

多項分布から Dirichlet 過程まで  
Version 0.03

村脇 有吾

最終更新: 2012 年 10 月 29 日

- 正規分布などは無視して多項分布一本に絞り込む
- かなり基本的なところから始める
- 図をほとんど書いていないので、適宜板書するかも
- Bayesian で、主に教師なし学習で使う話をやる
- Dirichlet 過程の説明は端折る
- inference の説明は省略するので、実装できるまでは至らない
- 実装したい人は Graham さんの [ノンパラメトリックベイズ入門](#)を参照のこと

確率変数  $x$  は  $K$  個のうちいずれかの値をとる。その確率は、例えば以下に従う。

$$\theta = (0.1, 0.6, 0.3)$$

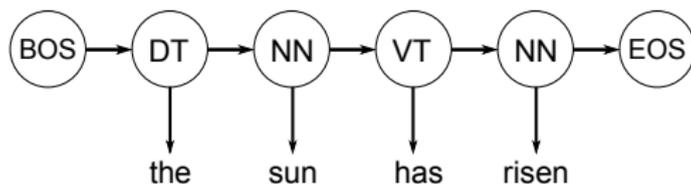
赤 黄 緑

多くのタスク設定では、真の  $\theta$  は不明であり、観測済みのデータ  $x = (x_1, \dots, x_N)$  から推定する。推定された  $\theta$  の主な用途は、新たな  $x'$  の確率を予測することである。

Bayes 推定では、 $\theta$  を積分消去してしまい、 $x$  から直接  $x'$  の確率を予測する。

- N-gram language model ( $V$  次元の語彙から一つの単語を生成)
- HMM POS tagging
- Probabilistic context-free grammar (PCFG)
- Topic model

## 例: Hidden Markov Model による POS tagging



品詞数を  $K$ 、単語数を  $V$  とする (簡単のため BOS と EOS は無視)。

遷移確率 (transition probability) は  $K$  個の多項分布  $(\theta_{DT}^{TRANS}, \theta_{NN}^{TRANS}, \dots)$  からなり、各多項分布は  $K$  値。

$$\theta_{DT}^{TRANS} = (0.21, 0.27, 0.09, \dots)$$

NN NNS ADJ

出力確率 (emission probability) は  $K$  個の多項分布  $(\theta_{DT}^{EMIT}, \theta_{NN}^{EMIT}, \dots)$  からなり、各多項分布は  $V$  値。

$$\theta_{NN}^{EMIT} = (0.012, 0.002, 0.005, \dots)$$

sun rose cat

- 多項分布を導入
- パラメータの最尤推定
- 事前確率の導入と MAP 推定
- Bayes 推定 (← 結局使うのはこれ!)

$\theta$  が既知とする。

確率変数  $x \in \{1, \dots, K\}$  は以下に従う。

$$p(x = k|\theta) = \theta_k$$

ここで  $0 \leq \theta_k \leq 1$ ,  $\sum_{k=1}^K \theta_k = 1$ 。

例えば、 $\theta = (0.1, 0.6, 0.3)$  について、 $p(x = 2|\theta) = \theta_2 = 0.6$ 。

$N$  回の観測  $\mathbf{x} = (x_1, \dots, x_N)$  を得る確率は以下に従う。

$$p(\mathbf{x}|\boldsymbol{\theta}) = \prod_{i=1}^N \theta_{x_i} = \prod_{k=1}^K \theta_k^{n_k}$$

ここで  $n_k$  は  $\mathbf{x}$  における値  $k$  の観測数 ( $\sum_{k=1}^K n_k = N$ )。

$p(\mathbf{x}|\boldsymbol{\theta})$  は  $\mathbf{x}$  内の順序に依存しない。 $N$  個の観測の総和のみに依存している (十分統計量)。

