

# トピックモデルの応用： 時系列データ

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”

- 購買データや科学論文など、時間変化をそもそも内包するデータは多数存在します
- 従って、時系列(時間変化)データ内のトピックの解析も多数試みが行なわれています

# 時系列データは数多く存在します

動画像



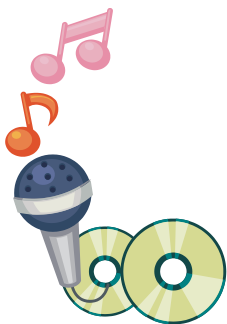
脳波・生体信号



購買履歴・市場インデックス



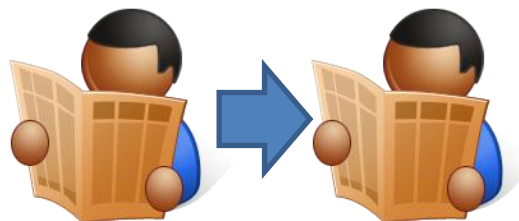
音楽・音響信号



新聞・ニュース記事

01/15

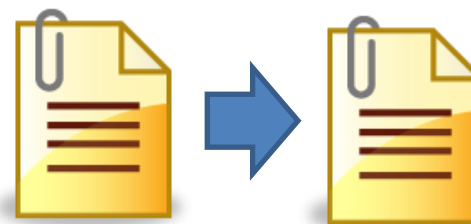
01/16



科学論文・特許

2012

2013



😊 トピックモデルの有用性が明白になったため、これら時系列データへトピックモデルを応用した研究が多数発表されています

# 時系列データモデリングのキモ： どこに時間依存性を入れるか？

- マルコフ性：前の時刻に依存して現在の時刻の状態が変化する
- 多くの時系列データでは、モデルのどの部分にマルコフ性のアイデアを導入するか、がポイントとなります
- これはトピックモデルの時系列データモデルでも同様です

# **Dynamic Topic Model**

**[Blei & Lafferty, 2006]**

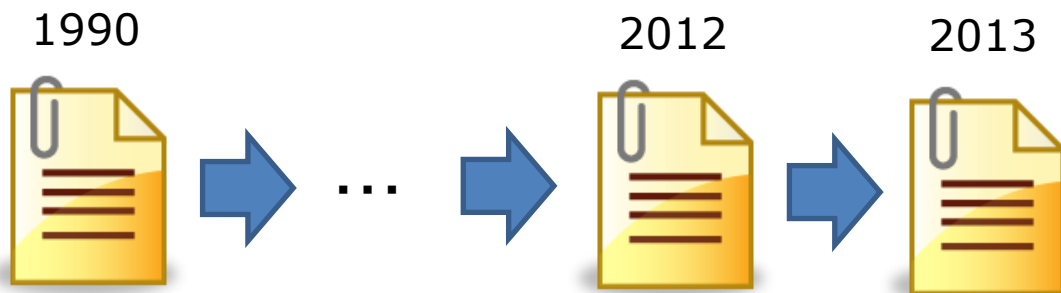
Blei and Lafferty,  
“Dynamic Topic Models”,  
in Proc. ICML, 2006.

# トピックモデルの大前提の仮定: **exchangeability**

- 簡単にいうと: 「各文書  $d$ , 各単語  $w$  のインデックスはただのシンボルで順番や名前には意味が無い」
- これのおかげで各種モデル推論が簡単になっています

# 文書コレクションの時系列データを考えると、これは問題です

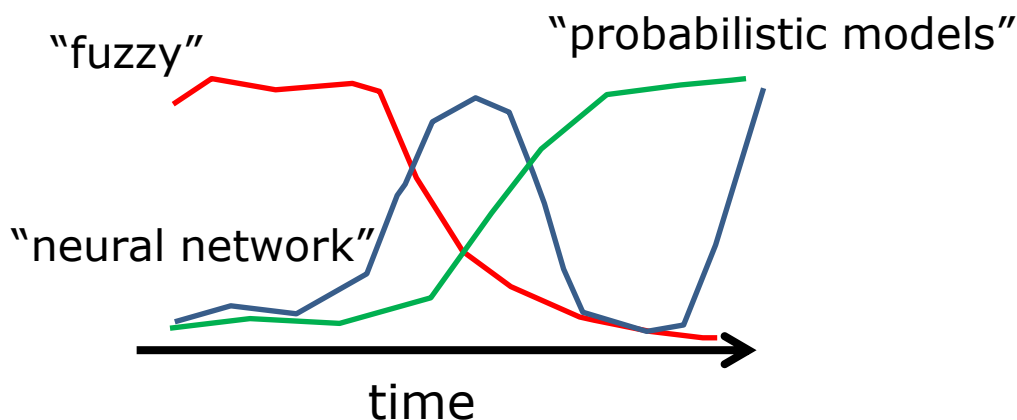
- 新聞記事は昨日までの報道の流れを汲んでいます
- 論文は先行研究の作った技術トレンドにのっています
- すなわち、文書  $d$  は一般には exchangeable ではありません！



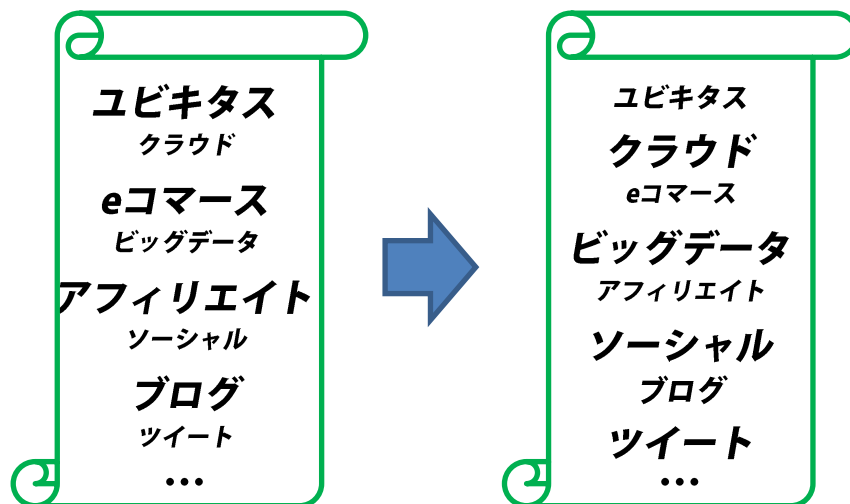
# 注目する時間依存性：トピック

- 1: 話題(topic)には流行り廃りがあります
- 2: トピックの中での言葉づかいも変化します
- これら2種類の「トピックの変化」を解析するモデルを考えたい

トピックの流行り廃り



トピックの中での言葉づかい



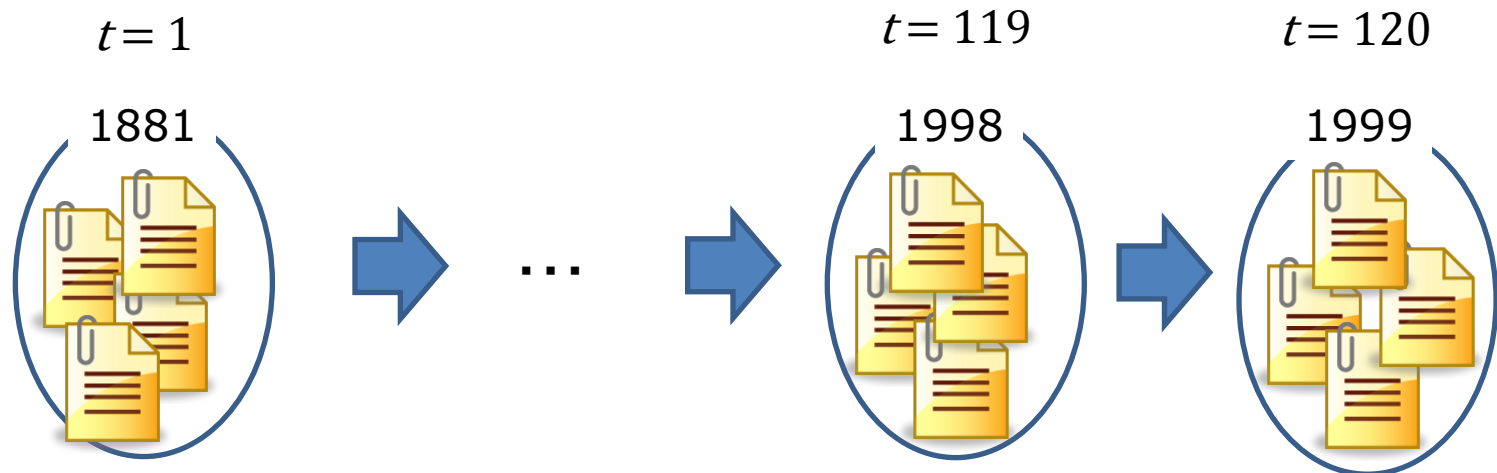


# 提案法: Dynamic Topic Models

- 😊 非常に有名な時系列トピックモデルです
- 科学誌ScienceのOCRデータを用いて、科学論文の時系列トピック解析を行います
- topic proportionとtopic-word proportionに時間マルコフ性を入れたものです
- 推論は非常に難しいです

