

P5-31 Contrastive Divergence Learning に対する新しい解釈とその理論解析

前田 新一[†], 石井 信[‡]
^{†‡}京大 情報学

[†]ichi@sys.i.kyoto-u.ac.jp, [‡]ishii@i.kyoto-u.ac.jp

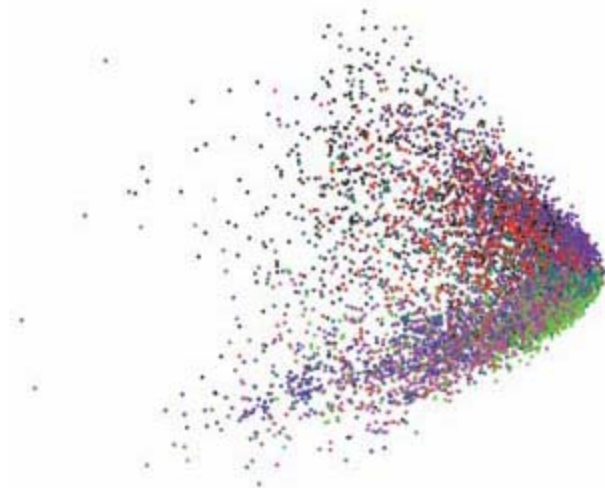
関連キーワード

ボルツマンマシン, Contrastive Divergence Learning,
詳細釣り合い条件

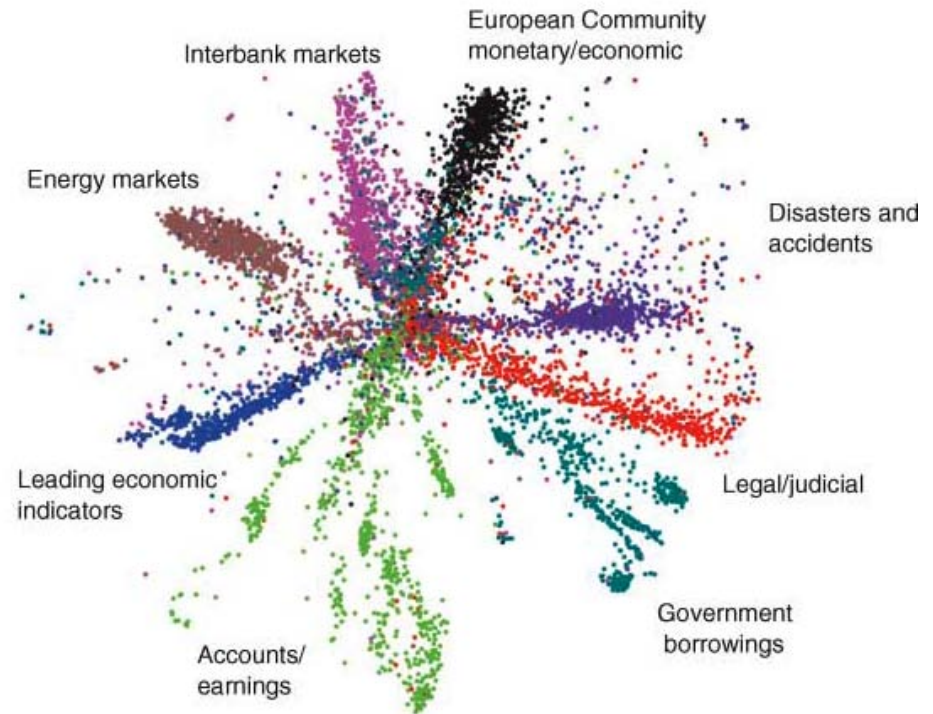
動機

なぜContrastive Divergence Learningだとうまくいく？

Hintonら(2006)は高次元パラメータをもつ階層ボルツマンマシンを大規模データを用いて学習させることに成功した。
この学習アルゴリズムにContrastive Divergence Learningが用いられており、Hintonらはこれが重要なファクターであることを示唆したが、このContrastive Divergence Learningの理論背景は明確ではなかった



