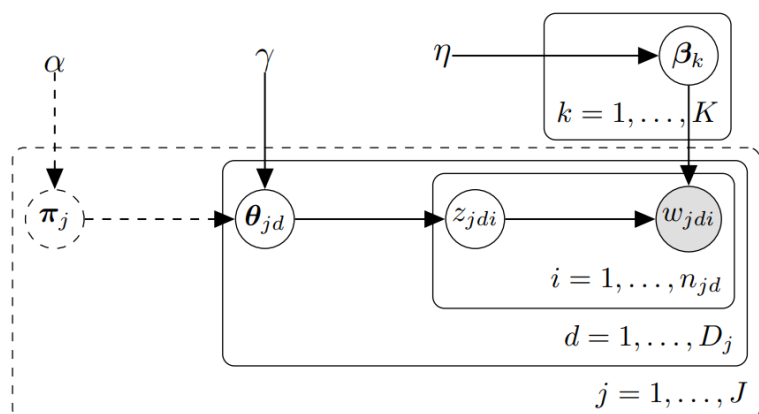
امضاهای جهش (mutational signatures)  $K$ :  $\beta_{K \times V}$  امضا جهش، هر کدام یک توزیع بر روی واژگان  $V$

نوع سرطان  $j$ : دارای  $D_j$  ژنوم سرطانی

ژنوم سرطانی  $d$  با نوع سرطان  $j$ : دارای  $n_{jd}$  جهش

جهش: index: واژگان متناظر با نوع جهش را دارد

hyperparameters:  $h = (\eta, \alpha, \gamma) \in (0, \infty)^3$



- $\beta_k \stackrel{\text{iid}}{\sim} \text{Dir}_V(\eta, \dots, \eta)$ , for topic  $k = 1, \dots, K$
- $\pi_j \stackrel{\text{iid}}{\sim} \text{Dir}_K(\alpha, \dots, \alpha)$ , for collection  $j = 1, \dots, J$
- $\theta_{jd} \stackrel{\text{iid}}{\sim} \text{Dir}_K(\gamma\pi_{j1}, \dots, \gamma\pi_{jK})$ , for document  $jd$
- $z_{jdi} \stackrel{\text{iid}}{\sim} \text{Mult}_K(\theta_{jd})$ , for each word  $w_{jdi}$
- $w_{jdi} \stackrel{\text{ind}}{\sim} \text{Mult}_V(\beta_{z_{jdi}})$

$$\beta = \begin{bmatrix} 0 \leq \leq 1 \\ \end{bmatrix}_{K \times V} \quad \pi = \begin{bmatrix} 0 \leq \leq 1 \\ \end{bmatrix}_{J \times K} \quad \theta = \begin{bmatrix} 0 \leq \leq 1 \\ \end{bmatrix}_{D_j \times K}$$

$$\alpha = \begin{bmatrix} \\ \end{bmatrix}_K \quad \eta = \begin{bmatrix} \\ \end{bmatrix}_V \quad z = \begin{bmatrix} 1 \leq \leq K \\ \end{bmatrix}_{D_j \times n_{jd}} \quad w = \begin{bmatrix} 1 \leq \leq V \\ \end{bmatrix}_{D_j \times n_{jd}}$$

توزیع پیشین:

$$p_h(z | \theta, \pi, \beta) p_h(\theta | \pi) p_h(\pi) p_h(\beta),$$

$$= \prod_{j=1}^J \left( \frac{\Gamma(K\alpha)}{\Gamma(\alpha)^K} \prod_{k=1}^K \pi_{jk}^{\alpha-1} \right) \left[ \prod_{j=1}^J \prod_{d=1}^{D_j} \prod_{k=1}^K \theta_{j dk}^{n_{j dk}} \right]$$

$$\left[ \prod_{j=1}^J \prod_{d=1}^{D_j} \left( \frac{\Gamma(\gamma)}{\prod_{k=1}^K \Gamma(\gamma \pi_{jk})} \prod_{k=1}^K \theta_{j dk}^{\gamma \pi_{jk} - 1} \right) \right] \left[ \prod_{k=1}^K \left( \frac{\Gamma(V\eta)}{\Gamma(\eta)^V} \prod_{v=1}^V \beta_{kv}^{\eta-1} \right) \right].$$

$$n_{j dk} = \sum_{i=1}^{n_{jd}} z_{j dik} \text{ تعداد جهش از ژنوم } d \text{ با نوع سرطان } z \text{ که به امضای جهش } k \text{ اختصاص یافته است.}$$

$$n_{j.k} = \sum_{d=1}^{D_j} \sum_{i=1}^{n_{jd}} z_{j dik} \text{ تعداد جهش در همه ژنوم‌هایی با نوع سرطان } z \text{ که به امضای جهش } k \text{ تخصیص یافته است.}$$

تابع درست‌نمایی:

$$p(w | z, \theta, \pi, \beta) = \prod_{j=1}^J \prod_{d=1}^{D_j} \prod_{i=1}^{n_{jd}} \prod_{k: z_{j dik}=1} \prod_{v=1}^V \beta_{kv}^{w_{j div}}$$

$$= \prod_{j=1}^J \prod_{d=1}^{D_j} \prod_{k=1}^K \prod_{v=1}^V \prod_{i \in S_{j dk}} \beta_{kv}^{w_{j div}}$$

$$= \prod_{j=1}^J \prod_{d=1}^{D_j} \prod_{k=1}^K \prod_{v=1}^V \beta_{kv}^{\sum_{i \in S_{j dk}} w_{j div}}$$

$$= \prod_{j=1}^J \prod_{d=1}^{D_j} \prod_{k=1}^K \prod_{v=1}^V \beta_{kv}^{m_{j dk v}},$$

$S_{jdk} = \{i: 1 \leq i \leq n_{jd} \text{ and } z_{jdk} = 1\}$  شماره جهش از ژنوم  $d$  با نوع سرطان  $j$  که امضای جهش آن  $k$  است.

