

# 自然言語処理のための Gibbs サンプルング

## Version 0.03

村脇 有吾

最終更新: 2012 年 11 月 5 日

- 多項分布から Dirichlet 過程までの続き
- 例によって図は板書する方向
- Gibbs サンプリングの理論的背景は端折る
- 引き続き例は LDA と Unigram 単語分割モデル
- 目標は Type-based サンプリングの理解

Dirichlet-多項分布の生成確率は以下で求まる (Dirichlet 過程の場合もほとんど同じ)。

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

LDA や Unigram 単語分割モデルは、Dirichlet(過程)-多項分布の組み合わせにすぎない。

## 前回の復習 2/2

$D$  個の文書、 $K$  個のトピックからなる LDA で、単語列  $W$  とトピック割り当て  $Z$  の生成確率は、 $D + K$  個の Dirichlet-多項分布の組み合わせで表される。

$$p(W, Z | \alpha, \beta) = \prod_{i=1}^D \left( \frac{\Gamma(\alpha_{(\cdot)})}{\Gamma(n_{i,(\cdot)}^{(\cdot)} + \alpha_{(\cdot)})} \prod_{k=1}^K \frac{\Gamma(n_{i,(\cdot)}^k + \alpha_k)}{\Gamma(\alpha_k)} \right) \\ \times \prod_{k=1}^K \left( \frac{\Gamma(\beta_{(\cdot)})}{\Gamma(n_{(\cdot),(\cdot)}^k + \beta_{(\cdot)})} \prod_v \frac{\Gamma(n_{(\cdot),v}^k + \beta_v)}{\Gamma(\beta_v)} \right)$$

Unigram 単語分割モデルの場合、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

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

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。 $G_0$  は任意の文字列からなる単語に適切な確率を与える分布。

モデルを設定し、観測データが与えられたとき、モデルに一番フィットする隠れ状態を求めたい。例えば、Unigram 単語分割の場合、 $\operatorname{argmax}_W p(W|C, \alpha, G_0) = \operatorname{argmax}_W p(C, W|\alpha, G_0)$  を求めたい。どうすればよいか。

- 解析的に求める (できるなら教えて!)
- $W$  を総当りで試して  $\max$  を求める (組合せ爆発により実際上不可能)
- 近似的に探索 (Gibbs サンプルングはこれ)

- ① 適当に  $W$  を初期化 ( $W_1$  とする)
- ②  $W_i$  を局所的に変更した  $W_{i+1}$  を確率的に得る

$p(W_i|C, \alpha, G_0)$  は上がったり下がったりするが、長期的には上がっていく。

確率分布からサンプルを得る。例えば、サイコロを振って出た目を記録する。実装的には、rand 関数の戻り値  $r$  ( $0 \leq r < 1$ ) を適当に加工して目的の値を得る。

例えば  $\theta = (0.1, 0.6, 0.3)$  からのサンプル  $x$  は以下の手続きで得られる。

$$x = \begin{cases} 1 & 0 \leq r < 0.1 \\ 2 & 0.1 \leq r < 0.7 \\ 3 & 0.7 \leq r < 1 \end{cases}$$

Dirichlet-多項分布の予測確率は以下で与えられる (復習):

$$p(x' = k' | \mathbf{x}, \boldsymbol{\alpha}) = \frac{n_{k'}(\mathbf{x}) + \alpha_{k'}}{n(\mathbf{x}) + \alpha_{(\cdot)}}$$

いま  $\mathbf{x} = (\text{赤}, \text{赤}, \text{黄}, \text{赤}, \text{赤})$ ,  $\boldsymbol{\alpha} = (\alpha_{\text{赤}}, \alpha_{\text{黄}}, \alpha_{\text{緑}}) = (1, 1, 1)$  とすると、

$$p(x' = \text{赤} | \mathbf{x}, \boldsymbol{\alpha}) = \frac{4 + 1}{5 + 3} = 0.625$$

$$p(x' = \text{黄} | \mathbf{x}, \boldsymbol{\alpha}) = \frac{1 + 1}{5 + 3} = 0.25$$

$$p(x' = \text{緑} | \mathbf{x}, \boldsymbol{\alpha}) = \frac{0 + 1}{5 + 3} = 0.125$$