参考: 従来法(Latent Semantec Analysis)による結果



学習データセット

Reuters Corpus Vol. II
2000語 103のトピック
302,207個の記事

テストデータセット

Reuters Corpus Vol. II
2000語 103のトピック
100,000個の記事

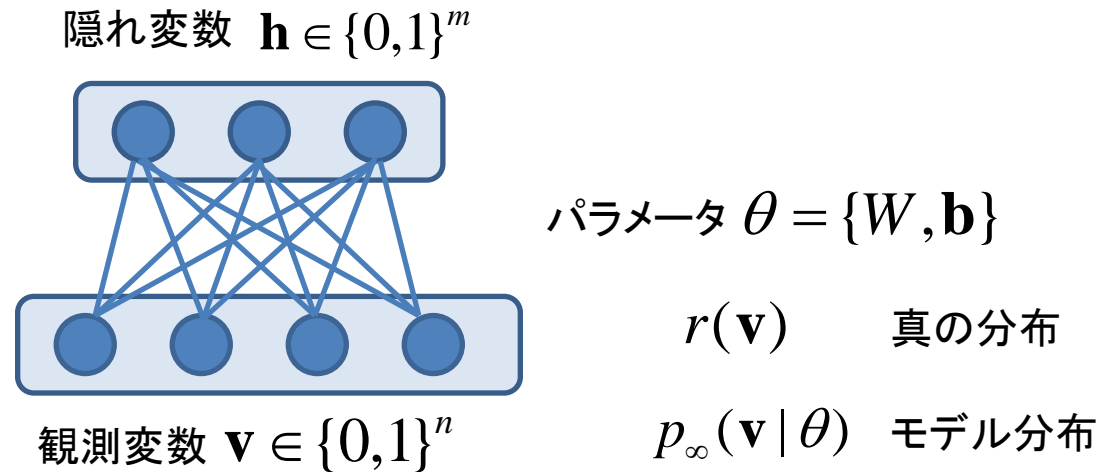
ネットワーク構成

encoder: 2000-500-250-125-2 (=約116万個のパラメータ数)
decoder: encoderと対称に配置

(Hinton and Salakhutdinov, Science, 2006 より引用)

Restricted Boltzmann Machine (RBM)

RBM = 層内に結合をもたない二層のボルツマンマシン



RBMの尤度

$$p_{\infty}(\mathbf{v} | \theta) = \frac{1}{Z(\theta)} \sum_{\mathbf{h}} \exp\left(-\sum_{ij} w_{ij} h_i v_j - \sum_j b_j v_j - \sum_i b_i h_i\right)$$

ただし、 $Z(\theta)$ は正規化定数

←評価が困難

RBMの条件付き分布

$$p(\mathbf{h} | \mathbf{v}, \theta) = \prod_{i=1}^m p(h_i | \mathbf{v}, \theta) \quad p(h_i = 1 | \mathbf{v}, \theta) = \frac{1}{1 + \exp\left(\sum_j w_{ij} v_j + b_i\right)} \equiv g\left(\sum_j w_{ij} v_j + b_i\right)$$
$$p(\mathbf{v} | \mathbf{h}, \theta) = \prod_{j=1}^n p(v_j | \mathbf{h}, \theta) \quad p(v_j = 1 | \mathbf{h}, \theta) = \frac{1}{1 + \exp\left(\sum_i w_{ij} h_i + b_j\right)} \equiv g\left(\sum_i w_{ij} h_i + b_j\right)$$

←評価が容易

Contrastive Divergence Learning (CDL)

(Hinton, 2002)

CDLのコスト関数(とされるもの)

$$CD(\theta) = KL[Q_0(\mathbf{v}) | Q_\infty(\mathbf{v} | \theta)] - KL[Q_1(\mathbf{v}) | Q_\infty(\mathbf{v} | \theta)]$$

$$\text{ただし、} Q_t(\mathbf{v} | \theta) \equiv \sum_{\mathbf{v}'} p_G(\mathbf{v} | \mathbf{v}', \theta) Q_{t-1}(\mathbf{v}' | \theta)$$

$$p_G(\mathbf{v} | \mathbf{v}', \theta) \equiv \sum_{\mathbf{h}} p(\mathbf{v} | \mathbf{h}', \theta) p(\mathbf{h}' | \mathbf{v}', \theta) \quad Q_0(\mathbf{v} | \theta) \equiv r(\mathbf{v}) \text{ (真の分布)}$$

確率勾配法による学習アルゴリズム \doteq CDL

$$\theta_{t+1} = \theta_t + \Delta \theta$$

$$\Delta \theta \propto \frac{\partial}{\partial \theta} CD(\theta)$$

$\theta = w_{ij}$ としたとき

$$\Delta w_{ij} \propto \left\langle h_i v_j \right\rangle_{r(\mathbf{v}) p(\mathbf{h} | \mathbf{v}, \theta)} - \left\langle h_i v_j \right\rangle_{Q_1(\mathbf{v} | \theta) p(\mathbf{h} | \mathbf{v}, \theta)} - \frac{\partial H(Q_1(\mathbf{v} | \theta_t))}{\partial w_{ij}}$$

