

混合分布学習の基礎から最先端

第1部: 混合分布の基礎 (藤巻)

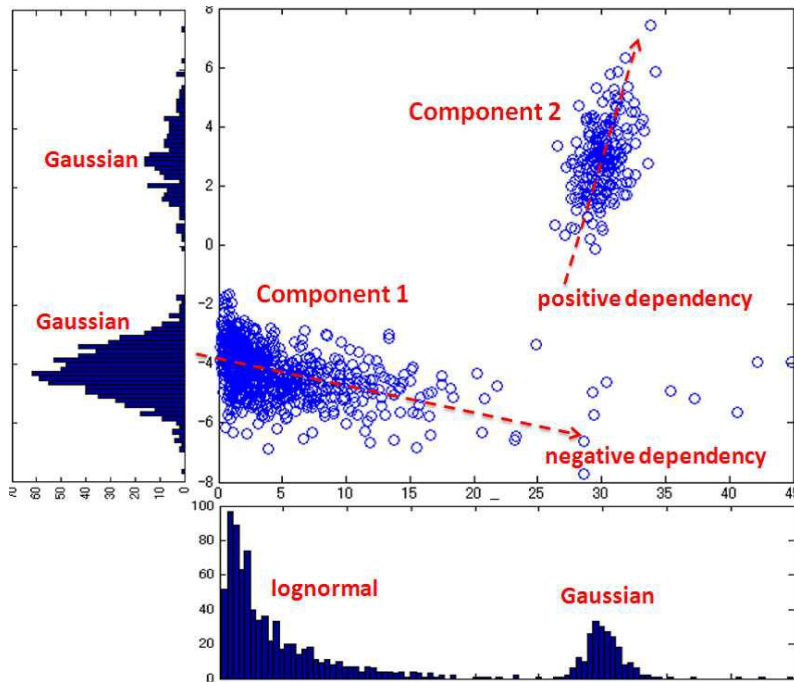
第2部: ノンパラメトリックベイズ法 (佐藤)

第3部: 異種混合学習と因子化漸近ベイズ (藤巻)

Ryohei Fujimaki (NEC Laboratories America)
joint Tutorial with
Issei Sato (University of Tokyo)

NEC 基幹技術フォーラム
2013/4/4

第3部 異種混合学習と因子化漸近ベイズ(藤 巻)



第3部の構成

- 異種混合学習とモデル選択
- 因子化情報量基準と因子化漸近ベイズ推論
- 因子化漸近ベイズ推論の性質
- 実験
- 発展モデルへ

分析案件を思い返してみると。。。

■ 自動車センサデータ

- 車種、走行環境(天候、道路状況。)、走行状態(加減速、停止、カーブ。)、、、

■ 健康診断データ

- 生活習慣、治療・通院状況、遺伝的リスク、、

■ サーバートラフィック分析

- サーバー特性、運用(データ移動、サーバー追加。)、周期(平日、週末、祝日。)、、、

■ 金融リスク分析

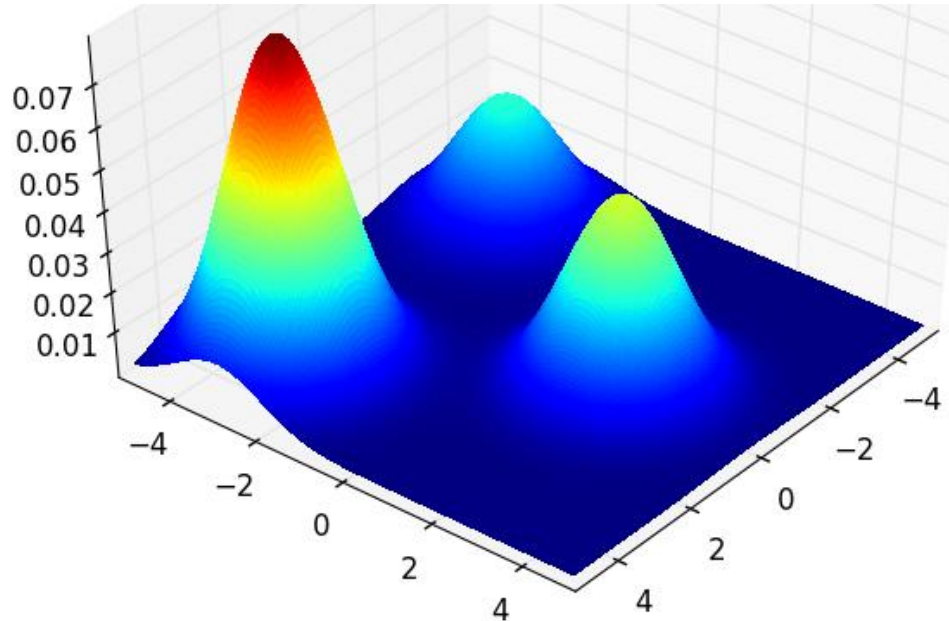
- リスク(高頻度低リスク、低頻度高リスク)、銀行タイプ(メガバンク、地銀。)、、、

• データは複雑怪奇に混ざっている。うまくモデルを選択しないと沈没

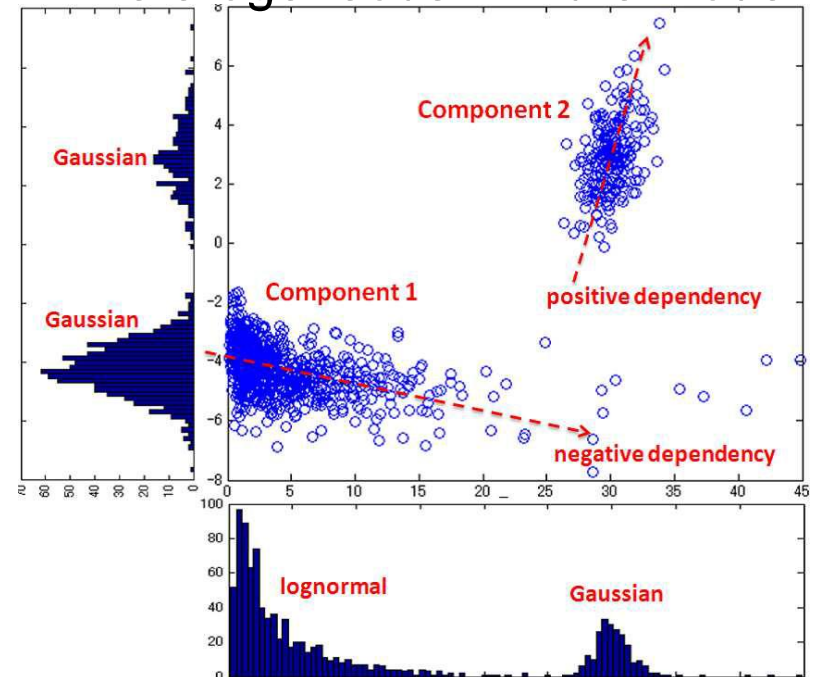
Mixture Distributions

- One of the most basic models in machine learning
 - Density estimation, clustering, anomaly detection, feature learning, ...

Gaussian Mixture Model



Heterogeneous Mixture Model



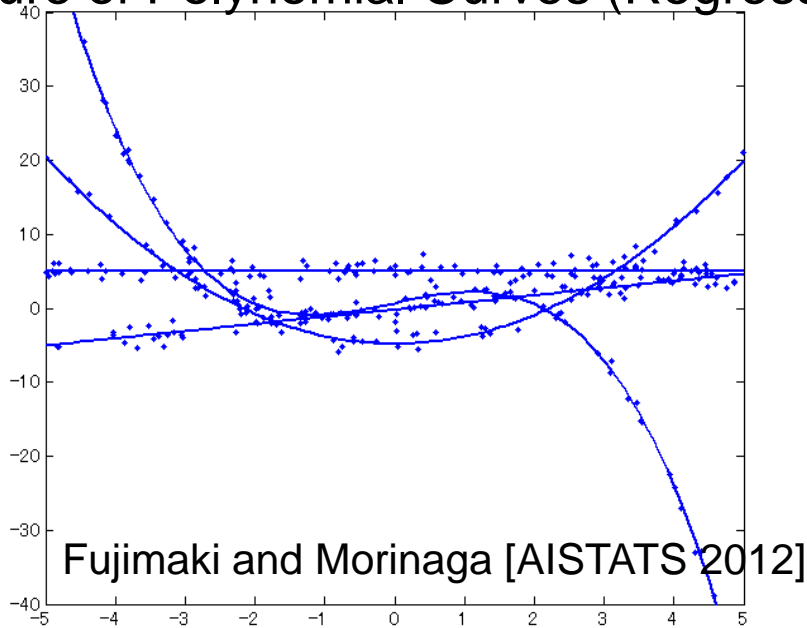
Fujimaki et al. [KDD2011]

- Learning Problem
 - How many distributions are mixed?
 - What are their parameters?
 - (What types of distributions are mixed?)

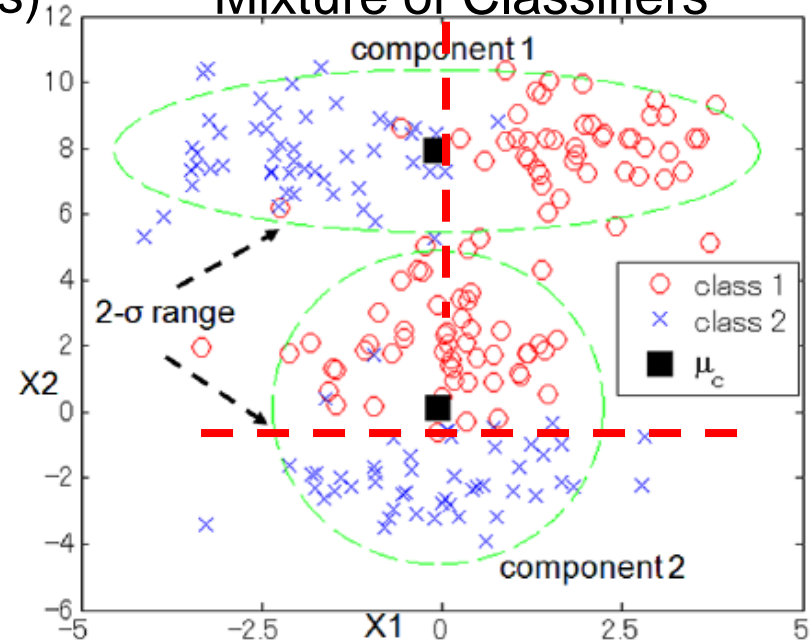
Mixture Experts

- Basic ensemble learners
 - Multimodality in classification, regression, etc.