したがって  $\theta = (0.625, 0.25, 0.125)$  からサンプルを得ると、それがこの予測分布からのサンプルになっている。



Dirichlet(過程)-多項分布の組み合わせ  $p(\mathbf{Z}) = p(z_1, \dots, z_N)$  から  $\mathbf{Z}$  を直接サンプリングするのは難しい。ただし、 $z_i$  を  $p(z_i|\mathbf{Z}^{-i})$  からサンプリングするのは容易。ここで  $\mathbf{Z}^{-i}$  は  $\mathbf{Z}$  から  $z_i$  を除いたもの。そこで、 $p(z_i|\mathbf{Z}^{-i})$  からのサンプリングを用いて、 $\mathbf{Z}$  のサンプルを得たい。これを実現するのが Gibbs サンプリング。

- 1  $\mathbf{Z}$  を適当に初期化
- 2 各ステップ  $\tau = 1, \dots, T$  で以下を行う
  - $z_i$  を  $z_i \sim p(z_i|\mathbf{Z}^{-i})$  からのサンプルで更新
  - 新たな  $\mathbf{Z}$  を用いて同様に  $z_{i+1}$  を更新
  - これをすべての  $i = 1, \dots, N$  に対して行う

これを十分に繰り返すと  $\mathbf{Z}_\tau$  は真の確率分布  $p(\mathbf{Z})$  からのサンプルとなる。理論的背景については PRML 11.3 を参照。

Gibbs サンプリングのキモは、 $p(z_i|Z^{-i})$  から容易にサンプリングできることである。Dirichlet-多項分布 (とその組み合わせ) の場合、交換可能性がなりたつことから容易にサンプリングできる。 $(z_1, \dots, z_{i-1}, z_i, z_{i+1}, \dots, z_N)$  を観測する確率は、順番を入れ換えた  $(z_1, \dots, z_{i-1}, z_{i+1}, \dots, z_N, z_i)$  を観測する確率と等しい。 $p(z_i|Z^{-i})$  は  $(z_1, \dots, z_{i-1}, z_{i+1}, \dots, z_N)$  を観測したもとの  $z_i$  の予測確率を表している。

$$\begin{aligned}
& p(z_i = k' | \mathbf{Z}^{-i}, \mathbf{W}, \alpha, \beta) \\
\propto & p(z_i = k', \mathbf{Z}^{-i}, \mathbf{W} | \alpha, \beta) \\
= & \prod_{j=1}^D \left( \frac{\Gamma(\alpha_{(\cdot)})}{\Gamma(n_{j,(\cdot)}^{(\cdot)}(\mathbf{Z}^{-i}) + I(w_i \in \mathbf{D}_j) + \alpha_{(\cdot)})} \right. \\
& \left. \prod_{k=1}^K \frac{\Gamma(n_{j,(\cdot)}^k(\mathbf{Z}^{-i}) + I(k' = k \& w_i \in \mathbf{D}_j) + \alpha_k)}{\Gamma(\alpha_k)} \right) \\
& \times \prod_{k=1}^K \left( \frac{\Gamma(\beta_{(\cdot)})}{\Gamma(n_{(\cdot),(\cdot)}^k(\mathbf{Z}^{-i}) + I(k' = k) + \beta_{(\cdot)})} \right. \\
& \left. \prod_v \frac{\Gamma(n_{(\cdot),v}^k(\mathbf{Z}^{-i}) + I(k' = k \& w_i = v) + \beta_v)}{\Gamma(\beta_v)} \right)
\end{aligned}$$

