

ノート：ディリクレ分布の  $\alpha$  のサンプリング\*

LDA のハイパーパラメータはトピック分布およびトピック-単語分布を生成するディリクレ分布のハイパーパラメータ  $\alpha, \eta$  ですが、これらはどうやって決めたらいいのでしょうか。

一般的にはトピック数を  $K$  として、 $\alpha = (50/K, \dots, 50/K)$ ,  $\eta = 0.01$  のようにヒューリスティックに決められることが多いのですが、これらのハイパーパラメータの影響を詳しく調べた研究[145]によると、これらもデータから学習した方がよいモデルとなることが示されています。とくに、非対称な  $\alpha = (\alpha_1, \alpha_2, \dots, \alpha_K)$  を推定することで、 $\alpha_k$  が大きいトピックに「が」「の」、英語なら “of”, “the” といった機能語<sup>\*28</sup> が集まり、他のトピックにこれらが混ざることが少なくなります。また、トピック数を増やしても不要なトピックは  $\alpha_k$  が小さく、自然に使われなくなることも報告されています。

$\alpha$  と  $\eta$  についての尤度関数は式(5.78)ですから、数学的にはこれは、このポリア分布のハイパーパラメータを推定することと等価です。ただし、このカウント  $n(i, k)$ ,  $m(k, v)$  はあくまで学習中の  $\mathbf{z}_1^N$  から計算される値ですから、式(3.61)で学習中に  $\alpha$  を最適化してしまうと、局所解に陥る危険があります。<sup>\*29</sup> よって、きちんと  $\alpha$  のベイズ推定を行うのが正しい方法でしょう。

ポリア分布の式

$$(5.79) \quad p(\mathbf{n}|\alpha) = \underbrace{\frac{\Gamma(\sum_k \alpha_k)}{\Gamma(N + \sum_k \alpha_k)}}_{(A)} \prod_{k=1}^K \underbrace{\frac{\Gamma(\alpha_k + n_k)}{\Gamma(\alpha_k)}}_{(B)}$$

は  $\alpha_k$  がガンマ関数  $\Gamma()$  の中に入っているため、そのままではサンプリングできる形になりません。しかし、巧妙な補助変数  $\theta, x$  を導入すると、ガンマ事後分布からサンプリングすることができます[60, Appendix C].<sup>\*30</sup>

\*28 もしこうした機能語を「ストップワード」として除いても、236 ページで説明したように、同様に意味の薄い語が必ず出現し、機能語とそれ以外は簡単に二分することができません。

\*29 すなわち、学習途中で  $\alpha$  を最適化してしまうと、学習アルゴリズム全体が Gibbs サンプリングではなく、245 ページのモンテカルロ EM アルゴリズムを行うことになってしまいます。

(A) ディリクレ分布で2次元の場合を考えると、**ベータ関数**

$$(5.80) \quad B(a, b) = \int_0^1 \theta^{a-1} (1-\theta)^{b-1} d\theta = \frac{\Gamma(a)\Gamma(b)}{\Gamma(a+b)}$$

が得られます。よって、

$$(5.81) \quad B(\sum_k \alpha_k, N) = \frac{\Gamma(\sum_k \alpha_k)\Gamma(N)}{\Gamma(\sum_k \alpha_k + N)}$$

ですから、(A) は補助変数  $\theta$  を積分消去した形で、

$$(5.82) \quad \frac{\Gamma(\sum_k \alpha_k)}{\Gamma(N + \sum_k \alpha_k)} = \frac{B(\sum_k \alpha_k, N)}{\Gamma(N)} = \frac{1}{\Gamma(N)} \int_0^1 \theta^{\sum_k \alpha_k - 1} (1-\theta)^{N-1} d\theta$$

と表すことができます。

(B) 110 ページのコラムでみたように、Pochhammer 関数  $\Gamma(\alpha+n)/\Gamma(\alpha)$  は

$$(5.83) \quad \frac{\Gamma(\alpha+n)}{\Gamma(\alpha)} = \alpha(\alpha+1)\cdots(\alpha+n-1) = \prod_{i=0}^{n-1} (\alpha+i)$$

と表すことができますが、これも  $\{0, 1\}$  の補助変数  $x_{ki}$  を周辺化した形で

$$(5.84) \quad \frac{\Gamma(\alpha_k + n_k)}{\Gamma(\alpha_k)} = \prod_{i=0}^{n_k-1} (\alpha_k + i) = \prod_{i=0}^{n_k-1} \sum_{x_{ki} \in \{0,1\}} \alpha_k^{x_{ki}} i^{1-x_{ki}}$$

と書くことができることに注意しましょう。

よって、 $\alpha_k$  がガンマ事前分布

$$(5.85) \quad p(\alpha_k) = \text{Ga}(a, b) = \frac{b^a}{\Gamma(a)} \alpha_k^{a-1} e^{-b\alpha_k}$$

に従っているとすると、式(5.82)、式(5.84)より、 $\alpha_k$  の事後分布は、ベイズの定理から

\*30 これは MCMC 法の基本テクニックの一つで、**補助変数法**とよばれています[146].