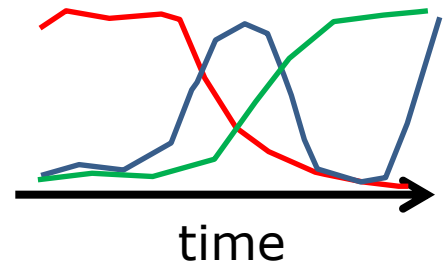
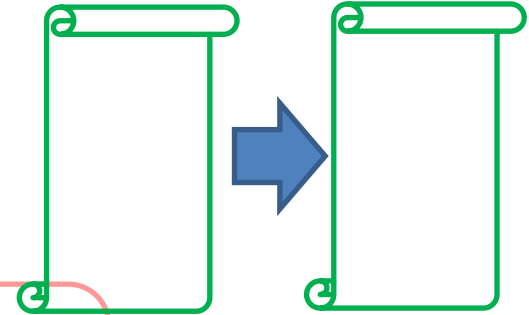
# 対象データ: Science誌

- 1880年にエジソンによって刊行された、非常に著名な科学論文誌
- OCRされた論文誌データ(JSTOR)を利用して、発行年度ごとの文書時系列データを作成



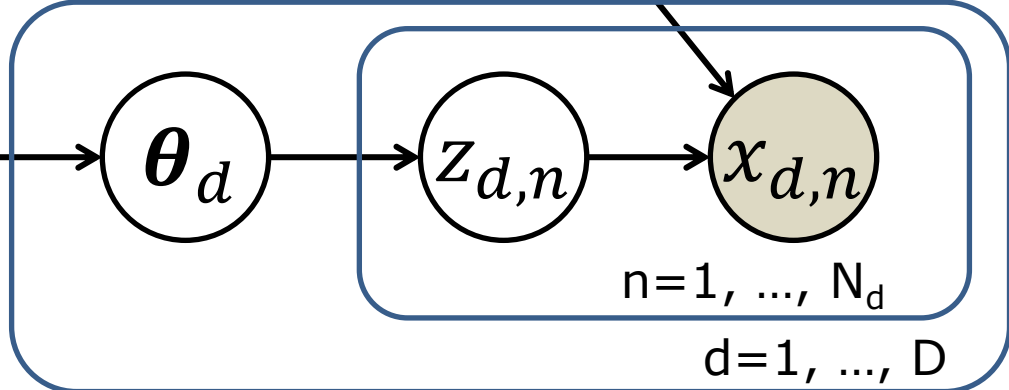
# 提案法のアイデア

- 以下の2点を時間発展させます
- $\alpha$ : トピックの流行り廃りを制御
- $\beta_k$ : トピックごとの単語分布



トピックの流行り廃り

$\alpha$

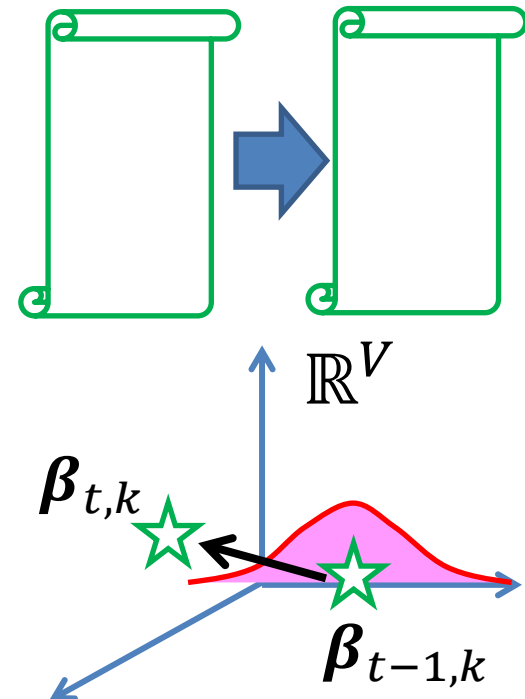
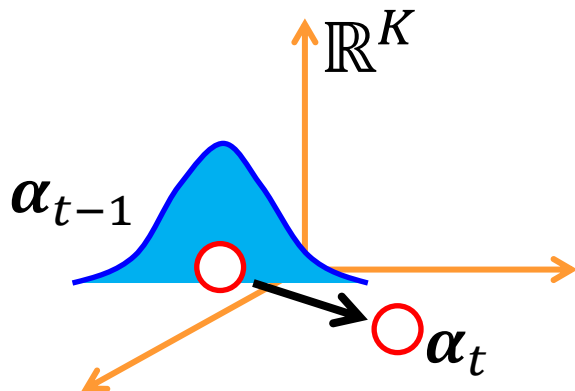
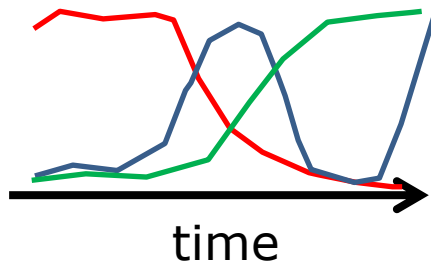


$\beta_k$   
 $k=1, \dots, K$

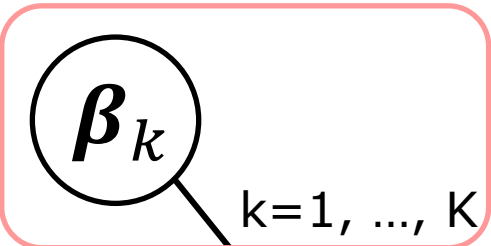
トピックの中での言葉づかい

# 時間発展のモデル： 正規ノイズでdrift

- 最も単純な時間発展モデルと言えます
- パラメータは前の時刻を中心に少しずつしか動かない、という想定です



LDA

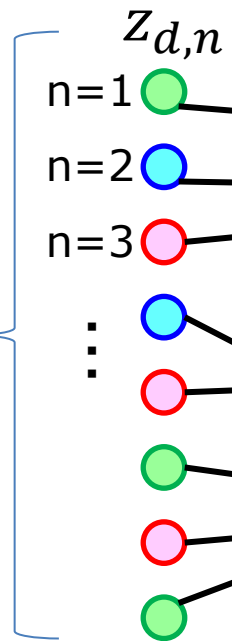
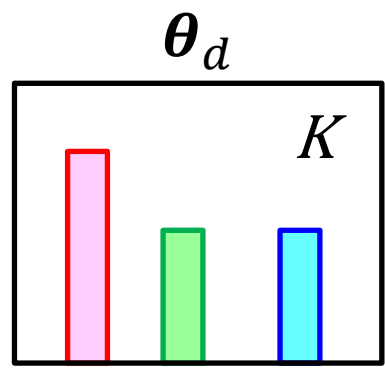
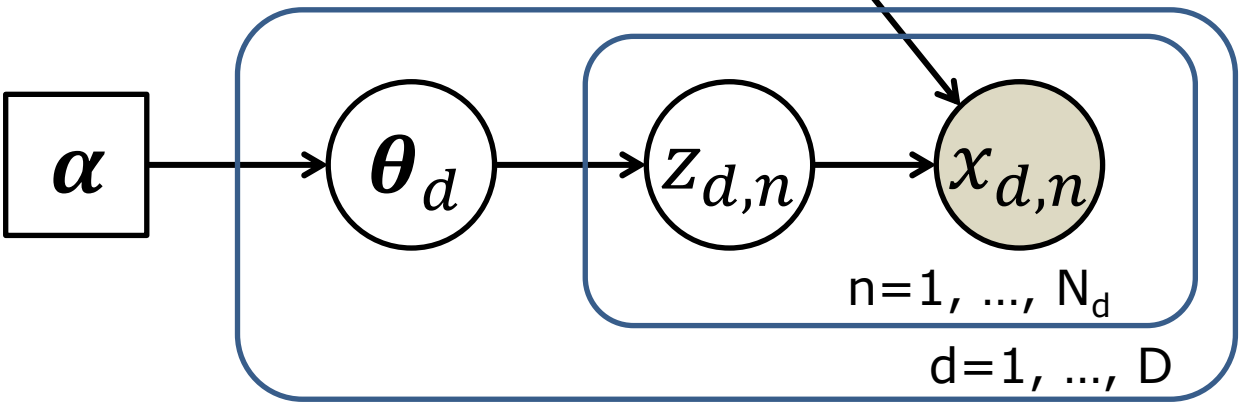


データ	.05
解析	.04
計算機	.03
...	...

リンク	.04
ソーシャル	.02
マイニング	.01
...	...

構造	.04
機械学習	.03
最適	.01
...	...

