

# ガウス過程に基づく連続空間トピックモデル

持橋大地

統計数理研究所 数理・推論研究系 学習推論グループ

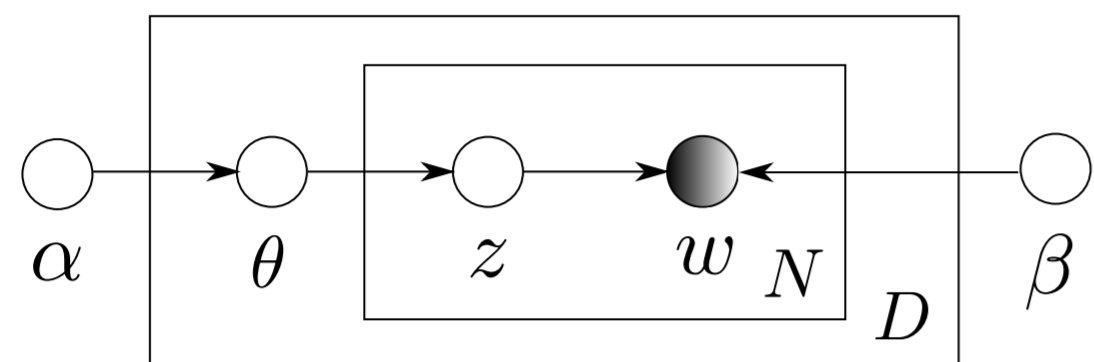
daichi@ism.ac.jp

## 1. 背景: トピックモデル

トピックモデル...文書(離散データ)の確率的生成モデル  
代表的なもの: Latent Dirichlet Allocation (Blei+ 2003)

生成プロセス: 各文書ごとに,

1. 潜在トピック分布  $\theta \sim \text{Dir}(\alpha)$  を生成
2. 各語ごとに,
  - (a) 潜在トピック  $z \sim \text{Mult}(\theta)$  を選択.
  - (b) 単語  $w \sim p(w|z) = \beta_{zw}$  を生成.



- $\theta$ : 多項分布
- $\beta$ : 潜在トピック別単語確率分布
- $w$ : 観測値(単語)

応用と拡張

- 非常に多くの分野で使われている  
- グラフィックス, ビジョン, バイオインフォマティクス, 推薦モデル, リンク解析, ネットワークモデル, ...
- 多数の改良の試み (sLDA, CTM, DHDP, Mixed Random Measures (Kim+ 2012), dHNRM (Chen+ 2012), ...)

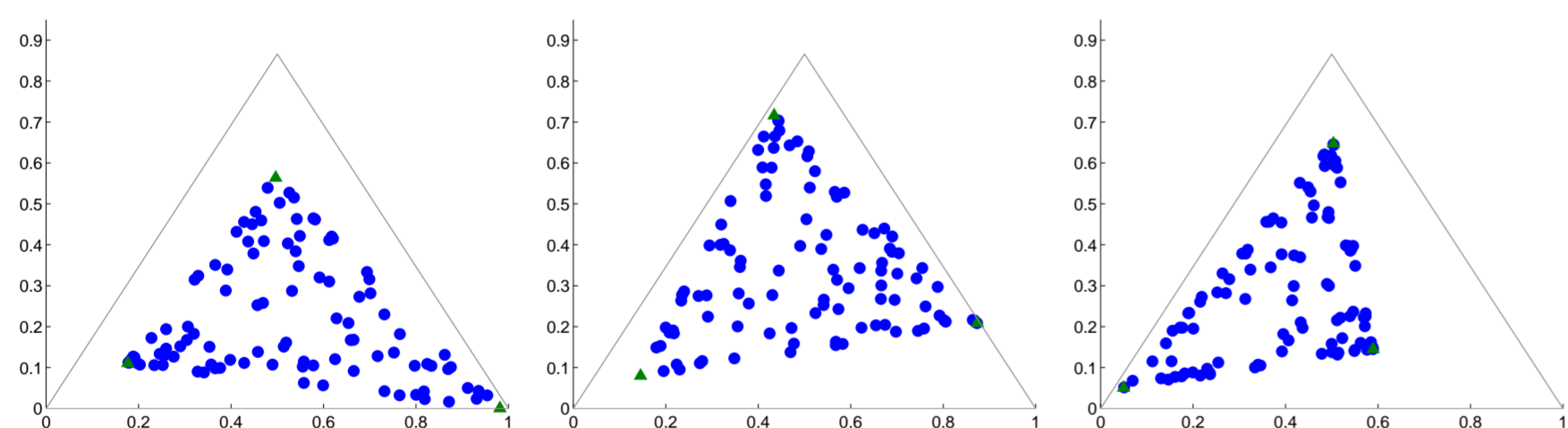
## 2. 問題点

- 混合モデルのため,  $\beta = p(w|z)$  は一度決まると固定  
- 語彙の情報を細かく制御できない = ラベルがトピックにしか影響しない (sLDA, dHNRM 等の拡張全て)

