

# 統計的機械学習理論と ボルツマン機械学習

山形大学 大学院理工学研究科  
安田 宗樹

# 統計的機械學習

# 統計的機械学習は何をするのか？

3 / 141

決定論的ニューラルネットワークによる機械学習

入出力関係をデータから学習する

入力と出力の決定論的な(入力から出力が一意に決まる)

関数関係を探し出すのが目的

---

統計的機械学習

データの確率的な生成規則を学習する

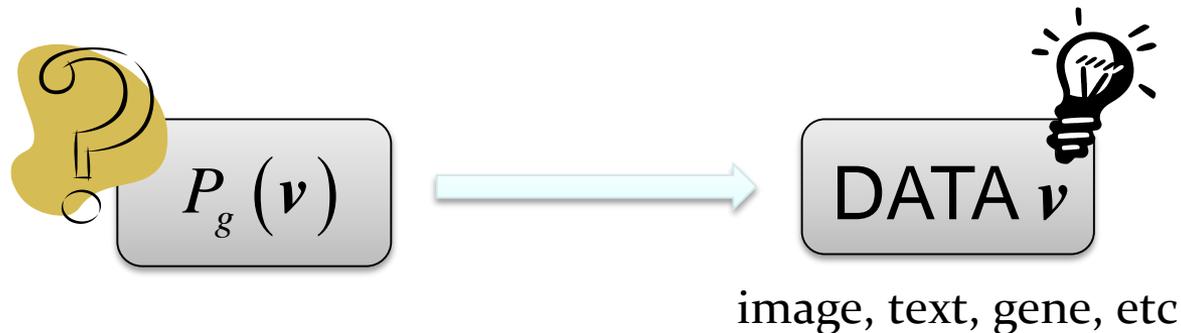
データパターンを生成する背景となっている

確率分布を探し出すのが目的

# なぜ確率を使うのか？

4 / 141

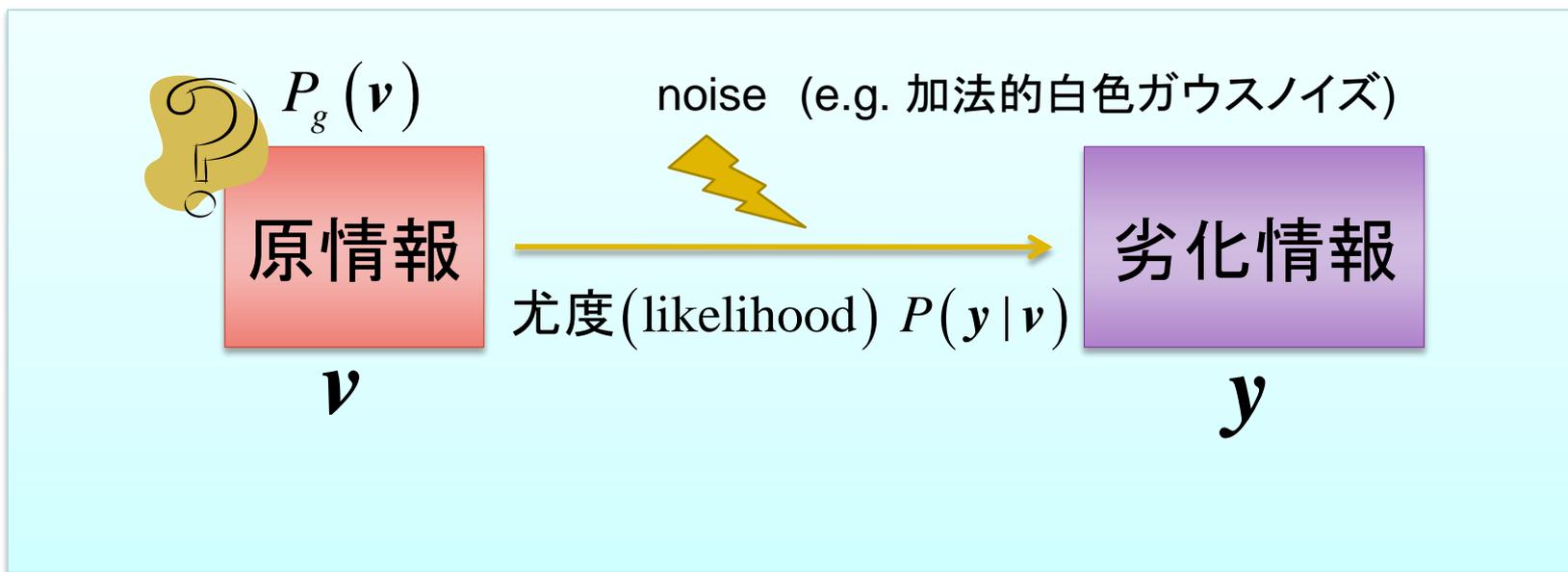
データパターンはとある確率規則に従って発生しているとする



何故確率が必要？

- 観測ノイズに関わる不可避の**不確実性**
- 知識不足による未知の因果に対する**不確実性**

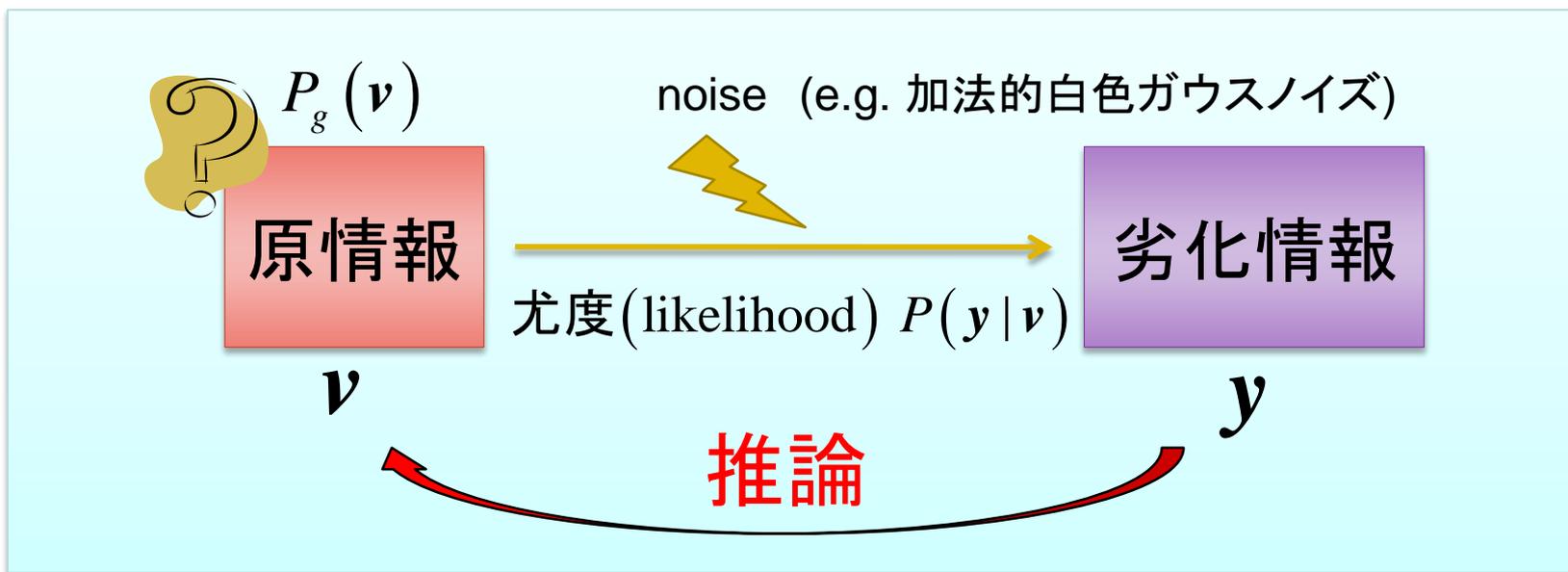
不確実な現象を表現する数学 = **確率理論**



原情報  $v$  が通信路ノイズによって劣化し

劣化情報  $y$  に変化した

我々は  $y$  しか知らない



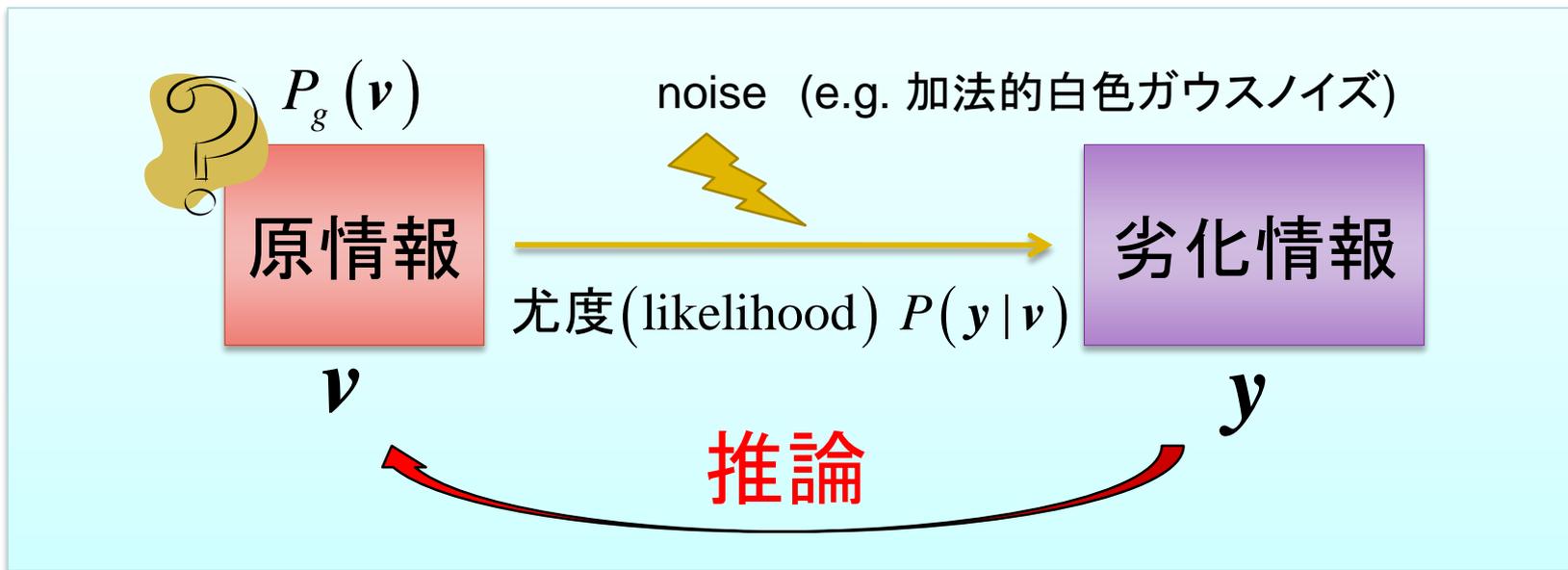
我々は  $y$  しか知らない



$y$  から  $v$  を推定したい!

# 事後確率を計算する

7 / 141



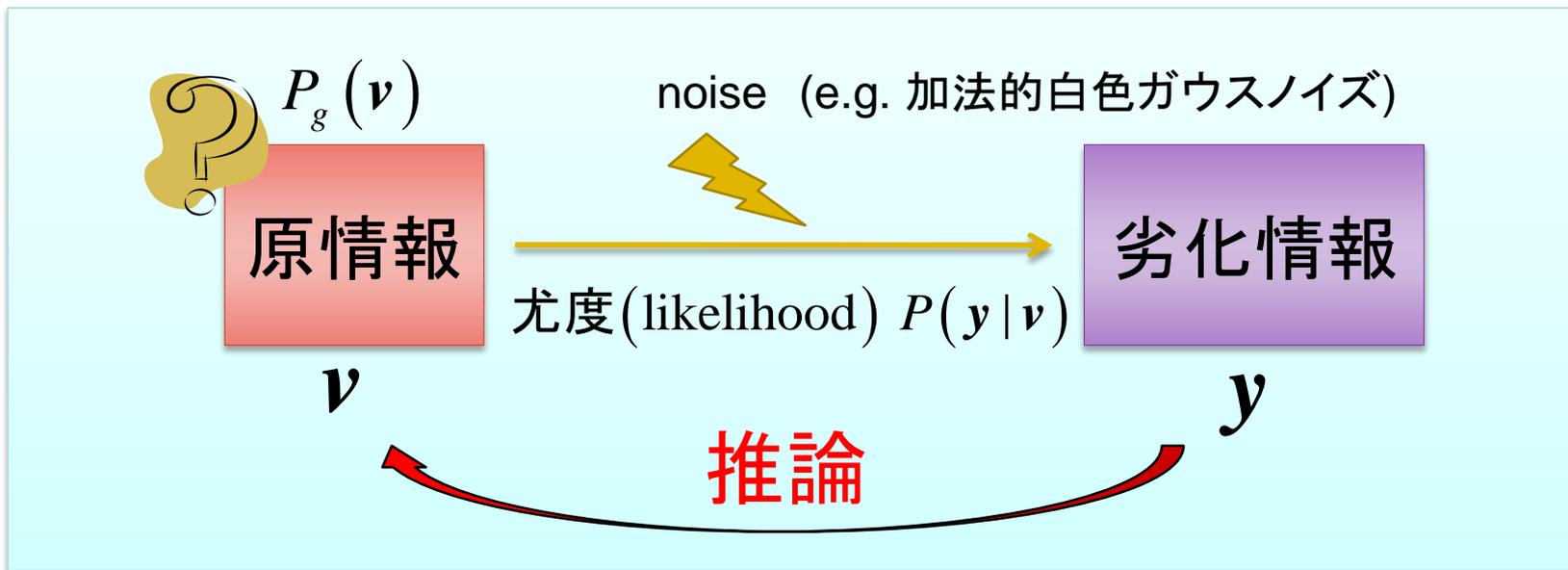
事後確率

ベイズの公式

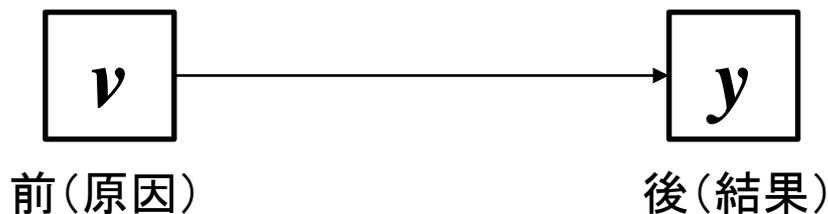
$$P(v|y) = \frac{P(v, y)}{P(y)} = \frac{P(y|v)P_g(v)}{P(y)}$$