$\beta_k$



**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

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、ソーシャルネットワークの履歴データを「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

$x_{d,n}$



# 生成モデル

for 時間  $t = 1, 2, \dots, T$

for theme (topic)  $k = 1, 2, \dots, K$

topic-word proportion drift

$$\boldsymbol{\beta}_{t,k} | \boldsymbol{\beta}_{t-1,k} \sim N(\boldsymbol{\beta}_{t-1,k}, \sigma^2 \mathbf{I})$$

topic proportion parameter drift

$$\boldsymbol{\alpha}_t | \boldsymbol{\alpha}_{t-1} \sim N(\boldsymbol{\alpha}_{t-1}, \delta^2 \mathbf{I})$$

for 文書  $d = 1, 2, \dots, D_t$

topic proportion

for 単語  $n = 1, 2, \dots, N_{t,d}$

topic-word assignment

word observation

for 時間  $t = 1, 2, \dots, T$

$$\boldsymbol{\alpha}_t | \boldsymbol{\alpha}_{t-1} \sim N(\boldsymbol{\alpha}_{t-1}, \delta^2 \mathbf{I}) \quad \boldsymbol{\beta}_{t,k} | \boldsymbol{\beta}_{t-1,k} \sim N(\boldsymbol{\beta}_{t-1,k}, \sigma^2 \mathbf{I})$$

for 文書  $d = 1, 2, \dots, D_t$

topic proportion  $\boldsymbol{\eta}_{t,d} | \boldsymbol{\alpha}_t \sim N(\boldsymbol{\alpha}_t, a^2 \mathbf{I})$

$$\boldsymbol{\theta}_{t,d} | \boldsymbol{\eta}_{t,d} = \pi(\boldsymbol{\eta}_{t,d})$$

for 単語  $n = 1, 2, \dots, N_d$

topic-word assignment

$$z_{t,d,n} | \boldsymbol{\theta}_{t,d} \sim \text{Multinomial}(\boldsymbol{\theta}_{t,d})$$

word observation

$$x_{d,n} | z_{d,n}, \{\boldsymbol{\beta}_{t,k}\} \sim \text{Multinomial}(\pi(\boldsymbol{\beta}_{t,z_{d,n}}))$$

$\pi$  は soft-max 関数  $\pi(\mathbf{v}) = \frac{\exp(v_k)}{\sum_l \exp(v_l)}$

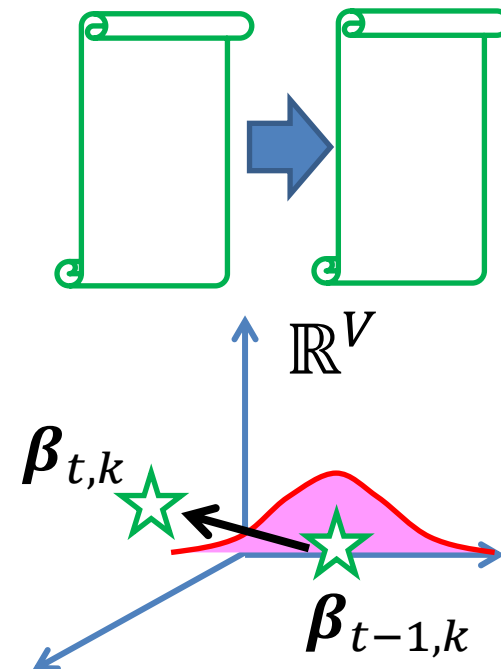
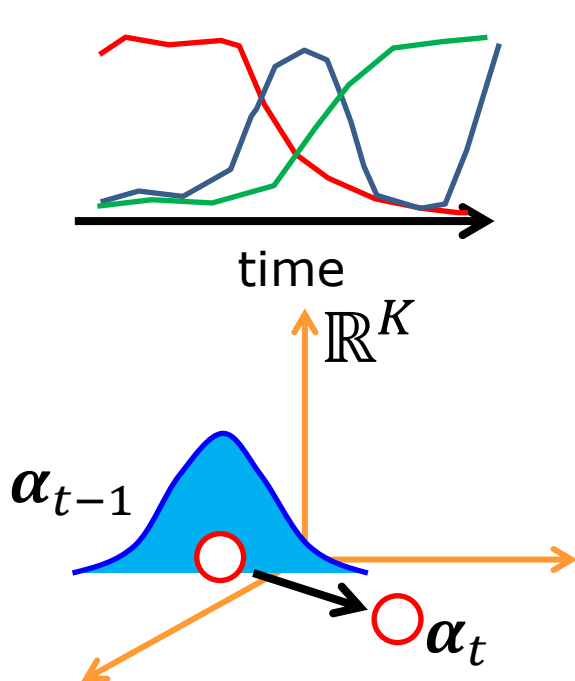


# 時間発展のモデル

- 正規分布を使って、1時刻のパラメータ遷移 (drift) をモデル化します

$$\alpha_t | \alpha_{t-1} \sim N(\alpha_{t-1}, \delta^2 I)$$

$$\beta_{t,k} | \beta_{t-1,k} \sim N(\beta_{t-1,k}, \sigma^2 I)$$



# トピックモデルへの適合

- 正規分布からは実数ベクトルが生成されるため、そのままでは多項分布(Multinomial)に使えません
- Soft-max関数を利用して変換します

$$\boldsymbol{\eta}_{t,d} | \boldsymbol{\alpha}_t \sim N(\boldsymbol{\alpha}_t, a^2 \mathbf{I}) \quad \text{時刻}t, \text{文書}d\text{のtopic proportion}$$

$$\boldsymbol{\theta}_{t,d} | \boldsymbol{\eta}_{t,d} = \pi(\boldsymbol{\eta}_{t,d}) \quad \text{Soft-max}$$

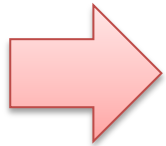
$$z_{t,d,n} | \boldsymbol{\theta}_{t,d} \sim \text{Multinomial}(\boldsymbol{\theta}_{t,d}) \quad \text{topic-word assign.}$$

$$\pi(\boldsymbol{v}) = \frac{\exp(v_k)}{\sum_l \exp(v_l)}$$

$$x_{d,n} | z_{d,n}, \{\boldsymbol{\beta}_{t,k}\} \sim \text{Multinomial}(\pi(\boldsymbol{\beta}_{t,z_{d,n}}))$$

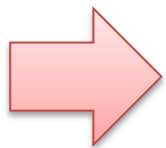
# 隠れ変数・パラメータの推定： 非常に難しくなります

- 原因1: Soft-max関数のため、共役性 (conjugate)を利用できません 😞



(collapsed) Gibbs samplingが不可能になるため、  
変分ベイズ法が候補になります

- 原因2: 時刻  $t$  が前時刻  $t-1$  に依存するため、時間依存性を考慮した推定が必要になります 😞



時間発展するパラメータは、時刻依存性を考慮して  
変分ベイズ法を構築する必要があります

# 時間発展パラメータの推定: 状態空間モデル解釈 [北川, 2005]

- 連続なパラメータの時間変化を追いかける定番の手法です
- DTMの時間発展部分も状態空間モデルとして解釈できます

一般の状態空間モデル

DTM(k, d, zなどを省略)

状態モデル  $y_t | y_{t-1} \sim f(y_{t-1}, \theta)$        $\beta_t | \beta_{t-1} \sim N(\beta_{t-1}, \sigma^2 I)$

観測モデル  $x_t | y_t \sim g(x_t, \varphi)$        $x_{t,n} | \beta_t \sim \text{Mult}(\pi(\beta_t))$

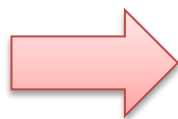
# 変分近似によるKalman filter

[Kalman, 1960]

- 状態モデル、観測モデルの双方が正規分布の場合, Kalman Filterを用いてexactな解が計算できます
- 変分事後分布として、観測モデルに正規分布を“強引に”仮定して推論します

$$\beta_t | \beta_{t-1} \sim N(\beta_{t-1}, \sigma^2 I)$$

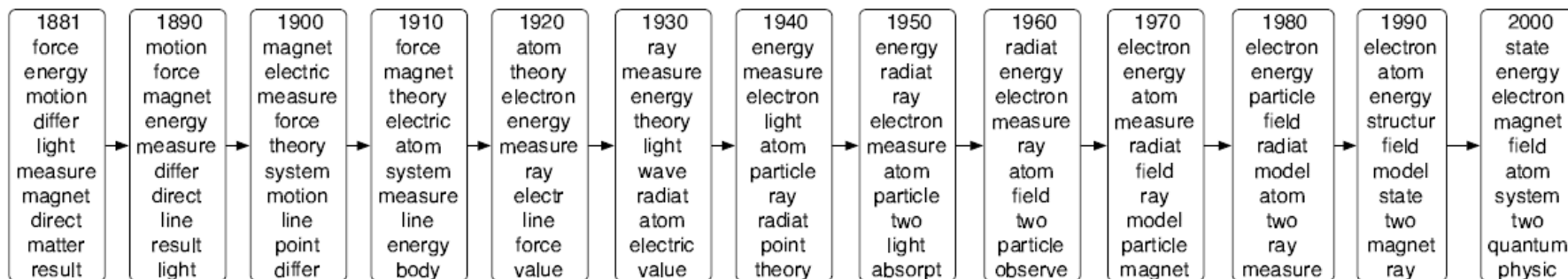
$$x_{t,n} | \beta_t \sim \text{Mult}(\pi(\beta_t))$$