実際の例

- 女性特有の言い回し, 語彙選択
- 会話文の文体 (トピック / 内容とは独立)
- ひらがな / 漢字ばかりの文章
- 著者の語彙的な癖 (-eous や re- が多い, など)

- 単語 Simplex の全領域をモデル化できない



↓  
単語 Simplex の全領域をモデル化し, 語彙に柔軟に対応する統計モデル.

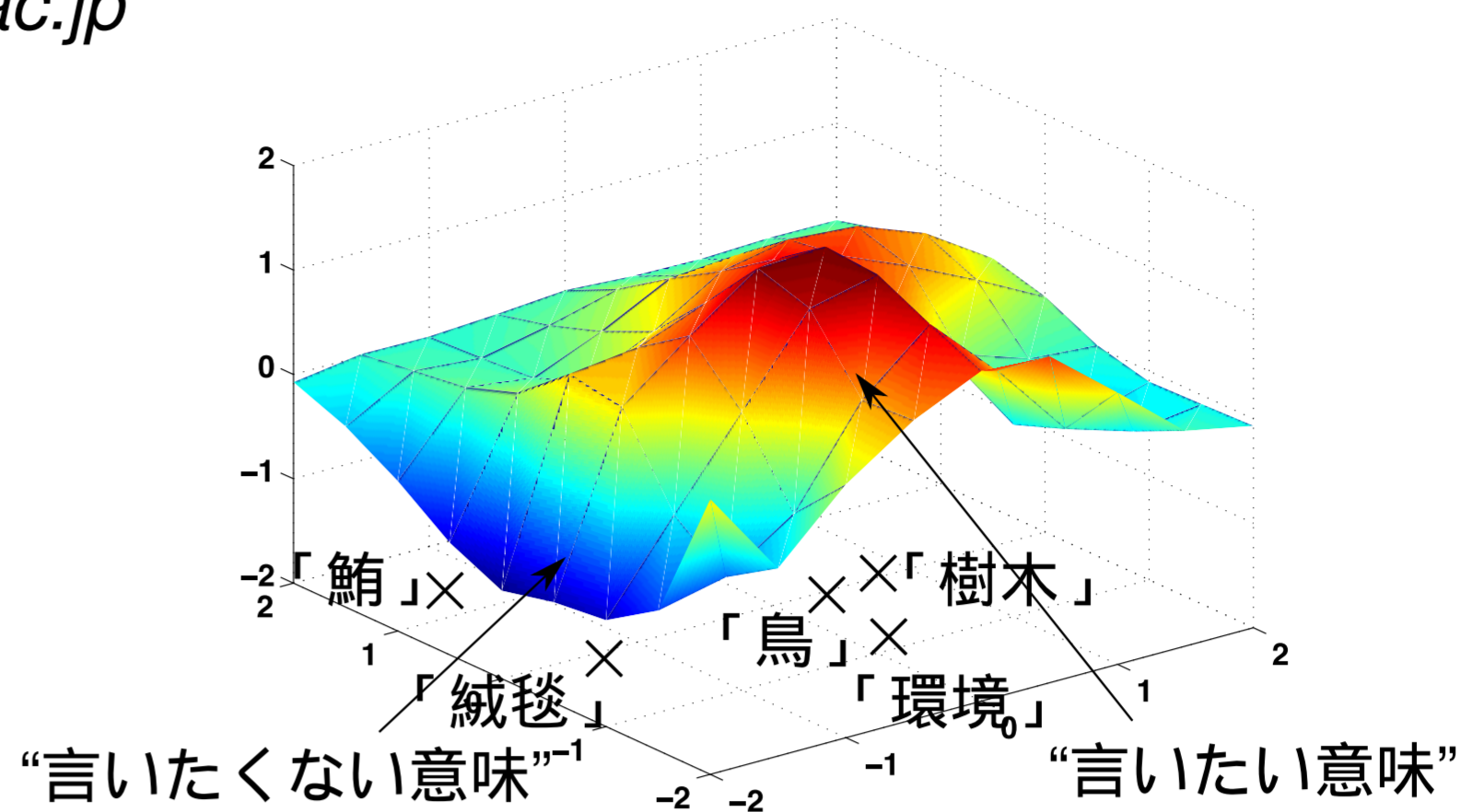
## 3. CSTM: Continuous Space Topic Model