تعداد جهش‌هایی از ژنوم  $d$  با نوع سرطان  $j$  که امضای جهش آن  $k$  است و نوع جهش آن  $v$  است.

توزیع پسین:

$$p(z, \theta, \pi, \beta | w) =$$

$$\left[ \prod_{k=1}^K \prod_{v=1}^V \beta_{kv}^{\sum_{j=1}^J \sum_{d=1}^{D_j} m_{jdkv} + \eta - 1} \right] \left[ \prod_{j=1}^J \prod_{d=1}^{D_j} \frac{\prod_{k=1}^K \theta_{jdk}^{n_{jdk} + \gamma \pi_{jk} - 1}}{\prod_{k=1}^K \Gamma(\gamma \pi_{jk})} \right] \left[ \prod_{j=1}^J \prod_{k=1}^K \pi_{jk}^{\alpha - 1} \right]$$

استنتاج تغییراتی:

در روش‌های تغییراتی، خانواده‌ای از کران‌های پایین برای توزیع پسین غیر قابل محاسبه در نظر می‌گیریم و از بین آنها نزدیک‌ترین کران را پیدا می‌کنیم. یک روش برای در نظر گرفتن خانواده‌ای کران‌های پایین پایین برای تقریب توزیع‌ها، استفاده از توزیع‌های پارامتری است. معمولاً این توزیع پارامتری با فرض استقلال بین متغیرهای پنهان ساده‌تر از توزیع پسین می‌شود. در نتیجه هدف شناسایی پارامترهایی است که کوچکترین کران پایین بدهد.

لگاریتم توزیع حاشیه‌ای:

$$\log p(w | \alpha, \eta, \gamma) = L(\lambda, \tau, \rho, \phi; \alpha, \eta, \gamma) + KL(q(\beta, \pi, \theta, z | \lambda, \tau, \rho, \phi) || p(\beta, \pi, \theta, z | w, \alpha, \eta, \gamma))$$

که در آن

$$L(\lambda, \tau, \rho, \phi; \alpha, \eta, \gamma) = \int \sum_z q(\beta, \pi, \theta, z) \log \left\{ \frac{p(\beta, \pi, \theta, z, w | \alpha, \eta, \gamma)}{q(\beta, \pi, \theta, z)} \right\} d\beta d\pi d\theta$$

$$KL(q, p) = - \int \sum_z q(\beta, \pi, \theta, z) \log \left\{ \frac{p(\beta, \pi, \theta, z, w | \alpha, \eta, \gamma)}{q(\beta, \pi, \theta, z)} \right\} d\beta d\pi d\theta$$

توزیع فاکتورسازی شده تغییراتی با پارامترهای تغییراتی  $\lambda, \tau, \rho, \phi$  :

$$q(\beta, \pi, \theta, z | \lambda, \tau, \rho, \phi) = [\prod_{k=1}^K q(\beta_k | \lambda_k)] [\prod_{j=1}^J q(\pi_j | \tau_j)] [\prod_{d=1}^{D_j} q(\theta_{j,d} | \rho_{j,d})] [\prod_{i=1}^{n_{j,d}} q(z_{j,d,i} | \phi_{j,d,i})]$$

توزیع مولفه‌های مستقل آن به صورت:

$$\beta_k \sim Dir_V(\lambda_k)$$

$$\pi_j \sim Dir_K(\tau_j)$$

$$\theta_{jd} \sim Dir_K(\rho_{jd})$$

$$z_{jdi} \sim Mult_K(\phi_{jdi})$$

---

**Algorithm 3:** Variational expectation maximization (VEM)

---

**Data:** Observed words  $w$  and document metadata

**Result:** Optimal variational parameters  $(\lambda, \tau, \rho, \phi)$

---

```
1 initialize  $(\lambda^{(0)}, \tau^{(0)})$ ;
2 while not converged do
    // Step 1: Expectation
3   for document  $d = 1, \dots, D_j, j = 1, \dots, J$  do
4       initialize  $\rho_{jdk}^{(0)} = \frac{\gamma \tau_{jk}}{\tau_j} + \frac{n_{jd}}{K}, k = 1, \dots, K$ ;
        // Variational updates for each document
5       while not converged do
6           for word  $w_{jdi}, i = 1, \dots, n_{jd}$  do
7               variational Multinomial update for  $\phi_{jdi}$  via (79);
8               variational Dirichlet update for  $\rho_{jd}$  via (82);
        // Variational updates for each topic
9       variational Dirichlet update for  $\lambda_k, k = 1, \dots, K$  via (85);
    // Step 2: Maximization
    // Constraint Newton updates for collection-level topic mixtures
10   for collection  $j = 1, \dots, J$  do
11       initialize  $(a_j^{(0)}, \omega_j^{(0)})$  based on the current  $\tau_j$ ;
12       while not converged do
13           constraint Newton update for  $\omega_j$  via (72);
14           Newton update for  $a_j$  via (76);
15       set  $\tau_j = a_j^{(\text{final})} * \omega_j^{(\text{final})}$ 
16   optimize hyperparameter  $h = (\alpha, \gamma, \eta)$ ;
```

---

روش تخمین پارامترهای تغییراتی  $\lambda, \tau, \rho, \phi$  با بیشینه کردن کران پایین  $L(\lambda, \tau, \rho, \phi; \alpha, \eta, \gamma)$  که با کمینه کردن معیار واگرایی کولبک-لیبلر بین توزیع تغییراتی  $q$  و توزیع پسین  $p$  معادل است، صورت می‌گیرد.