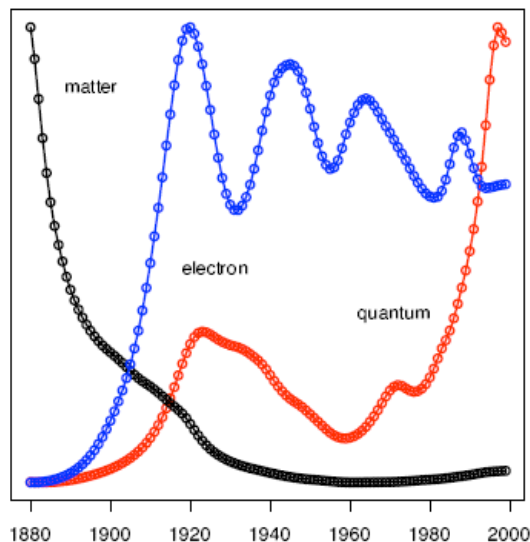
$$\beta_t | \beta_{t-1} \sim N(\beta_{t-1}, \sigma^2 I)$$

$$\hat{\beta}_t | \beta_t \sim N(\beta_t, \hat{v}_t I)$$

変分観測量

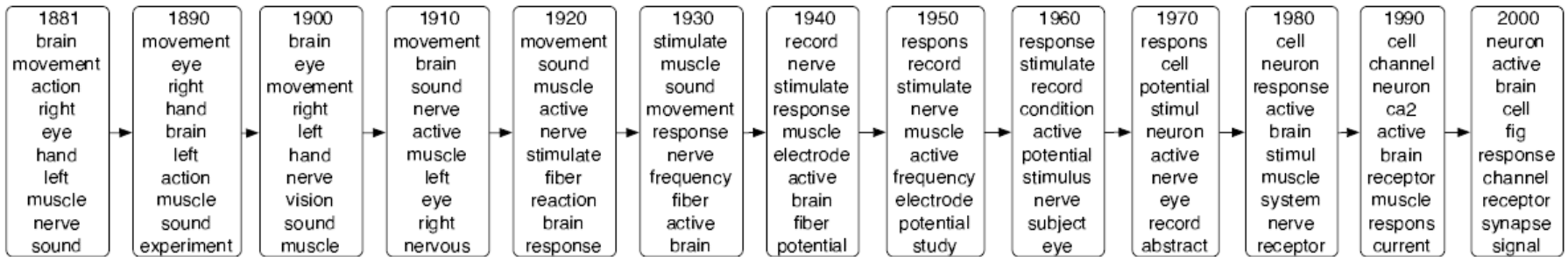


## "Atomic Physics"

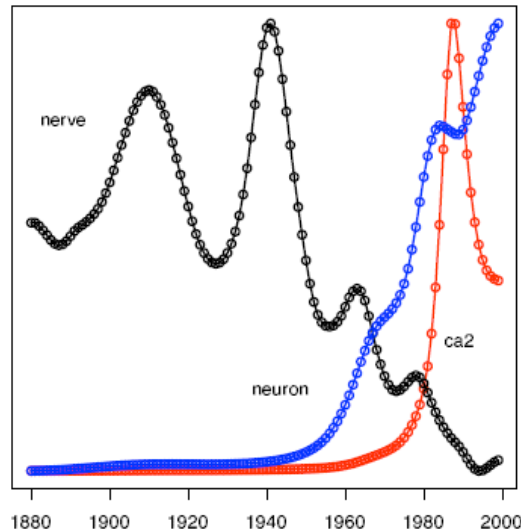


- 1881 On Matter as a form of Energy
- 1892 Non-Euclidean Geometry
- 1900 On Kathode Rays and Some Related Phenomena
- 1917 "Keep Your Eye on the Ball"
- 1920 The Arrangement of Atoms in Some Common Metals
- 1933 Studies in Nuclear Physics
- 1943 Aristotle, Newton, Einstein. II
- 1950 Instrumentation for Radioactivity
- 1965 Lasers
- 1975 Particle Physics: Evidence for Magnetic Monopole Obtained
- 1985 Fermilab Tests its Antiproton Factory
- 1999 Quantum Computing with Electrons Floating on Liquid Helium

[Blei & Lafferty, 2006]

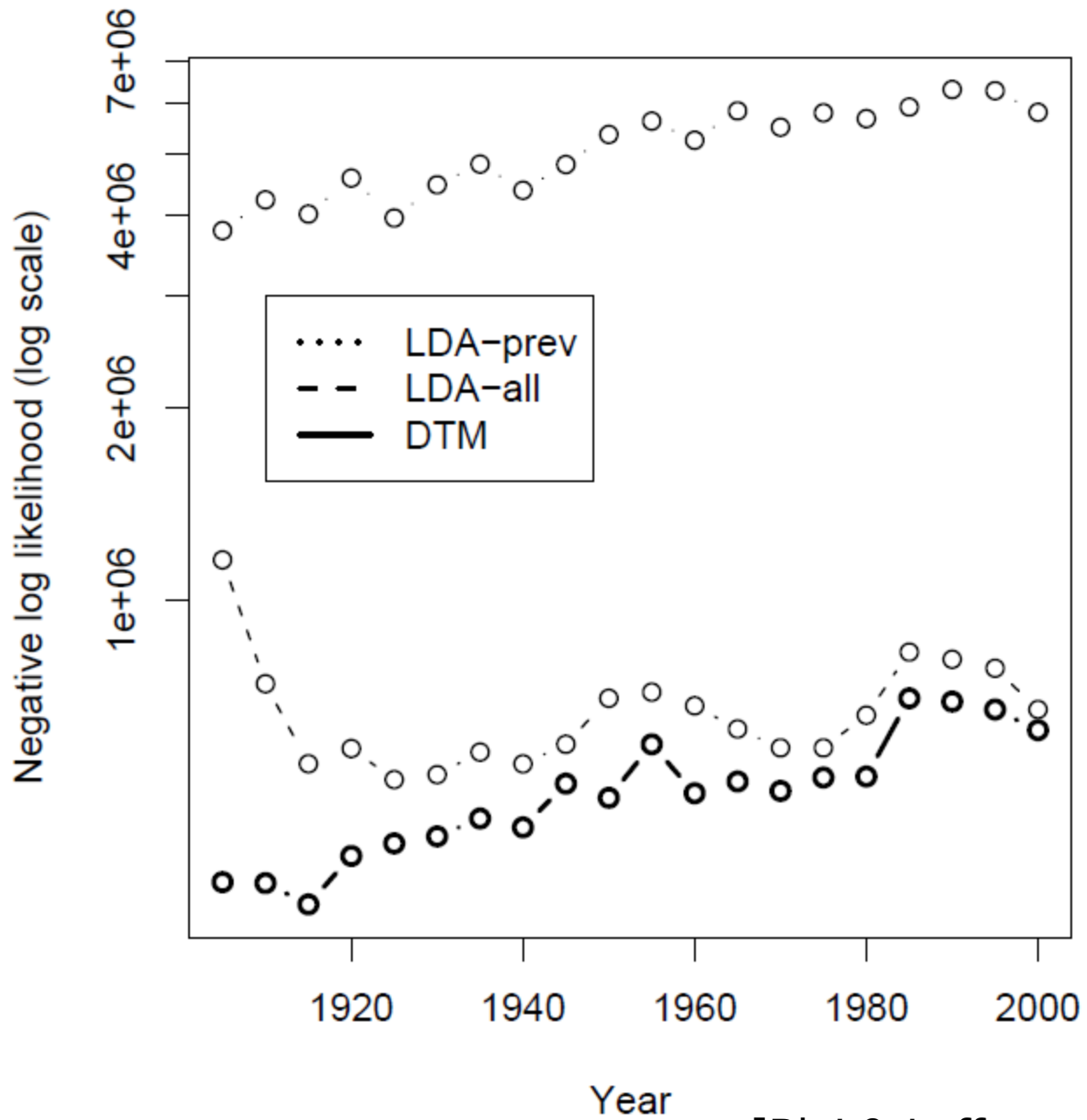


"Neuroscience"



- 1887 Mental Science
- 1900 Hemianopsia in Migraine
- 1912 A Defence of the ``New Phrenology''
- 1921 The Synchronal Flashing of Fireflies
- 1932 Myoesthesia and Imageless Thought
- 1943 Acetylcholine and the Physiology of the Nervous System
- 1952 Brain Waves and Unit Discharge in Cerebral Cortex
- 1963 Errorless Discrimination Learning in the Pigeon
- 1974 Temporal Summation of Light by a Vertebrate Visual Receptor
- 1983 Hysteresis in the Force-Calcium Relation in Muscle
- 1993 GABA-Activated Chloride Channels in Secretory Nerve Endings

[Blei & Lafferty, 2006]





# まとめ: Dynamic Topic Models

- トピックごとの単語分布、トピックの割合の二つを時間発展させたトピックモデルです
- 正規分布によるdriftで時間遷移を表現します
- 😊 非常に有名なので、時間モデルでは必ず押さえる必要がある論文です