$y$  を知った下での  $v$  の確率

# 事後確率 (用語の整理)



時間的には。。。

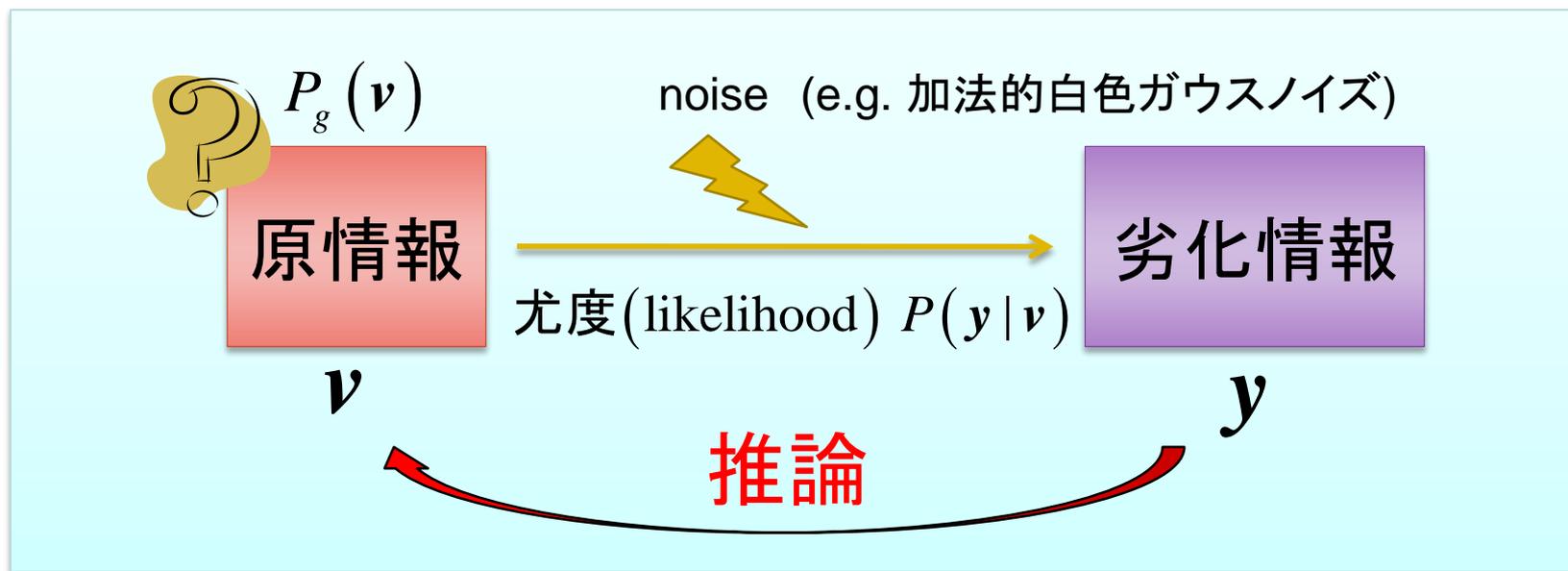


$$(\text{事後確率}) = P(v | y) = P(\text{原因} | \text{結果})$$

結果から原因を探るための定式化

# 最大事後確率 (MAP) 推定

9 / 141



$$\mathbf{v}^* = \arg \max_{\mathbf{v}} P(\mathbf{v}|\mathbf{y}) = \arg \max_{\mathbf{v}} P(\mathbf{y}|\mathbf{v}) P_g(\mathbf{v})$$

$\mathbf{y}$  から見て確率的にもっともあり得そうな  $\mathbf{v}$  が原情報だろうと考える

ベイズ的確率推定の基礎

$$P(\mathbf{v} | \mathbf{y}) \propto P(\mathbf{y} | \mathbf{v}) P_g(\mathbf{v})$$



$$\text{(事後確率)} \propto \text{(尤度)} \times \text{(事前確率)}$$

- 尤度 …… 原情報がどのような過程で変形を受けるかの確率
- 事前確率 …… 原情報がどのような確率規則で生成されるかの確率

事前確率  $P_g(\mathbf{v})$  : 真の情報の生成確率(情報生成の背後にある確率規則)

$$\begin{aligned} \text{(事前確率)} &= \text{(原情報生成の背後規則)} \\ &= \text{(生成モデル(generative model))} \end{aligned}$$

事前確率モデル = データの生成モデル(メカニズム)

- ✓ 適切に設計しなければ確率推論そのものがダメになる
- ✓ 一般的にはかなり難しい問題（メカニズム理解からのモデル構築）  
（ex. 遺伝子発現のメカニズム, 自然画像の生成メカニズム）

データ駆動型のモデル構築

- ✓ 生成モデルの詳細は分からないが、そこから発生したデータはもっている  
（例えば画像などは google で検索すればいくらでも手に入る）
- ✓ それらのデータは生成モデルに関する情報をもっているはず
- ✓ その情報を活かして生成モデルを探せないだろうか？

データ主体の生成モデル発掘技術 → 統計的機械学習

$N$  given ( $n$  dimensional) data points:

$$D = \underbrace{\{\mathbf{d}^{(1)}, \mathbf{d}^{(2)}, \dots, \mathbf{d}^{(N)}\}}_{N \text{ data points}}, \quad \mathbf{d}^{(\mu)} = \underbrace{\{d_1^{(\mu)}, d_2^{(\mu)}, \dots, d_n^{(\mu)}\}}_{n \text{ dimensional data point}}$$

$$d_i^{(\mu)} \in \{0, 1\}$$

$$\begin{array}{l} \boxed{0100101 \dots 001} \quad \mathbf{d}^{(1)} \\ \boxed{1011011 \dots 111} \quad \mathbf{d}^{(2)} \\ \vdots \\ \underbrace{\boxed{0010011 \dots 101}}_n \quad \mathbf{d}^{(N)} \end{array}$$

それぞれのデータはある生成モデルから独立に生成されているとする

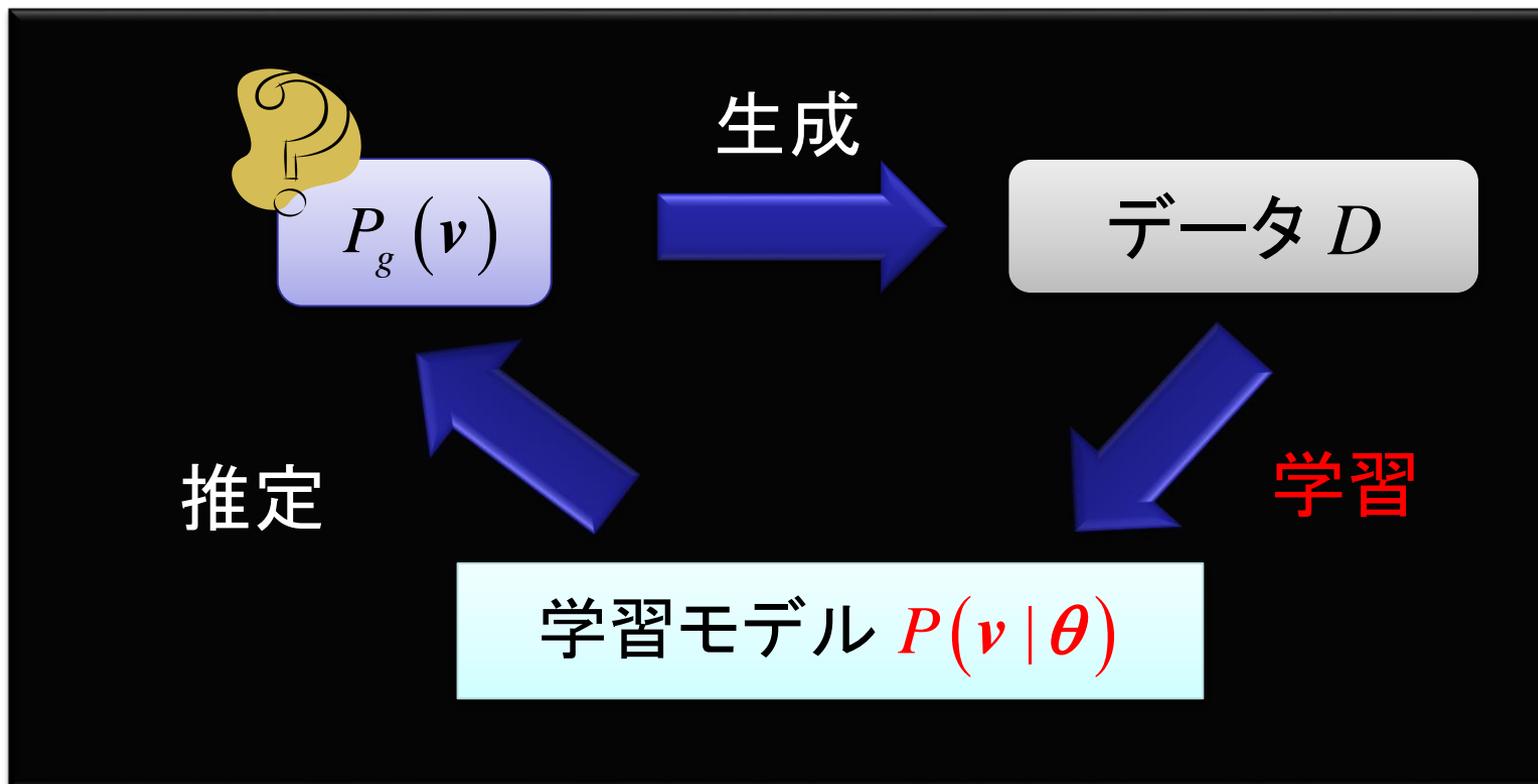
$$P_g(\mathbf{v}) \xrightarrow{\text{generate (i.i.d.)}} \mathbf{d}^{(1)}, \mathbf{d}^{(2)}, \dots, \mathbf{d}^{(N)}$$

$$\mathbf{v} = \{v_1, v_2, \dots, v_n\}$$

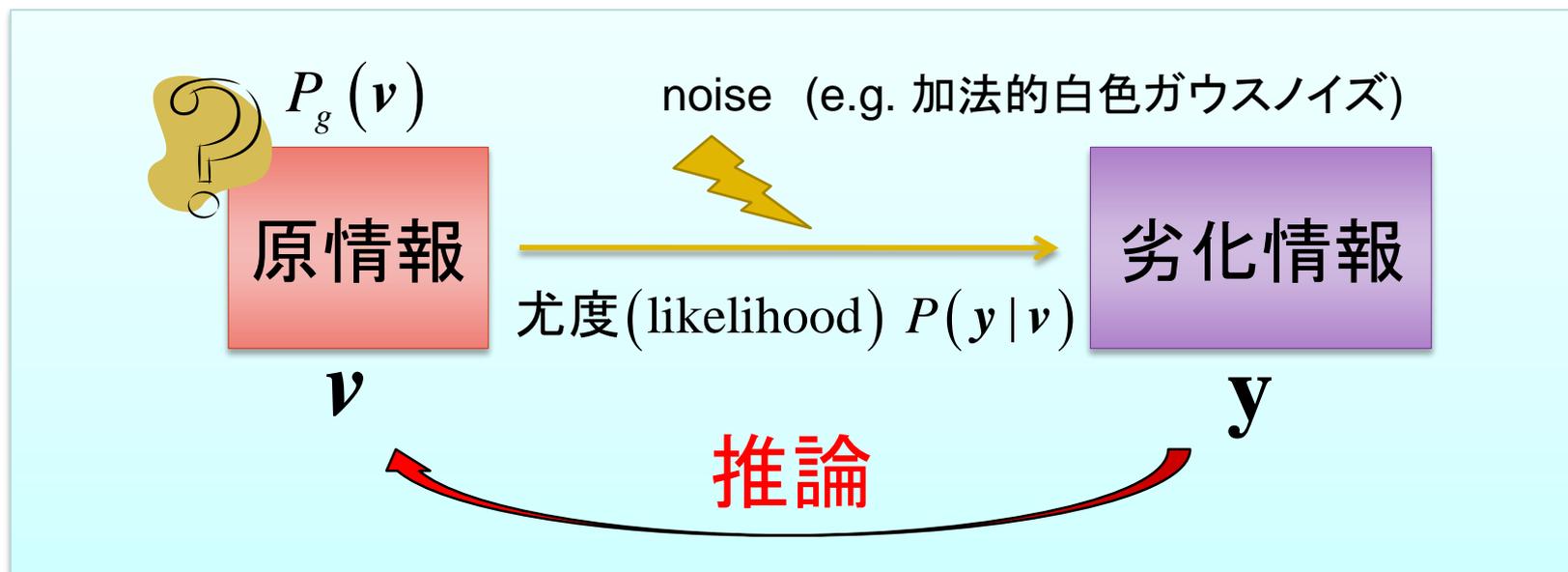
生成モデルは データ生成の確率規則を記述する

$$\mathbf{d}^{(\mu)} \text{ は確率 } p_\mu \text{ で生成される} \Leftrightarrow P_g(\mathbf{d}^{(\mu)}) = p_\mu$$

$P(v | \theta)$  ← パラメータ  $\theta$  をもったある学習モデルを仮定する



- データ  $D$  を利用して仮定した学習モデルを学習する
- 学習により(未知の)事前確率に学習モデルを近づける

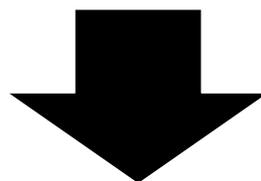


$$\mathbf{v}^* = \arg \max_{\mathbf{v}} P(\mathbf{y}|\mathbf{v}) \underbrace{P_g(\mathbf{v})}_{\text{生成モデル}} \approx \arg \max_{\mathbf{v}} P(\mathbf{y}|\mathbf{v}) \underbrace{P(\mathbf{v}|\boldsymbol{\theta})}_{\text{学習モデル}}$$

学習モデルで真の生成モデルを近似して  
原情報を推定する！

生成モデルを近似する学習モデルを仮定する必要がある

どのようなモデルを仮定すればよいか??