\uparrow 無視

$$\text{ここで } H(p(\mathbf{v})) \equiv -\sum_{\mathbf{v}} p(\mathbf{v}) \log p(\mathbf{v})$$

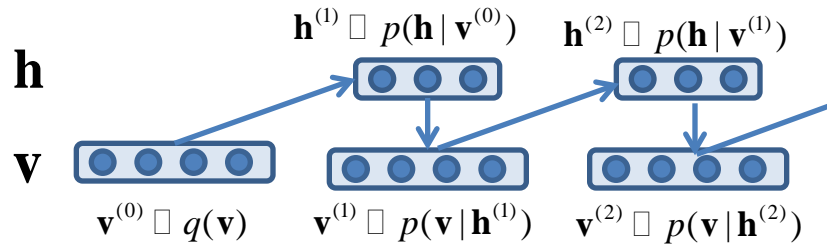
Gibbs Sampling とCDLの関係

尤度のパラメータ W による微分も計算困難

$$\Delta w_{ij} \propto \frac{\partial}{\partial w_{ij}} \sum_{\mathbf{v}} r(\mathbf{v}) \log p_{\infty}(\mathbf{v} | \theta) = \left\langle h_i v_j \right\rangle_{r(\mathbf{v}) p(h|\mathbf{v}, \theta)} - \left\langle h_i v_j \right\rangle_{p_{\infty}(\mathbf{v} | \theta) p(h|\mathbf{v}, \theta)} \quad \leftarrow \text{この項が問題}$$

$r(\mathbf{v})$: 真の分布

Gibbs Sampling



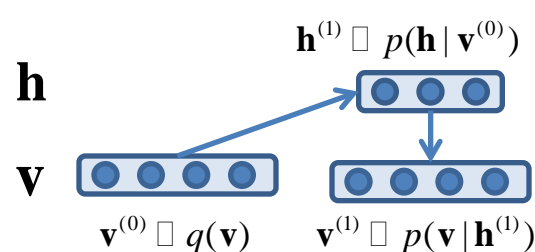
$$\mathbf{v}^{(\infty)} \sqsubset \lim_{t \rightarrow \infty} Q_t(\mathbf{v} | \theta) = p_{\infty}(\mathbf{v} | \theta)$$

$$\text{ただし、} Q_t(\mathbf{v} | \theta) \equiv \sum_{\mathbf{v}'} p_G(\mathbf{v} | \mathbf{v}', \theta) Q_{t-1}(\mathbf{v}' | \theta)$$

$$p_G(\mathbf{v} | \mathbf{v}', \theta) \equiv \sum_{\mathbf{h}} p(\mathbf{v} | \mathbf{h}', \theta) p(\mathbf{h}' | \mathbf{v}', \theta) \quad Q_0(\mathbf{v} | \theta) \equiv r(\mathbf{v})$$

マルコフ連鎖が平衡分布に
落ち着くまで時間がかかる

Contrastive Divergence Learning



$$\mathbf{v}^{(1)} \sqsubset Q_1(\mathbf{v} | \theta) \quad (\sqsubset p_{\infty}(\mathbf{v} | \theta)?)$$

$\lim_{t \rightarrow \infty} Q_t(\mathbf{v} | \theta)$ の代わりに $Q_1(\mathbf{v} | \theta)$ を用いる

早い
が、何を最適化したことになる？

Detailed Balance Learning (DBL)

$p(\mathbf{v} | \mathbf{v}', \theta)$: マルコフ連鎖

$p_{\infty}(\mathbf{v} | \theta)$: マルコフ連鎖の定常分布 (唯一の定常分布をもつと仮定)

コスト関数

隠れ変数なし $F(\theta, \bar{\theta}) = KL \left[p(\mathbf{v}' | \mathbf{v}, \bar{\theta}) r(\mathbf{v}) | p(\mathbf{v} | \mathbf{v}', \theta) r(\mathbf{v}') \right]$

隠れ変数あり $F(\theta, \bar{\theta}) = KL \left[p(\mathbf{v}', \mathbf{h} | \mathbf{v}, \bar{\theta}) r(\mathbf{v}) | p(\mathbf{v}, \mathbf{h} | \mathbf{v}', \theta) r(\mathbf{v}') \right]$

コスト関数の特徴

1. $F(\theta, \bar{\theta}) \geq 0$
2. $F(\theta, \theta) = 0$
となるのは
 $r(\mathbf{v}) = p_{\infty}(\mathbf{v} | \theta)$
のときのみ