# **Topic Tracking Model**

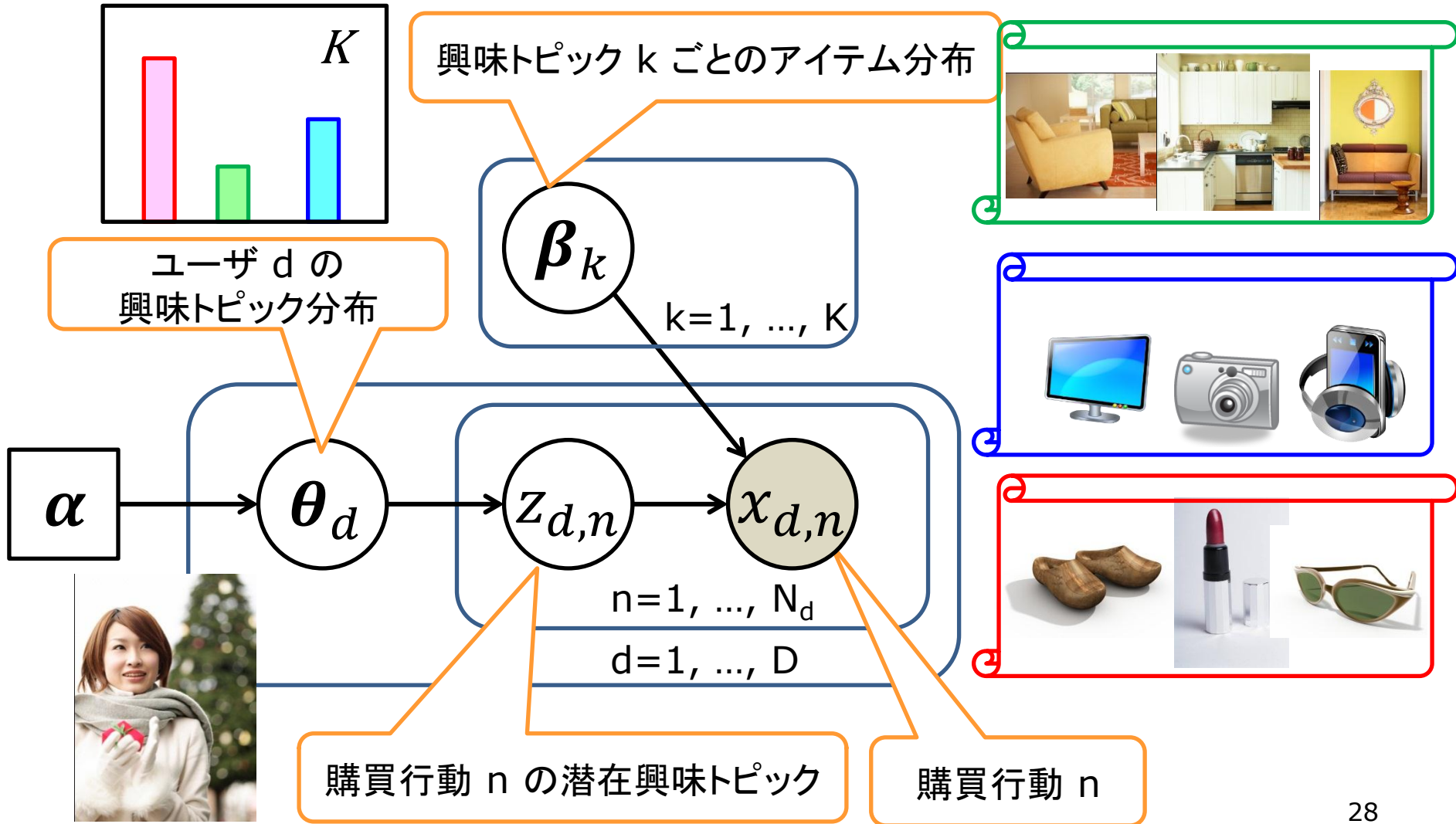
**[Iwata, 2009]**

Iwata et al,  
"Topic Tracking Model for Analyzing Consumer  
Purchase Behavior",  
in Proc. IJCAI, 2009.

# 購買履歴データへの トピックモデル応用

- PLSIなどのように、潜在変数モデルを使った購買履歴データのモデリングは多数存在します (e.g. [Jin, 2004])
- 当然、トピックモデルによる購買履歴データモデリングを考慮することもできます

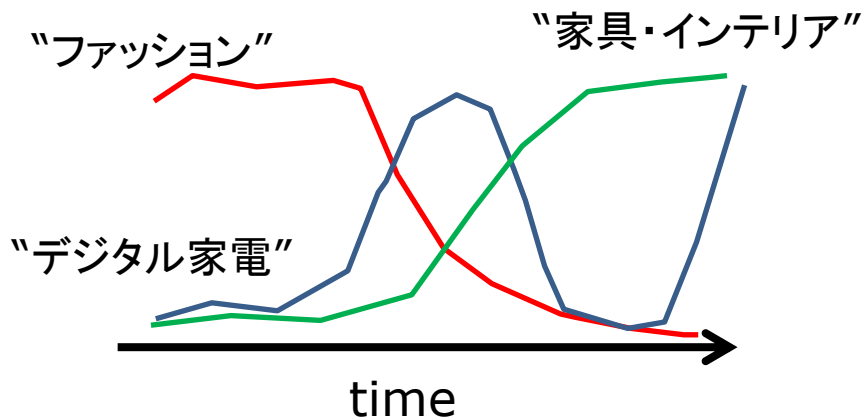
# 購買履歴データの 新しいトピック化



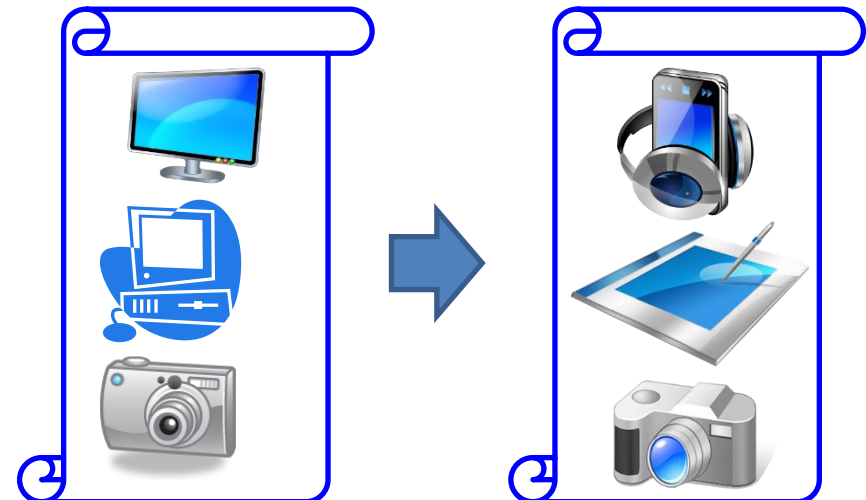
# 注目する時間依存性： ユーザ＋トピック

- 1: ユーザの興味は少しずつ変わります
- 2: 興味トピックの中でのアイテムの売れ筋も変化します

ユーザの中での興味トピック分布



興味トピックの中での売れ筋

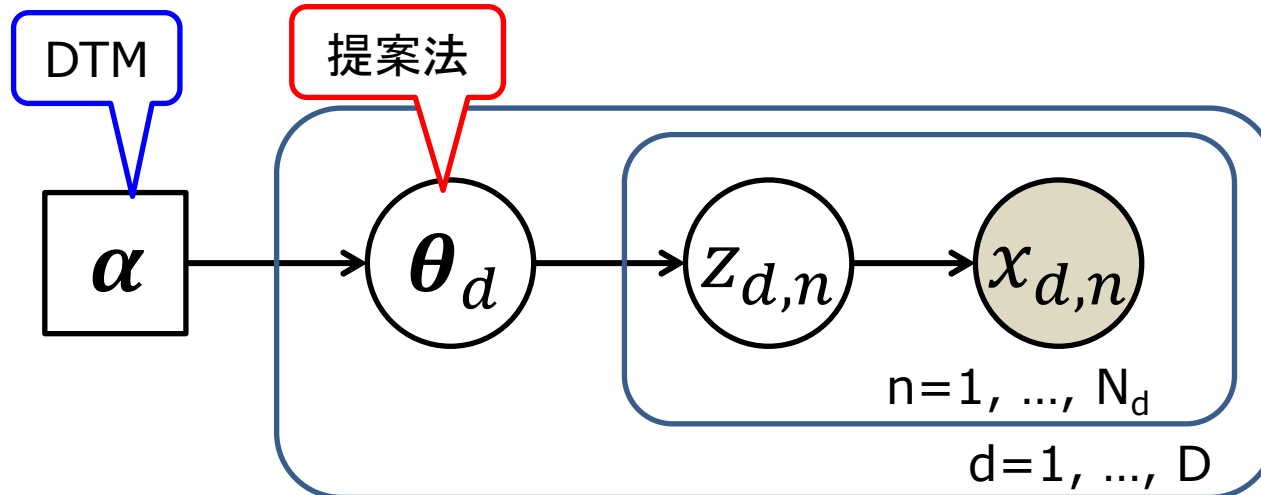


# 提案法: Topic Tracking Model

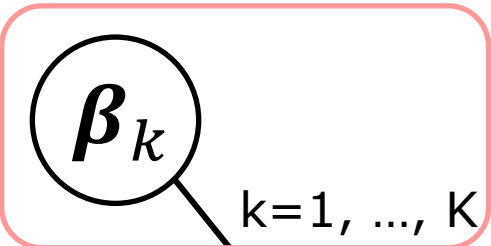
- Dynamic Topic Model(DTM)とはまた違う  
時系列トピックモデルです
- 文書(ユーザ)ごとのトピック分布と、トピック  
の単語(アイテム)分布が時間遷移します
- 推論はDTMに比べて少し簡単になるように  
工夫されています