کران پایین  $L(\lambda, \tau, \rho, \phi; \alpha, \eta, \gamma)$ :

$$\begin{aligned}
\mathcal{L}(q, p_h) = & \sum_{k=1}^K \mathbb{E}_{q_k} [\log p_\eta(\beta_k)] + \sum_{j=1}^J \mathbb{E}_{q_j} [\log p_\alpha(\pi)] + \sum_{j=1}^J \sum_{d=1}^{D_j} \mathbb{E}_{q_{jd}} [\log p_\gamma(\theta_{jd} | \pi_j)] \\
& + \sum_{j=1}^J \sum_{d=1}^{D_j} \sum_{i=1}^{n_{dj}} \mathbb{E}_{q_{jdi}} [\log p(z_{jdi} | \theta_{jd})] + \sum_{j=1}^J \sum_{d=1}^{D_j} \sum_{i=1}^{n_{dj}} \mathbb{E}_{q_{jdi}} [\log p(w_{jdi} | z_{jdi}, \beta)] \\
& - \sum_{k=1}^K \mathbb{E}_{q_k} [\log q(\beta_k | \lambda_k)] - \sum_{j=1}^J \mathbb{E}_{q_j} [\log q(\pi_j | \tau_j)] \\
& - \sum_{j=1}^J \sum_{d=1}^{D_j} \mathbb{E}_{q_{jd}} [\log q(\theta_{jd} | \rho_{jd})] - \sum_{j=1}^J \sum_{d=1}^{D_j} \sum_{i=1}^{n_{dj}} \mathbb{E}_{q_{jdi}} [\log q(z_{jdi} | \phi_{jdi})]
\end{aligned} \tag{63}$$

$$\rho_{jd.} = \sum_{k=1}^K \rho_{jdk}, \tau_{j.} = \sum_{k=1}^K \tau_{jk}, \lambda_{k.} = \sum_{v=1}^V \lambda_{kv}$$

میانگین مولفه‌های احتمالی به صورت زیر است.

$$\begin{aligned}
\beta_k &\sim \text{Dir}_V(\lambda_k) & E_q[\beta_{kv}] &= \lambda_{kv} / \lambda_{k.} \\
\pi_j &\sim \text{Dir}_K(\tau_j) & E_q[\pi_{jk}] &= \tau_{jk} / \tau_{j.} \\
\theta_{jd} &\sim \text{Dir}_K(\rho_{jd}) & E_q[\theta_{jdk}] &= \rho_{jdk} / \rho_{jd.}
\end{aligned}$$

میانگین لگاریتم مولفه‌های احتمالی به صورت زیر است.

$$\begin{aligned}
\beta_k &\sim \text{Dir}_V(\lambda_k) & \mathbb{E}_q[\log \beta_{kv} | \lambda_{kv}] &= \Psi(\lambda_{kv}) - \Psi(\lambda_{k.}) \\
\pi_j &\sim \text{Dir}_K(\tau_j) & \mathbb{E}_q[\log \pi_{jk} | \tau_{jk}] &= \Psi(\tau_{jk}) - \Psi(\tau_{j.}) \\
\theta_{jd} &\sim \text{Dir}_K(\rho_{jd}) & \mathbb{E}_q[\log \theta_{jdk} | \rho_{jdk}] &= \Psi(\rho_{jdk}) - \Psi(\rho_{jd.})
\end{aligned}$$

اگر  $\theta \sim \text{Dir}_K(\rho)$ ،  $E[\theta_k] = \rho_k / \rho$  و  $\alpha \in [0, \infty)$  در نظر بگیریم، عبارت غیر قابل محاسبه

$E[\log \Gamma(\alpha \theta_k)]$  با توجه به مقاله (Kim et al., 2013, Theorem 3.1) به صورت زیر قابل تخمین است:

$$\mathbb{E}[\log \Gamma(\alpha \theta_k)] \leq \log \Gamma(\alpha \mathbb{E}[\theta_k]) + \frac{\alpha}{\rho} (1 - \mathbb{E}[\theta_k]) + (1 - \alpha \mathbb{E}[\theta_k]) [\log \mathbb{E}[\theta_k] + \Psi(\rho) - \Psi(\rho_k)] \tag{65}$$

در نهایت هر میانگین از عبارت کران پایین به صورت زیر محاسبه میشود:

$$\begin{aligned}
\mathbb{E}_{q_k}[\log p_\eta(\boldsymbol{\beta}_k)] &= \log \Gamma(V\eta) - V \log \Gamma(\eta) + \sum_{v=1}^V (\eta - 1) [\Psi(\lambda_{kv}) - \Psi(\lambda_{k.})] \\
\mathbb{E}_{q_j}[\log p_\alpha(\boldsymbol{\pi})] &= \log \Gamma(K\alpha) - K \log \Gamma(\alpha) + \sum_{k=1}^K (\alpha - 1) [\Psi(\tau_{jk}) - \Psi(\tau_{j.})] \\
\mathbb{E}_{q_{jd}}[\log p_\gamma(\boldsymbol{\theta}_{jd} \mid \boldsymbol{\pi}_j)] &= \mathbb{E}_q[\log \Gamma(\gamma)] - \sum_{k=1}^K \mathbb{E}_q[\log \Gamma(\gamma\pi_{jk})] + \sum_{k=1}^K \mathbb{E}_q[(\gamma\pi_{jk} - 1) \log \theta_{jdk}] \\
&\geq \log \Gamma(\gamma) - \sum_{k=1}^K \left[ \log \Gamma(\gamma \mathbb{E}_q[\pi_{jk}]) + \frac{\gamma}{\tau_{j.}} (1 - \mathbb{E}_q[\pi_{jk}]) \right. \\
&\quad \left. + (1 - \gamma \mathbb{E}_q[\pi_{jk}]) [\log \mathbb{E}_q[\pi_{jk}] + \Psi(\tau_{j.}) - \Psi(\tau_{jk})] \right] \\
&\quad + \sum_{k=1}^K \left[ \gamma \mathbb{E}_q[\pi_{jk}] \mathbb{E}_q[\log \theta_{jdk}] - \mathbb{E}_q[\log \theta_{jdk}] \right] \\
&\geq \log \Gamma(\gamma) - \frac{\gamma}{\tau_{j.}} (K - 1) - (\gamma - K) \left[ \log \tau_{j.} - \Psi(\tau_{j.}) + \Psi(\rho_{jd.}) \right] \\
&\quad - \sum_{k=1}^K \left[ \log \Gamma\left(\frac{\gamma\tau_{jk}}{\tau_{j.}}\right) + \left(1 - \frac{\gamma\tau_{jk}}{\tau_{j.}}\right) [\log(\tau_{jk}) - \Psi(\tau_{jk}) + \Psi(\rho_{jdk})] \right] \\
\mathbb{E}_{q_{jdi}}[\log p(z_{jdi} \mid \boldsymbol{\theta}_{jd})] &= \sum_{k=1}^K \phi_{jdk} [\Psi(\rho_{jdk}) - \Psi(\rho_{jd.})] \\
\mathbb{E}_{q_{jdi}}[\log p(w_{jdi} \mid z_{jdi}, \boldsymbol{\beta})] &= \sum_{k=1}^K \sum_{v=1}^V \phi_{jdk} w_{jdiv} [\Psi(\lambda_{kv}) - \Psi(\lambda_{k.})] \\
\mathbb{E}_{q_k}[\log q(\boldsymbol{\beta}_k \mid \boldsymbol{\lambda}_k)] &= \log \Gamma(\lambda_{k.}) - \sum_{v=1}^V \log \Gamma(\lambda_{kv}) + \sum_{v=1}^V (\lambda_{kv} - 1) [\Psi(\lambda_{kv}) - \Psi(\lambda_{k.})] \\
\mathbb{E}_{q_j}[\log q(\boldsymbol{\pi}_j \mid \boldsymbol{\tau}_j)] &= \log \Gamma(\tau_{j.}) - \sum_{k=1}^K \log \Gamma(\tau_{jk}) + \sum_{k=1}^K (\tau_{jk} - 1) [\Psi(\tau_{jk}) - \Psi(\tau_{j.})] \\
\mathbb{E}_{q_{jd}}[\log q(\boldsymbol{\theta}_{jd} \mid \boldsymbol{\rho}_{jd})] &= \log \Gamma(\rho_{jd.}) - \sum_{k=1}^K \log \Gamma(\rho_{jdk}) + \sum_{k=1}^K (\rho_{jdk} - 1) [\Psi(\rho_{jdk}) - \Psi(\rho_{jd.})] \\
\mathbb{E}_{q_{jdi}}[\log q(z_{jdi} \mid \phi_{jdi})] &= \sum_{k=1}^K \phi_{jdk} \log \phi_{jdk}
\end{aligned}$$

(66)

به روز رسانی دیریکله تغییراتی برای انواع سرطان

عبارات شامل  $\tau_j$  را از کران پایین  $L$  جدا می‌کنیم:

$$\mathcal{L}_{[\tau_j]} = \mathbb{E}_{q_j}[\log p_{\alpha}(\boldsymbol{\pi})] - \mathbb{E}_{q_j}[\log q(\boldsymbol{\pi}_j | \tau_j)] + \sum_{d=1}^{D_j} \mathbb{E}_{q_{jd}}[\log p_{\gamma}(\boldsymbol{\theta}_{jd} | \boldsymbol{\pi}_j)]$$

این کران پایین یک فرم بسته برای به روز رسانی  $\tau_j$  ایجاد نمی‌کند. Kim et al. (2013) پیشنهاد می‌دهد که از به روز رسانی نیوتن با شرط‌های یکسان استفاده شود. برای اینکار لازم است  $\tau_j$  را به یک پارامتر مقیاس  $a_j$  و یک اندازه‌گیری پایه (base measure)  $w_j$  با شرط اینکه  $\sum_{k=1}^K w_{jk} = 1$  تجزیه کنیم. در نتیجه توزیع تغییراتی متناظر برای  $\boldsymbol{\pi}_j = \text{Dir}_K(a_j w_j)$  به صورت  $\boldsymbol{\pi}_j = \text{Dir}_K(a_j w_j)$  تعریف می‌شود. این تجزیه این امکان را می‌دهد تا از به روز رسانی نیوتن استفاده شود. کران پایین شامل عبارات  $\tau_j$  با  $a_j$  و  $w_j$  به صورت زیر با تعریف می‌شود:

$$\begin{aligned} \mathcal{L}_{[a_j w_j]} &= \sum_{k=1}^K \left[ (\alpha - a_j w_{jk}) [\Psi(a_j w_{jk}) - \Psi(a_j)] + \log \Gamma(a_j w_{jk}) \right] - \log \Gamma(a_j) \\ &\quad - \sum_{d=1}^{D_j} \left[ \frac{\gamma}{a_j} (K - 1) + (\gamma - K) [\log a_j - \Psi(a_j) + \Psi(\rho_{jd})] \right] \\ &\quad - \sum_{d=1}^{D_j} \sum_{k=1}^K \left[ \log \Gamma(\gamma w_{jk}) + (1 - \gamma w_{jk}) [\log(a_j w_{jk}) - \Psi(a_j w_{jk}) + \Psi(\rho_{jd})] \right] \end{aligned} \quad (67)$$