基本的なアイデア: 潜在単語空間上のガウス過程  $f \sim \text{GP}(0, K)$  によって, 「この文書で言いたいこと」を表現: 単語確率を, 下の積で Modify.

$$p(w|d) = e^f G_0(w)$$

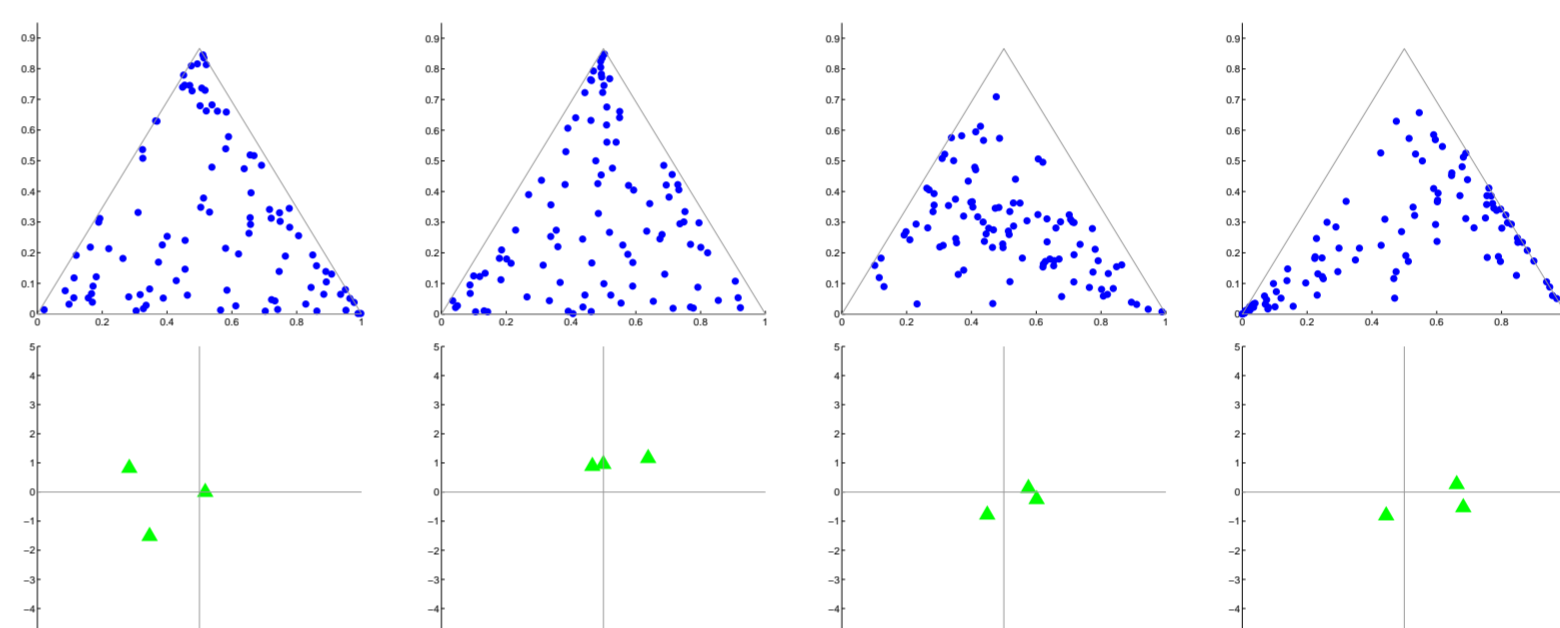
生成モデル

1. For  $w = 1 \dots W$ ,  $\phi(w) \sim N(0, I_d)$ . ( $w$  の潜在座標)
2. for  $d = 1 \dots D$ 
  - $f \sim \text{GP}(0, K)$  ;  $K(v, w) = \phi(v)^T \phi(w)$
  - $\alpha(w|d) = \alpha_0 G_0(w) \exp(f)$
  - $G \sim \text{DP}(\alpha(\cdot|d))$ ,  $\mathbf{w} \stackrel{\text{iid}}{\sim} G$ . (Polya 分布)



特徴:

- 語彙の確率を直接制御できる  
- 拡張: 文書ラベル  $\alpha(w|d) = \alpha_0 G_0(w) e^f e^{\theta^T c(d)}$   
- 拡張: 単語の特徴  $\alpha(w|d) = \alpha_0 G_0(w) e^f e^{\eta^T c(w)}$
- Product model ( $\leftrightarrow$  Mixture model)
- 単語 Simplex のほぼ全域をモデル化できる