もちろんそれは課題によって個別に設計されるべき

だが

**グラフィカルモデル**は中でもわりと一般的に使える  
可能な選択肢の一つ

# グラフィカルモデル

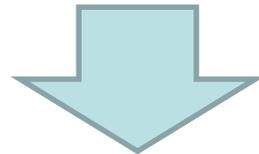
(生成モデル) = (データの確率的生成規則を記述する)

ふつう多次元 (e.g. 画像ならピクセルの個数分の次元)

- データというものはふつう独立的ではなく、データ要素間に何らかの関連性(確率的相関構造)をもつ
  - 出現しやすいデータパターンや出現しにくいパターンが存在する (A家の夕食がカレーのときは何故かB家の夕食が煮物であることが多い etc)
- データを確率的に扱うだけでなく、そのような関連性も扱うことが重要

データ要素間の関連性が意味のある情報の基本

生成モデルはデータ要素間に関連性を持たせたものであるべき



データ要素間の確率的関連性を明に取り入れるモデリング

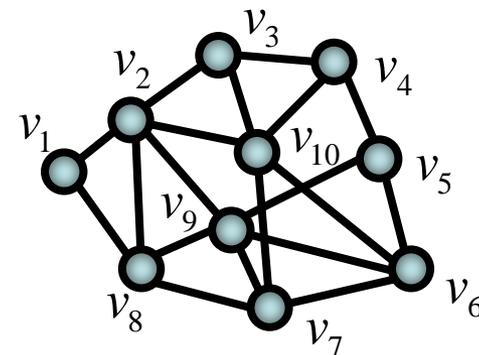
- マルコフ確率場 (Markov random field; MRF)
- ベイジアンネットワーク (Bayesian network; BN)

ともにグラフィカルモデルとよばれるモデル

# グラフィカルモデル？

19 / 141

ノードとリンクから構成されるグラフを作る



各ノードがデータの各要素(つまり各確率変数)に対応する

リンクは確率変数間の(直接の)関連性を表わす

関連性を直観的に  
導入することができる

作ったグラフ構造に則って確率モデルを作成する

## MRF

無向グラフ上の確率モデル

## BN

有向(非循環)グラフ上の確率モデル



## MRF と BN の関係

BN から MRF への変換 … 簡単

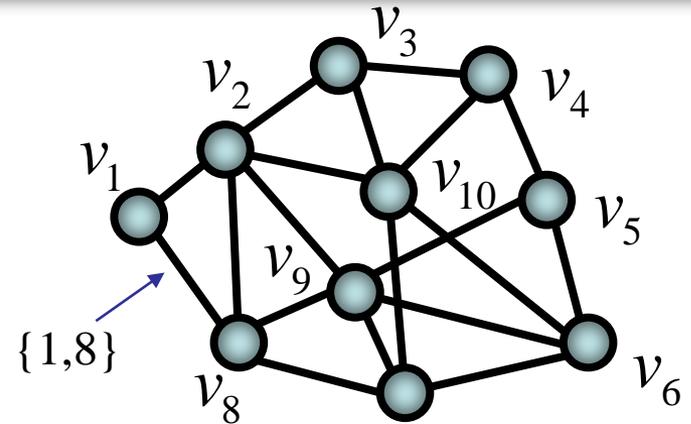
MRF から BN への変換 … NP-hard

# マルコフ確率場の導入

21 / 141

$G(V, E) \leftarrow$  無向グラフ

$V = \{1, 2, \dots, 10\} \leftarrow$  ノードラベルの集合



$E = \{\{1, 2\}, \{1, 8\}, \{2, 3\}, \{2, 8\}, \{2, 9\}, \dots\} \leftarrow$  リンクラベルの集合

- 各ノードが1つのデータ要素に対応  
(確率変数が割り当てられる)

$$\nu = \{\nu_1, \nu_2, \dots, \nu_{10}\}$$

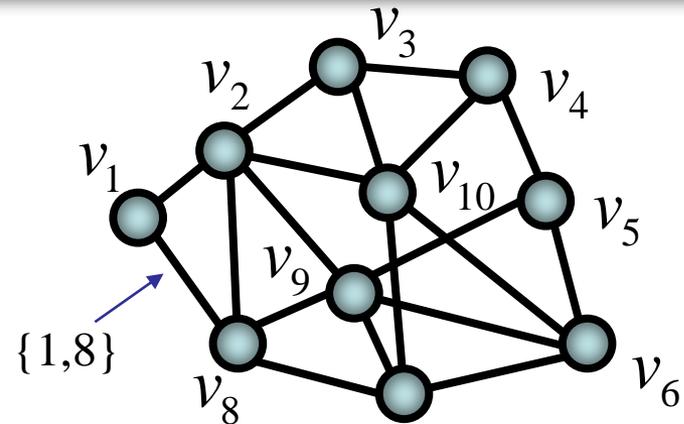
- リンクが変数同士の確率的な関連性を表現  
影響を直接及ぼすノード同士を結んでいる

# マルコフ確率場の導入

22 / 141

$G(V, E) \leftarrow$  無向グラフ

$V = \{1, 2, \dots, 10\} \leftarrow$  ノードラベルの集合



$E = \{\{1, 2\}, \{1, 8\}, \{2, 3\}, \{2, 8\}, \{2, 9\}, \dots\} \leftarrow$  リンクラベルの集合

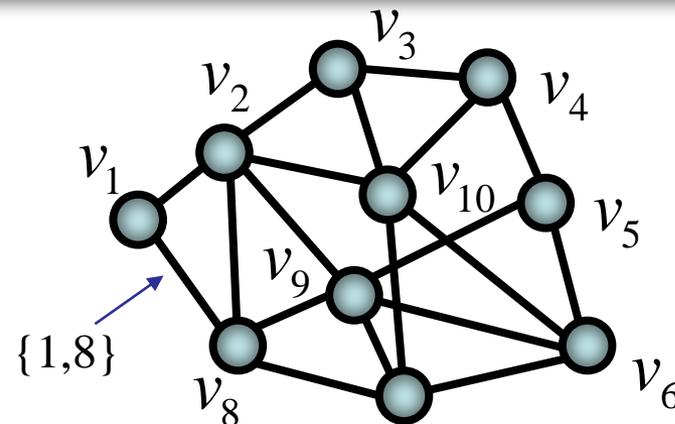
エネルギー関数(コスト関数)

$$C(\mathbf{v}) = \underbrace{-\sum_{i \in V} \Phi_i(v_i)}_{\text{各ノードごとのコスト}} - \underbrace{\sum_{\{i, j\} \in E} \Psi_{\{i, j\}}(v_i, v_j)}_{\text{各リンクごとのコスト}}$$

# マルコフ確率場の導入

$G(V, E) \leftarrow$  無向グラフ

$V = \{1, 2, \dots, 10\} \leftarrow$  ノードラベルの集合



$E = \{\{1, 2\}, \{1, 8\}, \{2, 3\}, \{2, 8\}, \{2, 9\}, \dots\} \leftarrow$  リンクラベルの集合

## エネルギー関数の和の意味

$$C(\nu) = \underbrace{-\sum_{i \in V} \Phi_i(\nu_i)}_{\text{各ノードごとのコスト}} - \underbrace{\sum_{\{i, j\} \in E} \Psi_{\{i, j\}}(\nu_i, \nu_j)}_{\text{各リンクごとのコスト}}$$

$\sum_{i \in V} \Phi_i(\nu_i) = \Phi_1(\nu_1) + \Phi_2(\nu_2) + \Phi_3(\nu_3) + \dots + \Phi_{10}(\nu_{10}) \leftarrow$  全てのノードに対する和

$\sum_{\{i, j\} \in E} \Psi_{\{i, j\}}(\nu_i, \nu_j) = \Psi_{\{1, 2\}}(\nu_1, \nu_2) + \Psi_{\{1, 8\}}(\nu_1, \nu_8) + \Psi_{\{2, 3\}}(\nu_2, \nu_3) + \dots \leftarrow$  全てのリンクに対する和

エネルギー関数(コスト関数)

$$C(\mathbf{v}) = \underbrace{-\sum_{i \in V} \Phi_i(v_i)}_{\text{各ノードごとのコスト}} - \underbrace{\sum_{\{i,j\} \in E} \Psi_{\{i,j\}}(v_i, v_j)}_{\text{各リンクごとのコスト}}$$

マルコフ確率場

$$P(\mathbf{v}) = \frac{1}{Z} \exp(-C(\mathbf{v}))$$

すべての変数の結合確率分布

規格化定数(分配関数)

確率分布の規格化条件を保証するための定数

$$\sum_{\mathbf{v}} P(\mathbf{v}) = 1$$

変数の実現可能な値のすべての組み合わせの和は必ず1になるという条件

すべての変数が0か1の2値を取る場合

$$\sum_{\mathbf{v}} (\dots) = \sum_{v_1=0,1} \sum_{v_2=0,1} \dots \sum_{v_n=0,1} (\dots)$$

$$\sum_{\mathbf{v}} P(\mathbf{v}) = \frac{1}{Z} \sum_{\mathbf{v}} \exp(-C(\mathbf{v})) = 1 \Leftrightarrow Z = \sum_{\mathbf{v}} \exp(-C(\mathbf{v}))$$

- 無向グラフのリンク構造に則ってエネルギー関数を設計
- エネルギー関数の各項の具体的な関数形は課題によって適宜設計する
- エネルギー関数の負を指数の肩に載せて規格化すると出来上がり！

$$P(\boldsymbol{v}) = \frac{1}{Z} \exp(-C(\boldsymbol{v}))$$

エネルギー関数が低い  $\boldsymbol{v}$  のとき確率が高くなる



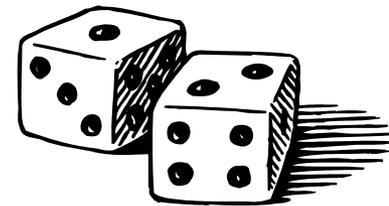
エネルギー関数を低くする  $\boldsymbol{v}$  のパターンが頻繁に生成される

# マルコフ確率場は結合確率分布

27 / 141

我々は今何を作ったのだろうか??

$$P(\mathbf{v}) = \frac{1}{Z} \exp(-C(\mathbf{v}))$$



例えばバイナリの確率変数が 3 個であるような状況を考える  $\mathbf{v} = \{v_i \in \{0,1\} \mid i = 1, 2, 3\}$

$\mathbf{v}$  は 000, 001, 010, ..., 111 の計  $2^3 = 8$  パターンをとる

そして

それぞれのパターンごとにそれぞれ MRF に従う確率値が割り当てられる

(MRF) = (8 面サイコロ)

ただし各面の出る確率は等確率でなく MRF が定義する確率に従う

エネルギー関数を低くするパターンの面がよく出るサイコロ

MRF はサイコロらしいのだが  
ではそのサイコロはどのように振ればよいのか？



MRF に従う疑似乱数の取得方法は？



多変数の確率分布においては

**ある変数の値が他の変数の値の確率に影響を及ぼす**ため

一般には**各変数ごとに独立に値を決めていくことができない**  
(ある変数の値を決める際に他の変数の値を気にする必要がある)

そのため疑似乱数のサンプリングが通常より難しい

MRF はサイコロらしいのだが  
ではそのサイコロはどのように振ればよいのか？



MRF に従う疑似乱数の取得方法は？

マルコフ連鎖モンテカルロ (Markov chain Monte Carlo; MCMC) 法

の中の

ギブス (Gibbs) サンプリング法

が有名

その他の変数すべてを条件とした  
ある確率変数  $v_i$  についての条件付き確率を計算する

$$P(v_i | \mathbf{v}_{-i}) = \frac{P(\mathbf{v})}{P(\mathbf{v}_{-i})} \propto \exp \left( \Phi_i(v_i) + \sum_{j \in \partial(i)} \Psi_{\{i,j\}}(v_i, v_j) \right)$$