برای بیشینه کردن  $L_{[a_j w_j]}$  براساس  $w_{jk}$ ، ابتدا عبارات شامل  $w_{jk}$  را جدا می‌کنیم:

$$\begin{aligned} \mathcal{L}_{[w_{jk}]} &= \Psi(a_j w_{jk}) \left[ \alpha + D_j - a_j w_{jk} - \gamma D_j w_{jk} \right] + \gamma w_{jk} \sum_{d=1}^{D_j} \Psi(\rho_{jd}) \\ &\quad + \log \Gamma(a_j w_{jk}) - D_j \left[ \log \Gamma(\gamma w_{jk}) + (1 - \gamma w_{jk}) \log(a_j w_{jk}) \right] \end{aligned} \quad (68)$$

مشتق اول  $g_{jk}$  و مشتق دوم  $h_{jk}$  آن:



$$\begin{aligned} \frac{\partial}{\partial \omega_{jk}} \mathcal{L}_{[\omega_{jk}]} &= a_j \Psi'(a_j \omega_{jk}) \left[ \alpha + D_j - a_j \omega_{jk} - \gamma D_j \omega_{jk} \right] - \gamma D_j \Psi(a_j \omega_{jk}) + \gamma \sum_{d=1}^{D_j} \Psi(\rho_{jd} k) \\ &\quad - D_j \left[ \gamma \Psi(\gamma \omega_{jk}) + \frac{1}{\omega_{jk}} - \gamma - \gamma \log(a_j \omega_{jk}) \right] \end{aligned} \quad (69)$$

$$\begin{aligned} \frac{\partial^2}{\partial \omega_{jk}^2} \mathcal{L}_{[\omega_{jk}]} &= a_j^2 \Psi''(a_j \omega_{jk}) \left[ \alpha + D_j - a_j \omega_{jk} - \gamma D_j \omega_{jk} \right] - a_j \Psi'(a_j \omega_{jk}) \left[ a_j + 2\gamma D_j \right] \\ &\quad - D_j \left[ \gamma^2 \Psi'(\gamma \omega_{jk}) - \frac{1}{\omega_{jk}^2} - \frac{\gamma}{\omega_{jk}} \right] \end{aligned} \quad (70)$$

ماتریس هسیان داده شده در (۷۰) قطری است. برای به دست آوردن مقدار به روز رسانی  $\Delta \omega_{jk}$  در هر مرحله نیوتن مجموعه معادلات خطی زیر را حل می‌کنیم:

$$\begin{bmatrix} \text{diag}(\mathbf{h}) & \mathbf{1} \\ \mathbf{1}^T & 0 \end{bmatrix} \begin{bmatrix} \Delta \omega_{jk} \\ u \end{bmatrix} = \begin{bmatrix} -\mathbf{g} \\ 0 \end{bmatrix}, \quad (71)$$

در نتیجه:

$$\Delta \omega_{jk} = \left\{ \frac{\sum_{k=1}^K \frac{g_{jk}}{h_{jk}}}{\sum_{k=1}^K \frac{1}{h_{jk}}} \right\} \begin{bmatrix} \frac{1}{h_{j1}} \\ \vdots \\ \frac{1}{h_{jK}} \end{bmatrix} - \begin{bmatrix} \frac{g_{j1}}{h_{j1}} \\ \vdots \\ \frac{g_{jK}}{h_{jK}} \end{bmatrix} \quad (72)$$

شرط  $\sum_{k=1}^K \Delta \omega_{jk} = 0$  همچنان برقرار است.

برای پیشینه کردن  $L_{[a_j \omega_j]}$  براساس  $a_j$ ، ابتدا عبارات شامل  $a_j$  را جدا می‌کنیم:

$$\begin{aligned} \mathcal{L}_{[a_j]} &= \sum_{k=1}^K [\Psi(a_j \omega_{jk}) - \Psi(a_j)] \left( \alpha + D_j - a_j \omega_{jk} - \gamma D_j \omega_{jk} \right) + \sum_{k=1}^K \log \Gamma(a_j \omega_{jk}) \\ &\quad - \log \Gamma(a_j) - \frac{\gamma D_j (K-1)}{a_j} \end{aligned} \quad (73)$$

مشتق اول  $g_j$  و مشتق دوم  $h_j$  آن:

$$\begin{aligned} \frac{\partial}{\partial a_j} \mathcal{L}_{[a_j]} &= \sum_{k=1}^K [\omega_{jk} \Psi'(a_j \omega_{jk}) - \Psi'(a_j)] \left( \alpha + D_j - a_j \omega_{jk} - \gamma D_j \omega_{jk} \right) \\ &\quad + (K-1) \gamma D_j a_j^{-2} \end{aligned} \quad (74)$$

$$\begin{aligned} \frac{\partial^2}{\partial a_j^2} \mathcal{L}_{[a_j]} &= \sum_{k=1}^K [\omega_{jk}^2 \Psi''(a_j \omega_{jk}) - \Psi''(a_j)] \left( \alpha + D_j - a_j \omega_{jk} - \gamma D_j \omega_{jk} \right) \\ &\quad - \sum_{k=1}^K \omega_{jk} [\omega_{jk} \Psi'(a_j \omega_{jk}) - \Psi'(a_j)] - 2(K-1) \gamma D_j a_j^{-3} \end{aligned} \quad (75)$$