コスト関数として適切な特徴をもつが、真の分布を含み直接、評価することは不可能

詳細釣り合い条件(Detailed Balance Condition)

隠れ変数なし

For any $\mathbf{v}, \mathbf{v}' \in S$

$$p(\mathbf{v}' | \mathbf{v}, \theta) r(\mathbf{v}) = p(\mathbf{v} | \mathbf{v}', \theta) r(\mathbf{v}')$$



$$r(\mathbf{v}) = p_{\infty}(\mathbf{v} | \theta)$$

隠れ変数あり

For any $\mathbf{v}, \mathbf{v}' \in \{0,1\}^n, \mathbf{h}' \in \{0,1\}^m$

$$p(\mathbf{v}', \mathbf{h} | \mathbf{v}, \theta) r(\mathbf{v}) = p(\mathbf{v}, \mathbf{h} | \mathbf{v}', \theta) r(\mathbf{v}')$$



$$r(\mathbf{v}) = p_{\infty}(\mathbf{v} | \theta)$$

Detailed Balance Learning (DBL)

- アルゴリズム -

DBLアルゴリズム

1. $t=1$ として θ_0, θ_1 に初期値を設定する.
2. $F(\theta_{t+1}, \theta_t) < F(\theta_t, \theta_t)$ を満たす θ_{t+1} を求める.
3. 終了条件を満たしたならばアルゴリズムを終了し, そうでないならば t を1増やし, ステップ2に進む.

$$F(\theta, \bar{\theta}) = \sum_{\mathbf{v}', \mathbf{v}} p(\mathbf{v}' | \mathbf{v}, \bar{\theta}) r(\mathbf{v}) \log p(\mathbf{v} | \mathbf{v}', \theta) + \text{const}$$

(ただし、 const はパラメータに依存しない定数)

より、 $F(\theta, \bar{\theta})$ は真の分布からのサンプルを用いたサンプル平均によって推定可能



ステップ2を解くためのコスト関数の(近似)評価は可能。


$p(\mathbf{v} | \mathbf{v}', \theta)$ がパラメータ微分可能であれば準ニュートン法などの最適化手法を適用可能

Detailed Balance Learning (DBL)

- アルゴリズム for RBM -

DBLアルゴリズム

1. $t=1$ として θ_0, θ_1 に初期値を設定する.
2. $F(\theta_{t+1}, \theta_t) < F(\theta_t, \theta_t)$ を満たす θ_{t+1} を求める.
3. 終了条件を満たしたならばアルゴリズムを終了し、そうでないならば t を1増やし、ステップ2に進む.


$$F(\theta, \bar{\theta}) = \sum_{\mathbf{v}', \mathbf{v}} p_G(\mathbf{v}' | \mathbf{v}, \bar{\theta}) r(\mathbf{v}) \log p(\mathbf{v} | \mathbf{v}', \theta) + \text{const}$$

確率勾配法を用いると、

$$w_{ij}^{(t+1)} = w_{ij}^{(t)} + \Delta w_{ij}$$

$$\Delta w_{ij} \propto \frac{\partial F(\theta, \theta_t)}{\partial w_{ij}} = \left\langle h_i v_j \right\rangle_{r(\mathbf{v}) p(\mathbf{h} | \mathbf{v}, \theta)} - \left\langle h_i v_j \right\rangle_{p_\infty(\mathbf{v} | \theta) p(\mathbf{h} | \mathbf{v}, \theta)}$$

CDLの学習則に一致

DBLの収束性

定理1

θ_t がDBLアルゴリズムに従って更新されたとする。

このとき、以下は θ_t が収束するための十分条件となる

$$F(\theta_{t+1}, \theta_t) < F(\theta_{t-1}, \theta_t) \text{ であるとき、}$$

$$F(\theta_t, \theta_{t+1}) < F(\theta_t, \theta_{t-1}) \text{ が成り立つ。}$$

上記の十分条件は、以下のKL擬距離の差分の対称性が成り立てば満たされる。

$$KL[q_t | p_{t+1}] - KL[q_t | p_t] > 0 \quad \Rightarrow \quad KL[p_{t+1} | q_t] - KL[p_t | q_t] > 0$$