Mixture of Polynomial Curves (Regressors)



Mixture of Classifiers

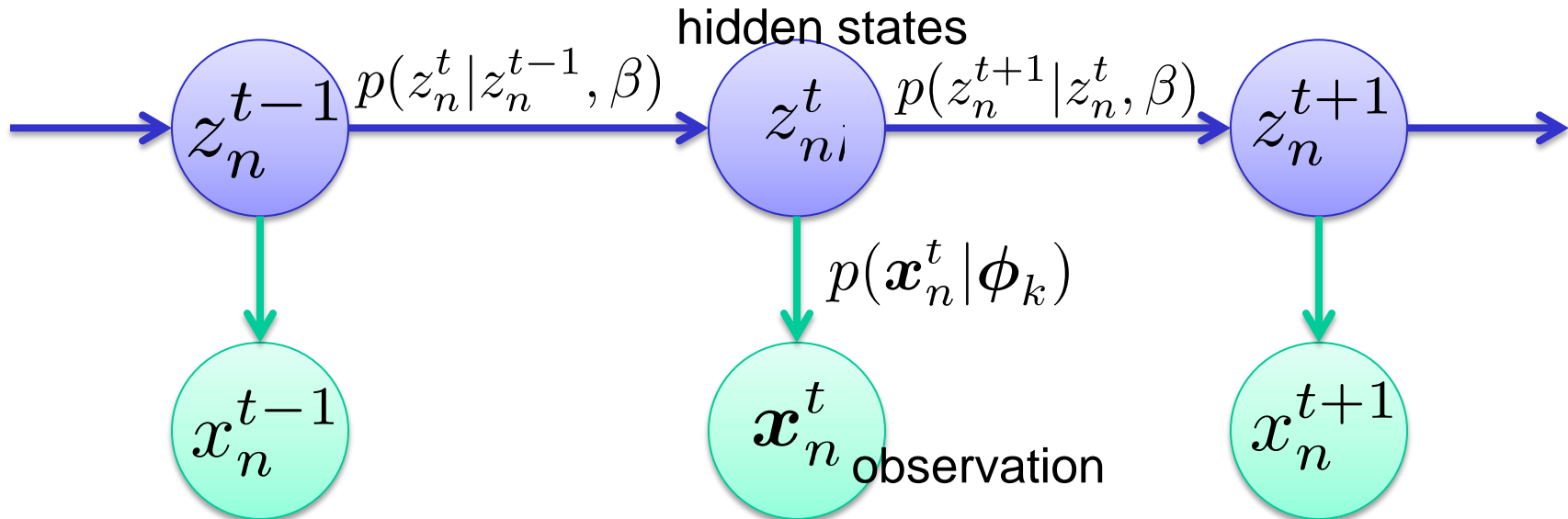


- Learning Problem

- How many experts are mixed?
- What are their parameters?
- What are appropriate levels of their regularizations? (or how many features does each expert use?)

Hidden Markov Models

- One of the most important models for sequential data.
 - time series analysis, speech recognition, text indexing, etc.

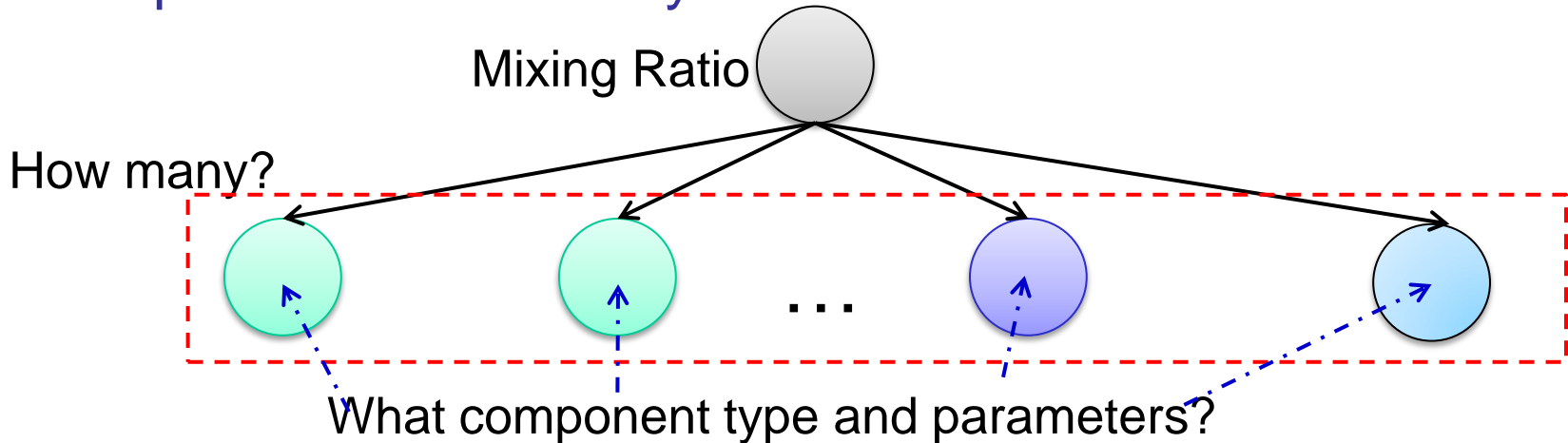


- Learning Problem

- How many hidden states are necessary?
- What are the transition and observation parameters?

Challenges

■ Computational Efficiency



■ Identifiability issue of the maximum likelihood (ML) estimators

- Ex. Mixture of two identical standard normal distributions

equivalent class (same likelihood)

$$\left[\begin{array}{l} 0.1 * N(0,1) + 0.9 * N(0,1) \\ 0.5 * N(0,1) + 0.5 * N(0,1) \\ \vdots \\ a * N(0,1) + b * N(0,1) \end{array} \right] = N(0,1)$$

Mapping is NOT 1-to-1 in likelihood space (non-identifiable)

- Generalization performance of a non-identifiable estimator is much worse than those of identifiable estimators (S.Watanabe 2000)

Bayesian Model Selection

- Bayesian Approximation Inference

$$M^* = \arg \max_M \log p(\mathbf{x}^N | M)$$

$$= \arg \max_M \log \int_{\boldsymbol{\theta}} p(\mathbf{x}^N | \boldsymbol{\theta}) p(\boldsymbol{\theta} | M) d\boldsymbol{\theta}$$

computationally/analytically intractable

- Bayesian marginalization automatically takes model complexity into account.

- Non-parametric Bayesian Modeling

- Plug-in “infinite prior” (Dirichlet process, hierarchical Dirichlet process, Pitman-Yaw process, etc.)
- Prior controls model complexity.

Bayes Information Criterion (BIC)

[Shwarz 1978]

- Bayes Information Criterion

parameter dimensionality

$$BIC(M|\mathbf{x}^N) = \log p(\mathbf{x}^N | \boldsymbol{\theta}_{ML}) + \frac{\mathcal{D}_M}{2} \log N$$

model complexity

- Non-regularity Issue and Non-Identifiability issues

Fisher information matrix is singular around ML estimator
 -> BIC's complexity loses theoretical justification
 (In fact, BIC over-estimates the complexity*)

ML estimator is not unique and non-identifiable
 -> Generalization performance becomes significantly worse**

*S. Watanabe (Neural Computation 2001)

**K. Yamazaki and S. Watanabe (Neural Networks 2003)

Variational Bayesian Inference

[MacKay 1997, Beal 2003]

- Maximize lower bound instead of marginal log-likelihood.

$$\log p(\mathbf{x}^N | M) \geq \sum_{\mathbf{z}^N} \int q(\mathbf{z}^N) q_{\theta}(\boldsymbol{\theta}) \log \frac{p(\mathbf{x}^N, \mathbf{z}^N | \boldsymbol{\theta}) p(\boldsymbol{\theta} | M)}{q(\mathbf{z}^N) q_{\theta}(\boldsymbol{\theta})} d\boldsymbol{\theta}$$

equality does NOT hold
in general

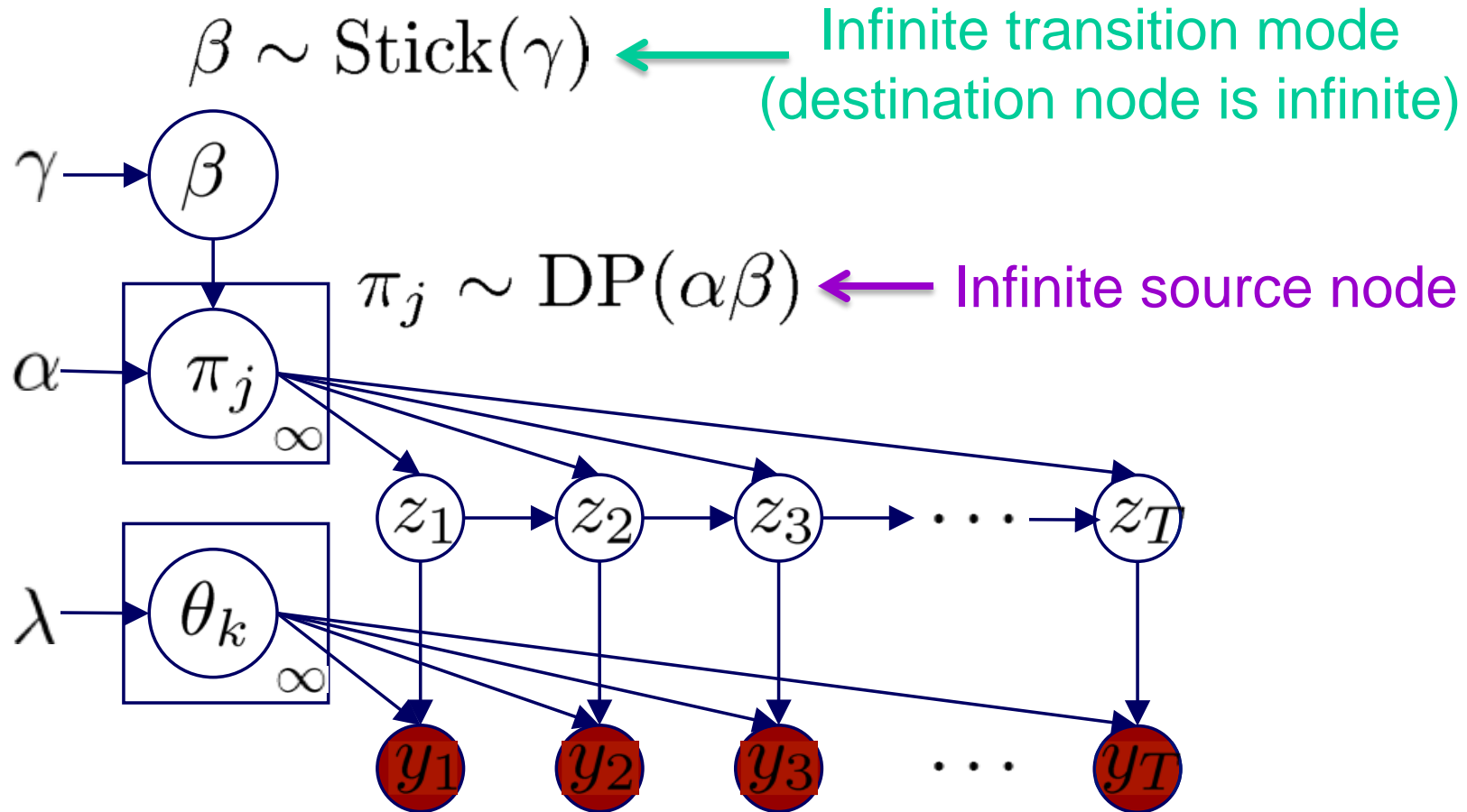
$$q(\mathbf{z}^N, \boldsymbol{\theta}) = q(\mathbf{z}^N) q_{\theta}(\boldsymbol{\theta})$$