# 提案法のアイデア: DTMとモデリングの観点が違います

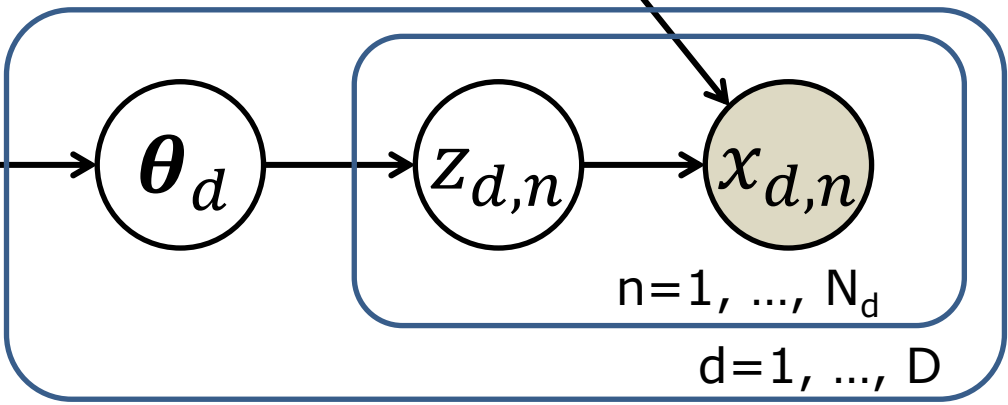
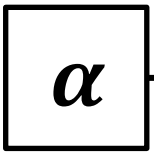
- DTM: 各年度での話題の隆盛が知りたい → トピック分布制御パラメータ  $\alpha$  を時間依存
- 提案法: ユーザの興味の変化が知りたい → 各ユーザ(文書)のトピック分布  $\theta$  を時間依存させる



LDA



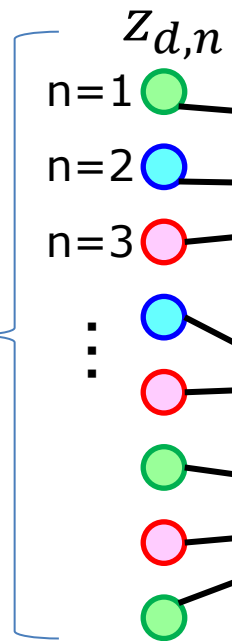
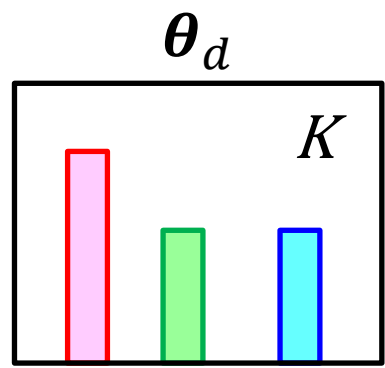
データ	.05
解析	.04
計算機	.03
...	...



リンク	.04
ソーシャル	.02
マイニング	.01
...	...

$\beta_k$

構造	.04
機械学習	.03
最適	.01
...	...



**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るため、**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**が大きな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**...

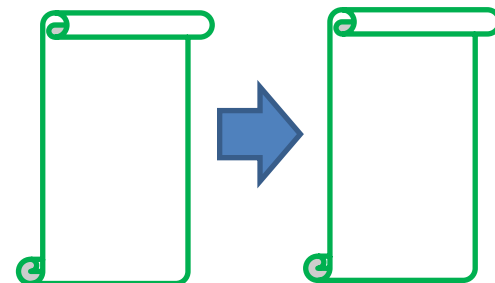
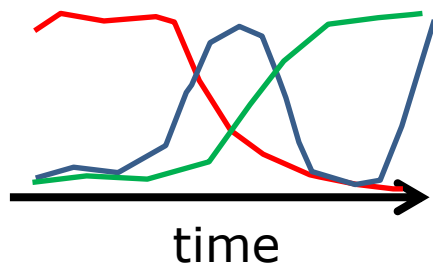
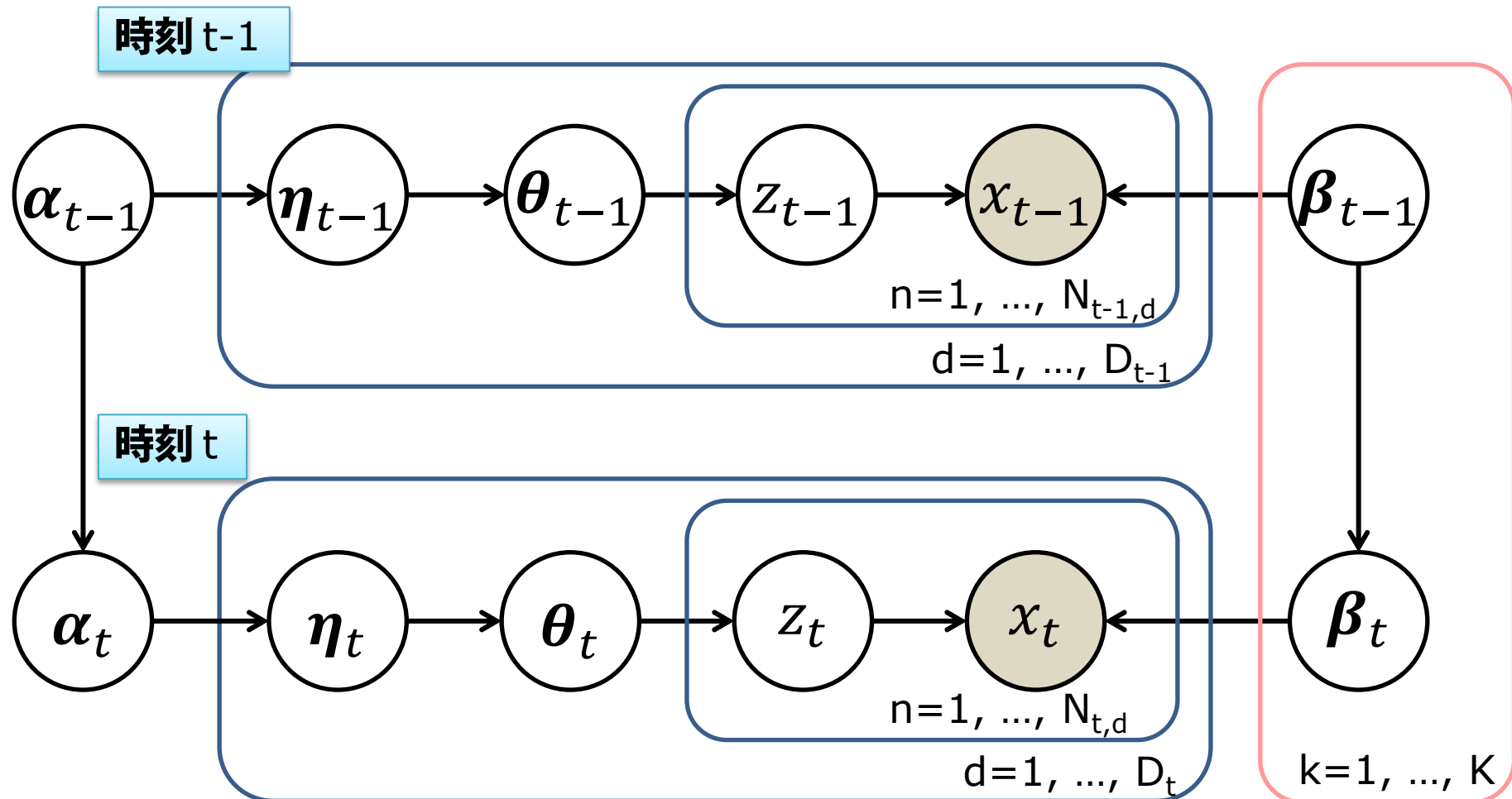
NTTコミュニケーション科学基礎研究所では、統計的・確率的基準のデータ解析技術に基づいた**データマイニング技術**の研究開発を行っています。多くの場合、**統計的機械学習**ではデータを数値化して取り扱います。本

顧客が、ある商品を何度購入したかという**データ**列をつくるのが可能です。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、**ソーシャルネットワーク**...