$\mathbf{v}_{-i} \leftarrow v_i$  以外のすべての変数     $\partial(i) \leftarrow$  ノード  $i$  とリンクでつながるノードの集合



$$P(v_i | \mathbf{v}_{-i}) = P(v_i | \mathbf{v}_{\partial(i)})$$

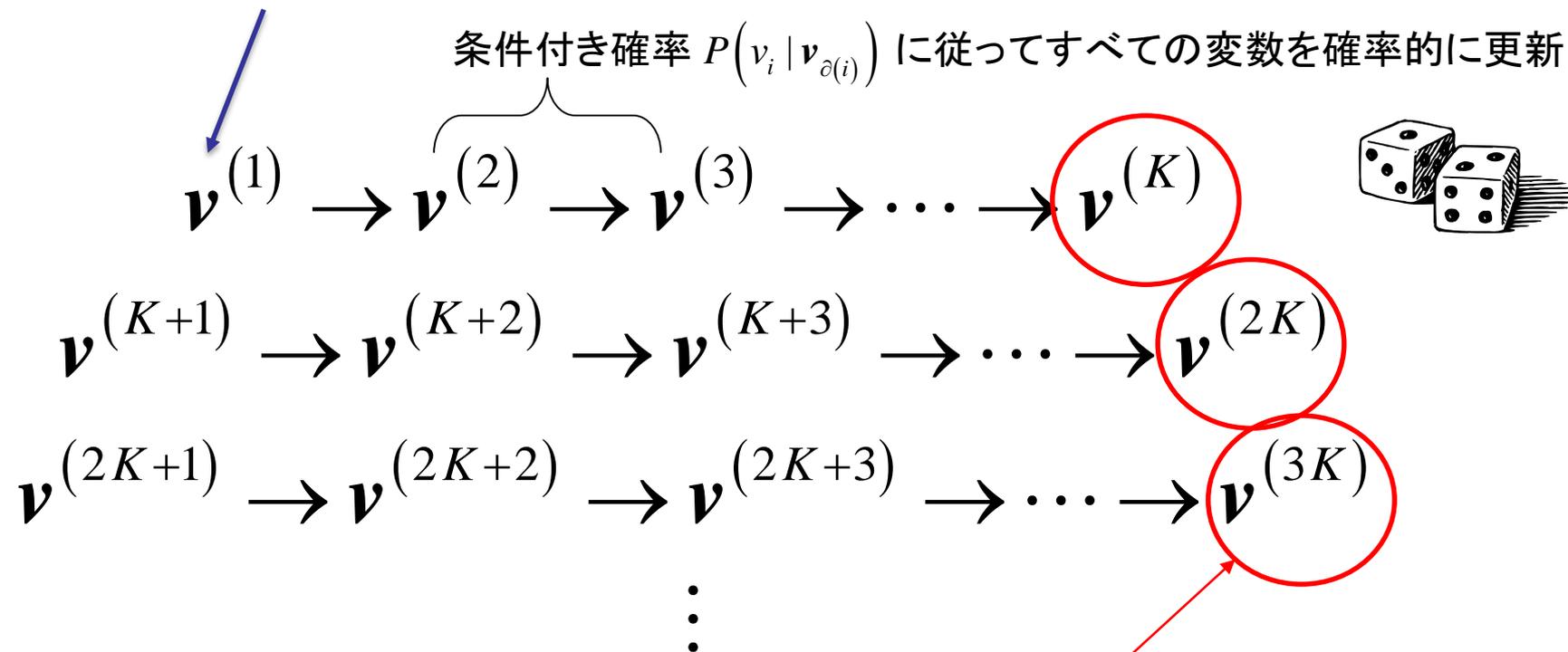
$\mathbf{v}_{\partial(i)} \leftarrow v_i$  とリンクでつながっている変数の集合

**リンクでつながっている変数のみにしか依存しない！**

# ギブスサンプリングの方法

31 / 141

全ての変数を適当に初期化

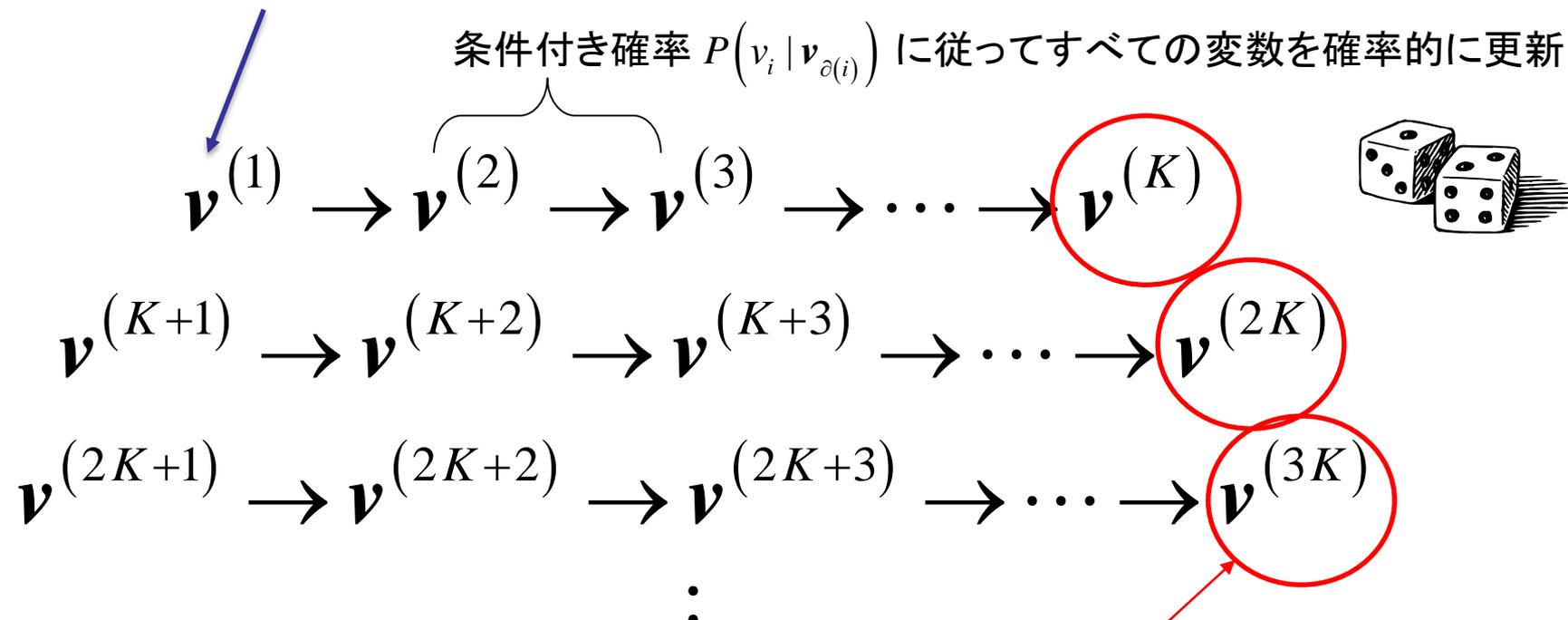


上記の手続きにより一定間隔 ( $K$  ステップ間隔) でサンプリングする

サンプル点ごとの相関を切って**独立なサンプリング**を実現するため

# ギブスサンプリングの方法

全ての変数を適当に初期化



初期の段階のサンプル点は出現確率が安定しないので捨てて  
(緩和期間)

十分回更新した後のサンプル点から採用する

# ボルツマンマシン ～MRFの基本モデル～

# ボルツマンマシン(Boltzmann machine; BM)とは？

34 / 141

- もっとも基礎的な MRF
- 連想記憶モデルであるホップフィールドモデルの確率的拡張モデル
- 相互結合型の確率的ニューラルネットワーク
- 最近のディープラーニング(深層学習)ブームで再燃焼

$\Phi_i(v_i) \rightarrow b_i v_i$     $\Psi_{\{i,j\}}(v_i, v_j) \rightarrow w_{ij} v_i v_j$  と単純な関数形で仮定する

## ボルツマンマシンのエネルギー関数

$$C(\mathbf{v}; \mathbf{b}, \mathbf{w}) = -\sum_{i \in V} b_i v_i - \sum_{\{i,j\} \in E} w_{ij} v_i v_j \quad v_i \in \{0,1\}$$

バイアスパラメータ

結合パラメータ

$$w_{ij} = w_{ji} \text{ (対称性)}$$

$$P(\mathbf{v} | \mathbf{b}, \mathbf{w}) = \frac{1}{Z(\mathbf{b}, \mathbf{w})} \exp \left( \sum_{i \in V} b_i v_i + \sum_{\{i, j\} \in E} w_{ij} v_i v_j \right)$$

$$v_i \in \{0, 1\}$$

$$C(\mathbf{v}; \mathbf{b}, \mathbf{w}) = - \sum_{i \in V} b_i v_i - \sum_{\{i, j\} \in E} w_{ij} v_i v_j$$

エネルギー関数が低いほど確率が高い

- もし  $b_i > 0$  ならば局所的には  $v_i = 1$  で逆に  $b_i < 0$  なら  $v_i = 0$  が確率的に高い
- もし  $w_{ij} > 0$  ならば  $v_i = 1$  &  $v_j = 1$  が局所的には確率が高い
- もし  $w_{ij} < 0$  ならば  $v_i = 1$  &  $v_j = 1$  でない方が局所的には確率が高い

$$P(\mathbf{v} | \mathbf{b}, \mathbf{w}) = \frac{1}{Z(\mathbf{b}, \mathbf{w})} \exp \left( \sum_{i \in V} b_i v_i + \sum_{\{i, j\} \in E} w_{ij} v_i v_j \right)$$

$$v_i \in \{0, 1\}$$

エネルギー関数中のリンクの項が

変数同士の複雑な関連性

を生み出している

ボルツマンマシン  
に対する  
統計的機械学習

$N$  given ( $n$  dimensional) data points:

$$D = \underbrace{\{\mathbf{d}^{(1)}, \mathbf{d}^{(2)}, \dots, \mathbf{d}^{(N)}\}}_{N \text{ data points}}, \quad \mathbf{d}^{(\mu)} = \underbrace{\{d_1^{(\mu)}, d_2^{(\mu)}, \dots, d_n^{(\mu)}\}}_{n \text{ dimensional data point}}$$

$$d_i^{(\mu)} \in \{0, 1\}$$

$$\begin{array}{l} \boxed{0100101 \dots 001} \quad \mathbf{d}^{(1)} \\ \boxed{1011011 \dots 111} \quad \mathbf{d}^{(2)} \\ \vdots \\ \underbrace{\boxed{0010011 \dots 101}}_n \quad \mathbf{d}^{(N)} \end{array}$$

それぞれのデータはある生成モデルから独立に生成されているとする

$$P_g(\boldsymbol{\nu}) \xrightarrow{\text{generate (i.i.d.)}} \mathbf{d}^{(1)}, \mathbf{d}^{(2)}, \dots, \mathbf{d}^{(N)}$$

観測データのセットを使って BM を設計する方法を考える

以降  $|V| = n$  を前提として話を進める

データの次元 ( $n$  次元)  
と BM の変数の数が同じ  
という意味

$N$  given ( $n$  dimensional) data points:

$$D = \underbrace{\{\mathbf{d}^{(1)}, \mathbf{d}^{(2)}, \dots, \mathbf{d}^{(N)}\}}_{N \text{ data points}}, \quad \mathbf{d}^{(\mu)} = \underbrace{\{d_1^{(\mu)}, d_2^{(\mu)}, \dots, d_n^{(\mu)}\}}_{n \text{ dimensional data point}}$$