$x = (x_1, \dots, x_N)$  の代わりに、観測数  $\mathbf{n} = (n_1, \dots, n_K)$  を得る確率を考えると、

$$\begin{aligned}\text{Multi}(\mathbf{n}|\boldsymbol{\theta}) &= \binom{N}{n_1 \dots n_K} \prod_{k=1}^K \theta_k^{n_k} \\ &= \frac{N!}{n_1! \dots n_K!} \prod_{k=1}^K \theta_k^{n_k} \\ &= \frac{\Gamma(N+1)}{\prod_{k=1}^K \Gamma(n_k+1)} \prod_{k=1}^K \theta_k^{n_k}\end{aligned}$$

となる。ここで  $\Gamma(n) = (n-1)!$ 。  $\mathbf{n}$  となるような  $x$  の組み合わせの数を  $p(x|\boldsymbol{\theta})$  にかけている。

自然言語処理の文献では  $\mathbf{n}$  ではなく  $x$  の分布をを多項分布とよぶ場合が多く、注意が必要 (Minka, 2003)。  $x$  の分布を Categorical distribution とよぶ流儀もある。

## $N$ 回の観測から $\theta$ を推定

$\theta$  が未知のとき、観測データ  $n$  から推定する。

まず尤度 (likelihood) 関数を以下のように定義。

$$L(\theta; n) = \text{Multi}(n|\theta)$$

以下を満たす  $\theta^{\text{ML}}$  を求める。これを最尤推定 (maximum likelihood estimation) とよぶ。

$$\theta^{\text{ML}} = \underset{\theta}{\operatorname{argmax}} L(\theta; n) = \underset{\theta}{\operatorname{argmax}} \text{Multi}(n|\theta)$$

これを解くと以下が得られる。

$$\theta_k^{\text{ML}} = \frac{n_k}{N}$$

これは Lagrange の未定乗数法を用いれば求まる (高村本 1.5.2 参照)。

低頻度のパラメータが信用出来ない。特に、観測回数が0のパラメータが出現すると、予測確率が0になってしまい都合が悪い。

観測データを大きくすれば大丈夫か？ No. 自然言語処理で用いられる多項分布は一般に sparse (参考: Zipf の法則)。常に低頻度のパラメータが存在する (Johnson, 2007)。

事前分布  $p(\theta)$  を導入する。これはデータを観測する前に (恣意的に) 仮定する  $\theta$  の分布。

Bayes の定理により

$$\begin{aligned} p(\theta|x) &= \frac{p(x|\theta)p(\theta)}{p(x)} \\ &\propto p(x|\theta)p(\theta) \end{aligned}$$

$p(\theta|x)$  を事後分布 (posterior distribution)、 $p(x|\theta)$  を尤度、 $p(\theta)$  を事前分布 (prior distribution) とよぶ。 $p(\theta|x)$  は、データを観測した後の  $\theta$  の分布と解釈できる。

多項分布に対しては、事前分布として Dirichlet 分布  $\text{Dir}(\theta|\alpha)$  を仮定すると計算上都合がよい。

$$p(\theta|n, \alpha) \propto \text{Multi}(n|\theta)\text{Dir}(\theta|\alpha)$$

ここで

$$\text{Dir}(\theta|\alpha) = \frac{\Gamma(A)}{\Gamma(\alpha_1) \cdots \Gamma(\alpha_K)} \prod_{k=1}^K \theta_k^{\alpha_k - 1}$$

ここで  $\alpha = (\alpha_1, \dots, \alpha_K)$ ,  $\alpha_k > 0$ ,  $A = \sum_{k=1}^K \alpha_k$ 。特に理由がなければ  $\alpha_i = \alpha_j$  とする。 $\alpha$  は事前に与えるパラメータであり、パラメータのパラメータであるためハイパーパラメータとよばれる。

$n_k$  が値  $k$  の観測上のカウントなのに対して、 $\alpha_k$  は事前に仮定されたカウント (擬似カウント) と解釈できる。ただし、 $\alpha_k$  は整数でなくてもよい。