Assume independency between
hidden states and parameters

- VB EM algorithm
 - VB Estep (VB forward-backward algorithm): maximization of the lower bound in terms of $q(\mathbf{z}^N)$
 - VB Mstep: maximization of the lower bound in terms of $q_{\theta}(\boldsymbol{\theta})$

Non-parametric Bayesian Method (infinite HMMs) [Beal 2002, van Gael 2008]

- Hierarchical Dirichlet Process Prior



- Efficient MCMC method : Beam Sampling

第3部の構成

- 異種混合学習とモデル選択
- 因子化情報量基準と因子化漸近ベイズ推論
- 因子化漸近ベイズ推論の性質
- 実験
- 発展モデルへ

Key Idea

- An alternative representation

$$\log p(\mathbf{x}^N | M) = \max_q \left\{ \sum_{\mathbf{z}^N} q(\mathbf{z}^N) \log \left(\frac{p(\mathbf{x}^N, \mathbf{z}^N | M)}{q(\mathbf{z}^N)} \right) \right\}$$

parameter is marginalized out

- “Factorized” Laplace Method

complete likelihood is regular

$$p(\mathbf{x}^N, \mathbf{z}^N | M) = \int \underbrace{p(\mathbf{z}^N | \boldsymbol{\alpha})}_{\text{Laplace approximation}} \prod_{c=1}^C \underbrace{p_c(\mathbf{x}^N | \mathbf{z}_c^N, \boldsymbol{\phi}_c)}_{\text{Laplace approximation}} p(\boldsymbol{\theta} | M) d\boldsymbol{\theta}$$

Laplace approximation

Laplace approximation

Each Laplace approximation is done around the ML estimator of complete likelihood

Key Idea

- An alternative representation

$$\log p(\mathbf{x}^N | M) = \max_q \left\{ \sum_{\mathbf{z}^N} q(\mathbf{z}^N) \log \left(\frac{p(\mathbf{x}^N, \mathbf{z}^N | M)}{q(\mathbf{z}^N)} \right) \right\}$$

parameter is marginalized out

- “Factorized” Laplace Method

complete ML estimator

$$p(\mathbf{x}^N, \mathbf{z}^N | M) \approx p(\mathbf{x}^N, \mathbf{z}^N | \bar{\boldsymbol{\theta}}) \frac{(2\pi)^{D_\alpha/2}}{N^{D_\alpha/2} |\bar{\mathcal{F}}_Z|^{1/2}} \times$$

$$\prod_{c=1}^C \frac{(2\pi)^{D_c/2}}{(\sum_{n=1}^N z_{nc})^{D_c/2} |\bar{\mathcal{F}}_c|^{1/2}}$$

Fisher information matrices

parameter dimensionalities

Factorized Information Criterion (FIC)

- FIC as an approximation of marginal log-likelihood

$$\log p(\mathbf{x}^N | M) \approx FIC(\mathbf{x}^N, M) \equiv \max_q \left\{ \mathcal{J}(q, \bar{\boldsymbol{\theta}}, \mathbf{x}^N) \right\}$$

$$\mathcal{J}(q, \bar{\boldsymbol{\theta}}, \mathbf{x}^N) = \sum_{\mathbf{z}^N} q(\mathbf{z}^N) \left(\log p(\mathbf{x}^N, \mathbf{z}^N | \bar{\boldsymbol{\theta}}) - \frac{\mathcal{D}_\alpha}{2} \log N \right)$$

complete ML estimator is not available

$$- \sum_{c=1}^C \left(\frac{\mathcal{D}_c}{2} \log \left(\sum_{n=1}^N z_{nc} \right) - \log q(\mathbf{z}^N) \right)$$

complexity of the c-th component is computationally intractable

Theorem 1 $FIC(\mathbf{x}^N, M)$ is asymptotically consistent with $\log p(\mathbf{x}^N | M)$.

Factorized Asymptotic Bayesian Inference (FAB)

- Lower bound of FIC

$$\begin{aligned}
 FIC(\mathbf{x}^N, M) &\geq \mathcal{G}(q, \tilde{q}, \theta, \mathbf{x}^N) \\
 &\equiv \sum_{\mathbf{z}^N} q(\mathbf{z}^N) \left(\log p(\mathbf{x}^N, \mathbf{z}^N | \theta) - \frac{\mathcal{D}_\alpha}{2} \log N \right. \\
 &\quad \left. - \sum_{c=1}^C \frac{\mathcal{D}_c}{2} \mathcal{L} \left(\sum_{n=1}^N z_{nc}, \sum_{n=1}^N \tilde{q}(z_{nc}) \right) - \log q(\mathbf{z}^N) \right)
 \end{aligned}$$

- FAB solves model selection problems by maximizing the lower bound of FIC

$$M^*, q^*, \theta^*, \tilde{q}^* = \arg \max_{M, q, \theta, \tilde{q}} \mathcal{G}(q, \tilde{q}, \theta, \mathbf{x}^N)$$

Brief Reminder of the EM algorithm for MMs

- Parameter estimation of latent variable models

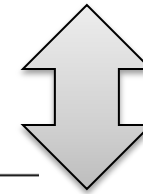
Maximization step (M-step):

Estimating the component parameters

$$\alpha_c^t \propto \sum_{n=1}^N q_{nc}^t$$

$$\phi_c^t = \arg \max_{\phi} \sum_{n=1}^N q_{nc}^t \log p(\mathbf{x}_n | \phi_c^{t-1})$$

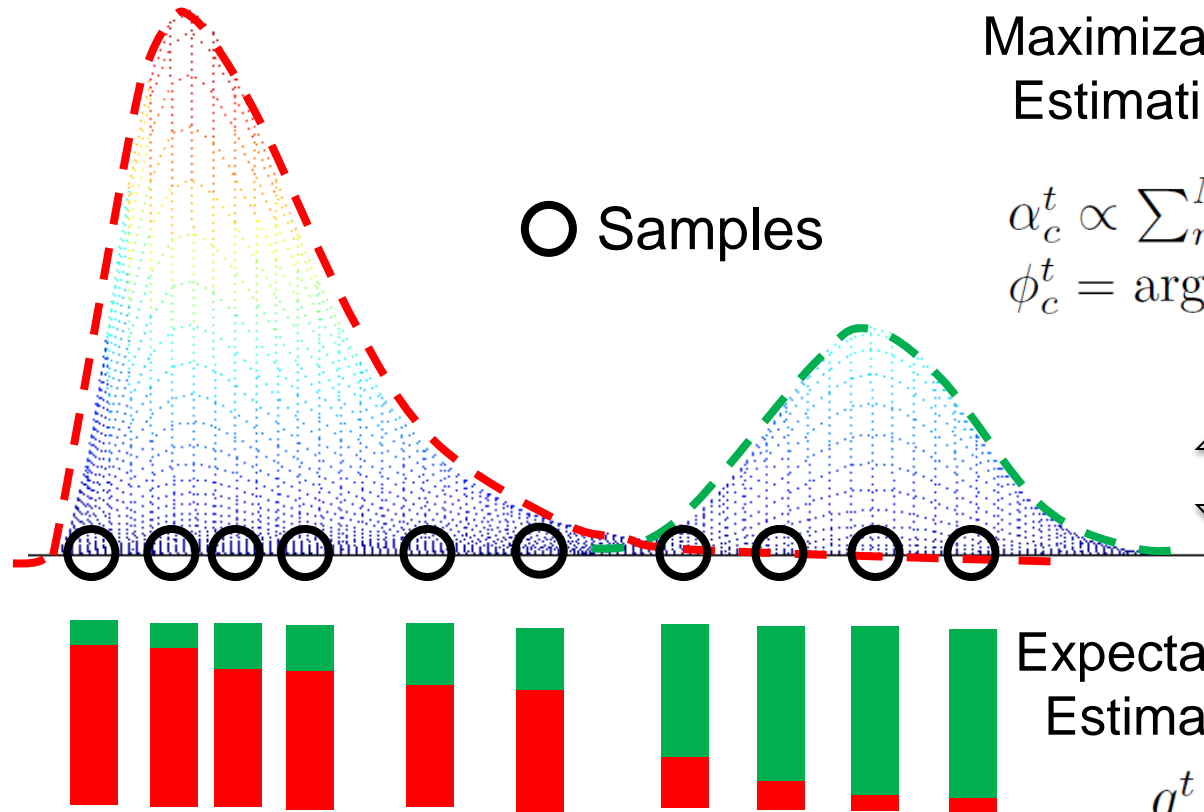
○ Samples



Expectation step (E-step):