$i$  はコーパス全体での通し番号。 $\mathbf{D}_j$  は  $j$  番目の文書の単語の集合。  
ほとんどの項は  $z_i$  と無関係なので無視できる。後半へ続く。

$$\begin{aligned}
 & p(z_i = k' | \mathbf{Z}^{-i}, \mathbf{W}, \boldsymbol{\alpha}, \boldsymbol{\beta}) \\
 \propto & \frac{\Gamma(n_{j,(\cdot)}^{(\cdot)}(\mathbf{Z}^{-i}) + \alpha_{(\cdot)})}{\Gamma(n_{j,(\cdot)}^{(\cdot)}(\mathbf{Z}^{-i}) + 1 + \alpha_{(\cdot)})} \times \frac{\Gamma(n_{j,(\cdot)}^{k'}(\mathbf{Z}^{-i}) + 1 + \alpha_{k'})}{\Gamma(n_{j,(\cdot)}^{k'}(\mathbf{Z}^{-i}) + \alpha_{k'})} \\
 & \times \frac{\Gamma(n_{(\cdot),(\cdot)}^{k'}(\mathbf{Z}^{-i}) + \beta_{(\cdot)})}{\Gamma(n_{(\cdot),(\cdot)}^{k'}(\mathbf{Z}^{-i}) + 1 + \beta_{(\cdot)})} \times \frac{\Gamma(n_{(\cdot),v'}^{k'}(\mathbf{Z}^{-i}) + 1 + \beta_{v'})}{\Gamma(n_{(\cdot),v'}^{k'}(\mathbf{Z}^{-i}) + \beta_{v'})} \\
 = & \frac{n_{j,(\cdot)}^{k'}(\mathbf{Z}^{-i}) + \alpha_{k'}}{n_{j,(\cdot)}^{(\cdot)}(\mathbf{Z}^{-i}) + \alpha_{(\cdot)}} \times \frac{n_{(\cdot),v'}^{k'}(\mathbf{Z}^{-i}) + \beta_{v'}}{n_{(\cdot),(\cdot)}^{k'}(\mathbf{Z}^{-i}) + \beta_{(\cdot)}}
 \end{aligned}$$

ここで  $v' = w_i$ ,  $j'$  は  $w_i$  の属す文書のインデックスとする。  
 第1項は文書  $D_j$  に属すトピックの Dirichlet-多項分布に  $k'$  を追加するときの予測確率、第2項はトピック  $k'$  に属す単語の Dirichlet-多項分布に  $v'$  を追加するときの予測確率となっている (実際には第1項の分母は  $k'$  によらないので無視できる)。

$$\frac{n_{j',(\cdot)}^{k'}(Z^{-i})+\alpha_{k'}}{n_{j',(\cdot)}^{(\cdot)}(Z^{-i})+\alpha_{(\cdot)}} \times \frac{n_{(\cdot),v'}^{k'}(Z^{-i})+\beta_{v'}}{n_{(\cdot),(\cdot)}^{k'}(Z^{-i})+\beta_{(\cdot)}}$$

- $z_i$  の値は、第 1 項と第 2 項の積を大きくするような  $k$  に (確率的に) なりやすい
- 第 1 項を大きくするには、文書  $D_j$  の中で頻出する  $k$  を選ばよ
- 第 2 項を大きくするには、 $w_i$  が頻出するトピック  $k$  を選ばよ
- サンプリングを繰り返すと、 $Z$  の値は、これらを全体的に満たすようになっていく

以下のデータを管理する。

- $W$  と現在の  $Z$
- $n_{j,(\cdot)}^k$  (文書  $D_j$  中のトピック  $k$  の割り当て数)
- $(n_{j,(\cdot)}^{(\cdot)})$  (文書  $D_j$  中の単語総数)
- $n_{(\cdot),v}^k$  (トピック  $k$  の単語  $v$  の頻度)
- $n_{(\cdot),(\cdot)}^k$  (トピック  $k$  の単語総数)

$z_i$  のサンプリングは以下の手順で行う。

- 1 現在の  $z_i$  に対応するカウントを decrement
- 2  $p(z_i|Z^{-i}, W, \alpha, \beta)$  (に比例する値) を  $K$  通り計算
- 3 乱数を振って新たな  $z_i$  を決め、 $Z$  を更新
- 4 新たな  $z_i$  に対応するカウントを increment