$x_{d,n}$

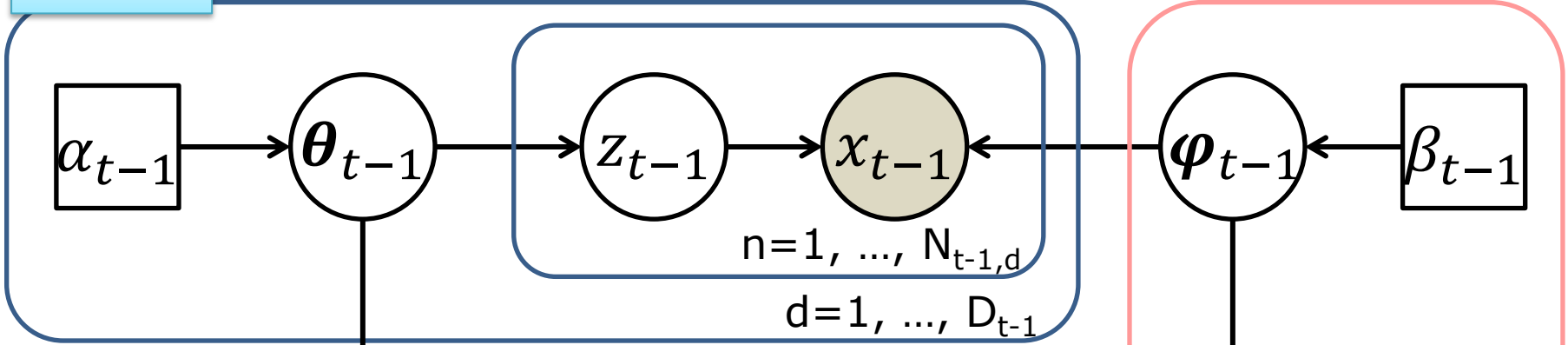


Dynamic Topic Model (添え字  $d, n, k$  は省略)

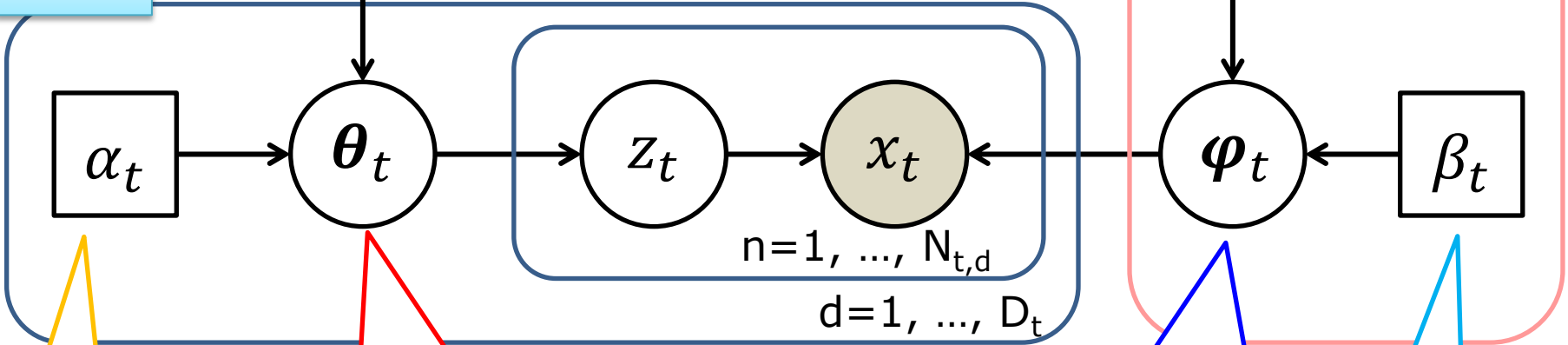


Topic Tracing Model (添え字  $d, n, k$  は省略)

時刻  $t-1$



時刻  $t$



Persistencey  
パラメータ

各ユーザ  $d$  の  
興味トピック分布が  
マルコフ依存

各トピック  $k$  の  
アイテム単語の分布が  
マルコフ依存

Persistencey  
パラメータ

# 生成モデル

for 時間  $t = 1, 2, \dots, T$

for 興味topic  $k = 1, 2, \dots, K$

topic-item word proportion parameter  $\beta_{t,k}$

for ユーザ  $d = 1, 2, \dots, D_t$

topic proportion parameter  $\alpha_{t,d}$

topic proportion

for 購買行動  $n = 1, 2, \dots, N_{t,d}$

topic-item word assignment

item word observation

for 時間  $t = 1, 2, \dots, T$

for 興味topic  $k = 1, 2, \dots, K$

topic-item word proportion evolution

$$\boldsymbol{\varphi}_{t,k} | \hat{\boldsymbol{\varphi}}_{t-1,k}, \beta_{t,k} \sim \text{Dir}(\beta_{t,k} \hat{\boldsymbol{\varphi}}_{t-1,k})$$

for ユーザ  $d = 1, 2, \dots, D_t$

topic proportion evolution

$$\boldsymbol{\theta}_{t,d} | \hat{\boldsymbol{\theta}}_{t-1,d}, \alpha_{t,d} \sim \text{Dir}(\alpha_{t,d} \hat{\boldsymbol{\theta}}_{t-1,d})$$

for 購買行動  $n = 1, 2, \dots, N_d$

topic-item word assignment

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

item word observation

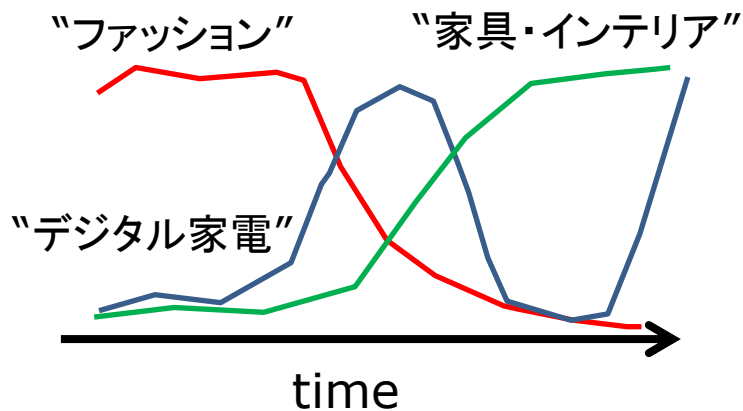
$$x_{t,d,n} | z_{t,d,n}, \{\boldsymbol{\varphi}_{t,k}\} \sim \text{Mult}(\boldsymbol{\varphi}_{t,z_{t,d,n}})$$

$\hat{\cdot}$  は"事後分布での期待値"を表す

# 興味トピック分布のモデル

- DTMと違い、ディリクレ分布を利用して時間発展をモデル化しています
- ユーザ、時間ごとに、興味トピックの持続度 (persistence) もモデル化します

平均  $\hat{\theta}_{t-1,d}$  のディリクレ分布



$$\theta_{t,d} | \hat{\theta}_{t-1,d}, \alpha_{t,d} \sim \text{Dir}(\alpha_{t,d} \hat{\theta}_{t-1,d})$$

$\hat{\theta}_{t-1,d}$  : 前時刻の事後分布期待値

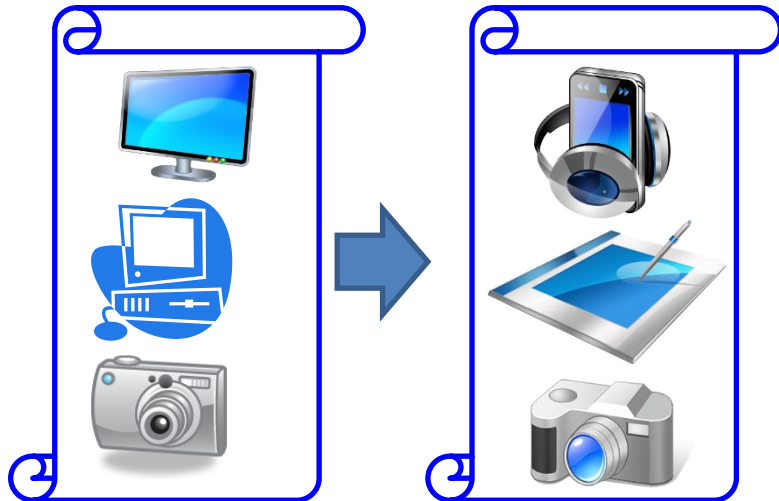
$\alpha_{t,d}$  : 持続度パラメータ

$\alpha$ 大 =  $\theta_t$ の分散小 → 小さな時間変化

$\alpha$ 小 =  $\theta_t$ の分散大 → 大きな時間変化

# トピック-アイテム（単語）分布のモデル

- 興味トピックと同様です



平均  $\hat{\varphi}_{t-1,d}$  のディリクレ分布

$$\varphi_{t,k} | \hat{\varphi}_{t-1,k}, \beta_{t,k} \sim \text{Dir}(\beta_{t,k} \hat{\varphi}_{t-1,k})$$

$\hat{\varphi}_{t-1,d}$ : 前時刻の事後分布期待値

$\beta_{t,d}$ : 持続度パラメータ

$\beta$ 大 =  $\phi_t$ の分散小  $\rightarrow$  小さな時間変化

$\beta$ 小 =  $\phi_t$ の分散大  $\rightarrow$  大きな時間変化

# 長期時間依存モデル

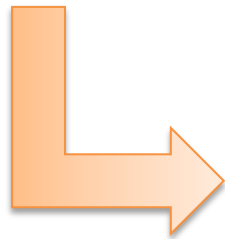
- 1時刻前に依存するだけでなく、数ステップ前までに依存する形への拡張も簡単です