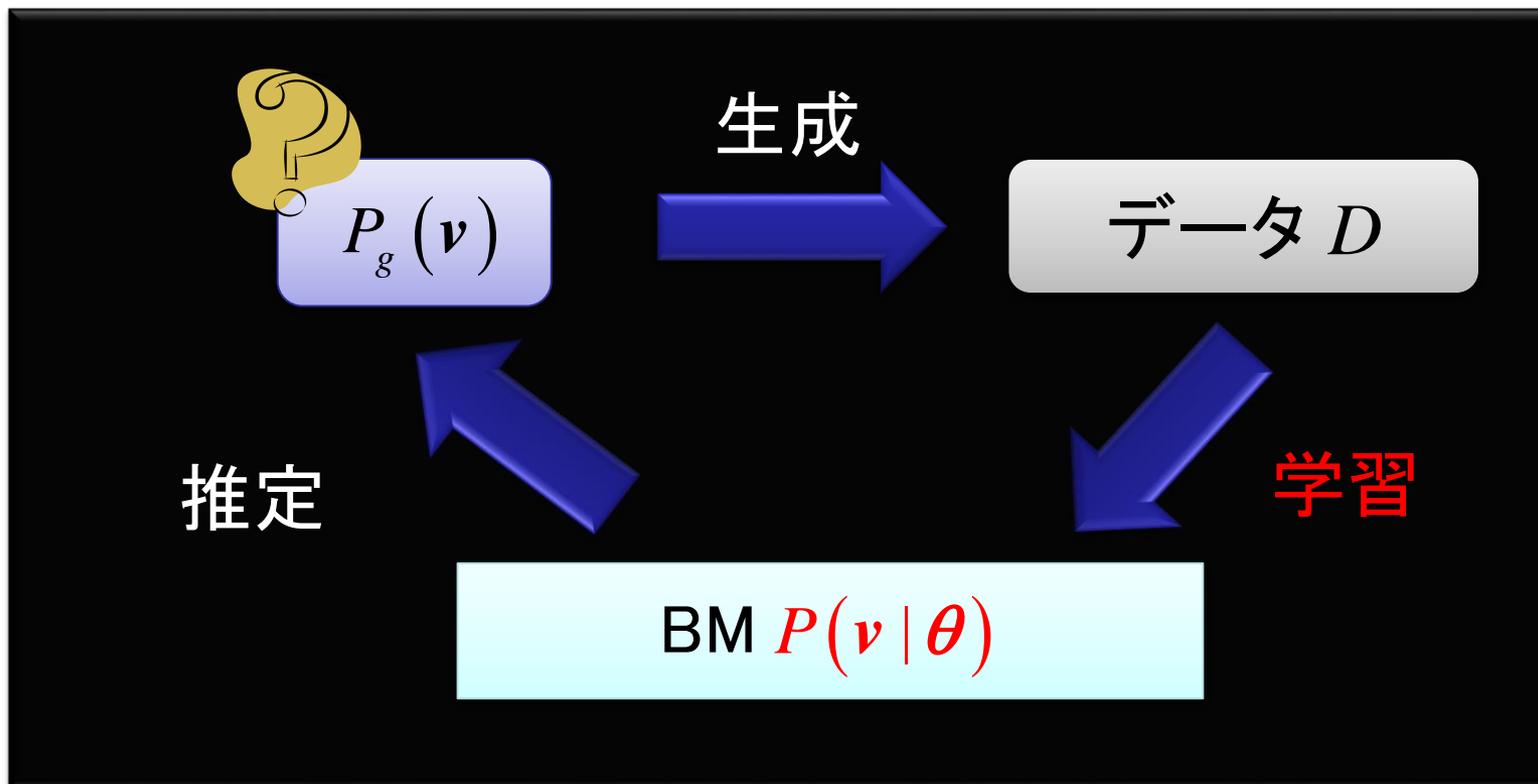
$$d_i^{(\mu)} \in \{0, 1\}$$

$$\begin{array}{l} \boxed{0100101 \dots 001} \quad \mathbf{d}^{(1)} \\ \boxed{1011011 \dots 111} \quad \mathbf{d}^{(2)} \\ \vdots \\ \underbrace{\boxed{0010011 \dots 101}}_n \quad \mathbf{d}^{(N)} \end{array}$$

$d_i^{(\mu)}$  を  $v_i$  に対応させる

データの要素一つ一つが  
各変数に対応している状況

$P(v | \theta) \leftarrow$  パラメータ  $\theta$  をもった **BM** を仮定する  $\theta = \{b, w\}$

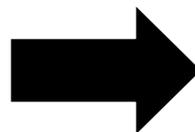


- データ  $D$  を利用して仮定した **BM** を学習する
- 学習により(未知の)生成モデルに **BM** を近づける

## 統計的機械学習の基本的考え方

### 最尤推定 (maximum likelihood estimation; MLE)

BM  $P(v | \theta)$



データ  $D$

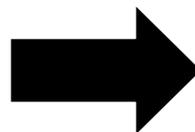
観測データは BM (学習モデル) から i.i.d. で生成されたと考える

実際は未知の事前確率から生成されているのだが  
事前確率を学習モデルで近似するという考え方からこのように考える

## 統計的機械学習の基本的考え方

### 最尤推定 (maximum likelihood estimation; MLE)

BM  $P(v | \theta)$



データ  $D$

観測データを生成するに

**もっとも尤もらしい** 確率になるようにパラメータを調整する

何故なら観測データはそのモデルから生成されたと考えるため

# 最尤推定 (尤度関数)

$$D = \underbrace{\{\mathbf{d}^{(1)}, \mathbf{d}^{(2)}, \dots, \mathbf{d}^{(N)}\}}_{N \text{ data points}}, \quad \mathbf{d}^{(\mu)} = \underbrace{\{d_1^{(\mu)}, d_2^{(\mu)}, \dots, d_n^{(\mu)}\}}_{n \text{ dimensional data point}}$$

観測データと同じ次元 ( $n$ 次元) の確率変数をもつ BM を用意する

$$P(\mathbf{v} | \theta) \quad \mathbf{v} = \{v_1, v_2, \dots, v_n\}$$



$$P(\mathbf{v} = \mathbf{d}^{(\mu)} | \theta) \leftarrow \text{MRFが } \mu \text{ 番目のデータを生成する確率}$$

BM が観測データセット  $D$  を生成する確率

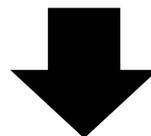
$$\underbrace{P(\mathbf{v} = \mathbf{d}^{(1)} | \theta) P(\mathbf{v} = \mathbf{d}^{(2)} | \theta) \dots P(\mathbf{v} = \mathbf{d}^{(N)} | \theta)}_{\text{各データは独立に生成されるので積}} = \prod_{\mu=1}^N P(\mathbf{v} = \mathbf{d}^{(\mu)} | \theta) = \text{(尤度関数)}$$

## 尤度関数

$$L_D(\theta) = \prod_{\mu=1}^N P(\mathbf{v} = \mathbf{d}^{(\mu)} | \theta)$$

仮定した学習モデルが  
実際に観測データを生成する確率

観測データセットに対して一番尤もらしいモデル



尤度関数を最大とするパラメータ値のときのモデル

※ 観測データセットを生成する確率が一番高いモデル

## 最尤推定

$$\theta^* = \arg \max_{\theta} L_D(\theta)$$



$P(\mathbf{v} | \theta^*)$  が観測データに対して一番尤もらしいモデル



(モデル設計の枠内で) 事前確率に一番近いモデル

$$P_g(\mathbf{v}) \approx P(\mathbf{v} | \theta^*)$$

$$l_D(\boldsymbol{\theta}) = \frac{1}{N} \ln L_D(\boldsymbol{\theta}) = \frac{1}{N} \sum_{\mu=1}^N \ln P(\boldsymbol{v} = \mathbf{d}^{(\mu)} | \boldsymbol{\theta})$$

尤度関数の対数をとったもの

積が和に変身するので計算的に扱い易いしオーバー(アンダー)フロー等にも強い

対数関数は単調増加関数なので

(尤度関数の最大点) = (対数尤度関数の最大点)

なので通常最尤推定をするときは

尤度関数よりもむしろ対数尤度関数を用いる

カルバック・ライブラー (Kullback-Leibler; KL) 情報量

2つの確率分布  $P$  と  $Q$  に対して

$$D(P \parallel Q) = \sum_{\mathbf{v}} P(\mathbf{v}) \ln \frac{P(\mathbf{v})}{Q(\mathbf{v})}$$

$$D(P \parallel Q) \geq 0 \leftarrow \text{非負}$$

$$D(P \parallel Q) = 0 \text{ iff } P(\mathbf{v}) = Q(\mathbf{v}) \leftarrow \text{2つの確率が一致するとき } 0$$

2つの確率分布間の距離の指標として使用される

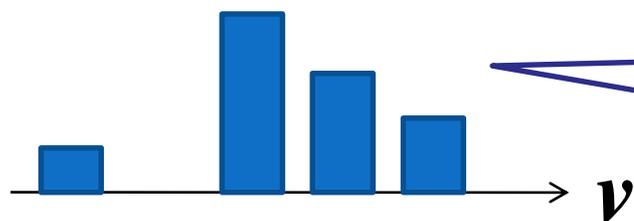
※ 距離の定義を満足しないので厳密な意味での距離ではない

$$\mathbf{D} = \underbrace{\{\mathbf{d}^{(1)}, \mathbf{d}^{(2)}, \dots, \mathbf{d}^{(N)}\}}_{N \text{ data points}}, \quad \mathbf{d}^{(\mu)} = \underbrace{\{d_1^{(\mu)}, d_2^{(\mu)}, \dots, d_n^{(\mu)}\}}_{n \text{ dimensional data point}}$$

$$H_D(\mathbf{v}) = \frac{1}{N} \sum_{\mu=1}^N \delta(\mathbf{v}, \mathbf{d}^{(\mu)}) \leftarrow \underbrace{\text{観測データの分布 (ヒストグラム)}}_{\text{経験分布 (empirical distribution)}}$$

クロネッカーのデルタ関数  $\delta(\mathbf{v}, \mathbf{d}^{(\mu)}) = \begin{cases} 1 & (\mathbf{v} = \mathbf{d}^{(\mu)}) \\ 0 & (\mathbf{v} \neq \mathbf{d}^{(\mu)}) \end{cases}$

※ 確率変数が連続の場合はディラックのデルタ関数になる



観測データが存在するところに  
1個ずつ積み上げていく

経験分布  $H_D(\mathbf{v})$  と学習モデル  $P(\mathbf{v}|\theta)$  との間のKL情報量

$$D(H_D \parallel P) = \sum_{\mathbf{v}} H_D(\mathbf{v}) \ln \frac{H_D(\mathbf{v})}{P(\mathbf{v}|\theta)}$$

$$D(H_D \parallel P) = \underbrace{\sum_{\mathbf{v}} H_D(\mathbf{v}) \ln H_D(\mathbf{v})}_{\theta \text{ に関係ない項}} - \underbrace{l_D(\theta)}_{\text{対数尤度関数}}$$

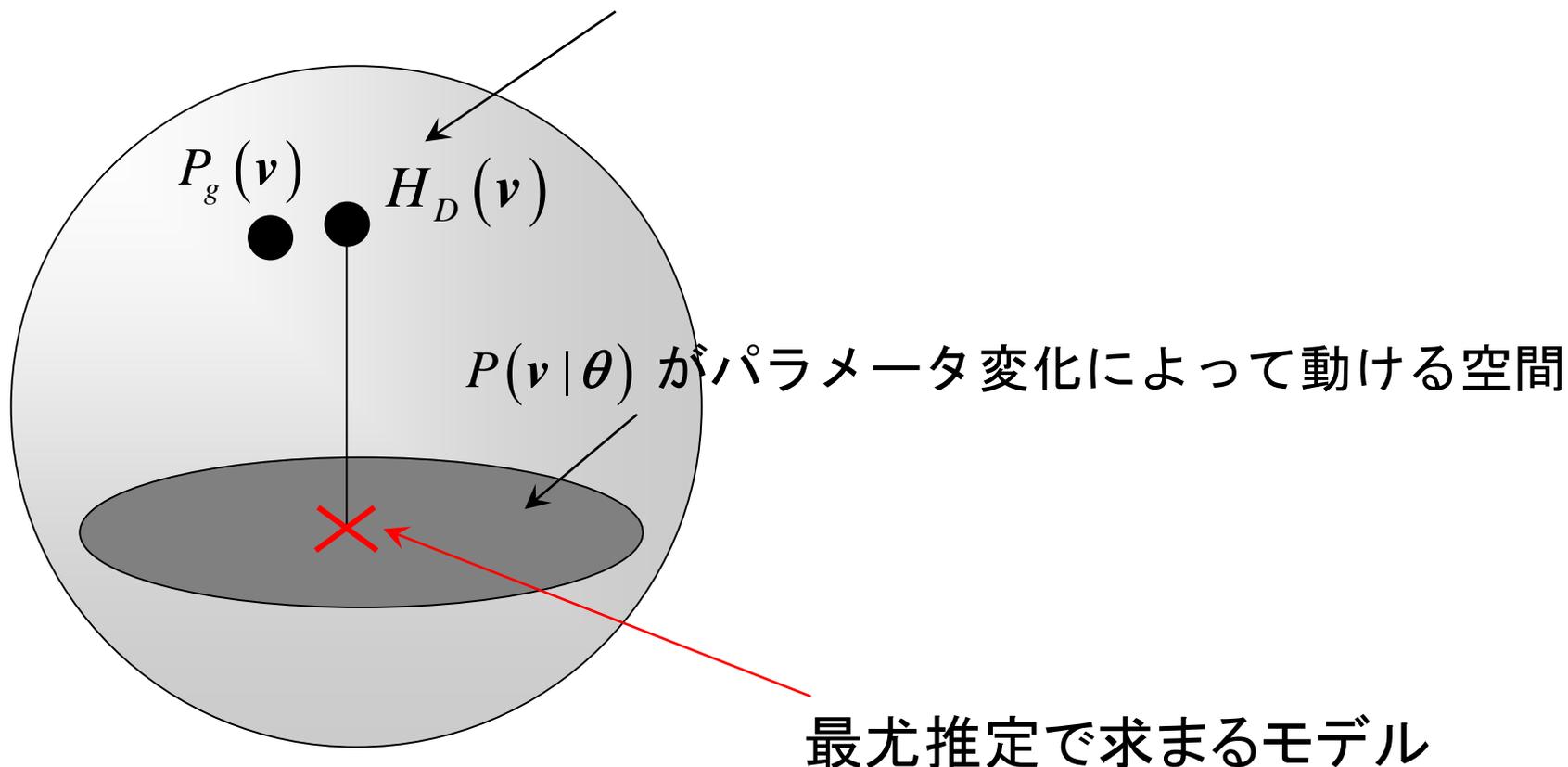


(最尤推定) = (KL情報量の最小化)



(最尤推定) = (観測データの分布に一番近い分布を見つけること)

事前確率 (生成モデル) と経験分布は比較的近いだろう



# ボルツマンマシンの最尤推定



$$P(\mathbf{v} | \boldsymbol{\theta}) = \frac{1}{Z(\boldsymbol{\theta})} \exp \left( \sum_{i \in V} b_i v_i + \sum_{\{i, j\} \in E} w_{ij} v_i v_j \right) \quad v_i \in \{0, 1\}$$



対数尤度関数  $\boldsymbol{\theta} = \{\mathbf{b}, \mathbf{w}\}$

$$l_D(\boldsymbol{\theta}) = \sum_{i \in V} b_i \mathbb{E}_D[v_i] + \sum_{\{i, j\} \in E} w_{ij} \mathbb{E}_D[v_i v_j] - \ln Z(\boldsymbol{\theta})$$

$$\mathbb{E}_D[\dots] \leftarrow \text{観測データの標本平均} \quad \left( \text{e.g. } \mathbb{E}_D[v_i] = \frac{1}{N} \sum_{\mu=1}^N d_i^{(\mu)}, \quad \mathbb{E}_D[v_i v_j] = \frac{1}{N} \sum_{\mu=1}^N d_i^{(\mu)} d_j^{(\mu)} \right)$$

勾配

$$\mathbb{E}[\dots | \boldsymbol{\theta}] = \sum_{\mathbf{v}} (\dots) P(\mathbf{v} | \boldsymbol{\theta})$$

$$\frac{\partial l_D(\boldsymbol{\theta})}{\partial b_i} = \mathbb{E}_D[v_i] - \mathbb{E}[v_i | \boldsymbol{\theta}], \quad \frac{\partial l_D(\boldsymbol{\theta})}{\partial w_{ij}} = \mathbb{E}_D[v_i v_j] - \mathbb{E}[v_i v_j | \boldsymbol{\theta}]$$

# モーメントマッチング方程式

$$\frac{\partial l_D(\boldsymbol{\theta})}{\partial b_i} = E_D[v_i] - E[v_i | \boldsymbol{\theta}], \quad \frac{\partial l_D(\boldsymbol{\theta})}{\partial w_{ij}} = E_D[v_i v_j] - E[v_i v_j | \boldsymbol{\theta}]$$

対数尤度関数を最大化すればよいので**勾配上昇法**

ところで最大点(勾配が0の点)では。。

指数分布族と呼ばれる分布族をモデルにすると一般的に表れる方程式

$$E_D[v_i] = E[v_i | \boldsymbol{\theta}], \quad E_D[v_i v_j] = E[v_i v_j | \boldsymbol{\theta}]$$

(データの標本平均) = (モデルの対応する期待値)

→ **モーメントマッチング方程式** !