Estimating expected assignments

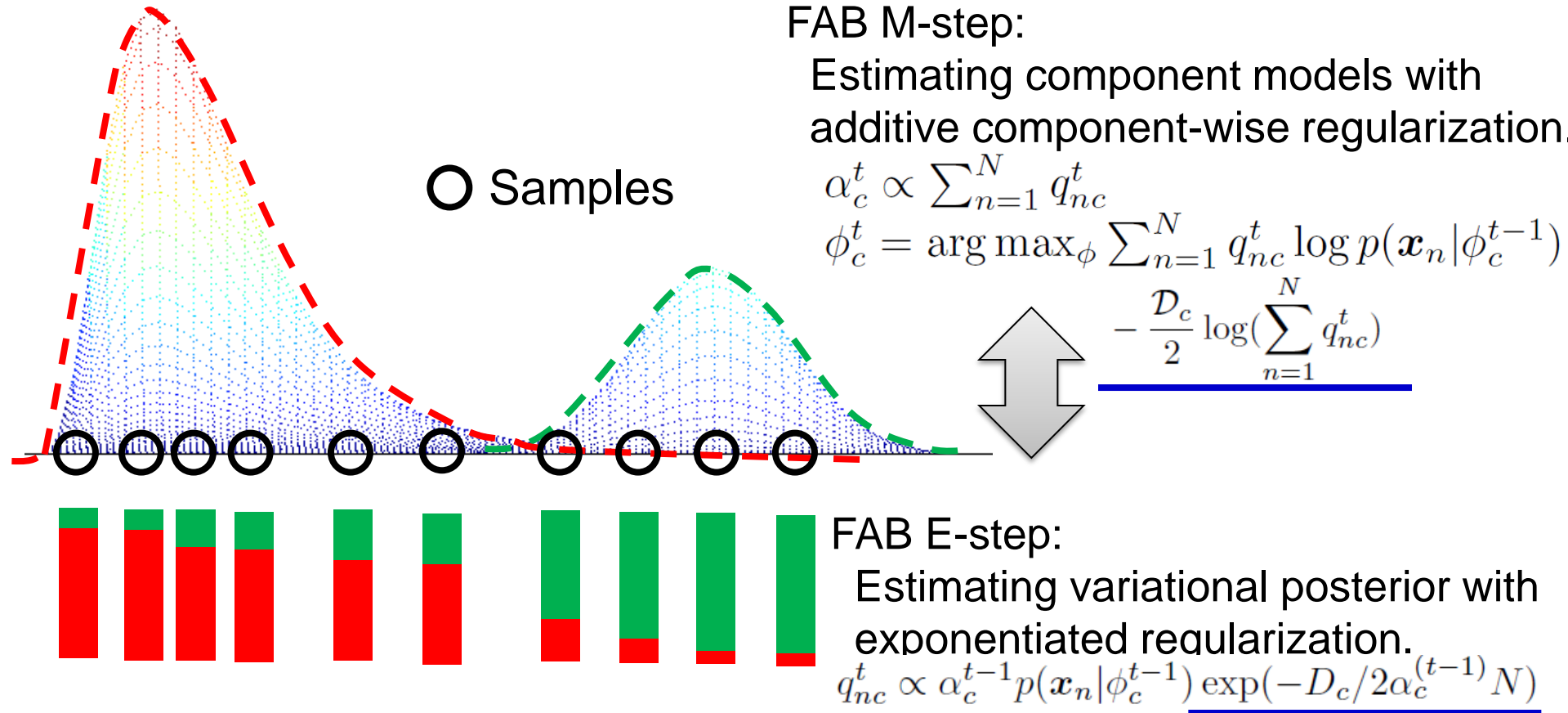
$$q_{nc}^t \propto \alpha_c^{t-1} p(\mathbf{x}_n | \phi_c^{t-1})$$



- Convergence to the ML estimator (non-identifiable)
- Additional model selection (e.g., CV, BIC, ...)

Factorized Asymptotic Bayesian Inference for MMs

- Simultaneous model selection and parameter estimation



- Identifiable estimation based on model-induced regularization
- No hand-determined model selection parameters.

Factorized Asymptotic Bayesian Inference (FAB)

- Gaussian Mixture Models

$$\min_{C, \alpha, \mu_c, \Sigma_c} \left\{ \sum_{n=1}^N -\log \sum_{c=1}^C \alpha_c \mathcal{N}(x_n | \mu_c, \Sigma_c) \right\}$$

- E step : update posterior of latent variables

$$q^{(t)}(z_{nc}) \propto \alpha_c^{(t-1)} p(\mathbf{x}_n | \phi_c^{(t-1)}) \exp\left(\frac{-\mathcal{D}_c}{2\alpha_c^{(t-1)} N}\right)$$

$$\sum_{c=1}^C q^{(t)}(z_{nc}) = 1 \quad \tilde{q} = q^{(t-1)}$$

Smaller/more complex component is more strongly regularized

- M step : update parameters

$$\alpha_c^{(t)} = \sum_{n=1}^N q^{(t)}(z_{nc}) / N \quad \phi_c^{(t)} = \arg \max_{\phi_c} \sum_{n=1}^N q^{(t)}(z_{nc}) \log p(\mathbf{x}_n | \phi_c)$$

Factorized Asymptotic Bayesian Inference (FAB)

- Mixture of regressors

$$\min_{C, \alpha, w_c, \sigma_c} \left\{ \sum_{n=1}^N -\log \sum_{c=1}^C \alpha_c P(y_n | w_c^T x_n, \sigma_c^2) + \sum_{c=1}^C \lambda_c r(w_c) \right\}$$

- E step : update posterior of latent variables

$$q^{(t)}(z_{nc}) \propto \alpha_c^{(t-1)} p(y_n | \phi_c^{(t-1)}, x_n) \exp\left(\frac{-\mathcal{D}_c}{2\alpha_c^{(t-1)} N}\right)$$

- M step : update parameters

$$\alpha_c^{(t)} = \sum_{n=1}^N q^{(t)}(z_{nc}) / N$$

$$\phi_c^{(t)} = \arg \max_{\phi_c} \sum_{n=1}^N q^{(t)}(z_{nc}) \log p(y_n | \phi_c, x_n) - \frac{\mathcal{D}_c}{2} \log \left(\sum_{n=1}^N q^{(t)}(z_{nc}) \right)$$

Component regularization
is automatically determined

第3部の構成

- 異種混合学習とモデル選択
- 因子化情報量基準と因子化漸近ベイズ推論
- 因子化漸近ベイズ推論の性質
- 実験
- 発展モデルへ

Automatic Component Shrinkage

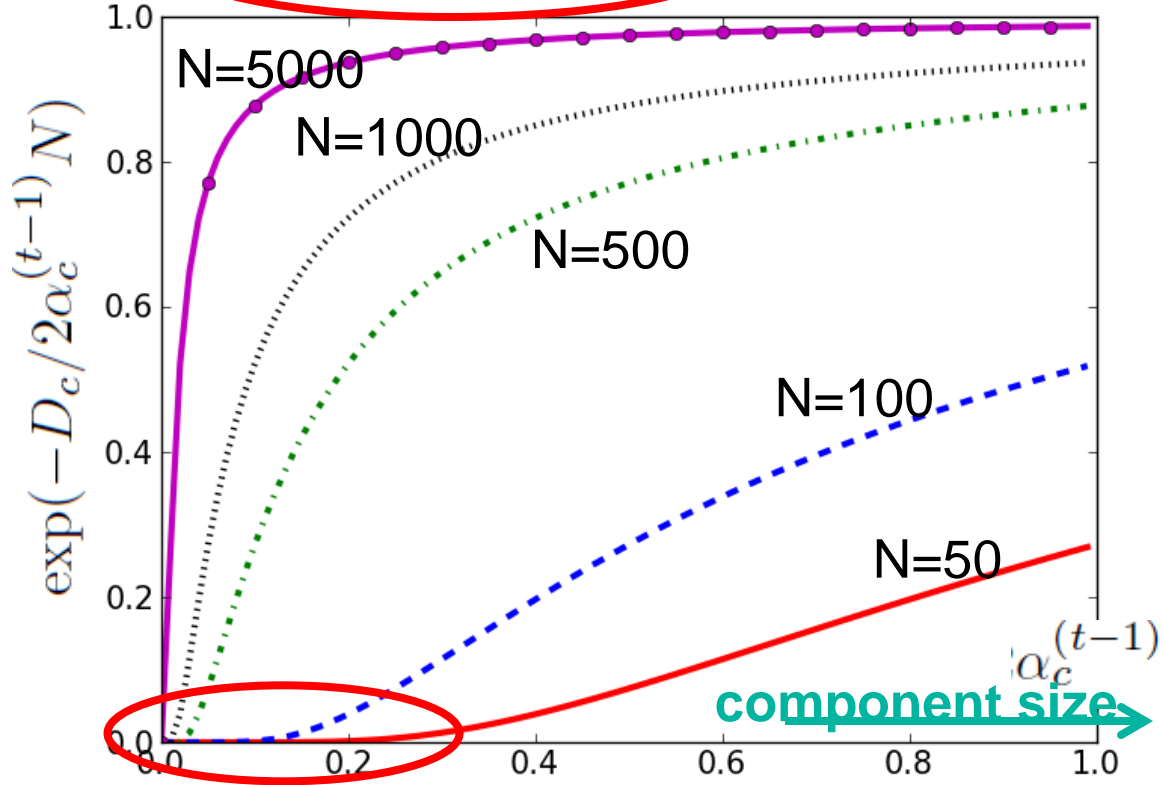
Mechanism

- FAB-unique exponentiated regularization in E-steps

$$q_{nc}^t \propto \alpha_c^{t-1} p(\mathbf{x}_n | \phi_c^{t-1}) \exp(-D_c / 2\alpha_c^{(t-1)} N)$$

Gaussian Mixture Models
D=10 (Dc = 65)