Dirichlet 分布と多項分布 (再掲) の式の類似に注意。

$$\text{Multi}(\mathbf{n}|\boldsymbol{\theta}) = \frac{\Gamma(N+1)}{\prod_{k=1}^K \Gamma(n_k+1)} \prod_{k=1}^K \theta_k^{n_k}$$

事後分布  $p(\boldsymbol{\theta}|\mathbf{n}, \boldsymbol{\alpha}) \propto \text{Multi}(\mathbf{n}|\boldsymbol{\theta})\text{Dir}(\boldsymbol{\theta}|\boldsymbol{\alpha})$  を求める (計算省略)。

$$\begin{aligned} p(\boldsymbol{\theta}|\mathbf{n}, \boldsymbol{\alpha}) &= \frac{\Gamma(N+A)}{\Gamma(n_1+\alpha_1)\cdots\Gamma(n_K+\alpha_K)} \prod_{k=1}^K \theta_k^{n_k+\alpha_k-1} \\ &= \text{Dir}(\boldsymbol{\theta}|\mathbf{n}+\boldsymbol{\alpha}) \end{aligned}$$

事後分布は事前分布と同じ形になる。この性質を共役性 (conjugacy) とよぶ。

## MAP (最大事後確率) 推定

事後確率を最大化するような  $\theta$  の推定を最大事後確率 (maximum a posteriori) 推定とよぶ。

$$\begin{aligned}\theta^{\text{MAP}} &= \underset{\theta}{\operatorname{argmax}} p(\theta|\mathbf{n}) \\ &= \underset{\theta}{\operatorname{argmax}} p(\mathbf{n}|\theta)p(\theta)\end{aligned}$$

Dirichlet-多項分布の場合は以下のようになる。

$$\begin{aligned}\theta^{\text{MAP}} &= \underset{\theta}{\operatorname{argmax}} \operatorname{Dir}(\theta|\mathbf{n} + \alpha) \\ \theta_k^{\text{MAP}} &= \frac{n_k + \alpha_k - 1}{\sum_{k=1}^K (n_k + \alpha_k - 1)} = \frac{n_k + \alpha_k - 1}{N + A - K}\end{aligned}$$

最尤推定  $\theta_k^{\text{ML}} = \frac{n_k}{N}$  と比較すると、加算スムージングが行われていることが確認できる。

Bayes 推定では、MAP 推定のように  $\theta$  を一つに決めるのではなく、 $\theta$  の分布を保持したまま、積分消去して利用する。

$$p(\mathbf{x}) = \int_{\theta} p(\mathbf{x}|\theta)p(\theta) d\theta$$

新たな  $x'$  の生成確率 (予測確率) を以下で求める。

$$\begin{aligned} p(x'|\mathbf{x}) &= \int_{\theta} p(x'|\theta)p(\theta|\mathbf{x}) d\theta \\ &= \int_{\theta} p(x'|\theta) \frac{p(\mathbf{x}|\theta)p(\theta)}{p(\mathbf{x})} d\theta \end{aligned}$$