例：分布 p_t, p_{t+1}, q_t が異なる平均、等しい共分散行列を持つガウス分布のときこの対称性に関する条件は常に満たされる。

DBLとCDLの関係

定理2

マルコフ連鎖 $p(\mathbf{v}|\mathbf{v}',\theta)$ の定常分布 $p_\infty(\mathbf{v}|\theta)$ が
詳細釣り合い条件を満たすとき、以下が成り立つ。

$$F(\theta, \theta) = CD(\theta) + KL[Q_1(\mathbf{v}|\theta) | r(\mathbf{v})]$$

系

$CD(\theta) \geq 0, KL[Q_1(\mathbf{v}|\theta) | r(\mathbf{v})] \geq 0$ より

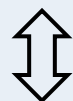
$$F(\theta, \theta) \leq CD(\theta), F(\theta, \theta) \leq KL[Q_1(\mathbf{v}|\theta) | r(\mathbf{v})]$$

詳細釣り合い条件の成立条件

定理3

$p(\mathbf{v}, \mathbf{h} | \theta)$ の条件付き分布 $p(\mathbf{v} | \mathbf{h}, \theta)$ と $p(\mathbf{h} | \mathbf{v}, \theta)$ から
構成されるマルコフ連鎖 $p_G(\mathbf{v} | \mathbf{v}', \theta) = \sum_{\mathbf{h}} p(\mathbf{v} | \mathbf{h}, \theta) p(\mathbf{h} | \mathbf{v}', \theta)$
が唯一の定常分布をもつとき、以下が成り立つ。

$$\text{For any } \mathbf{v}, r(\mathbf{v}) = \sum_{\mathbf{h}} p(\mathbf{v}, \mathbf{h} | \theta)$$



For any $\mathbf{v}, \mathbf{v}', \mathbf{h}$

$$p(\mathbf{v}' | \mathbf{h}, \theta) p(\mathbf{h} | \mathbf{v}, \theta) r(\mathbf{v}) = p(\mathbf{v} | \mathbf{h}, \theta) p(\mathbf{h} | \mathbf{v}', \theta) r(\mathbf{v}')$$