$$(5.86) \quad p(\boldsymbol{\alpha}|\mathbf{n}) \propto p(\mathbf{n}|\boldsymbol{\alpha})p(\boldsymbol{\alpha}) = \prod_{k=1}^K \frac{b^a}{\Gamma(a)} \alpha_k^{a-1} e^{-b\alpha_k} \times \\ \frac{1}{\Gamma(N)} \int_0^1 \theta^{\sum_k \alpha_k - 1} (1-\theta)^{N-1} d\theta \times \prod_{k=1}^K \left( \prod_{i=0}^{n_k-1} \sum_{x_{ki} \in \{0,1\}} \alpha_k^{x_{ki}} i^{1-x_{ki}} \right)$$

となります。これは  $\theta$  と  $x_{ki}$  を周辺化した形ですから、同時確率は

$$(5.87) \quad p(\boldsymbol{\alpha}, \theta, x|\mathbf{n}) \propto \prod_{k=1}^K \frac{b^a}{\Gamma(a)} \alpha_k^{a-1} e^{-b\alpha_k} \times \frac{1}{\Gamma(N)} \theta^{\sum_k \alpha_k - 1} (1-\theta)^{N-1} \\ \times \prod_{k=1}^K \left( \prod_{i=0}^{n_k-1} \alpha_k^{x_{ki}} i^{1-x_{ki}} \right)$$

です。よって各  $\alpha_k$  の事後分布は、式(5.87)から  $\alpha_k$  に関する項を抜き出すと、補助変数  $\theta, x_{ki}$  が与えられた下では

$$(5.88) \quad p(\alpha_k|\mathbf{n}, \theta, x) \propto p(\alpha_k, \theta, x|\mathbf{n}) = \alpha_k^{a-1} \cdot e^{-b\alpha_k} \cdot \theta^{\alpha_k} \cdot \prod_{i=0}^{n_k-1} \alpha_k^{x_{ki}} \\ = \alpha_k^{a + \sum_{i=0}^{n_k-1} x_{ki} - 1} e^{-(b - \log \theta) \alpha_k} \\ \sim \text{Ga}(a + \sum_{i=0}^{n_k-1} x_{ki}, b - \log \theta)$$

となることがわかります。同様にして補助変数  $\theta, x_{ki}$  についても対応する項を抜き出すと、その分布は

$$(5.89) \quad p(\theta|\mathbf{n}, \boldsymbol{\alpha}, y) \sim \text{Be}(\sum_k \alpha_k, N)$$

$$(5.90) \quad p(x_{ki}|\mathbf{n}, \boldsymbol{\alpha}, \theta) \propto \alpha_k^{x_{ki}} i^{1-x_{ki}} \propto \left( \frac{\alpha_k}{\alpha_k + i} \right)^{x_{ki}} \left( \frac{i}{\alpha_k + i} \right)^{1-x_{ki}} \\ \sim \text{Bernoulli} \left( \frac{\alpha_k}{\alpha_k + i} \right)$$

となります。この分布から補助変数  $\theta, x_{ki}$  をサンプリングしてから、式(5.88)を使って  $\alpha_k$  をサンプリングします。実際には、式(5.79)の  $p(\mathbf{n}|\boldsymbol{\alpha})$  は  $N$  個の文書について存在しますから、同様に計算すると、 $N_i$  を文書  $i$  の長さとして

$$(5.91) \quad \begin{cases} p(\theta_i | n, \boldsymbol{\alpha}, x) \sim \text{Be}(\sum_k \alpha_k, N_i) \\ p(x_{ikj} | n, \boldsymbol{\alpha}, \theta) \sim \text{Bernoulli}\left(\frac{\alpha_k}{\alpha_k + j}\right) \\ p(\alpha_k | n, \theta, x) \sim \text{Ga}\left(a + \sum_{i=1}^N \sum_{j=0}^{n(i,k)-1} x_{ikj}, b - \sum_{i=1}^N \log \theta_i\right) \end{cases}$$

で,  $\boldsymbol{\alpha}$  をガンマ事後分布からサンプリングすることができます.  $\square$

$\eta$  についても, 式(5.78)から同様にして計算すれば, 補助変数  $\phi_k, y_{kvj}$  を導入することで

$$(5.92) \quad \begin{cases} p(\phi_k | m, \eta, y) \sim \text{Be}(V\eta, \sum_{v=1}^V m(k, v)) \\ p(y_{kvj} | m, \eta, \phi) \sim \text{Bernoulli}\left(\frac{\eta}{\eta + j}\right) \\ p(\eta | m, \phi, y) \sim \text{Ga}\left(a + \sum_{k=1}^K \sum_{v=1}^V \sum_{j=0}^{m(k,v)-1} y_{kvj}, b - \sum_{k=1}^K \log \phi_k\right) \end{cases}$$

で, ガンマ事後分布からサンプリングを行えることがわかります.  $\square$