## 4. 計算と推論

GP  $f$  は直接表現が難しい 補助変数  $u \sim N(0, I_d)$  のとき,  $f = \Phi u$  の分布は  $u$  を消去して

$$f|\Phi \sim N(0, \Phi^T \Phi) = N(0, K).$$

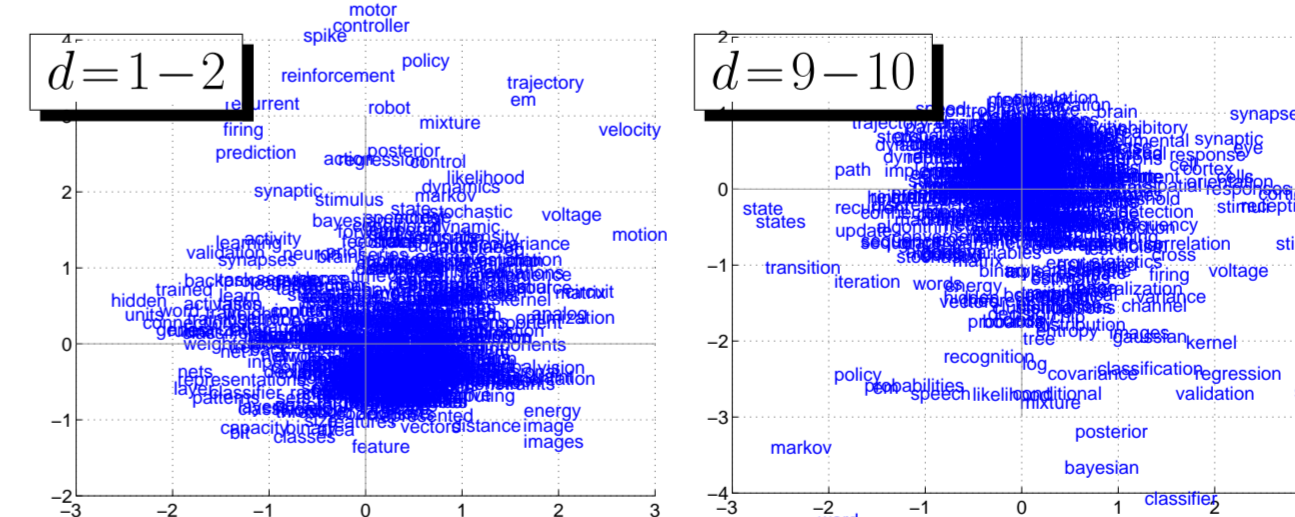
よって, 文書  $d$  毎に  $f = \Phi u_d$  として  $u_d$  と  $\Phi$  を MCMC で更新.  $\alpha_0 \sim \text{Ga}(a_0, b_0)$ ,  $\eta_{kv} \sim N(0, I_d)$ ,  $\theta_{cv} \sim \text{Exp}(\beta)$ .

## 5. 実験結果

Perplexity

| データ  | CSTM    | CSTM(+lex) | LDA      |
|------|---------|------------|----------|
| NIPS | 1410.96 | -          | 1648.3   |
| KOS  | 1632.35 | -          | 1730.7   |
| 毎日新聞 | 473.954 | 473.273    | 507.3855 |

Latent Embedding (NIPS,  $d=10$ )



Covariate modeling

| 文書   | $e^\eta$ | テキスト 2364:   | $\eta$ (ひらがな) | テキスト 4580: |  |
|------|----------|--|---------------|------------|--|
| 2364 | 1.498    | #日公開の映画ではウォンカーウエイ監督の花様年華がよねんかがかんヌ国際映画祭最優秀男優トニレオン高等技術院賞受賞のかくかくたる戦果をあげての香港凱旋がいせんだあまりにも古風な映画でカーウエイ監督ファンはびっくりするかも日本映画は連弾がはじけるおもしろさ人間の肩もけこう見せるデニエロのくせ者ぶりが楽しめるミートザベアレンツ小粒でも... | 4580          | 1.720      | 文 小森香折 こもりかおり 絵 広瀬 悠 ひろせげん ゆうが押し入れをかけたげいているとこらへどうぞという父とうさんの声こえがきこえてきました押し入れのむこうほうはならないの部屋へやですゆうは押し入れにもぐりこんで耳みみをあててました女のお客さやくさんとあいつているのがきこえてきますあのおうへがみさまはこれがおすきだとうかがいまして... |
| 4597 | 1.471    |  | 9961          | 1.501      |  |
| 442  | 1.440    |  | 5238          | 1.494      |  |
| 4633 | 1.433    |  | 7420          | 1.470      |  |
| 1520 | 1.422    |  | 8375          | 1.452      |  |

## Future Work

Random Poisson intensity として, HDP/HPY の拡張  $\rightarrow$  PCFG,  $n$  グラムモデル等の潜在空間埋め込み拡張