複雑な関数の最大化問題(最小化問題)における

常套手段は

## 勾配法

を利用したコンピュータによる

繰り返し計算による数値的最適化法

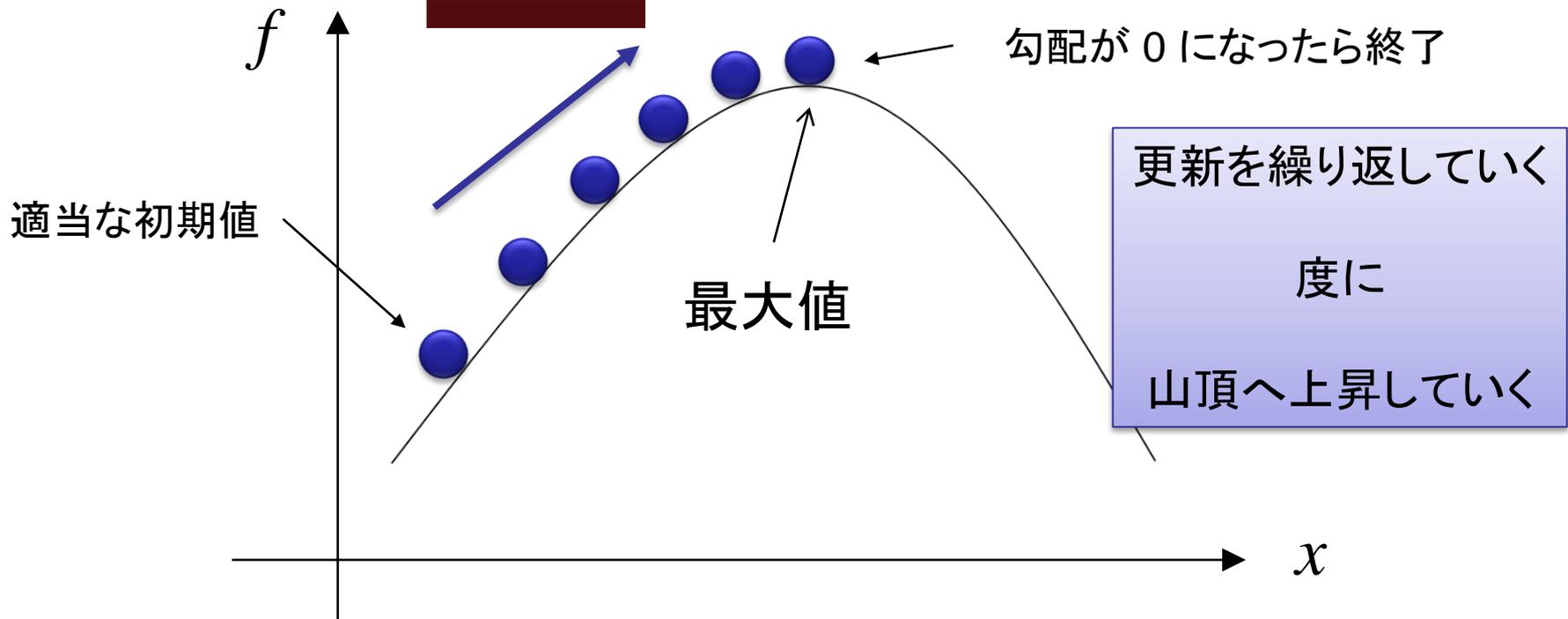
山に登らせる  
更新式

$$x^{\text{new}} \leftarrow x^{\text{old}} + \varepsilon \frac{\partial f(x^{\text{old}})}{\partial x}$$

勾配  
(導関数)



$\varepsilon$ : 小さな正の数 (学習率)



というわけで

コンピュータを使って勾配上昇法

$$b_i^{\text{new}} \leftarrow b_i^{\text{old}} + \varepsilon \left( \mathbb{E}_D [v_i] - \mathbb{E} [v_i | \boldsymbol{\theta}^{\text{old}}] \right)$$

$$w_{ij}^{\text{new}} \leftarrow w_{ij}^{\text{old}} + \varepsilon \left( \mathbb{E}_D [v_i v_j] - \mathbb{E} [v_i v_j | \boldsymbol{\theta}^{\text{old}}] \right)$$

BM の対数尤度関数は  
パラメータに関して凹関数なので

原理的には勾配法で最適な学習解を得ることができる

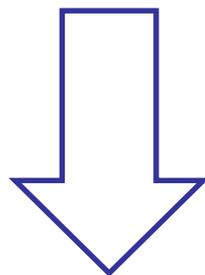
# マルコフ確率場 の欠点？

# 何故今になってマルコフ確率場？

57 / 141

MRF はコンピュータビジョン(CV)等のパターン認識分野での  
スタンダードモデルになり始めている

MRF 自体はかなり昔から知られたモデル



何故今になってマルコフ確率場が浮上してきたのか？  
何故今まではほとんど利用されてきてこなかったのか？

確率モデルにおいて期待値計算は基本中の基本！

変数の数を  $n$  とする

$$E[v_i | \theta] = \sum_{\mathbf{v}} v_i P(\mathbf{v} | \theta)$$

$$\sum_{\mathbf{v}} (\dots) = \sum_{v_1} \sum_{v_2} \dots \sum_{v_n} \leftarrow \text{実現可能なすべて値の組み合わせに関する和}$$

変数が 0 or 1 のバイナリだとすると

この和は  $2^n$  個の項の和になる

指数関数的な計算量増加はかなり危険！

仮に 1 秒間に 1 億回の足し算が可能であるPCがあるとする

$2^n$  回の演算に必要な時間は...



同様なロジックで  
規格化定数  $Z$  も  
そもそも計算できない



$n$	time
10	約 0.000001 秒
30	約 0.18 分
50	約 130 日
70	約 37万4千 年
100	約 4千億 年

指数関数的な計算量増加はかなり危険！

組み合わせ爆発の問題により

MRF の実装は

非現実的と考えられていた

BMの**近似学習アルゴリズム**は多数知られている

確率伝搬法 (belief propagation) [Yasuda & Horiguchi, 06]

コントラストティブ・ダイバージェンス (contrastive divergence; CD) [Hinton, 02]

疑似最尤法 (pseudo likelihood estimation; PLE) [Basag, 75]

複合最尤法 (composite likelihood estimation) [Lindsay, 88]

スコア・マッチング (score matching) [Hyvärinen, 05]

最小確率流 (minimum probability flow) [Sohl-Dickstein *et. al.*, 11]

拡張モンテカルロ積分法 (spatial Monte Carlo integration) [Yasuda, 15]

etc.

## 確率変数が連続の実数値の場合の MRF

(連続値のボルツマンマシン)

$$P(\mathbf{v} | \boldsymbol{\lambda}, \mathbf{b}, \mathbf{w}) = \frac{1}{Z(\boldsymbol{\lambda}, \mathbf{b}, \mathbf{w})} \exp\left(-\frac{1}{2} \sum_{i \in V} \lambda_i v_i^2 + \sum_{i \in V} b_i v_i + \sum_{\{i, j\} \in E} w_{ij} v_i v_j\right) \quad (\lambda_i > 0)$$

一般形

$v_i \in (-\infty, +\infty) \leftarrow$  確率変数が実数値

規格化定数が解析的に計算できる！

$$Z(\boldsymbol{\lambda}, \mathbf{b}, \mathbf{w}) = \exp\left(\frac{1}{2} \mathbf{b}^T \mathbf{C}^{-1} \mathbf{b}\right) \sqrt{(2\pi)^{|V|} \det \mathbf{C}^{-1}}$$

$$C_{ij} = \begin{cases} \lambda_i & (i = j) \\ -w_{ij} & (\{i, j\} \in E) \\ 0 & (\text{otherwise}) \end{cases}$$

ガウシアングラフィカルモデル  
(Gaussian graphical model; GGM)

GGM は多次元のガウス分布

$$P(\mathbf{v} | \boldsymbol{\lambda}, \mathbf{b}, \mathbf{w}) = \frac{1}{\sqrt{(2\pi)^{|\mathbf{v}|} \det \mathbf{C}^{-1}}} \exp\left(-\frac{1}{2}(\mathbf{v} - \mathbf{m})^T \mathbf{C}(\mathbf{v} - \mathbf{m})\right)$$

$$\mathbf{m} = \mathbf{C}^{-1}\mathbf{b}$$

期待値もとても簡単に計算可能

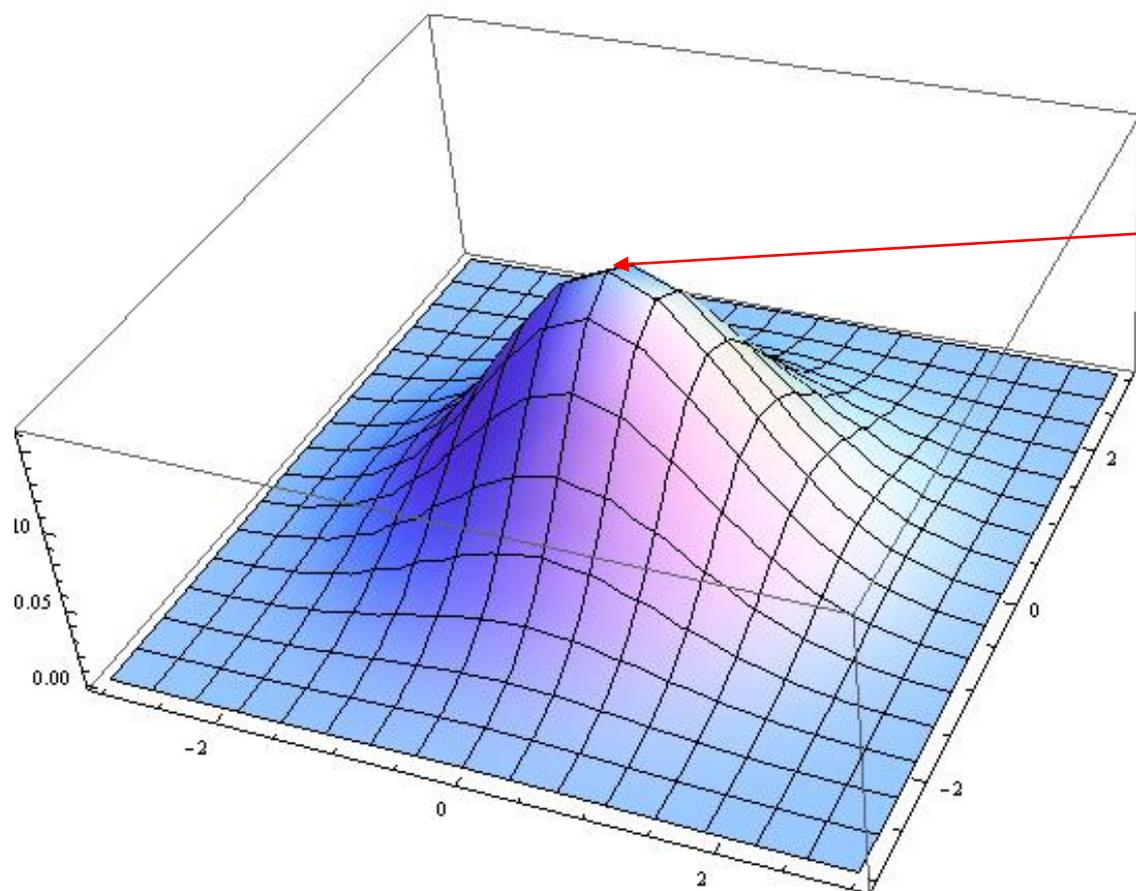
$$E[v_i | \boldsymbol{\theta}] = \int_{-\infty}^{\infty} v_i P(\mathbf{v} | \boldsymbol{\lambda}, \mathbf{b}, \mathbf{w}) d\mathbf{v} = m_i$$

$$E[v_i^2 | \boldsymbol{\theta}] = \int_{-\infty}^{\infty} v_i^2 P(\mathbf{v} | \boldsymbol{\lambda}, \mathbf{b}, \mathbf{w}) d\mathbf{v} = C_{ii}^{-1} + m_i^2$$

$$E[v_i v_j | \boldsymbol{\theta}] = \int_{-\infty}^{\infty} v_i v_j P(\mathbf{v} | \boldsymbol{\lambda}, \mathbf{b}, \mathbf{w}) d\mathbf{v} = C_{ij}^{-1} + m_i m_j$$