つまり、モデル分布が真の分布を含むならば
真の分布で詳細釣り合い条件を成り立たせるパラメータ θ が存在する

計算機実験

モデル分布が真の分布を**表現可能**な場合

真の分布：RBM（観測変数4次元, 隠れ変数**3**次元）（パラメータは乱数で決定）

モデル分布：RBM（観測変数4次元, 隠れ変数**3**次元）（初期パラメータは乱数で決定）

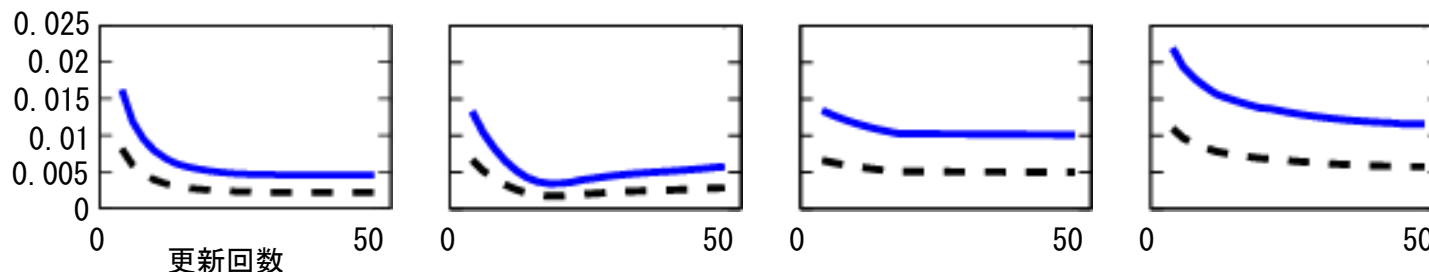
サンプル数：1000点

学習則：DBL（ステップ2の実行には準ニュートン法を用いた）

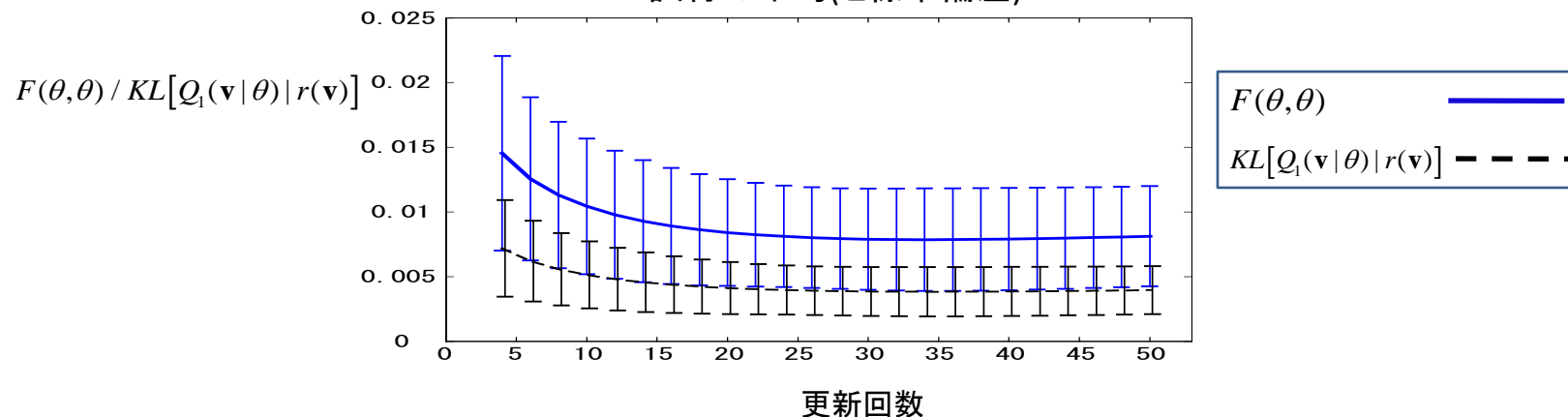
× 50試行

ランダムに選んだ4試行の結果

$F(\theta, \theta) / KL[Q_1(\mathbf{v} | \theta) | r(\mathbf{v})]$



50試行の平均(と標準偏差)

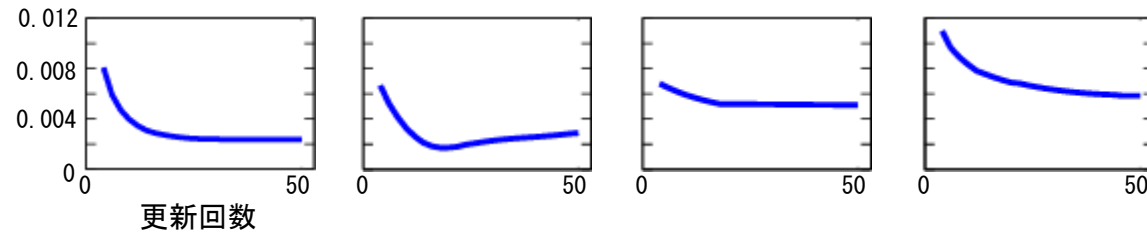


計算機実験

モデル分布が真の分布を**表現可能**な場合

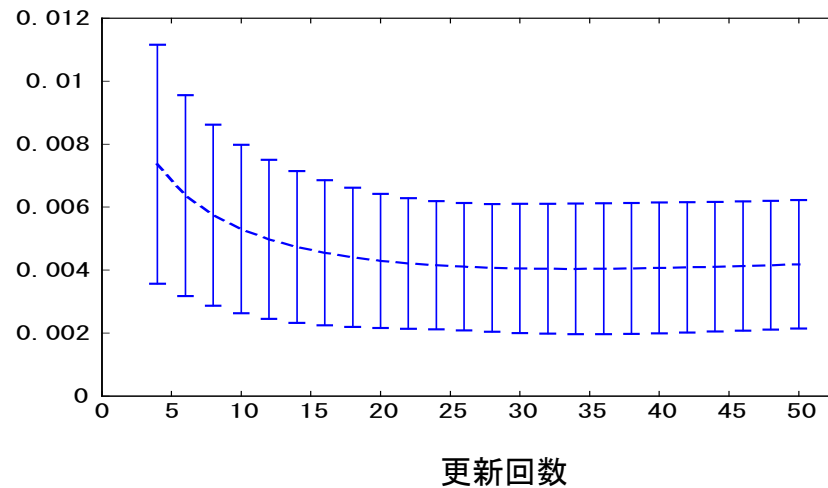
ランダムに選んだ4試行の結果

$$KL[r(\mathbf{v}) | p_{\infty}(\mathbf{v} | \theta)]$$



50試行の平均(と標準偏差)

