

トピックモデルの応用： イントロダクション

NTT コミュニケーション科学基礎研究所

石黒 勝彦

2013/01/15-16 統計数理研究所 会議室1

このスライドの“トピック”

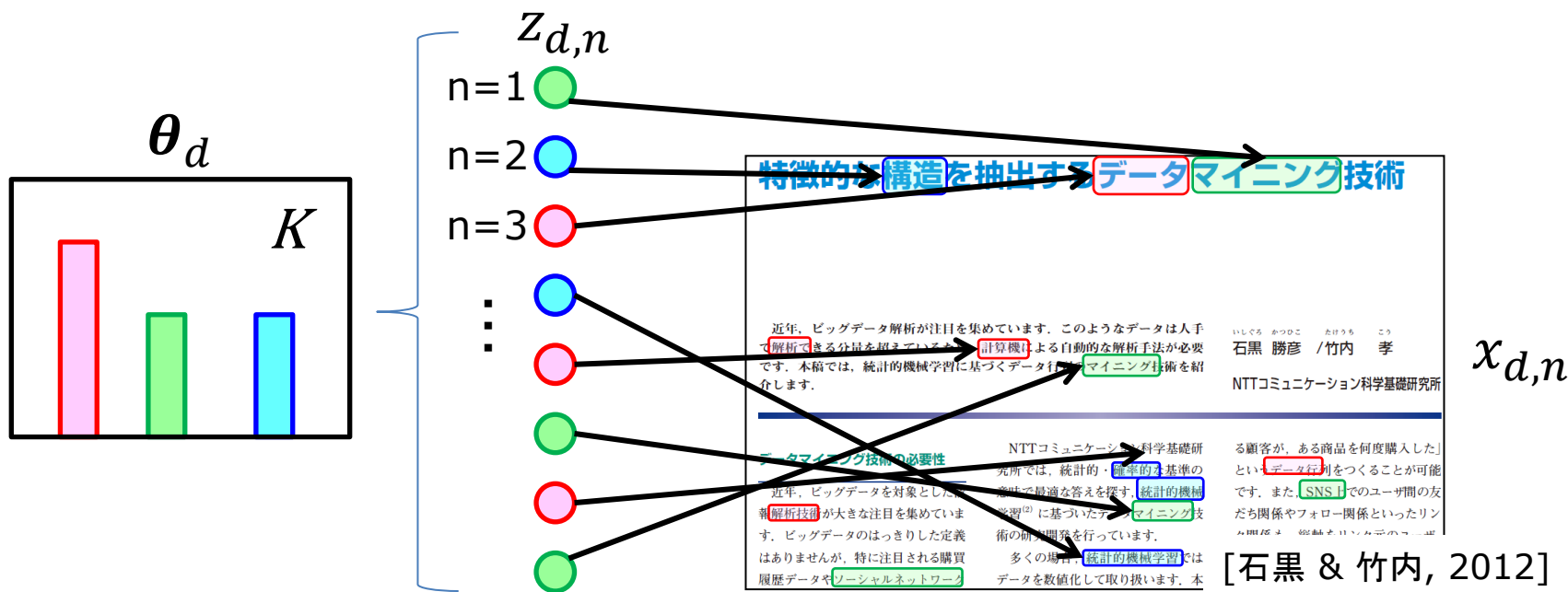
- まずは本日の講義のイントロダクションです
- LDAの復習
 - 使用する変数, モデルなど
 - 解法: Gibbs sampler, VB-EM
- 今日の講義全体に関する注意
 - notation policyなど
 - 参考文献

本日の予定

- LDAの拡張手法と応用について紹介します
 - 基本的には、様々な論文の説明になります
- 午前：トピックモデル(LDA)の拡張モデル
- 午後：トピックモデルの各種ドメインデータへの応用