D: dimensionality of \mathbf{x}_n
Dc : dimensionality of $\phi_c^{(t-1)}$

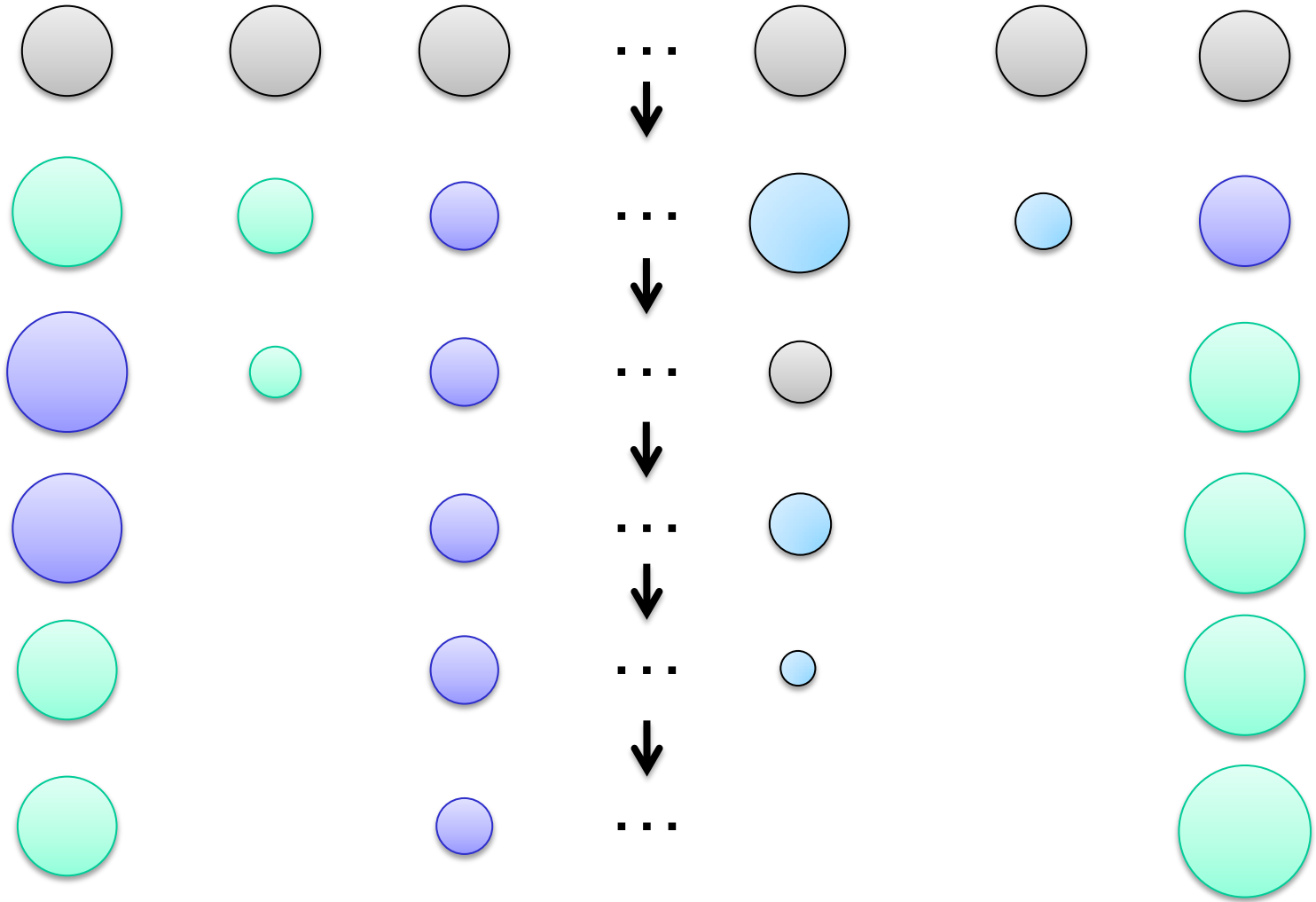


Small Components are shrunk through the EM iterations.

Automatically controls the number of latent variables (complexity)

How FAB works?

Random
initialization

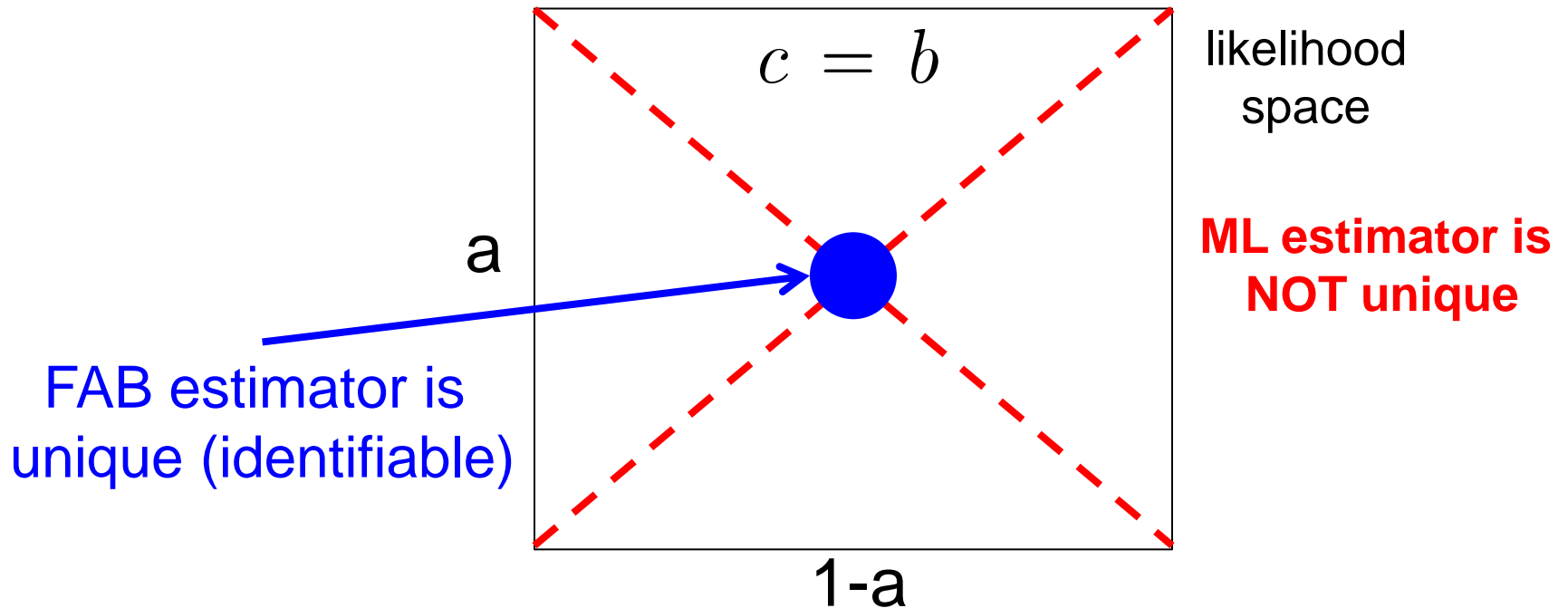


FAB EM-iteration

BIC and Difficulties of Model Selection for MMs

- Non-regularity Issue and Non-Identifiability issues


$$p(X|a, b, c) = ap(X|b) + (1 - a)p(X|c)$$



Identifiability Analysis on FAB



- Convergence point

$$\sum_{n=1}^N q^*(z_{nc}) = \alpha_c^* N = \sum_{n=1}^N \alpha_c^* p(\mathbf{x}_n | \phi_c^*) \exp\left(\frac{-\mathcal{D}_c}{2\alpha_c^* N}\right) / Z_n$$


$$N = \sum_{n=1}^N p(\mathbf{x}_n | \phi_c^*) \exp\left(\frac{-\mathcal{D}_c}{2\alpha_c^* N}\right) / Z_n$$

- What happen if $b = c$ $p(X|a, b, c) = ap(X|b) + (1 - a)p(X|c)$

$$\sum_{n=1}^N p(\mathbf{x}_n | \phi_1^*) \exp\left(\frac{-\mathcal{D}_1}{2\alpha_1^* N}\right) / Z_n = \sum_{n=1}^N p(\mathbf{x}_n | \phi_2^*) \exp\left(\frac{-\mathcal{D}_2}{2\alpha_2^* N}\right) / Z_n$$



$$\exp\left(\frac{-\mathcal{D}_1}{2\alpha_1^* N}\right) = \exp\left(\frac{-\mathcal{D}_2}{2\alpha_2^* N}\right)$$
 
$$\alpha_1^* = \alpha_2^*$$

- Components having the same parameter have the same size
- The maximum entropy distribution (in terms of qZ) is selected in the equivalent class

Non-Identifiability (EM case)

- Convergence point (EM case)

$$\sum_{n=1}^N q^*(z_{nc}) = \alpha_c^* N = \sum_{n=1}^N \alpha_c^* p(\mathbf{x}_n | \phi_c^*) \quad \blacksquare / Z_n$$

 $N = \sum_{n=1}^N p(\mathbf{x}_n | \phi_c^*) \quad \blacksquare / Z_n$

- What happen if $b = c$ $p(X|a, b, c) = ap(X|b) + (1 - a)p(X|c)$

$$\sum_{n=1}^N p(\mathbf{x}_n | \phi_1^*) \quad \blacksquare / Z_n = \sum_{n=1}^N p(\mathbf{x}_n | \phi_2^*) \quad \blacksquare / Z_n$$

- No "identifiable" conclusion

Comparison with BIC (HMM)

- FAB lower bound

$$\begin{aligned}
 FIC_{LB}^{(i)}(\mathbf{x}^N, M) = & \sum_{n=1}^N \sum_{t=1}^{T_n} \log \zeta_n^{t(i)} + \underbrace{\sum_{n,t=1}^{N, T_n} \log \Delta^t - \frac{\mathcal{D}_\alpha}{2} \log N}_{\text{red line}} \\
 & - \underbrace{\sum_{k=1}^K \left(\frac{\mathcal{D}_{\beta_k}}{2} (\log(\sum_{n,t}^{N, T_n-1} \tilde{q}(z_{nk}^t)) - 1) + \frac{\mathcal{D}_{\phi_k}}{2} (\log(\sum_{n,t}^{N, T_n} \tilde{q}(z_{nk}^t)) - 1) \right)}_{\text{red line}}
 \end{aligned}$$

- BIC (no theoretical justification)

$$BIC = \sum_{n,t=1}^{N, T_n} \log \zeta_n^t(ML) - \underbrace{\frac{\mathcal{D}}{2} \log \sum_{n=1}^N T_n}_{\text{red line}}$$

- Stochastic complexities of HMMs are theoretically shown to be much smaller than the BIC's complexity term (Yamazaki and Watanabe, 2005)

Comparison with VB (HMM)

- FAB lower bound

$$\begin{aligned}
 FIC_{LB}^{(i)}(\mathbf{x}^N, M) &= \sum_{n=1}^N \sum_{t=1}^{T_n} \log \zeta_n^{t(i)} + \sum_{n,t=1}^{N,T_n} \log \Delta^t - \frac{\mathcal{D}_{\alpha}}{2} \log N \\
 &\quad - \sum_{k=1}^K \left(\frac{\mathcal{D}_{\beta_k}}{2} \left(\log \left(\sum_{n,t}^{N,T_n-1} \tilde{q}(z_{nk}^t) \right) - 1 \right) + \frac{\mathcal{D}_{\phi_k}}{2} \left(\log \left(\sum_{n,t}^{N,T_n} \tilde{q}(z_{nk}^t) \right) - 1 \right) \right)
 \end{aligned}$$

- VB lower bound

$$\begin{aligned}
 \mathcal{F}_{VB} &= \sum_{n,t=1}^{N,T_n} \log \zeta_n^t(VB) + \int d\alpha q(\alpha) \log \frac{p(\alpha)}{q(\alpha)} \\
 &\quad + \int d\beta q(\beta) \log \frac{p(\beta)}{q(\beta)} + \sum_{k=1}^K \int d\phi_k q(\phi_k) \log \frac{p(\phi_k)}{q(\phi_k)}
 \end{aligned}$$