$$KL[r(\mathbf{v}) | p_{\infty}(\mathbf{v} | \theta)]$$



計算機実験

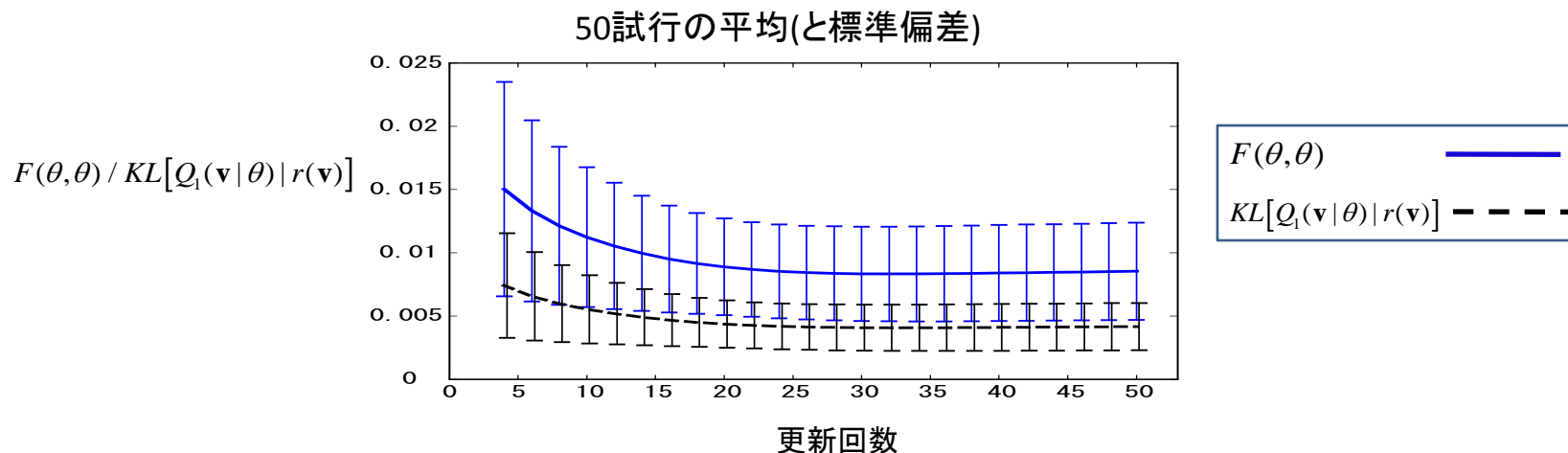
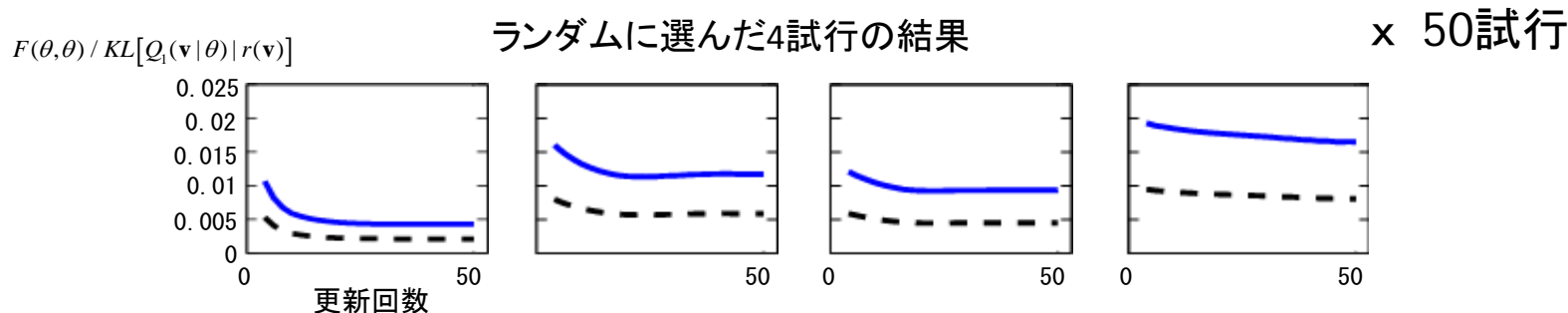
モデル分布が真の分布を**表現不可能**な場合