一般に積分は困難だが、Dirichlet-多項分布の場合は解析的に解ける。

$$\begin{aligned}
p(\mathbf{n}|\boldsymbol{\alpha}) &= \int_{\boldsymbol{\theta}} \text{Multi}(\mathbf{n}|\boldsymbol{\theta}) \text{Dir}(\boldsymbol{\theta}|\boldsymbol{\alpha}) d\boldsymbol{\theta} \\
&= \int_{\boldsymbol{\theta}} \left( \frac{\Gamma(N+1)}{\prod_{k=1}^K \Gamma(n_k+1)} \prod_{k=1}^K \theta_k^{n_k} \right) \left( \frac{\Gamma(A)}{\Gamma(\alpha_1) \cdots \Gamma(\alpha_K)} \prod_{k=1}^K \theta_k^{\alpha_k-1} \right) d\boldsymbol{\theta} \\
&= \frac{\Gamma(N+1)}{\prod_{k=1}^K \Gamma(n_k+1)} \frac{\Gamma(A)}{\Gamma(\alpha_1) \cdots \Gamma(\alpha_K)} \int_{\boldsymbol{\theta}} \prod_{k=1}^K \theta_k^{n_k+\alpha_k-1} d\boldsymbol{\theta} \\
&= \frac{\Gamma(N+1)}{\prod_{k=1}^K \Gamma(n_k+1)} \frac{\Gamma(A)}{\prod_{k=1}^K \Gamma(\alpha_k)} \frac{\prod_{k=1}^K \Gamma(n_k+\alpha_k)}{\Gamma(N+A)} \\
&= \frac{\Gamma(N+1)}{\prod_{k=1}^K \Gamma(n_k+1)} \frac{\Gamma(A)}{\Gamma(N+A)} \prod_{k=1}^K \frac{\Gamma(n_k+\alpha_k)}{\Gamma(\alpha_k)}
\end{aligned}$$

$p(n|\alpha)$  の最初の項は組み合わせ数。これを除くと  $x$  の生成確率となる。

$$p(\mathbf{x}|\alpha) = \frac{\Gamma(A)}{\Gamma(N+A)} \prod_{k=1}^K \frac{\Gamma(n_k + \alpha_k)}{\Gamma(\alpha_k)}$$

なお、 $x$  の生成過程を以下のように表記する。

$$x_n | \theta \sim \theta$$

$$\theta | \alpha \sim \text{Dir}(\alpha)$$

$x_n | \theta \sim \text{Multi}(\theta)$  と表記することもある。

$x'$  の予測確率は以下で与えられる ( $\theta^{\text{MAP}}$  と比較せよ)。

$$\begin{aligned}
 p(x' = k' | \mathbf{x}, \alpha) &= \frac{p(x' = k', \mathbf{x} | \alpha)}{p(\mathbf{x} | \alpha)} \\
 &= \frac{\frac{\Gamma(A)}{\Gamma(N+1+A)} \prod_{k=1}^K \frac{\Gamma(n_k + I(k=k') + \alpha_k)}{\Gamma(\alpha_k)}}{\frac{\Gamma(A)}{\Gamma(N+A)} \prod_{k=1}^K \frac{\Gamma(n_k + \alpha_k)}{\Gamma(\alpha_k)}} \\
 &= \frac{\Gamma(N+A)}{\Gamma(N+1+A)} \frac{\Gamma(n_{k'} + 1 + \alpha_{k'})}{\Gamma(n_{k'} + \alpha_{k'})} \\
 &= \frac{n_{k'} + \alpha_{k'}}{N + A}
 \end{aligned}$$

ここで  $I(\text{statement})$  は statement が真なら 1、そうでなければ 0。

$$\begin{aligned}
 p(\mathbf{x}_{1:N}|\boldsymbol{\alpha}) &= p(x_N|\mathbf{x}_{1:N-1}, \boldsymbol{\alpha})p(\mathbf{x}_{1:N-1}|\boldsymbol{\alpha}) \\
 &= \frac{n_{x_N}(\mathbf{x}_{1:N-1}) + \alpha_{x_N}}{N - 1 + A} p(\mathbf{x}_{1:N-1}|\boldsymbol{\alpha}) \\
 &= \frac{n_{x_N}(\mathbf{x}_{1:N-1}) + \alpha_{x_N}}{N - 1 + A} p(x_{N-1}|\mathbf{x}_{1:N-2}, \boldsymbol{\alpha})p(\mathbf{x}_{1:N-2}|\boldsymbol{\alpha}) \\
 &= \frac{n_{x_N}(\mathbf{x}_{1:N-1}) + \alpha_{x_N}}{N - 1 + A} \frac{n_{x_{N-1}}(\mathbf{x}_{1:N-2}) + \alpha_{x_{N-1}}}{N - 2 + A} p(\mathbf{x}_{1:N-2}|\boldsymbol{\alpha}) \\
 &= \frac{\prod_{k=1}^K \prod_{i=1}^{n_k(\mathbf{x}_{1:N})} i - 1 + \alpha_k}{\prod_{i=1}^N i - 1 + A} \\
 &= \frac{\Gamma(A)}{\Gamma(N + A)} \prod_{k=1}^K \frac{\Gamma(n_k(\mathbf{x}_{1:N}) + \alpha_k)}{\Gamma(\alpha_k)}
 \end{aligned}$$