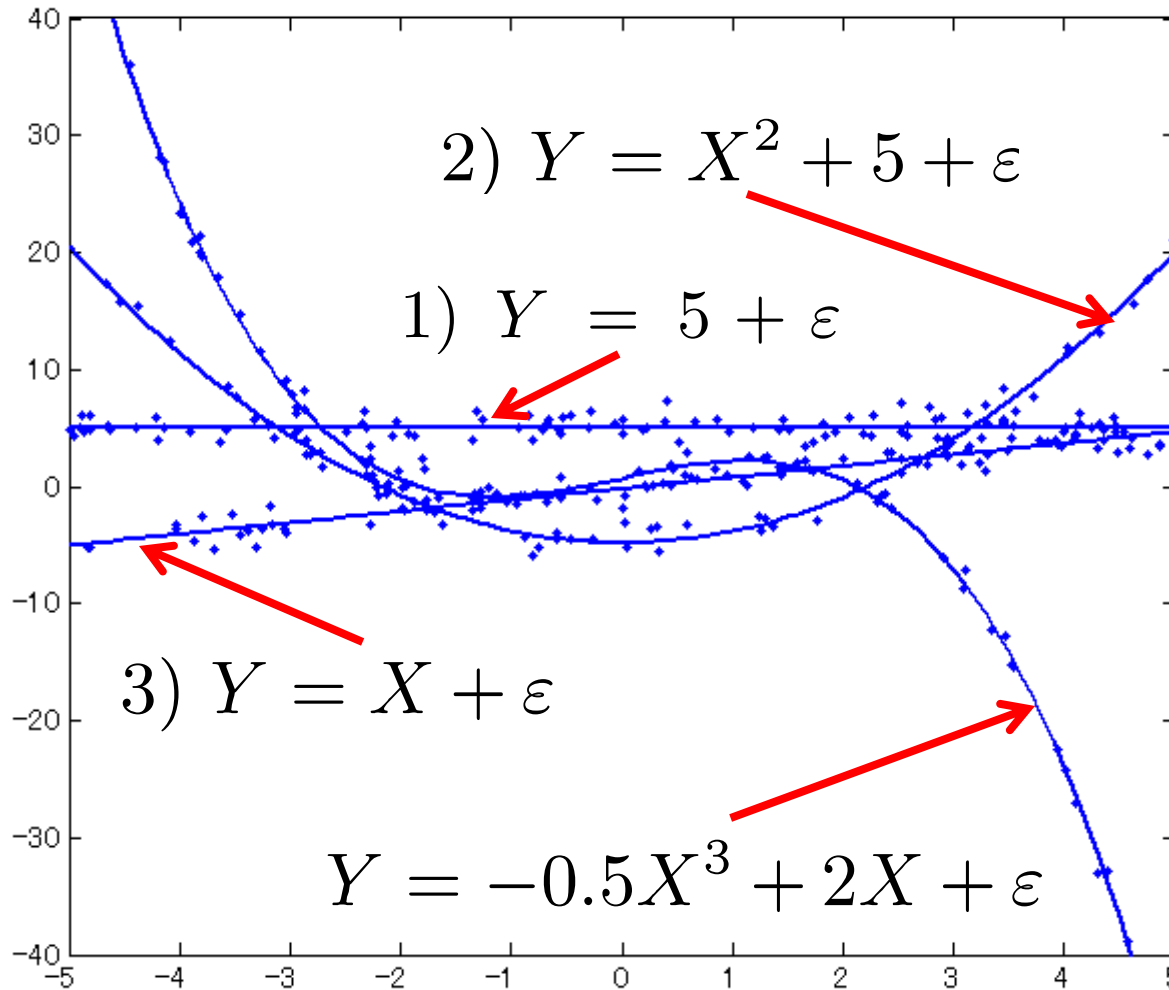
- Only FAB explicitly takes dependency between hidden states and model parameter into account.

第3部の構成

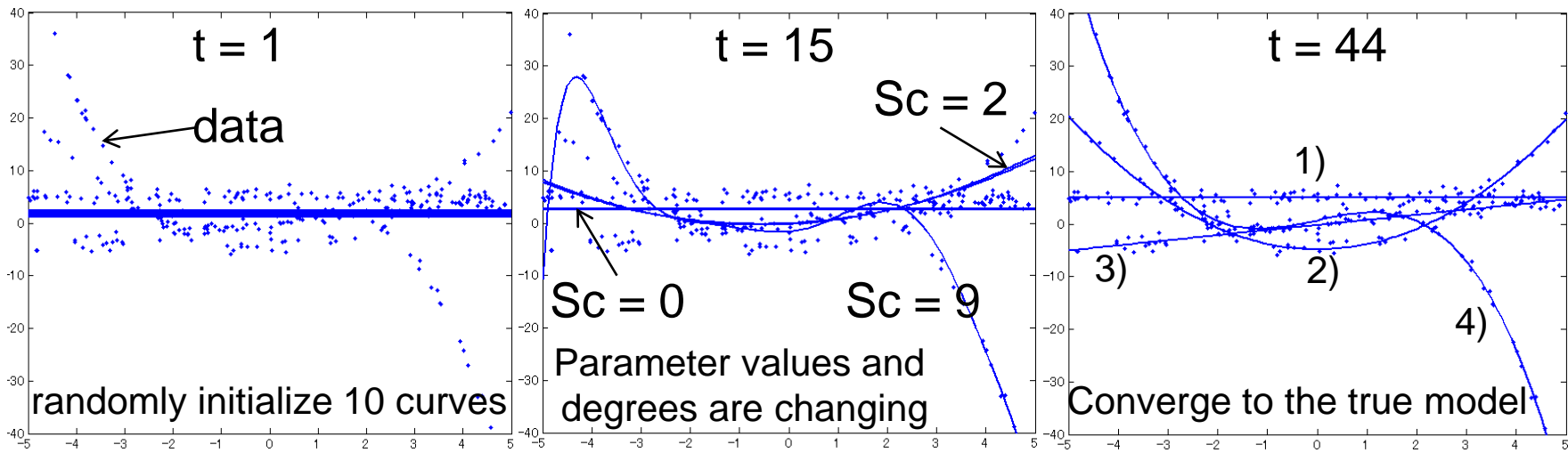
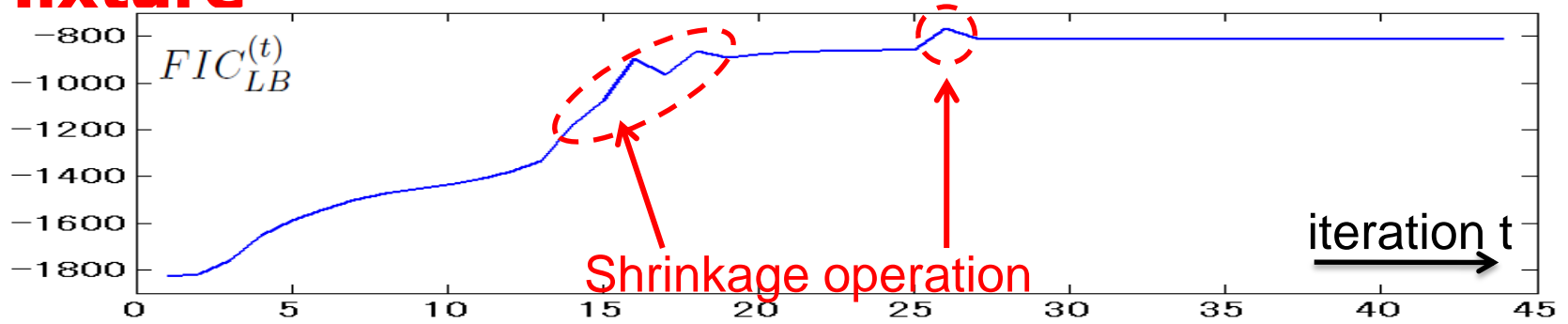
- 異種混合学習とモデル選択
- 因子化情報量基準と因子化漸近ベイズ推論
- 因子化漸近ベイズ推論の性質
- 実験
- 発展モデルへ

MM: Illustrative Polynomial Curve Mixture

- Artificially generate 4 curves



MM: Illustrative Polynomial Curve Mixture



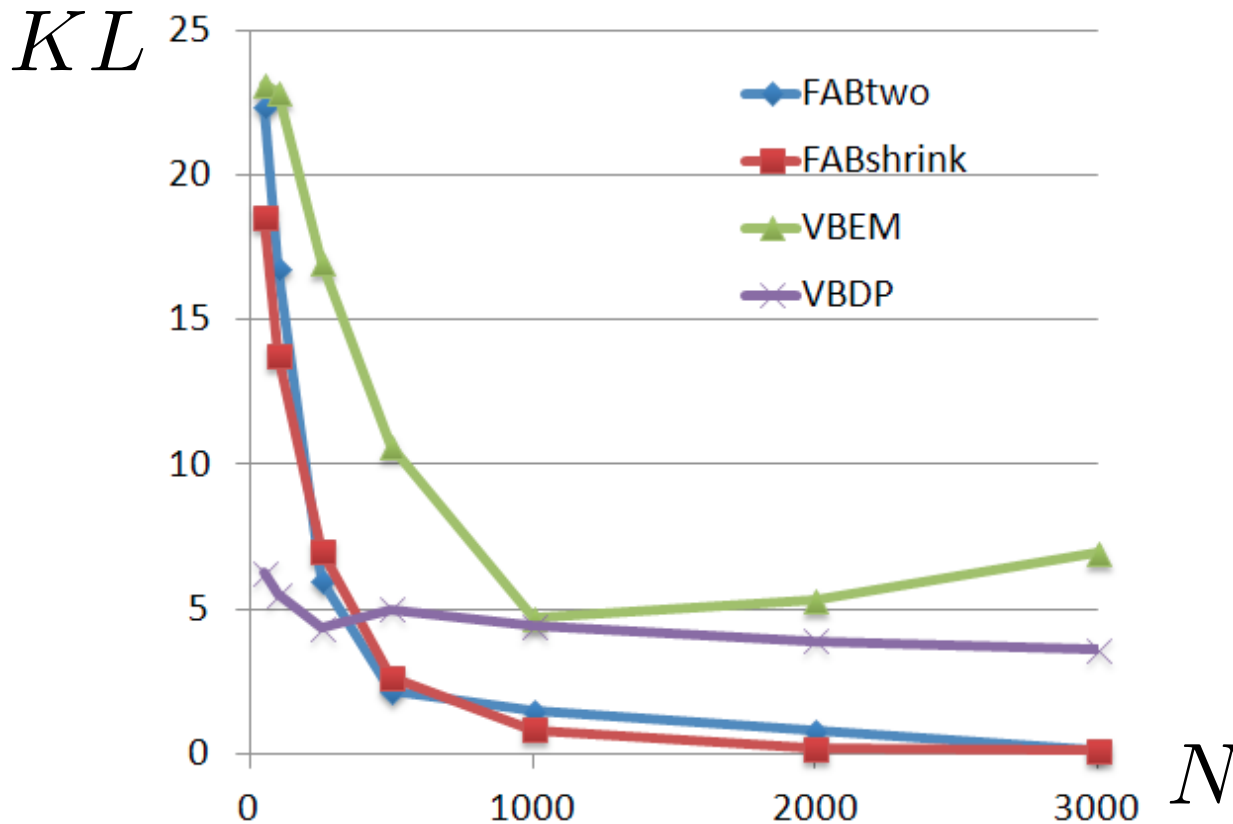
polynomial curve with Gaussian noise

$$p(\boxed{Y} | X, \boldsymbol{\theta}, \mathbf{S}) = \sum_{c=1}^C \alpha_c p(Y | X, \boxed{\phi_c}, \boxed{S_c})$$