真の分布：RBM（観測変数4次元, 隠れ変数**5**次元）（パラメータは乱数で決定）

モデル分布：RBM（観測変数4次元, 隠れ変数**3**次元）（初期パラメータは乱数で決定）

サンプル数：1000点

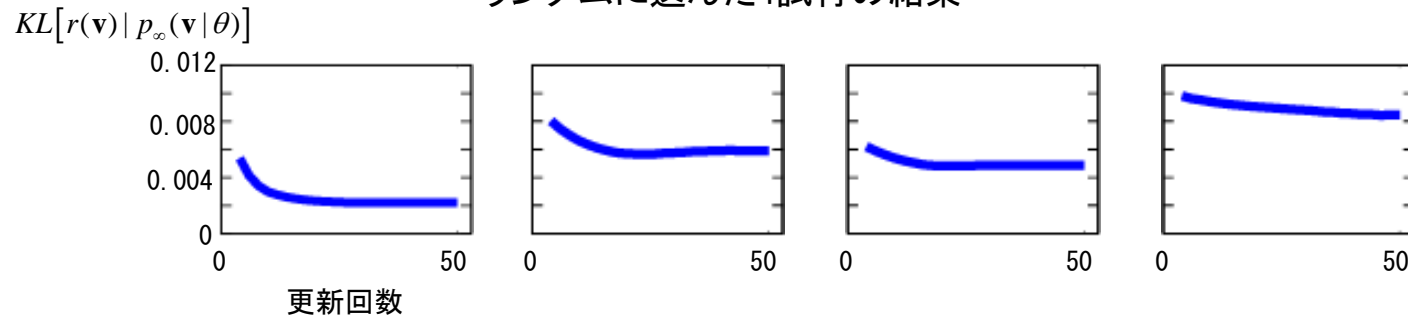
学習則：DBL（ステップ2の実行には準ニュートン法を用いた）



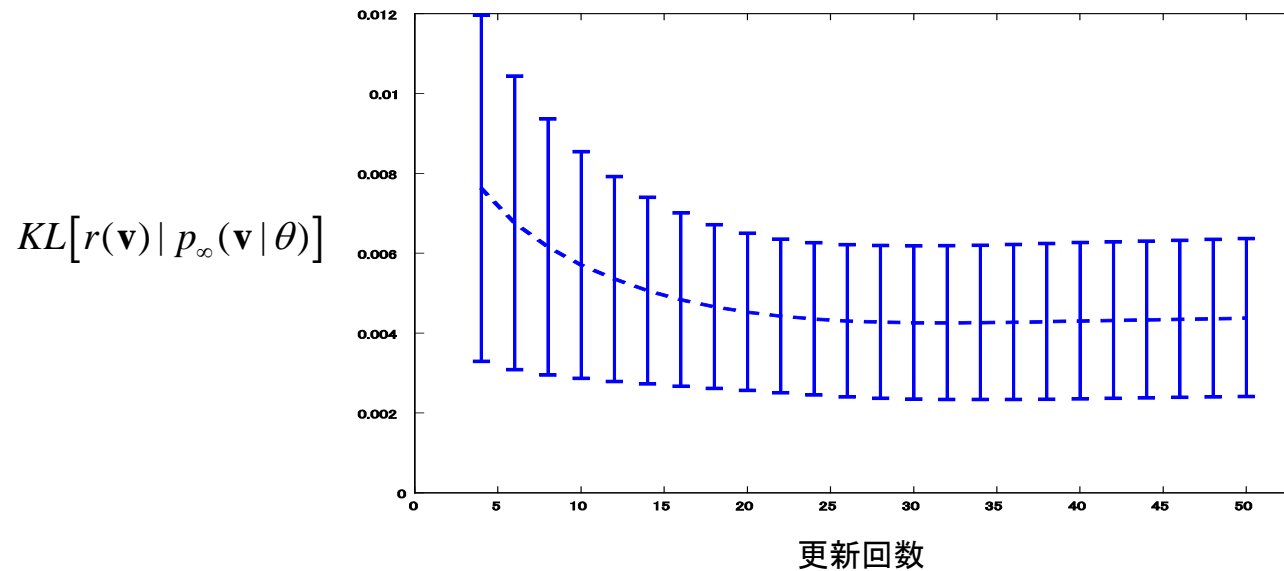
計算機実験

モデル分布が真の分布を表現不可能な場合

ランダムに選んだ4試行の結果



50試行の平均(と標準偏差)



まとめ

1. マルコフ連鎖の定常分布の 新しい学習法(Detailed Balance Learning)を提案した。

(ボルツマンマシンを含むギプスサンプリングの定常分布
としてに表現されるモデル分布の学習に幅広く適用可能)

2. Detailed Balance Learningによって Contrastive Divergence Learningを説明 (定理2)した。

(一般に用いられるCDLの学習則は、DBLを確率勾配法で解いたもの)

(CDLのコスト関数が明らかになったことにより、より高速な学習アルゴリズムの構築可能?)

3. Detailed Balance Learningの収束条件(定理1)、 最適化後の分布が真の分布と一致する条件(定理3)を求めた。

(つまり、一般に用いられているCDLの収束条件、
最適化後の分布が真の分布に一致する条件を求めることができた。)

RBMとCDLの説明



DBLの説明



DBLの理論解析



計算機実験