$C_{ij}^{-1} \leftarrow$  行列  $C$  の逆行列  $C^{-1}$  の  $(i, j)$  成分

$$P(\mathbf{v} | \lambda, \mathbf{b}, \mathbf{w}) = \frac{1}{\sqrt{(2\pi)^{|\mathbf{v}|} \det \mathbf{C}^{-1}}} \exp\left(-\frac{1}{2}(\mathbf{v} - \mathbf{m})^T \mathbf{C}(\mathbf{v} - \mathbf{m})\right)$$



ピーク点は平均ベクトル

$$\mathbf{m} = \mathbf{C}^{-1}\mathbf{b}$$

**確率最大の点も  
簡単に分かる！**

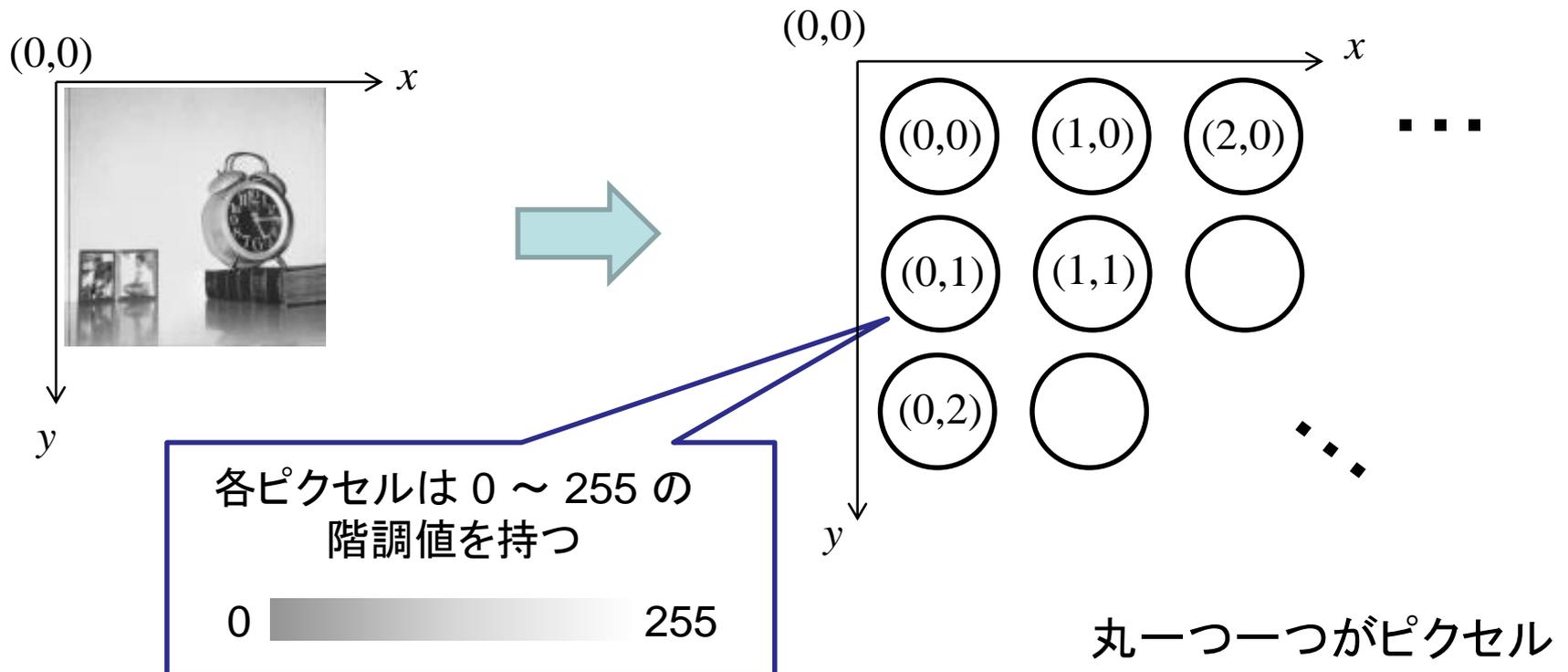
ガウシアングラフィカルモデル  
ならば

特別な近似なしに

期待値計算と最大値計算  
が厳密に可能である

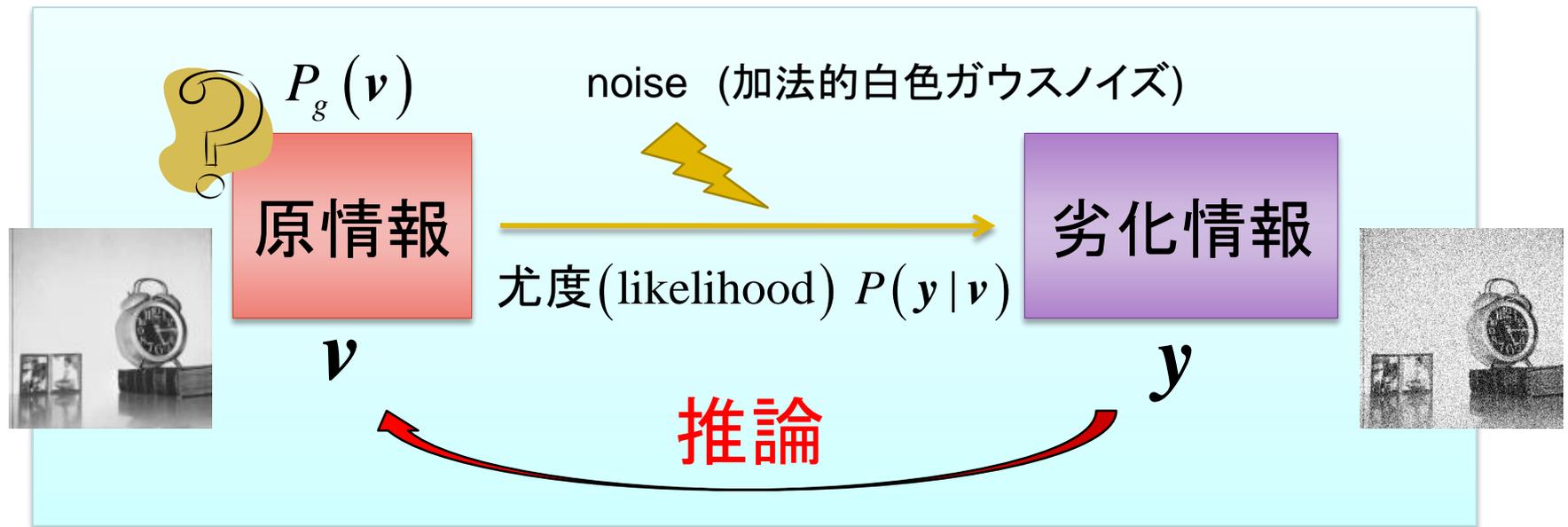
# ベイズ推定とMRF の応用例

- 碁盤の目(正方格子)状に並べられた**ピクセル**(pixel)とよばれる画素の集まりで表現される
- 各ピクセルには**階調値**(輝度値)とよばれる輝度(色)が割り当てられる



# MRFによる画像ノイズ除去の枠組み

68 / 141



劣化した画像から原画像をベイズの枠組みにより推定する

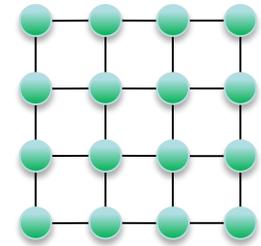
尤度はガウスノイズで設計することとするとしても

さて事前確率はどうしましょう？

→ MRF で設計することとしましょう

# MRFで画像の事前確率を作る

- 各ピクセルを各確率変数に対応させる



→ 確率変数はピクセルの個数分ある(超高次元)

- 確率変数のとる値は階調値にする  $v_i \in \{0, 1, 2, \dots, 255\}$
- ピクセルは正方格子状にならんでいるので同様に**正方格子のグラフ**上に定義する(各ノードが各ピクセル)
- リンクは最近接のピクセル対にのみあるとする

$$C(\mathbf{v}) = \underbrace{-\sum_{i \in V} \Phi_i(v_i)}_{\text{各ノードごとのコスト}} - \underbrace{\sum_{\{i, j\} \in E} \Psi_{\{i, j\}}(v_i, v_j)}_{\text{各リンクごとのコスト}}$$

エネルギー関数をどのように設計すればよいか??

$$C(\mathbf{v}) = \underbrace{-\sum_{i \in V} \Phi_i(v_i)}_{\text{各ノードごとのコスト}} - \underbrace{\sum_{\{i,j\} \in E} \Psi_{\{i,j\}}(v_i, v_j)}_{\text{各リンクごとのコスト}}$$

画像に対する我々の**先見知識**を基に設計する

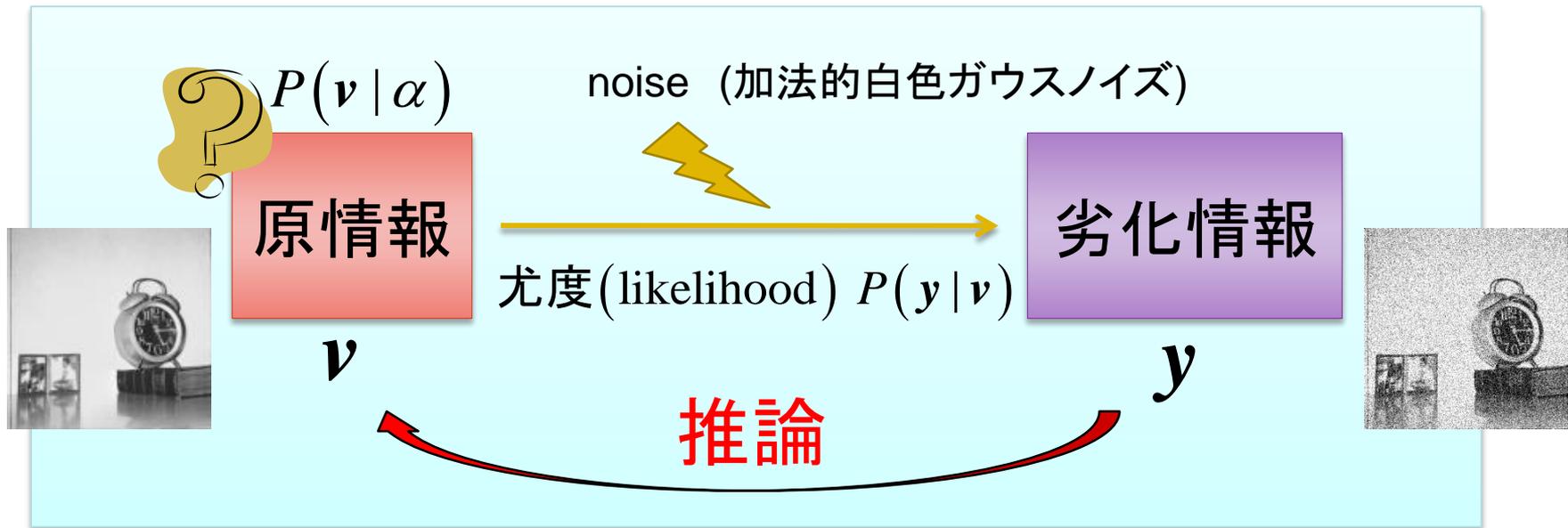
空間的に近くにあるピクセル同士は同じ色を取りやすいと仮定

リンクコスト:  $\Psi_{\{i,j\}}(v_i, v_j) \cdots v_i$  と  $v_j$  の値が近い程大きくなるようにする

例  $\Psi_{\{i,j\}}(v_i, v_j) = -\alpha |v_i - v_j|$ ,  $\Psi_{\{i,j\}}(v_i, v_j) = -\alpha (v_i - v_j)^2$  ( $\alpha > 0$ )

座標ごとの一貫した個性は特にないと仮定

$$\Phi_i(v_i) = 0$$



$$P_g(v) \xrightarrow{\text{MRFでモデル化}} P(v|\alpha) = \frac{1}{Z(\alpha)} \exp\left(-\alpha \sum_{\{i,j\} \in E} |v_i - v_j|\right)$$

ノイズは標準偏差  $\sigma$  のガウスノイズでピクセルごとに独立にのるとすると

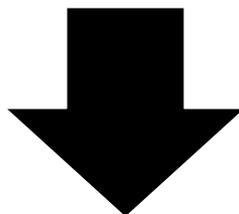
$$(\text{尤度}) = P(y|v, \sigma^2) = \prod_{i \in V} \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{(y_i - v_i)^2}{2\sigma^2}\right)$$