مقدار به روز رسانی  $a_j$  در هر مرحله نیوتن:

$$\Delta a_j = -h_j^{-1} g_j. \quad (76)$$

به روز رسانی دیریکله تغییراتی برای جهش‌ها

برای به روز رسانی پارامتر تغییراتی  $\phi_{jdk}$  (احتمال اینکه کلمه  $z_{di}$  ام توسط امضای جهش  $k$  ایجاد شده باشد) با بیشینه کردن کران پایین  $L$  با شرط اینکه  $\sum_{k=1}^K \phi_{jdk} = 1$ ، تمام عبارات شامل  $\phi_{jdk}$  را از  $L$  جدا می‌کنیم:

$$\begin{aligned} \mathcal{L}_{[\phi_{jdk}]} &= \phi_{jdk} \left[ [\Psi(\rho_{jdk}) - \Psi(\rho_{jd.})] + [\Psi(\lambda_{kv}) - \Psi(\lambda_{k.})] - \log \phi_{jdk} \right] \\ &\quad + \mu_{jdi} \left[ \sum_{k=1}^K \phi_{jdk} - 1 \right] \end{aligned} \quad (77)$$

از عبارت بالا نسبت به  $\phi_{jdk}$  مشتق می‌گیریم:

$$\frac{\partial}{\partial \phi_{jdk}} \mathcal{L}_{[\phi_{jdk}]} = [\Psi(\rho_{jdk}) - \Psi(\rho_{jd.})] + [\Psi(\lambda_{kv}) - \Psi(\lambda_{k.})] - \log \phi_{jdk} - 1 + \mu_{jdi} \quad (78)$$

مشتق بالا را مساوی صفر قرار داده و مقدار بیشینه  $\phi_{jdk}$  را بدست می‌آوریم:

$$\phi_{jdk} \propto \exp \left( [\Psi(\rho_{jdk}) - \Psi(\rho_{jd.})] + [\Psi(\lambda_{kv}) - \Psi(\lambda_{k.})] \right) \quad (79)$$

به روز رسانی دیریکله تغییراتی برای ژنوم‌ها

کران پایین  $L$  براساس  $\rho_{jdk}$  بیشینه می‌کنیم. عبارات شامل  $\rho_{jdk}$  در  $L$ :

$$\begin{aligned}\mathcal{L}_{[\rho_{jd}]} = & \sum_{k=1}^K \left( \frac{\gamma \tau_{jk}}{\tau_{j.}} + \sum_{i=1}^{n_{jd}} \phi_{jdik} - \rho_{jdk} \right) [\Psi(\rho_{jdk}) - \Psi(\rho_{jd.})] \\ & - \log \Gamma(\rho_{jd.}) + \sum_{k=1}^K \log \Gamma(\rho_{jdk})\end{aligned}\quad (80)$$

از عبارت بالا نسبت به  $\rho_{jdk}$  مشتق می‌گیریم:

$$\frac{\partial}{\partial \rho_{jdk}} \mathcal{L}_{[\rho_{jd}]} = [\Psi'(\rho_{jdk}) - \Psi'(\rho_{jd.})] \left( \frac{\gamma \tau_{jk}}{\tau_{j.}} + \sum_{i=1}^{n_{jd}} \phi_{jdik} - \rho_{jdk} \right) \quad (81)$$

مشتق بالا را مساوی صفر قرار داده و مقدار بیشینه  $\rho_{jdk}$  را بدست می‌آوریم:

$$\rho_{jdk} = \frac{\gamma \tau_{jk}}{\tau_{j.}} + \sum_{i=1}^{n_{jd}} \phi_{jdik} \quad (82)$$

به روز رسانی دیریکله تغییراتی برای امضاهای جهش

کران پایین  $L$  براساس  $\lambda_{kv}$  بیشینه می‌کنیم. عبارات شامل  $\lambda_{kv}$  در  $L$ :

$$\begin{aligned}\mathcal{L}_{[\lambda_k]} = & \sum_{k=1}^K \sum_{v=1}^V \left( \eta - \lambda_{kv} + \sum_{j=1}^J \sum_{d=1}^{D_j} \sum_{i=1}^{n_{jd}} \phi_{jdik} w_{jdiv} \right) [\Psi(\lambda_{kv}) - \Psi(\lambda_{k.})] \\ & - \sum_{k=1}^K [\log \Gamma(\lambda_{k.}) - \sum_{v=1}^V \log \Gamma(\lambda_{kv})]\end{aligned}\quad (83)$$

از عبارت بالا نسبت به  $\lambda_{kv}$  مشتق می‌گیریم:

$$\frac{\partial}{\partial \lambda_{kv}} \mathcal{L}_{[\lambda_k]} = [\Psi'(\lambda_{kv}) - \Psi'(\lambda_{k.})] \left( \eta - \lambda_{kv} + \sum_{j=1}^J \sum_{d=1}^{D_j} \sum_{i=1}^{n_{jd}} \phi_{jdik} \right) \quad (84)$$

مشتق بالا را مساوی صفر قرار داده و مقدار بیشینه  $\lambda_{kv}$  را بدست می‌آوریم:

$$\lambda_{kv} = \eta + \sum_{j=1}^J \sum_{d=1}^{D_j} \sum_{i=1}^{n_{jd}} \phi_{jdik} \quad (85)$$

بهینه‌سازی هایپر پارامترها  $\eta, \alpha, \gamma$

مانند الگوریتم تغییراتی EM برای LDA (Blei et al., 2003) ، E-step در cLDA پارامترهای تغییراتی براساس عبارت گفته شده در بالا به روزرسانی می‌شوند. از کران پایین بهینه شده  $L(q^*, p)$  به عنوان تقریب قابل محاسبه برای لگاریتم درست‌نمایی حاشیه‌ای  $\log p(w)$  می‌توانیم استفاده کنیم. در M-step می‌توانیم هایپر پارامترهای  $\eta, \alpha, \gamma$  را با بیشینه کردن کران پایین بهینه شده براساس  $\eta, \alpha, \gamma$  به روز رسانی کنیم.