ここで  $\mathbf{x}_{1:N} = (x_1, \dots, x_N)$ ,  $n_k(\mathbf{x}) =$  number of times  $k$  appears in  $\mathbf{x}$ 。

## 交換可能性 (Exchangeability)

$p(x|\alpha)$  は  $x$  内の順序に依存せず、 $N$  個の観測の総和のみに依存している。つまり、 $x$  内の任意の要素ペア  $x_i, x_j$  を置換して作った  $x'$  を考えた時、

$$p(x'|\alpha) = p(x|\alpha)$$

がなりたつ。この性質を交換可能性という。交換可能性のおかげで inference が簡単に実現できるのだが、それを説明するにはまずサンプリングを説明しないといけない。

赤、黄、緑の3種類のいずれかの値を取る  $x$  を考える。いま (赤, 赤, 黄) を観測したとき、この3個の生成確率は以下で求まる。

$$p(x|\alpha) = \frac{0 + \alpha_{\text{赤}}}{0 + A} \frac{1 + \alpha_{\text{赤}}}{1 + A} \frac{0 + \alpha_{\text{黄}}}{2 + A}$$

この結果は観測する順序を変えても変わらない。

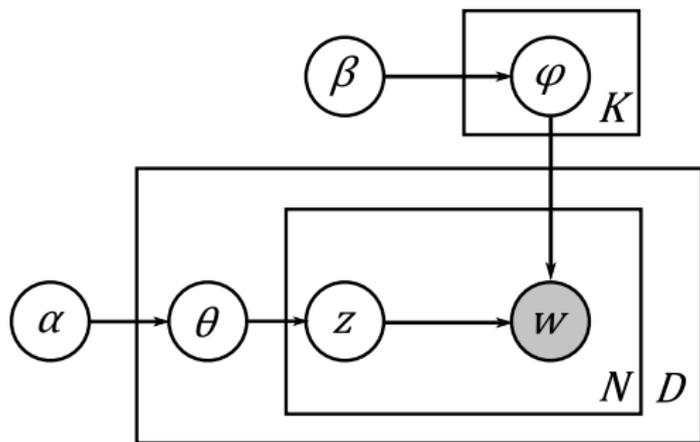
## 例題: LDA (Latent Dirichlet allocation)

(Blei+, 2002 のオリジナルではなく、Griffiths+, 2004 のモデル)

$D$  個の文書が与えられ、文書  $i$  について  $N_i$  個の単語が観測されるとする。LDA では、各文書は  $K$  個のトピックの混合であると仮定する。単語  $w_{i,j}$  はトピック  $k$  に属しており、文書  $i$  は  $K$  個のトピックの混合分布  $\theta_i$  を持つ。トピック  $k$  は、それに対応する単語分布  $\varphi_k$  を持っており、 $w_{i,j}$  は単語分布  $\varphi_k$  から生成されたものである。

## 例題: LDA (Latent Dirichlet allocation)

生成過程は有向グラフで表現できる。



プレートは繰り返し、色付きのノードは観測変数を表す。  
 $\alpha$ と $\beta$ のような所与のパラメータは黒点で表す流儀もある(その方が良さそう)。