target variable weight vector and noise level degree of curve

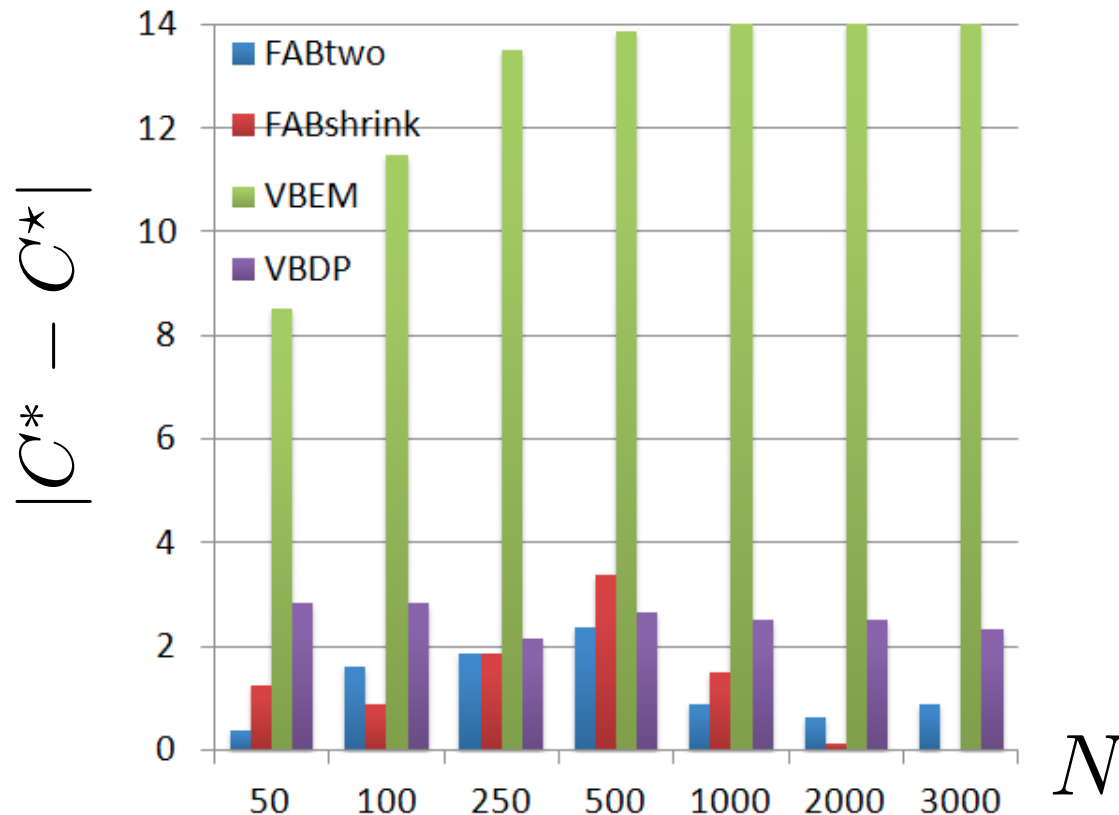
MM: Artificial GMM

- Comparison with VB and Non-parametric Bayes
 - VBEM (Bishop 2005): VBEM algorithm with Dirichlet prior (model selection via variational free energy)
 - VBDP (Blei and Jordan 2006): VB Dirichlet process mixture model (model selection via Dirichlet process prior)



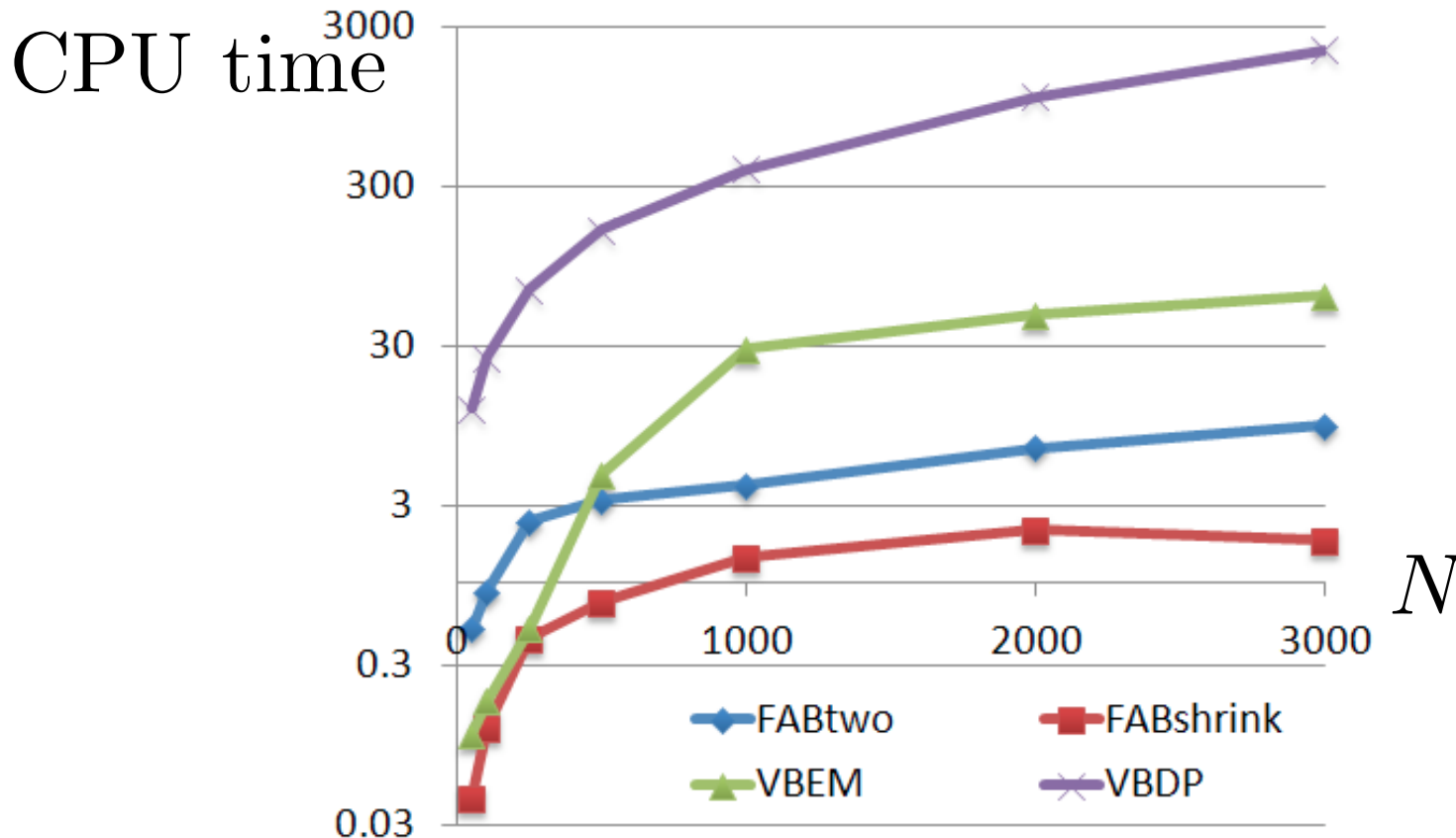
MM: Artificial GMM

- Comparison with VB and Non-parametric Bayes
 - VBEM (Bishop 2005): VBEM algorithm with Dirichlet prior (model selection via variational free energy)
 - VBDP (Blei and Jordan 2006): VB Dirichlet process mixture model (model selection via Dirichlet process prior)



MM: Artificial GMM

- Comparison with VB and Non-parametric Bayes
 - VBEM (Bishop 2005): VBEM algorithm with Dirichlet prior (model selection via variational free energy)
 - VBDP (Blei and Jordan 2006): VB Dirichlet process mixture model (model selection via Dirichlet process prior)



HMMs: E-Book Dataset

- Six books in different topics (<http://www.gutenberg.org>)
 - Alice's Adventures in Wonderland
 - The Art of War
 - The Metamorphosis
 - The Republic
 - The United States Declaration of Independence
 - The Adventures of Sherlock Holmes
- The first 5000 characters for training and the following 5000 characters for testing.
- Comparison methods
 - VB: variational Bayesian HMMs (Mackay 1997, Beal 2003)
 - iHMM: infinite HMMs with beam sampler (van Gael et al. 2008)