1ステップ前からの依存関係モデル

$$\varphi_{t,k} | \hat{\varphi}_{t-1,k}, \beta_{t,k} \sim \text{Dir}(\beta_{t,k} \hat{\varphi}_{t-1,k})$$

$$\theta_{t,d} | \hat{\theta}_{t-1,d}, \alpha_{t,d} \sim \text{Dir}(\alpha_{t,d} \hat{\theta}_{t-1,d})$$

Lステップ前からの依存関係モデル



$$\varphi_{t,k} | \hat{\varphi}_{t-1,k}, \beta_{t,k} \sim \text{Dir} \left( \sum_{l=1}^L \beta_{t,k,l} \hat{\varphi}_{t-l,k} \right)$$

$$\theta_{t,d} | \hat{\theta}_{t-1,d}, \alpha_{t,d} \sim \text{Dir} \left( \sum_{l=1}^L \alpha_{t,d,l} \hat{\theta}_{t-l,d} \right)$$

# 隠れ変数・パラメータの推定

- 😊 非常に簡単な逐次推定アルゴリズムが導出できます
  - 正規分布やsoft-maxがないため！！
  - LDAのGibbs, VB (EM) を導出したことがある方にとっては自明な解が得られます
- 😞 ただし、DTMのように系列としての最適解は得られません



Table 1: Average  $N$ -best accuracies (%) over time. The digit in the bracket is the standard deviation.

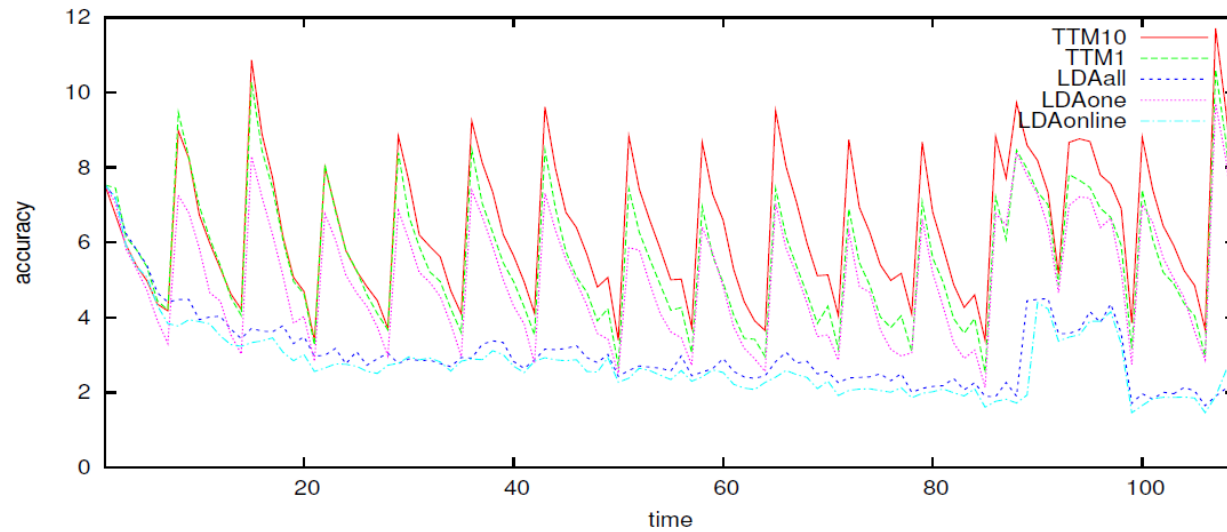
(a) movie

$N$	LDAall	LDAonline	LDAone	TTM1	TTM10
1	1.21 (0.61)	1.08 (0.54)	1.91 (0.78)	2.22 (0.91)	<b>2.46</b> (0.92)
2	2.18 (0.79)	2.00 (0.78)	3.52 (1.22)	3.99 (1.33)	<b>4.47</b> (1.36)
3	3.06 (1.04)	2.81 (1.02)	5.04 (1.64)	5.60 (1.75)	<b>6.35</b> (1.85)
4	3.90 (1.27)	3.56 (1.24)	6.24 (1.90)	6.82 (2.01)	<b>7.82</b> (2.15)
5	4.70 (1.51)	4.26 (1.44)	7.37 (2.20)	7.92 (2.26)	<b>9.20</b> (2.42)

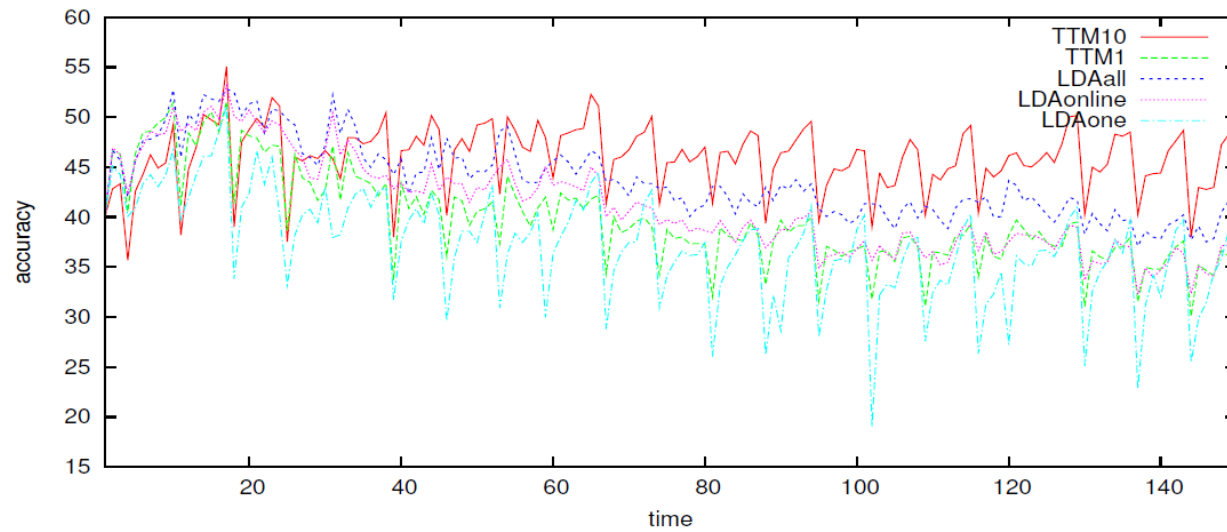
(b) cartoon

$N$	LDAall	LDAonline	LDAone	TTM1	TTM10
1	27.0 (3.3)	26.0 (3.5)	24.8 (4.5)	26.8 (4.2)	<b>30.5</b> (3.4)
2	37.3 (3.6)	35.1 (4.2)	32.4 (4.9)	34.2 (4.5)	<b>39.9</b> (3.5)
3	43.7 (3.9)	41.1 (4.8)	37.2 (5.3)	39.8 (4.6)	<b>45.9</b> (3.3)
4	48.5 (4.0)	45.8 (5.1)	40.9 (5.3)	44.5 (4.6)	<b>50.6</b> (3.2)
5	52.4 (4.2)	49.6 (5.4)	44.1 (5.4)	48.5 (4.6)	<b>54.4</b> (3.0)

[Iwata, 2009]



(a) movie



(b) cartoon

[Iwata, 2009]

Figure 3: Three-best accuracies (%) for each day.

# まとめ: Topic Tracking model

- ユーザ(文書)ごとのトピック分布、トピックの単語分布を時間発展させたトピックモデル
- Dirichletで時間遷移を表現したことで、非常に簡単に解を導出できます

# その他の時系列データ応用

- Wang and McCallum, “Topics over Time: A Non-Markov Continuous-Time model of Topical Trends”, in Proc. KDD, 2006.
- Iwata et al., “Sequential Modeling of Topic Dynamics with Multiple Timescales”, ACM Trans. on Knowledge Discovery from Data. Vol. 5(4). pp. 19:1-19:27, 2012.
- Pruteanu-Malinici, et al., “Hierarchical Bayesian Modeling of Topics in Time-Stamped Documents”, IEEE Trans. PAMI, Vol. 32(6), pp.996-1011, 2010.

# 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.
- [Blei & Lafferty, 2006], Blei and Lafferty, “Dynamic Topic Models”, in Proc. ICML, 2006.
- [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.
- [北川, 2005] 北川, “時系列解析入門”, 岩波書店, 2005.
- [Kalman, 1960] Kalman, “A New Approach to Linear Filtering and Prediction Problems”, Journal of Basic Engineering, 1960.
- [Iwata, 2009] Iwata et al, “Topic Tracking Model for Analyzing Consumer Purchase Behavior”, Proc. in IJCAI, 2009.
- [Jin, 2004] Jin et al, “Web Usage Mining based on Probabilistic Latent Semantic Anlysis”, Proc. in KDD, 2004.