## 例題: LDA (Latent Dirichlet allocation)

以下の生成過程で表される。

- ① 各トピック  $k \in \{1, \dots, K\}$  について、単語分布  $\varphi_k \sim \text{Dir}(\beta)$  を生成
- ② 文書  $i \in \{1, \dots, D\}$  について、トピック分布  $\theta_i \sim \text{Dir}(\alpha)$  を生成
- ③ 文書  $i$  内の単語  $j \in \{1, \dots, N_i\}$  について
  - ① トピック  $z_{i,j} \sim \text{Multi}(\theta_i)$  を生成
  - ② 単語  $w_{i,j} \sim \text{Multi}(\varphi_{z_{i,j}})$  を生成

以後、コーパス中の単語列を  $W$ 、対応するトピック列を  $Z$  で表す。いま  $W$  が観測されたとき、 $Z$  を推定したい。

## 例題: LDA (Latent Dirichlet allocation)

まず同時分布を考える。

$$p(\mathbf{W}, \mathbf{Z}, \boldsymbol{\theta}, \boldsymbol{\varphi} | \boldsymbol{\alpha}, \boldsymbol{\beta}) = \prod_{k=1}^K p(\boldsymbol{\varphi}_k | \boldsymbol{\beta}) \times \prod_{i=1}^D \left( p(\boldsymbol{\theta}_i | \boldsymbol{\alpha}) \prod_{j=1}^{N_i} p(z_{i,j} | \boldsymbol{\theta}_i) p(w_{i,j} | \boldsymbol{\varphi}_{z_{i,j}}) \right)$$

$\boldsymbol{\theta}_i$  と  $\boldsymbol{\varphi}_k$  は積分消去できる。

$$\begin{aligned} p(\mathbf{W}, \mathbf{Z} | \boldsymbol{\alpha}, \boldsymbol{\beta}) &= \int_{\boldsymbol{\theta}} \int_{\boldsymbol{\varphi}} p(\mathbf{W}, \mathbf{Z}, \boldsymbol{\theta}, \boldsymbol{\varphi} | \boldsymbol{\alpha}, \boldsymbol{\beta}) d\boldsymbol{\theta} d\boldsymbol{\varphi} \\ &= \prod_{i=1}^D \left( \frac{\Gamma(\alpha_{(\cdot)})}{\Gamma(n_{i,(\cdot)}^{(\cdot)} + \alpha_{(\cdot)})} \prod_{k=1}^K \frac{\Gamma(n_{i,(\cdot)}^k + \alpha_k)}{\Gamma(\alpha_k)} \right) \\ &\quad \times \prod_{k=1}^K \left( \frac{\Gamma(\beta_{(\cdot)})}{\Gamma(n_{(\cdot),(\cdot)}^k + \beta_{(\cdot)})} \prod_v \frac{\Gamma(n_{(\cdot),v}^k + \beta_v)}{\Gamma(\beta_v)} \right) \end{aligned}$$

ここで  $n_{i,v}^k$  は文書  $i$  における単語  $v$  がトピック  $k$  を取る回数。  
 $D + K$  個の Dirichlet-多項分布の組み合わせであることが確認できる。

## 例題: LDA (Latent Dirichlet allocation)

文書からトピックの、トピックからの単語の逐次生成を考える。  
交換可能性が成り立つことも容易に確かめられる。

文書 ID	トピック	単語	生成確率
1	政治	オバマ	$\frac{0+\alpha_{\text{政治}}}{0+\alpha_{(\cdot)}} \times \frac{0+\beta_{\text{オバマ}}}{0+\beta_{(\cdot)}}$
	経済	FRB	$\frac{0+\alpha_{\text{経済}}}{1+\alpha_{(\cdot)}} \times \frac{0+\beta_{\text{FRB}}}{0+\beta_{(\cdot)}}$
	政治	選挙	$\frac{1+\alpha_{\text{政治}}}{2+\alpha_{(\cdot)}} \times \frac{0+\beta_{\text{選挙}}}{1+\beta_{(\cdot)}}$
	政治	オバマ	$\frac{2+\alpha_{\text{政治}}}{3+\alpha_{(\cdot)}} \times \frac{1+\beta_{\text{オバマ}}}{2+\beta_{(\cdot)}}$
2	経済	バーナンキ	$\frac{0+\alpha_{\text{経済}}}{0+\alpha_{(\cdot)}} \times \frac{0+\beta_{\text{バーナンキ}}}{1+\beta_{(\cdot)}}$
	経済	FRB	$\frac{1+\alpha_{\text{経済}}}{1+\alpha_{(\cdot)}} \times \frac{1+\beta_{\text{FRB}}}{2+\beta_{(\cdot)}}$