サンプリングで得られるのは、あくまで確率分布からのサンプルであり、明示的に最大値を求めているわけではない。LDA の場合、 $p(Z|W, \alpha, \beta)$  からのサンプルであり、この値を小さくするような  $Z$  も (非常に小さな確率で) 得られる。ただ、LDA などの確率分布は非常に peaky であり、単にサンプリングを繰り返しているだけで極大値付近に滞留することが多い。

真面目に最大値を近似的に求めるための拡張として焼きなまし法がある (詳細省略)。

サンプリングがうまく行っているかを確認するには、以下の値を追跡する。

$$p(W, Z|\alpha, \beta) = \prod_{i=1}^D \left( \frac{\Gamma(\alpha_{(\cdot)})}{\Gamma(n_{i,(\cdot)}^{(\cdot)} + \alpha_{(\cdot)})} \prod_{k=1}^K \frac{\Gamma(n_{i,(\cdot)}^k + \alpha_k)}{\Gamma(\alpha_k)} \right) \\ \times \prod_{k=1}^K \left( \frac{\Gamma(\beta_{(\cdot)})}{\Gamma(n_{(\cdot),(\cdot)}^k + \beta_{(\cdot)})} \prod_v \frac{\Gamma(n_{(\cdot),v}^k + \beta_v)}{\Gamma(\beta_v)} \right)$$

この値が下がり続けていたりしたらバグっている可能性が高い。

実装上のヒント:  $p(W, Z|\alpha, \beta)$  は非常に小さな値なので、そのままでは underflow する。代わりに  $\log p(W, Z|\alpha, \beta)$  を求める。

$\log \Gamma$  は C の `lgamma` で実装されている。Perl の場合は外部ライブラリを使う (e.g. `Math::Cephes` の `lgam`)。



準備として、単語列  $W = (w_1, w_2, \dots)$  と等価な 2 値変数の列  $B = (b_1, b_2, b_3, \dots)$  を考える。各変数は文字の間の分割状態に対応し、1 なら境界あり、0 なら境界なしを表す。例えば、単語列 (京都, 大学, だ) は  $(0, 1, 0, 1)$  に対応する (文頭、文末は固定なので無視)。

Unigram 単語分割モデルのサンプリングは  $B$  に対して行う。例えば、2 番目の値を flip して  $B = (0, 0, 0, 1)$  とすると、 $W = (\text{京都大学}, \text{だ})$  になる。

$b_i$  のサンプリングは、 $W$  の対応する局所領域を 2 単語 (京都, 大学) とするか、1 単語 (京都大学) とするかに対応する。以下ではサンプリング箇所を  $\cdot$  で示し、京都 $\cdot$ 大学 と表記する。

$b_i$  に対応する局所領域を  $a.b$  (e.g. 京都) としたとき、分割しない確率は以下のように、 $W^{-ab}$  を観測したもとでの単語  $ab$  の予測確率に比例する。

$$\begin{aligned}
 & P(b_i = 0 | \mathbf{B}^{-i}, \mathbf{C}, \alpha, G_0) \\
 = & P(ab | \mathbf{W}^{-ab}, \mathbf{C}, \alpha, G_0) \\
 \propto & P(ab, \mathbf{W}^{-ab}, \mathbf{C} | \alpha, G_0) \\
 = & \frac{\Gamma(\alpha)}{\Gamma(n_{(\cdot)}(\mathbf{W}^{-ab}) + 1 + \alpha)} \\
 & \times \prod_{l \in \text{voc}(\mathbf{W}^{-ab}) \cup \{ab\}} \frac{\Gamma(n_l(\mathbf{W}^{-ab}) + I(l = ab) + \alpha G_0(l))}{\Gamma(\alpha G_0(l))} \\
 \propto & \frac{\Gamma(n_{(\cdot)}(\mathbf{W}^{-ab}) + \alpha)}{\Gamma(n_{(\cdot)}(\mathbf{W}^{-ab}) + 1 + \alpha)} \times \frac{\Gamma(n_{ab}(\mathbf{W}^{-ab}) + 1 + \alpha G_0(ab))}{\Gamma(n_{ab}(\mathbf{W}^{-ab}) + \alpha G_0(ab))} \\
 = & \frac{n_{ab}(\mathbf{W}^{-ab}) + \alpha G_0(ab)}{n_{(\cdot)}(\mathbf{W}^{-ab}) + \alpha}
 \end{aligned}$$