عبارات شامل هر هایپر پارامتر را از کران پایین بهینه شده جدا می‌کنیم:

$$\mathcal{L}_{[\alpha]} = J \log \Gamma(K\alpha) - JK \log \Gamma(\alpha) + \sum_{j=1}^J \sum_{k=1}^K \alpha [\Psi(\tau_{jk}) - \Psi(\tau_{j.})] \quad (86)$$

$$\mathcal{L}_{[\eta]} = K \log \Gamma(V\eta) - KV \log \Gamma(\eta) + \sum_{k=1}^K \sum_{v=1}^V \eta [\Psi(\lambda_{kv}) - \Psi(\lambda_{k.})] \quad (87)$$

$$\begin{aligned} \mathcal{L}_{[\gamma]} = & \sum_{j=1}^J D_j [\log \Gamma(\gamma) - \frac{\gamma}{\tau_{j.}} (K-1)] - \gamma \sum_{j=1}^J \sum_{d=1}^{D_j} [\log \tau_{j.} - \Psi(\tau_{j.}) + \Psi(\rho_{jd.})] \\ & - \sum_{j=1}^J \sum_{d=1}^{D_j} \sum_{k=1}^K \left[ \log \Gamma\left(\frac{\gamma \tau_{jk}}{\tau_{j.}}\right) - \frac{\gamma \tau_{jk}}{\tau_{j.}} [\log(\tau_{jk}) - \Psi(\tau_{jk}) + \Psi(\rho_{jdk})] \right] \end{aligned} \quad (88)$$

مشتق اول و دوم هر کدام را محاسبه می‌کنیم:

$$\frac{\partial}{\partial \alpha} \mathcal{L}_{[\alpha]} = JK [\Psi(K\alpha) - \Psi(\alpha)] + \sum_{j=1}^J \sum_{k=1}^K [\Psi(\tau_{jk}) - \Psi(\tau_{j.})] \quad (89)$$

$$\frac{\partial^2}{\partial \alpha^2} \mathcal{L}_{[\alpha]} = JK^2 \Psi'(K\alpha) - JK \Psi'(\alpha) \quad (90)$$

$$\frac{\partial}{\partial \eta} \mathcal{L}_{[\eta]} = KV [\Psi(V\eta) - \Psi(\eta)] + \sum_{k=1}^K \sum_{v=1}^V [\Psi(\lambda_{kv}) - \Psi(\lambda_{k.})] \quad (91)$$

$$\frac{\partial^2}{\partial \eta^2} \mathcal{L}_{[\eta]} = KV^2 \Psi'(V\eta) - KV \Psi'(\eta) \quad (92)$$

$$\begin{aligned} \frac{\partial}{\partial \gamma} \mathcal{L}_{[\gamma]} = & \sum_{j=1}^J D_j [\Psi(\gamma) - \frac{1}{\tau_{j.}} (K-1)] - \sum_{j=1}^J \sum_{d=1}^{D_j} \left[ \log \tau_{j.} - \Psi(\tau_{j.}) + \Psi(\rho_{jd.}) \right] \\ & - \sum_{j=1}^J \sum_{d=1}^{D_j} \sum_{k=1}^K \frac{\tau_{jk}}{\tau_{j.}} \left[ \Psi\left(\frac{\gamma \tau_{jk}}{\tau_{j.}}\right) - [\log(\tau_{jk}) - \Psi(\tau_{jk}) + \Psi(\rho_{jdk})] \right] \end{aligned} \quad (93)$$

$$\frac{\partial^2}{\partial \gamma^2} \mathcal{L}_{[\gamma]} = \sum_{j=1}^J D_j \Psi'(\gamma) - \sum_{j=1}^J \sum_{d=1}^{D_j} \sum_{k=1}^K \frac{\tau_{jk}^2}{\tau_{j.}^2} \Psi'\left(\frac{\gamma \tau_{jk}}{\tau_{j.}}\right) \quad (94)$$

با استفاده از مشتق‌های بدست آمده می‌توانیم  $\eta$ ,  $\alpha$ ,  $\gamma$  با روش نیوتن به روز رسانی کنیم (Blei et al., 2003; Minka, 2000):

روش‌های Newton-Raphson برای یک هسیان با ساختار خاص (Blei et al., 2003):

این روش برای تخمین حداکثر درست‌نمایی توزیع دیریکله مورد استفاده قرار می‌گیرد. این تکنیک یک نقطه ثابت از یک تابع را با تکرار عبارت زیر پیدا می‌کند:

$$\alpha_{\text{new}} = \alpha_{\text{old}} - H(\alpha_{\text{old}})^{-1} g(\alpha_{\text{old}})$$

که در آن  $g(\alpha)$  و  $H(\alpha)$  به ترتیب گرادیان و ماتریس هسیان در نقطه  $\alpha$  است. این الگوریتم به دلیل محاسبه ماتریس معکوس به زمان  $O(N^3)$  نیاز دارد.