求めたいのは

$$\operatorname{argmax}_Z p(Z|W, \alpha, \beta) = \operatorname{argmax}_Z p(W, Z|\alpha, \beta)$$

探索空間は  $K^{|W|}$  通り。すべての組み合わせを調べるのは現実的でない。そこで、Gibbs サンプリングを用いた random walk による探索がよく用いられる (詳細省略)。

ただし、

- ① サンプリング自体は、 $\operatorname{argmax}$  を求めるアルゴリズムではない。
- ② サンプリングが簡単に実現できるのは、交換可能性がなりたつから。

## Dirichlet 過程 (Dirichlet process)

LDA のトピック数  $K$  は事前に人手で与える。しかし、最適な  $K$  は事前にはわからず、データから自動推定したい。

そこで導入されるのが Dirichlet 過程である。Dirichlet 過程も多項分布と共役であり、事後分布も Dirichlet 過程となる。

$$\begin{aligned}x|G &\sim G \\ G|\alpha, G_0 &\sim \text{DP}(\alpha, G_0)\end{aligned}$$

ここで  $G_0$  は基底測度。確率分布と思えばよい。Dirichlet 分布の  $\alpha_k$  が  $\alpha G_0(k)$  に相当する。ただし、Dirichlet 過程においては、 $K$  は無限である。もっとも、有限の  $N$  個のデータが属すトピック数  $K'$  は有限 ( $K' \leq N$ ) である。よって、無限のトピック数のうち、 $K'$  個のトピックのみを明示的に扱い、残りのトピックは遅延評価すれば手続き的に扱える。

$G$  を積分消去する (Blackwell and MacQueen, 1973)。

$$p(x' = k' | \mathbf{x}, \alpha, G_0) = \int_G p(x' = k' | G) P(G | \mathbf{x}, \alpha, G_0) dG$$

$$\begin{aligned} x' | \mathbf{x}, \alpha, G_0 &\sim \frac{1}{N + \alpha} \sum_{j=1}^N \delta(x_j) + \frac{\alpha}{N + \alpha} G_0 \\ &\sim \sum_k \frac{n_k}{N + \alpha} \delta(k) + \frac{\alpha}{N + \alpha} G_0 \end{aligned}$$

ここで  $\delta(k)$  は  $k$  にすべての確率質量を置く分布。

Dirichlet 過程は中華料理店過程により解釈できる。無限個のテーブルを持ち、各テーブルに無限の客を収容できる料理店を考える。新たに入ってきた客は、既に客がいるテーブルに客の数に比例した確率で座る。あるいは、新たなテーブルに  $\alpha$  に比例する確率で座る。

$$p(x_i = l_k | \mathbf{x}, \alpha) = \begin{cases} \frac{n_k}{N + \alpha} & 1 \leq k \leq K(\mathbf{x}) \\ \frac{\alpha}{N + \alpha} & k = K(\mathbf{x}) + 1 \end{cases}$$

ここで  $K(\mathbf{x})$  は、 $\mathbf{x}$  が座っているテーブルの総数。なお、最初の客は必ず新たなテーブルに座る。

トピック・モデルの場合、 $k$  はトピックのインデックスを表し、単語分布  $\varphi_k$  に対応付けられている。 $G_0$  を一様な連続分布とすれば、各テーブルがトピックに対応する。ただし、このままではトピック・モデルは機能しない。各文書の  $G_i$  間でトピックを共有するには階層 Dirichlet 過程が必要となる (Teh et al., 2006)。以下では、より簡単な unigram 単語分布を例題とする。