E-book Dataset

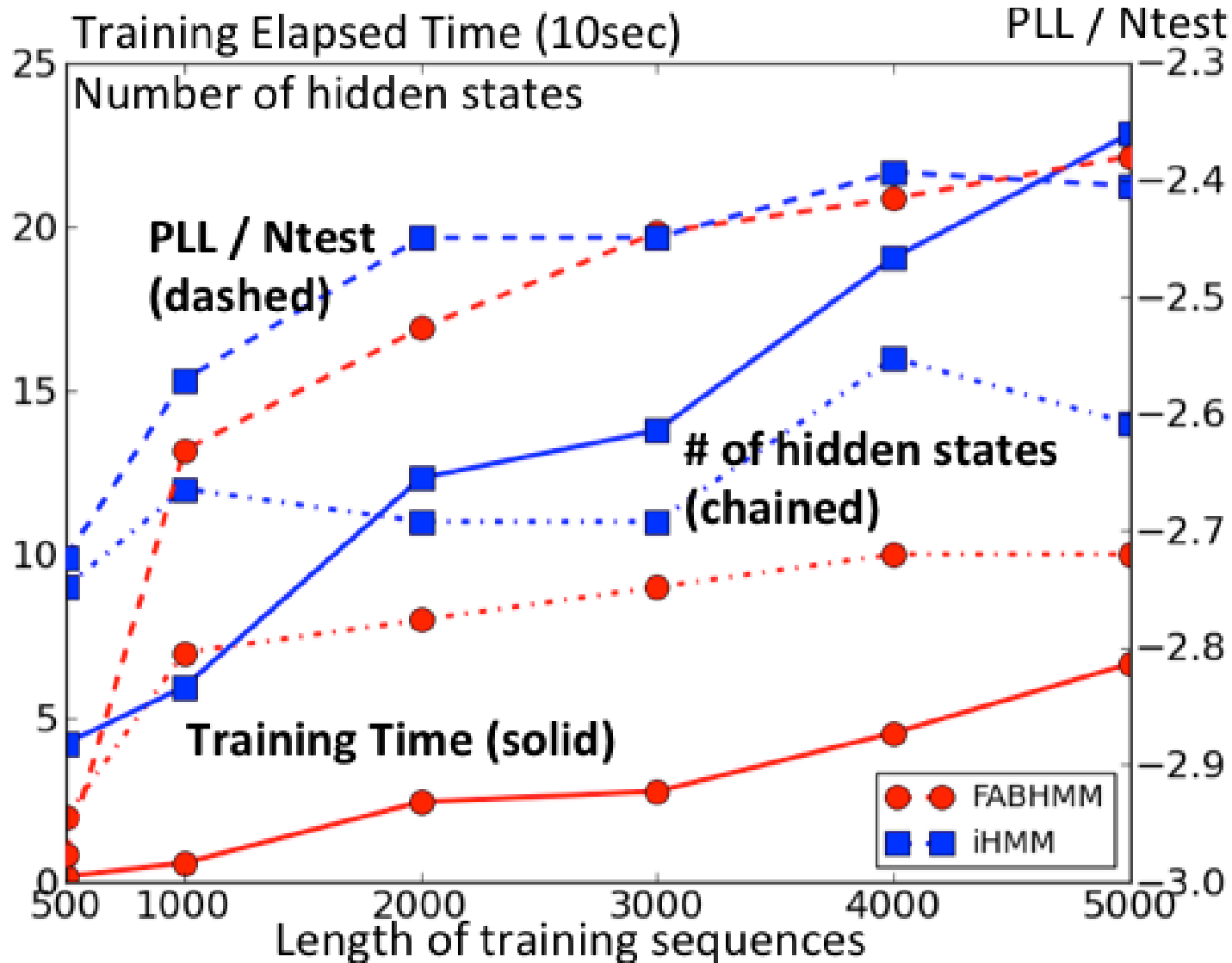
- Predictive Log-likelihood and Elapsed Time

Table 2. Estimated number of hidden states K , training time (sec), and PLLs on the ebook data sets.

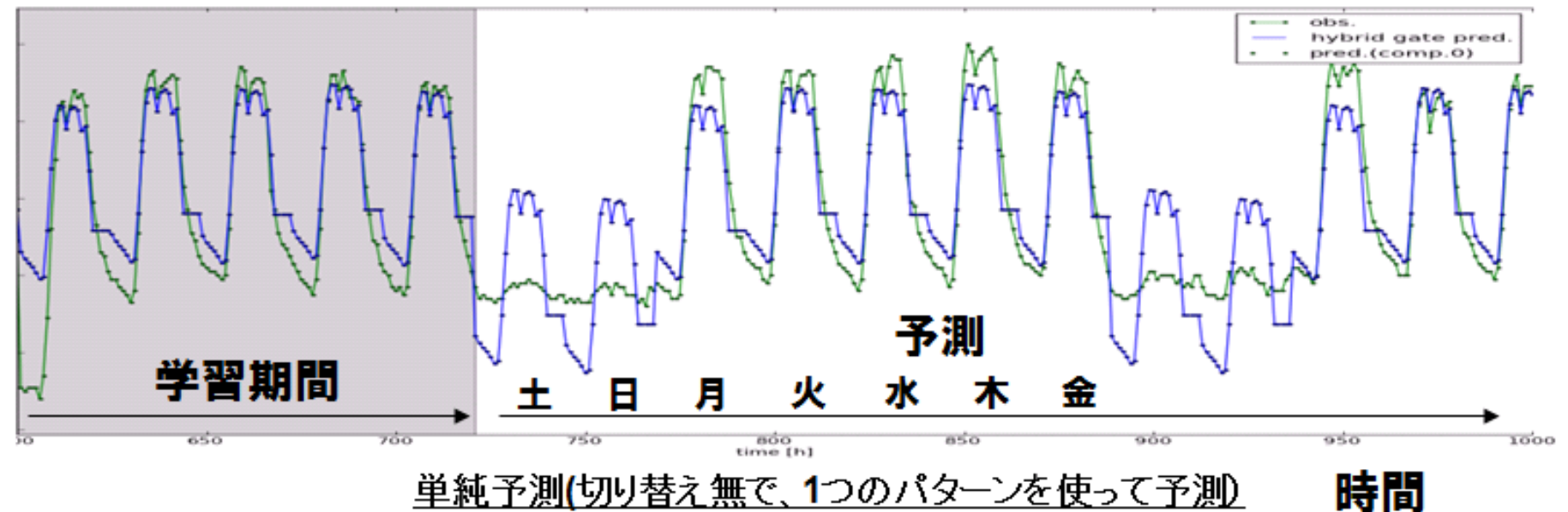
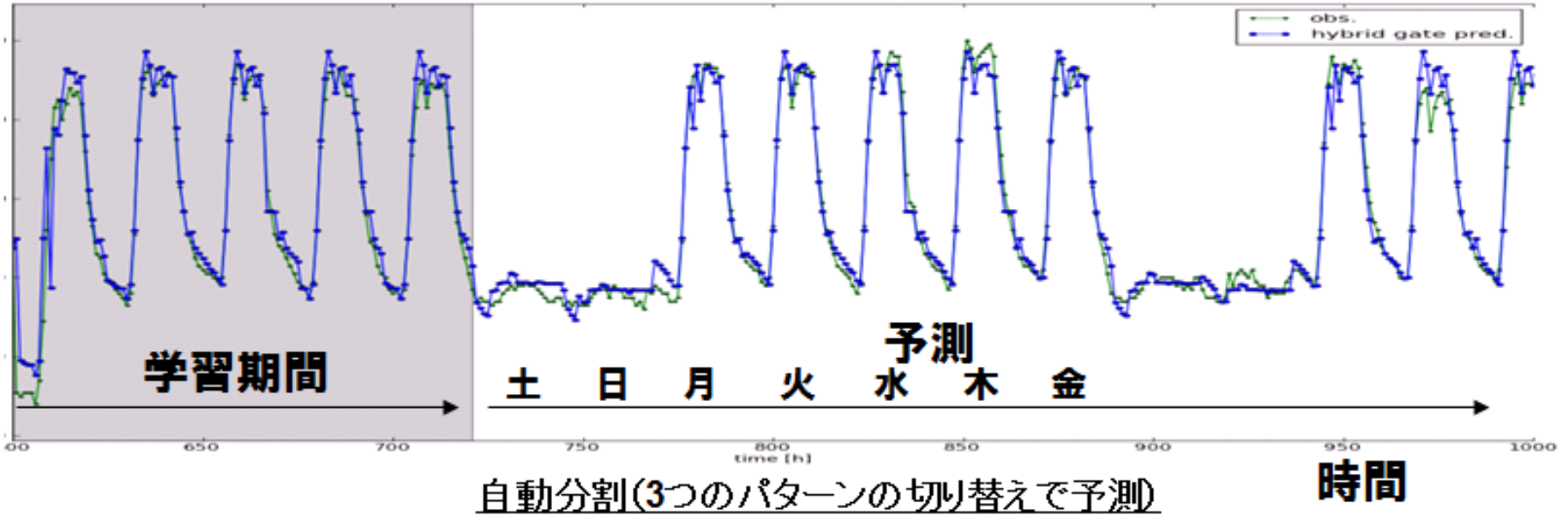
DATA	FABHMM			iHMM			VBHMM		
	K	TIME	PLL	K	TIME	PLL	K	TIME	PLL
ALICE	9	64.2	-2.57	13	234	-2.54	15	122.3	-2.73
KAFKA	10	47.0	-2.36	11	218	-2.40	12	104.2	-2.62
PLATO	10	66.6	-2.38	14	228	-2.41	8	144.9	-2.63
SHERL	11	72.3	-2.58	12	227	-2.52	19	98.0	-2.75
SUNZI	10	63.2	-2.56	14	228	-2.52	14	110.8	-2.72
DOI	10	73.8	-2.98	12	232	-2.75	11	94.6	-2.76

E-book Dataset

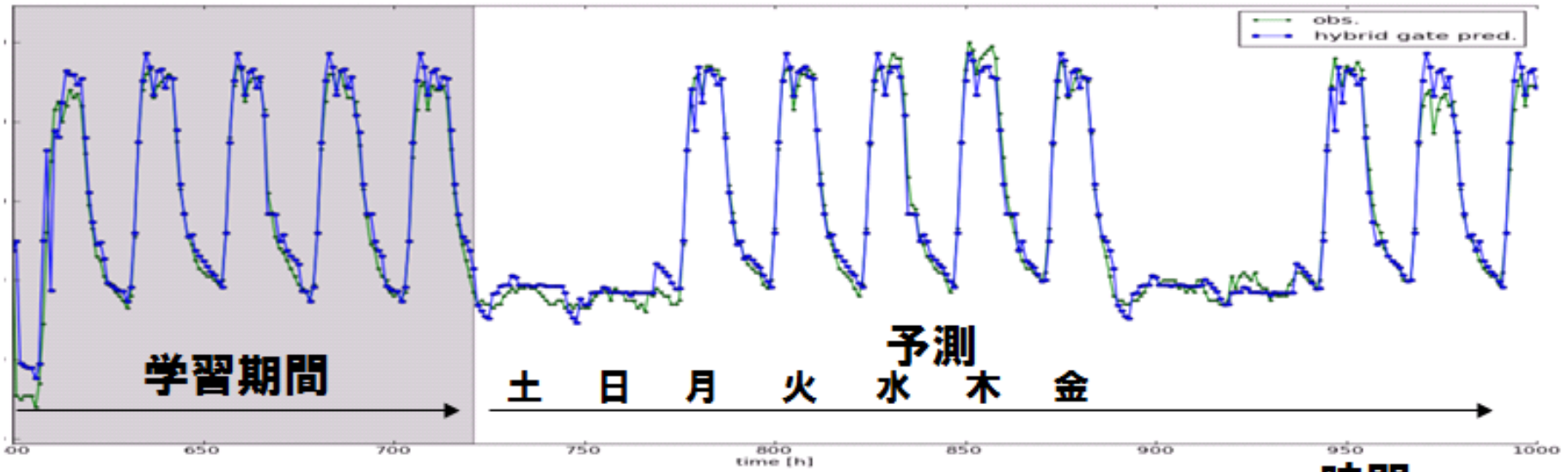
- Predictive Log-likelihood and Elapsed Time



Application to Building Power Prediction



Application to Building Power Prediction



自動分割(3つのパターンの切り替えで予測)