اگر ماتریس هسیان فرم زیر را داشته باشد:

$$H = \text{diag}(h) + \mathbf{1}z\mathbf{1}^T,$$

ماتریس معکوس به صورت زیر قابل محاسبه است:

$$H^{-1} = \text{diag}(h)^{-1} - \frac{\text{diag}(h)^{-1} \mathbf{1} \mathbf{1}^T \text{diag}(h)^{-1}}{z^{-1} + \sum_{j=1}^k h_j^{-1}}$$

با ضرب گرادیان در آن:

$$(H^{-1}g)_i = \frac{g_i - c}{h_i} \quad c = \frac{\sum_{j=1}^k g_j / h_j}{z^{-1} + \sum_{j=1}^k h_j^{-1}}.$$

این عبارت تنها به  $2k$  مقدار  $g_i$  و  $h_i$  بستگی دارد و در نتیجه پیچیدگی زمانی خطی دارد.

### نکات پیاده‌سازی:

مقدار اولیه  $\lambda$  و  $\tau$  در مقاله (George et al., 2019) به صورت زیر قرار گرفته است:

$$\lambda_{kv} = \text{uniform}(0,1) + 1/V$$

$$\tau_{jk} = \alpha + n_j / K$$

$\alpha$  به صورت ورودی گرفته می‌شود.  $n_j$  تعداد جهش‌های همه ژنوم‌هایی که نوع سرطان  $j$  را دارند. مقدار  $\tau_{jk}$  نرمال‌سازی می‌شود تا جمع برابر ۱ شود.

مقادیر اولیه در مقاله (Matsutani et al., 2019) به صورت  $\text{rnd}$ - قرار گرفته است و نرمال‌سازی می‌شوند.

من در پیاده‌سازی مقادیر اولیه را به صورت  $\text{rnd}$ - دادم و نرمال‌سازی کردم.

مقدار اولیه  $a_j$  و  $w_{jk}$  به صورت زیر قرار گرفته است:

$$a_j = \tau_j.$$

$$w_{jk} = \tau_{jk} / \tau_j.$$

برای چک کردن همگرایی کل الگوریتم کل مولفه‌های کران پایین L (66) را در نظر می‌گیریم و اگر مجموع آنها از یک threshold کمتر شود، تکرار را متوقف می‌کنیم.

برای چک کردن همگرایی  $\rho_{jdk}$  تمام مولفه‌های کران پایین مرتبط با ژنوم‌ها شامل مولفه ۱،۶، ۲،۷ در (66) را در نظر می‌گیریم و اگر مجموع آنها از یک threshold کمتر شود، تکرار را متوقف می‌کنیم.

برای چک کردن همگرایی  $\tau_{jk}$  تمام مولفه‌های کران پایین مرتبط با انواع سرطان شامل مولفه ۳، ۲،۷ در (66) را در نظر می‌گیریم و  $a_{jk}w_j$  را با  $\tau_{jk}$  جابه‌جا می‌کنیم و اگر مجموع آنها از یک threshold کمتر شود، تکرار را متوقف می‌کنیم.

برای به روز رسانی هایپرپارامترها مقاله (George et al., 2019) از عبارات زیر استفاده می‌کند:

$$\alpha^{new} = \exp(\log(\alpha^{old}) - \frac{f'(\alpha)}{f''(\alpha) \times \alpha + f'(\alpha)})$$

$$\eta^{new} = \exp(\log(\eta^{old}) - \frac{f'(\eta)}{f''(\eta) \times \eta + f'(\eta)})$$

و مقدار  $\gamma$  به روزرسانی نشده است.

برای به روز رسانی هایپرپارامترها مقاله (Matsutani et al., 2019) از عبارات زیر استفاده می‌کند:

$$\alpha_k = \frac{\sum_{s=1}^S [\Psi(E[n_{s,k} + \hat{\alpha}_k]) - \Psi(\hat{\alpha}_k)] \hat{\alpha}_k}{\sum_{s=1}^S [\Psi(n_s + \sum_{k=1}^K \hat{\alpha}_k) - \Psi(\sum_{k=1}^K \hat{\alpha}_k)]}$$

$$\beta_v = \frac{\sum_{k=1}^K [\Psi(E[n_{k,v} + \hat{\beta}_v]) - \Psi(\hat{\beta}_v)] \hat{\beta}_v}{\sum_{k=1}^K [\Psi(\sum_{v=1}^V E[n_{k,v}] + \hat{\beta}_v) - \Psi(\sum_{v=1}^V \hat{\beta}_v)]}$$

که به این عبارات در مقاله (Minka, 2000) در قسمت توزیع چندجمله‌ای - دیریکله اشاره شده است.

من برای به روز رسانی  $\eta$ ,  $\alpha$ ,  $\gamma$  از روش Newton-Raphson که در بالا گفته شد استفاده کردم. این روش در مقاله (Blei et al., 2003) و (Minka, 2000) در قسمت توزیع دیریکله اشاره شده است.

1. Clint P. George , Wei Xia , and George Michailidis. Analyses of Multi-collection Corpora via Compound Topic Modeling, Machine Learning, Optimization, and Data Science. LOD 2019. Lecture Notes in Computer Science, vol 11943. Springer, Cham, 2020  
(<https://arxiv.org/pdf/1907.01636.pdf>)  
(<https://github.com/clintpgeorge/clda>)
2. Taro Matsutani, Yuki Ueno, Tsukasa Fukunaga, and Michiaki Hamada. Discovering novel mutation signatures by latent Dirichlet allocation with variational Bayes inference, Bioinformatics, 35(22), 4543–4552, 2019
3. Thomas Minka. Estimating a dirichlet distribution, 2000a.
4. Thomas P. Minka. Beyond Newton's method. Technical report, Microsoft, 2000b
5. David M. Blei, Andrew Y. Ng, and Michael I. Jordan. Latent Dirichlet allocation. Journal of Machine Learning Research, 3:993–1022, 2003.