## 2 段階生成モデル

文を単語へ分割することを考える (いわゆる単語分割)。考えられる文の分割、それに応じた単語ラベルの総数は組み合わせ爆発を起こす。そこで、任意の文字列に対して 0 より大きな確率を与える単語確率が必要となる。こうしたモデルを Dirichlet 過程を用いて実現する。

以下のように、文字列から単語を生成すれば、任意の文字列からなる単語に適当な確率を割り振れる。

$$G_0(l = c_1, \dots, c_M) = p_{\#}(1 - p_{\#})^{M-1} \prod_{i=1}^M p(c_i)$$

ここで  $p_{\#}$  は単語境界の確率。この  $G_0$  を基底分布とする単語分布を考える (unigram モデル)。

## 2 段階生成モデル

CRP の各テーブルに単語ラベル  $l$  を関連付ける必要がある。これは、 $G_0$  による単語ラベルの生成 (generator) と、CRP による客のテーブルへの割り当て (adaptor) という 2 段階の生成によって実現される。

新たな客が来たとき以下のいずれかが行われる。

- $\frac{n_k}{N+\alpha}$  の確率で既存のテーブル  $k$  に座る。このとき、割り当て済みの単語ラベル  $l_k$  が使われる。
- $\frac{\alpha}{N+\alpha}$  の確率で新たなテーブルに座る。このとき、 $G_0$  に従って新たなラベル  $l_{K(x+1)}$  を生成し、これをテーブルに割り当てる。

一般に複数のテーブルに同じラベルが割り当てられることに注意。

## 2 段階生成モデル

単語列  $w_{1:N}$  およびそのテーブル割り当て  $z_{1:N}$  が与えられたとき、予測確率は以下で与えられる。

$$\begin{aligned} p(w_i = l | w_{1:N}, z_{1:N}, \alpha, G_0) &= \sum_{k=1}^{K(z_{1:N})} I(l_k = l) \frac{n_k(z_{1:N})}{N + \alpha} + \frac{\alpha}{\alpha + N} G_0(l) \\ &= \frac{n_l(w_{1:N}) + \alpha G_0(l)}{N + \alpha} \end{aligned}$$

ここで  $n_l(w)$  は、 $w$  におけるラベル  $l$  の出現総数。つまり、実は予測確率はテーブル割り当てに関係なく

$$p(w_i = l | w_{1:N}, \alpha, G_0) = \frac{n_l(w_{1:N}) + \alpha G_0(l)}{N + \alpha}$$

となる (階層 Dirichlet 過程や、(階層) Pitman-Yor 過程を用いた場合はテーブル割り当てに依存する)。

分割なしのコーパス  $C$ 、これに対応する単語列  $W$  の生成確率は、Dirichlet-多項分布の生成確率と同様に以下で表される。

$$p(C, W | \alpha, G_0) = \frac{\Gamma(\alpha)}{\Gamma(n_{(\cdot)}(W) + \alpha)} \prod_{l \in \text{voc}(W)} \frac{\Gamma(n_l(W) + \alpha G_0(l))}{\Gamma(\alpha G_0(l))}$$

ここで  $\text{voc}(W)$  は  $W$  に含まれる単語ラベルの集合を返す関数。

$\text{argmax}_W p(W|C, \alpha, G_0) = \text{argmax}_W p(C, W|\alpha, G_0)$  を求めたい。すべての  $W$  を調べるのは不可能なので、Gibbs サンプリングによる randomized search が用いられる。

- Gibbs サンプルングの理論的背景、中身 & 実装まわり
- ブロック・サンプルング、特に type-based のサンプルング
- Stick-breaking 過程による Dirichlet 過程の解釈
- 変分 Bayes (variational Bayes)
- 階層 Dirichlet 過程
- Dirichlet 過程の一般化である Pitman-Yor 過程