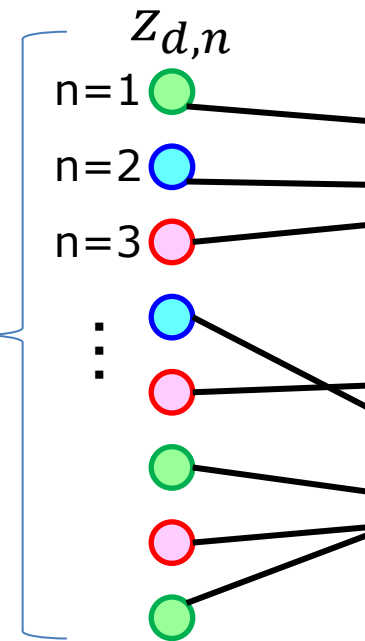
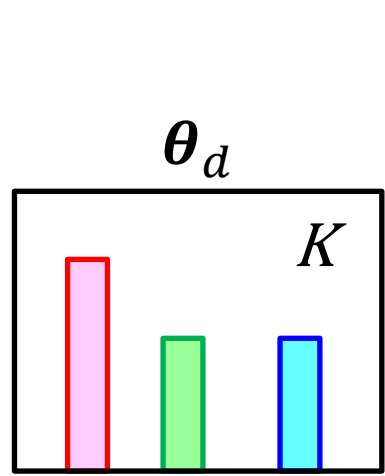
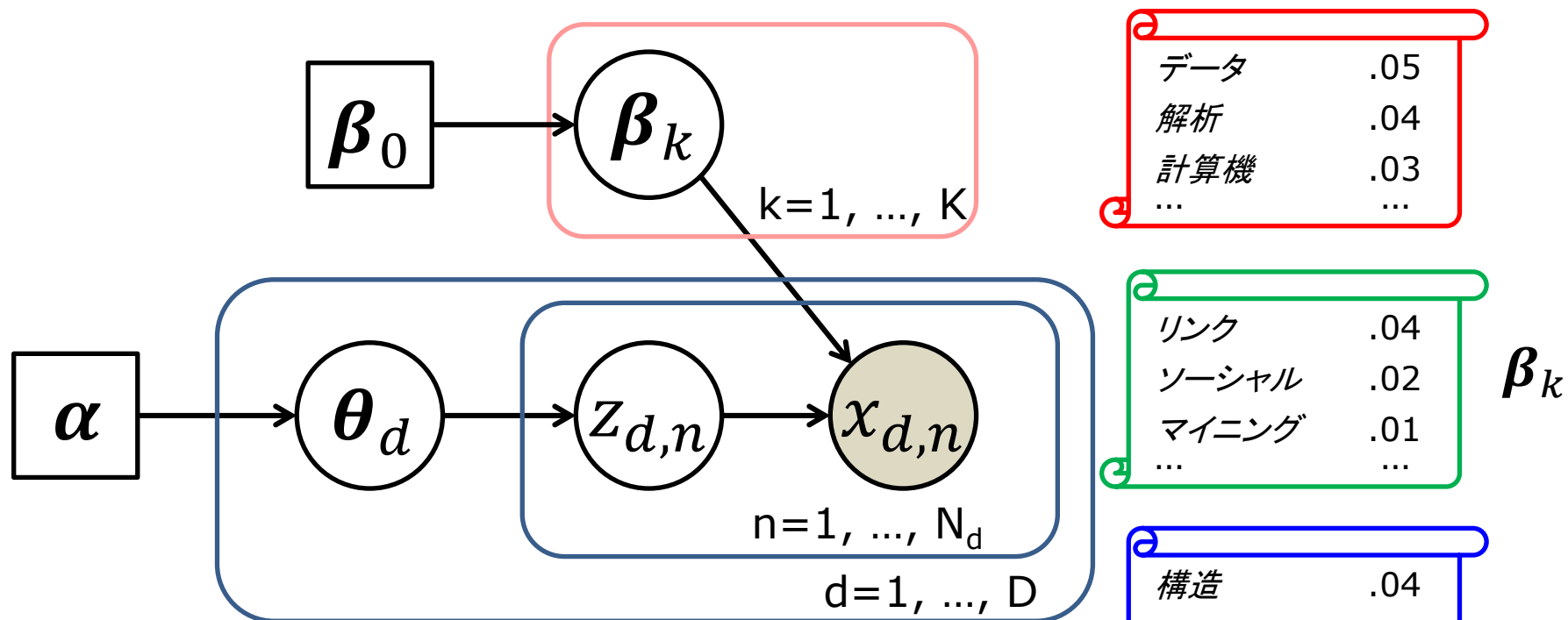
Latent Dirichlet Allocation [Blei, 2003]

- 様々な離散データに隠された潜在的なトピックを推定するベイジアンモデル



インデックス、定数、変数

- 文書インデックス d
- トピックインデックス k
- 単語インデックス n
- 文章数 D
- トピック数 K
- 文書 d 中の単語数 N_d
- 単語の種類 V
- 観測された単語 x
- 単語のtopic assignment z
- 文書のtopic proportion θ
- トピックのword proportion β
- θ の事前分布パラメータ α
- β の事前分布パラメータ β_0



特徴的な「構造」を抽出する「データマイニング」技術

近年、ビッグデータ解析が注目を集めています。このようなデータは人手で解析できる分量を超えています。計算機による自動的な解析手法が必要です。本稿では、統計的機械学習に基づくデータマイニング技術を紹介いたします。

NTTコミュニケーション科学基礎研究所

石黒 勝彦 / 竹内 孝

近年、ビッグデータを対象としたデータマイニング技術の重要性が注目を集めています。ビッグデータのはっきりした定義はありませんが、特に注目される購買履歴データをソーシャルネットワーク