分割する確率も同様だが、 $a = b$  のときに注意が必要。

$$\begin{aligned}
 & P(b_i = 1 | \mathbf{B}^{-i}, \mathbf{C}, \alpha, G_0) \\
 = & P(a, b | \mathbf{W}^{-ab}, \mathbf{C}, \alpha, G_0) \\
 \propto & \frac{n_a(\mathbf{W}^{-ab}) + \alpha G_0(a)}{n_{(\cdot)}(\mathbf{W}^{-ab}) + \alpha} \times \frac{n_b(\mathbf{W}^{-ab}) + I(a = b) + \alpha G_0(b)}{n_{(\cdot)}(\mathbf{W}^{-ab}) + 1 + \alpha}
 \end{aligned}$$

$\mathbf{W}^{-ab}$  を観測したもとの  $a$  を観測し、次に  $(\mathbf{W}^{-ab}, a)$  を観測したもとの  $b$  を観測する予測確率となっている (LDA の場合と同様、第 1 項の分母は無視できる)。

$$\frac{n_{ab}(W^{-ab})+\alpha G_0(ab)}{n_{(\cdot)}(W^{-ab})+\alpha} \text{ vs. } \frac{n_a(W^{-ab})+\alpha G_0(a)}{n_{(\cdot)}(W^{-ab})+\alpha} \times \frac{n_b(W^{-ab})+I(a=b)+\alpha G_0(b)}{n_{(\cdot)}(W^{-ab})+1+\alpha}$$

- 分母に注目すると、1以下の数を2つかける(分割する)より、1つだけの(分割しない)方が大きくなる。極端な話、1文を1単語にすれば掛け算せずですんで大きな値になる。
- 分子の  $n_x(W^{-ab})$  は、 $x$  が  $W^{-ab}$  の中で出てきた回数だから、これを大きくしようと思ったら、なるべく区切った方がよい。極端な話、1文字1単語にすれば、再利用されまくる。
- Unigram 単語分割モデルは両者のバランスをとっている。

以下のデータを管理する。

- 現在の単語列  $W$
- $n_l$  (各単語の頻度)
- $n_{(\cdot)}$  (全単語数)

$z_i$  のサンプリングは以下の手順で行う。

- ① 現在の  $b_i$  に対応するカウントを decrement ( $b_i = 0$  のときは 1 単語、 $b_i = 1$  のときは 2 単語)
- ②  $P(b_i | \mathbf{B}^{-i}, C, \alpha, G_0)$  を 2 通り計算
- ③ 乱数を振って新たな  $b_i$  を決め、 $W$  を更新
- ④ 新たな  $b_i$  に対応するカウントを increment

普通の Gibbs サンプリングは局所的な探索を繰り返すので、山と山の間には深い谷があるとなかなか移れない。

例えば、Unigram 単語分割モデルにおいて、 $W$  中で a.b (e.g. 京都) が 100 回出現し、すべてマージすると  $p(W|C, \alpha, G_0)$  が最大値になるとする。いま、a.b が 97 回分割されているとする。(a, b) をサンプリングすると、97 - 1 回の出現の影響により  $P(b_i = 1 | \dots) \approx 1$  となり、ほぼ必ず分割される。したがって、すべてを分割する状態に収束し、すべてマージする状態に移る可能性は 0 に近い。つまり極所解にとらわれてしまう。

Type-based サンプリング (Liang et al., 2010) では、この 100 個の出現をまとめてサンプリングすることにより、谷を一気に飛び越える。

$a.b$  をサンプルングするとき、単語列  $a, b$  および  $ab$  は同じ型 (type) に属す。同じ型に属すものの集合を  $S$  とする。Type-based サンプルングでは  $S$  をまとめてサンプルングする。

$S$  に対応する 2 値変数の集合を  $b_S$  とする。 $b_S$  に対応する  $|S|$  個の局所領域を  $W$  から除いたものを  $W^{-S}$  とする。

Type-based サンプルングのキモは、 $b_S$  に属す各  $b_i$  が交換可能なことである。久しぶりに (サンプルではなく) 観測数の分布を考え、以下の 2 段階でサンプルングを行う。

- ① 最初に観測数の分布から分割数  $k$  ( $0 \leq k \leq |S|$ ) をサンプル
- ② 次に  $|S|$  個の状態に対して  $k$  個の分割をランダムに割り当て

分割数  $k$  の分布は以下の通り (ただし  $a \neq b$ )。

$$\begin{aligned}
 & P(k|\mathbf{B}^{-S}, \mathbf{C}, \alpha, G_0) \\
 = & P(k|\mathbf{W}^{-S}, \mathbf{C}, \alpha, G_0) \\
 \propto & \binom{|\mathbf{S}|}{k} \frac{\Gamma(n_{(\cdot)}(\mathbf{W}^{-S}) + \alpha)}{\Gamma(n_{(\cdot)}(\mathbf{W}^{-S}) + |\mathbf{S}| + k + \alpha)} \times \frac{\Gamma(n_a(\mathbf{W}^{-S}) + k + \alpha G_0(a))}{\Gamma(n_a(\mathbf{W}^{-S}) + \alpha G_0(a))} \\
 & \times \frac{\Gamma(n_b(\mathbf{W}^{-S}) + k + \alpha G_0(b))}{\Gamma(n_b(\mathbf{W}^{-S}) + \alpha G_0(b))} \times \frac{\Gamma(n_{ab}(\mathbf{W}^{-S}) + |\mathbf{S}| - k + \alpha G_0(ab))}{\Gamma(n_{ab}(\mathbf{W}^{-S}) + \alpha G_0(ab))}
 \end{aligned}$$

$a = b$  のときは、第3,4項

$$\frac{\Gamma(n_a(\mathbf{W}^{-S}) + k + \alpha G_0(a))}{\Gamma(n_a(\mathbf{W}^{-S}) + \alpha G_0(a))} \times \frac{\Gamma(n_b(\mathbf{W}^{-S}) + k + \alpha G_0(b))}{\Gamma(n_b(\mathbf{W}^{-S}) + \alpha G_0(b))}$$

を以下で置き換える

$$\frac{\Gamma(n_a(\mathbf{W}^{-S}) + 2k + \alpha G_0(a))}{\Gamma(n_a(\mathbf{W}^{-S}) + \alpha G_0(a))}$$



分割数の分布からのサンプルングにより分割数  $k$  が求まる。次に  $S$  に対して  $k$  を割り振る。どのような分割の割り当ても等確率なので、ランダムに割り当てる。

特殊な場合には注意が必要である。例えば  $a_1, a_2, a_3$  ( $a_1 = a_2 = a_3$ ) について、 $a_1.a_2$  をサンプリングすることを考える。このとき、 $a_2.a_3$  は  $a_1.a_2$  と同じ型に属す。しかし、 $a_2$  を共有しているため、 $a_1.a_2$  と  $a_2.a_3$  は交換可能ではない。仕方がないので、 $a_1.a_2$  をサンプリングするときには、 $a_2.a_3$  はサンプリング対象から外す。

以下のデータを管理する。

- 現在の単語列  $W$
- $n_l$  (各単語の頻度)
- $n_{(\cdot)}$  (全単語数)
- **型を管理するデータ構造**

サンプリングは以下の手順で行う。

- ①  $b_i$  について、同じ型で、衝突のないものの集合  $S$  を集める
- ②  $S$  に対応するカウントを decrement
- ③  $P(k|\mathbf{B}^{-S}, C, \alpha, G_0)$  を  $|S| + 1$  通り計算
- ④ 乱数を振って  $k$  を決め、分割をランダムに割り振る
- ⑤ 割り当てられた分割に応じて  $W$  を更新
- ⑥ 割り当てられた分割に応じてカウントを increment