NTTコミュニケーション科学基礎研究所では、統計的・確率的基準のデータ解析技術に基づいたデータマイニング技術の研究開発を行っています。多くの場合、統計的機械学習ではデータを数値化して取り扱います。本

顧客が、ある商品を何度購入した」とい「データ」列をつくることが可能です。また「SNS」でのユーザー間の友だち関係やフォロー関係といったリンク関係も、総称として「ソーシャルネットワーク

$x_{d,n}$

[石黒 & 竹内, 2012]

生成モデル

for 文書 $d = 1, 2, \dots, D_t$

topic proportion

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

for 単語 $n = 1, 2, \dots, N_d$

topic-word assignment

$$z_{d,n} | \boldsymbol{\theta}_d \sim \text{Mult}(\boldsymbol{\theta}_d)$$

word observation

$$x_{d,n} | z_{d,n}, \{\boldsymbol{\beta}_k\} \sim \text{Mult}(\boldsymbol{\beta}_{z_{d,n}})$$

for トピック $k = 1, 2, \dots, K$

topic-word proportion

$$\boldsymbol{\beta}_k | \boldsymbol{\beta}_0 \sim \text{Dir}(\boldsymbol{\beta}_0)$$

topic-word proportionは事前分布を
仮定しない場合も多々あります

Gibbs sampler

- 最も正確な解を得ることができる解法です
- 各文書 d の単語 n を一つずつトピック k に割り当てていきます

$$p(z_{d,n} = k | x_{d,n} = w, \mathbf{X}_{\neg(d,n)}, \mathbf{Z}_{\neg(d,n)}, \boldsymbol{\alpha}, \boldsymbol{\beta}_0) \propto \frac{m_{dk} + \alpha_k}{\sum_{k'} (m_{dk'} + \alpha_{k'})} \frac{m_{kw} + \beta_{0,w}}{\sum_{w'} (m_{kw'} + \beta_{0,w'})}$$

文書 d から
トピック k が
生成される確率

トピック k から
単語 w が
生成される確率

m_{dk} : 文書 d 内でトピック k にアサインされた観測量のカウント ($x_{d,n}$ を除く)

m_{kw} : 文書全体でトピック k にアサインされた観測量のうち
単語 w だった観測量のカウント ($x_{d,n}$ を除く)

変分ベイズ法(VB-EM)

- 文書データの周辺化尤度をJensen不等式で下から押さえて、それを最大化する変分事後分布 $q()$ を求める手法です

$$\begin{aligned}\log p(\mathbf{X}) &= \iint \log p(\mathbf{X}, \mathbf{Z}, \boldsymbol{\theta}) d\mathbf{Z}d\boldsymbol{\theta} \\ &\geq \iint q(\mathbf{Z}, \boldsymbol{\theta}|\mathbf{X}) \log \frac{p(\mathbf{X}, \mathbf{Z}, \boldsymbol{\theta})}{q(\mathbf{Z}, \boldsymbol{\theta}|\mathbf{X})} d\mathbf{Z}d\boldsymbol{\theta}\end{aligned}$$

$$q(z_{d,n} = k) \propto \beta_{k,x_{d,n}} \exp\left(\Psi(\hat{\theta}_{d,k})\right)$$

$$\hat{\theta}_{d,k} = \alpha_k + \sum_n q(z_{d,n} = k)$$

Notationについて

- 直観的な理解のために、できる限り同じ意味を持つ量は同じ名前の変数・インデックスで表します
- また、一部のモデルについてはよりわかりやすく等価なモデルで説明します
- したがって、原論文とは変数の名前やグラフィカルモデルの形が違う可能性があります

全体を通じての重要な参考文献

- Blei et al, “Latent Dirichlet Allocation”,
Journal of Machine Learning
Research, Vol. 3, pp. 993-1022,
2003.
 - LDAの原論文です

全体を通じての重要な参考文献

- ビショップ、"パターン認識と機械学習", 丸善出版, 2012
 - "PRML本": 和書の中では、現状、機械学習に関するもっとも良い本の一つです



引用及び参考文献

- [Blei, 2003] Blei et al, “Latent Dirichlet Allocation”, Journal of Machine Learning Research, Vol. 3, pp. 993-1022, 2003.
- [PRML] Blei and Lafferty, “A Correlated Topic Model of Science”, The Annals of Applied Statistics, Vol. 1(1), pp. 17-35, 2007.
- [石黒 & 竹内, 2012] 石黒, 竹内, “特徴的な構造を抽出するデータマイニング技術”, NTT技術ジャーナル, Vol. 24, No. 9, 2012.