(事後確率)  $\propto$  (尤度)  $\times$  (事前確率)

$$P(\mathbf{v} | \mathbf{y}, \alpha, \sigma^2) \propto P(\mathbf{y} | \mathbf{v}, \sigma^2) P(\mathbf{v} | \alpha) \propto \exp \left( - \sum_{i \in V} \frac{(y_i - v_i)^2}{2\sigma^2} - \alpha \sum_{\{i, j\} \in E} |v_i - v_j| \right)$$

$x$  ... 原画像 (確率変数)

$y$  ... 劣化画像 (データ)



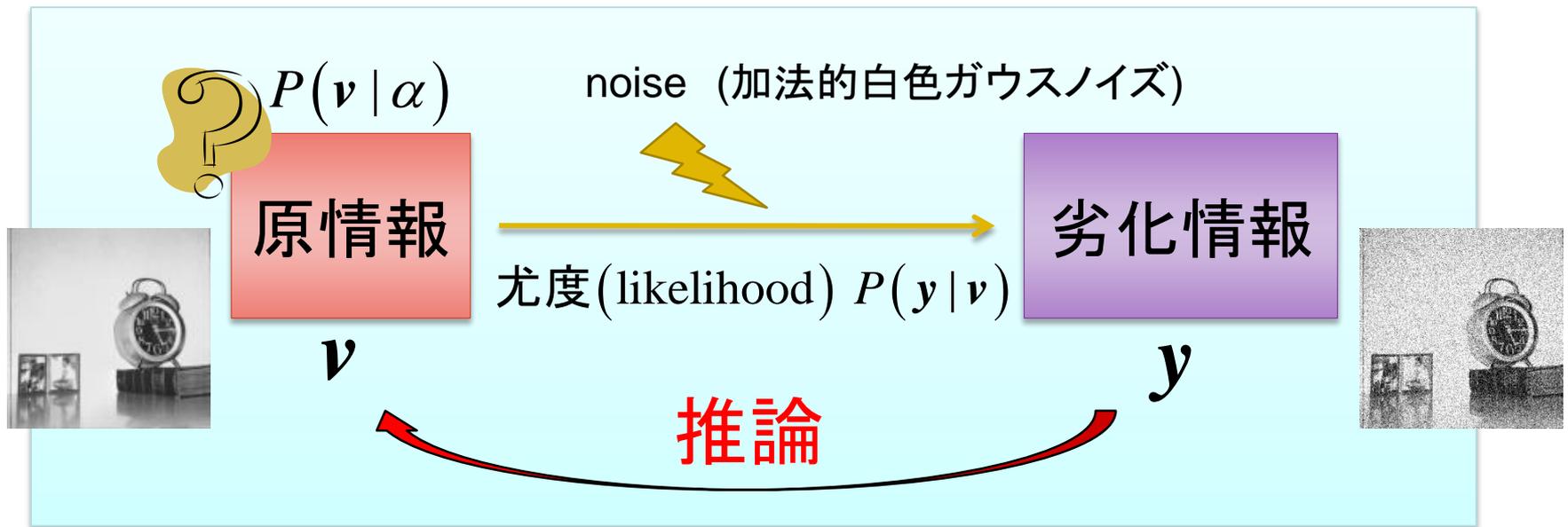
$$P(\mathbf{v} | \mathbf{y}, \alpha, \sigma^2) = \frac{1}{Z(\alpha, \sigma^2)} \exp \left( - \sum_{i \in V} \frac{(y_i - v_i)^2}{2\sigma^2} - \alpha \sum_{\{i, j\} \in E} |v_i - v_j| \right)$$

ノードのコスト  
(データ項)

リンクのコスト  
(平滑化項)

事後確率も MRF になる

# MRFによるノイズ除去フィルタ

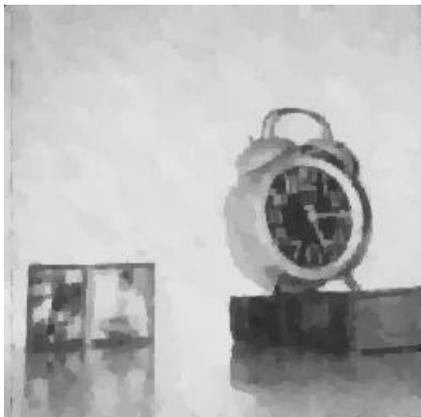
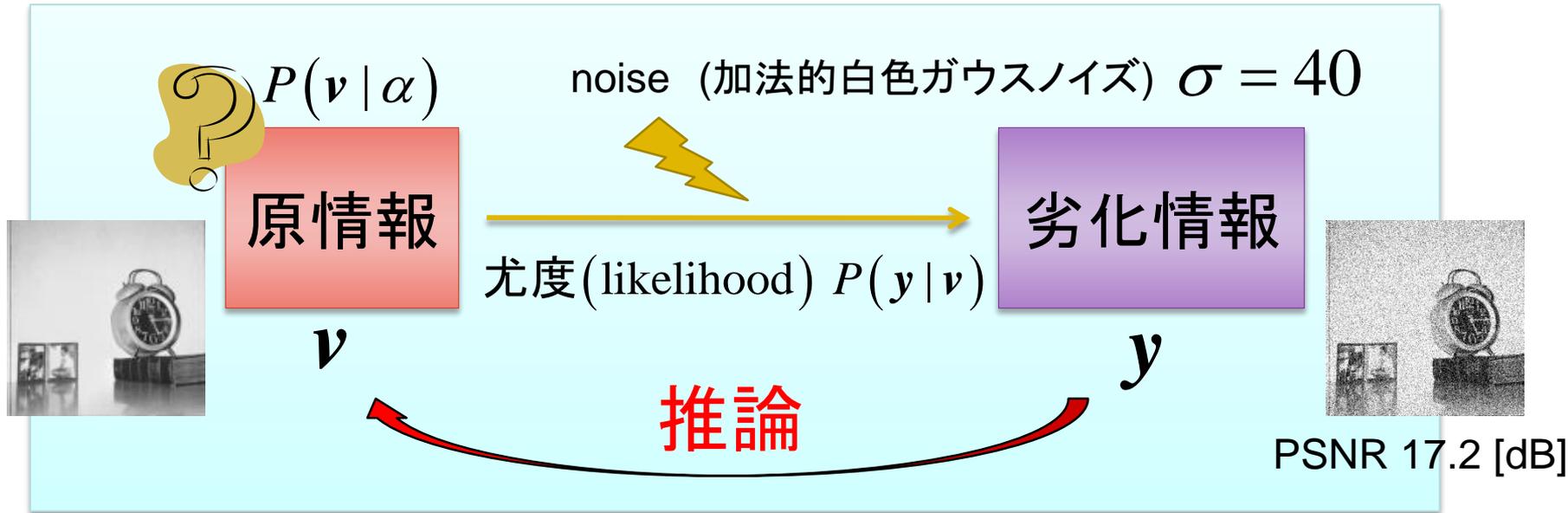


最大事後確率 (MAP) 推定

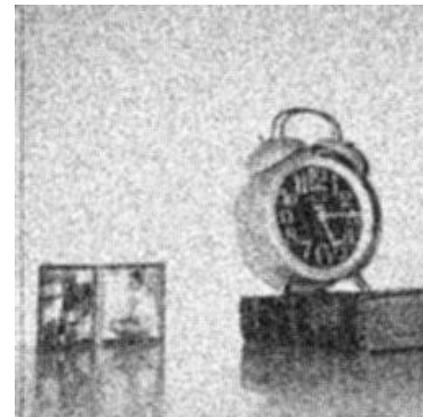
$$\left( \text{推定画像 } \mathbf{v}^* \right) = \arg \max_{\mathbf{v}} \underbrace{P(\mathbf{v} | \mathbf{y}, \alpha, \sigma^2)}_{\text{事後確率}}$$

max-product アルゴリズムにより最大化を実行

# MRFによるノイズ除去フィルタ



MRF PSNR 25.6 [dB]

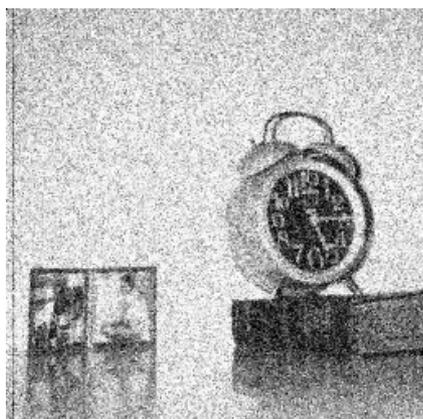


平滑化フィルタ PSNR 24.2 [dB]

$$P(\mathbf{v} | \mathbf{y}, \alpha, \sigma^2) = \frac{1}{Z(\alpha, \sigma^2)} \exp \left( - \sum_{i \in V} \frac{(y_i - v_i)^2}{2\sigma^2} - \alpha \sum_{\{i, j\} \in E} |v_i - v_j| \right)$$

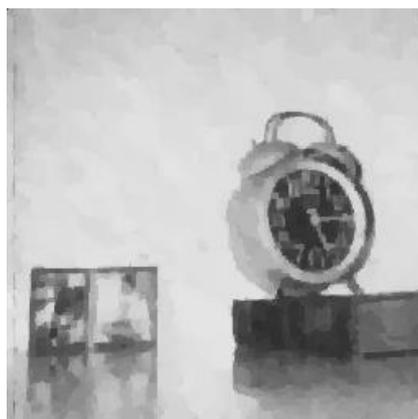
$\alpha$  ← データ項と平滑化項の重要度を調整するパラメータ

$\alpha = 0.002$



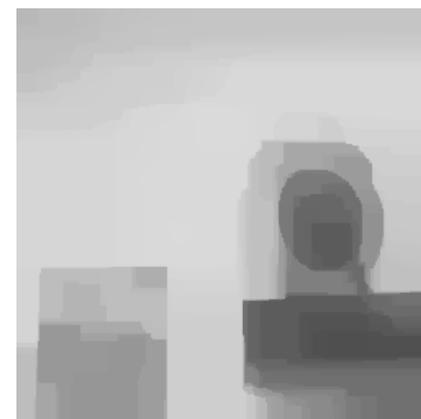
PSNR 19.0 [dB]

$\alpha = 0.02$



PSNR 25.6 [dB]

$\alpha = 0.2$



PSNR 18.8 [dB]

パラメータの値によって結果はかなり異なる

良い結果を得るためには  
パラメータの値を良いものにしないといけない  
しかし我々はその値を事前には知り得ない  
どうすればよいただろうか??

先程の例はノイズの分散も既知としたものである

しかし実際にはそれすら知り得ないこともある

$\alpha$  と  $\sigma^2$  の値をなんらかの方法で推定したい

# 結合確率から出発する

原画像  $\mathbf{v}$  と劣化画像  $\mathbf{y}$  の結合確率は

$$P(\mathbf{v}, \mathbf{y} | \alpha, \sigma^2) = P(\mathbf{y} | \mathbf{v}, \sigma^2) P(\mathbf{v} | \alpha)$$
$$= \left\{ \prod_{i \in V} \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{(y_i - v_i)^2}{2\sigma^2}\right) \right\} \times \left\{ \frac{1}{Z(\alpha)} \exp\left(-\alpha \sum_{\{i,j\} \in E} |v_i - v_j|\right) \right\}$$

$$P(\mathbf{v}, \mathbf{y} | \alpha, \sigma^2)$$

データとして  
観測されない変数

隠れ変数  
(潜在変数)

データ(劣化画像)として  
観測されている変数

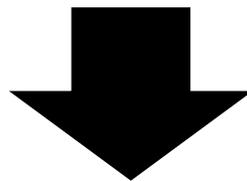
可視変数  
(観測変数)

決めたい  
パラメータ

データとして観測されない変数(隠れ変数)を含む

確率モデルのパラメータの値を

観測したデータ(劣化画像)を用いて決定したい！



expectation-maximization (EM) アルゴリズム

ある時刻  $t$  において次の **Q-関数** を定義する

$$Q(\theta_{t+1} | \theta_t, y) = \sum_x \underbrace{P(x | y, \theta_t)}_{\text{事後確率}} \ln \underbrace{P(x, y | \theta_{t+1})}_{\text{結合確率}}$$

$\theta_t = \{\alpha_t, \sigma_t^2\}$  ← パラメータをまとめて記述する

$\theta_t$  を既知とするとこの関数は  $\theta_{t+1}$  の関数となる

- (1)  $t = 0$  としてパラメータ  $\theta_0$  の値を適当に初期化
- (2)  $Q$ -関数  $Q(\theta_{t+1} | \theta_t, y)$  を計算する ← Eステップ
- (3)  $Q$ -関数  $Q(\theta_{t+1} | \theta_t, y)$  を  $\theta_{t+1}$  に関して最大化する ← Mステップ

$$\theta_{t+1}^{\max} = \arg \max_{\theta_{t+1}} Q(\theta_{t+1} | \theta_t, y)$$

- (4)  $\theta_t = \theta_{t+1}^{\max}$  となっていれば  $\theta_{t+1}^{\max}$  を結果として出力して終了  
そうでなければ  
 $\theta_t \leftarrow \theta_{t+1}^{\max}$  と値を更新して時刻を一つ進め ( $t \leftarrow t + 1$ ) て(2)に戻る

# ノイズ除去の場合のEステップ

$$Q(\boldsymbol{\theta}_{t+1} | \boldsymbol{\theta}_t, \mathbf{y}) = -\frac{1}{2\sigma_{t+1}^2} \sum_{i \in V} \mathbf{E}_{\text{pos}} \left[ (y_i - v_i)^2 \right]_t \\ - \alpha_{t+1} \sum_{\{i,j\} \in E} \mathbf{E}_{\text{pos}} \left[ |v_i - v_j| \right]_t - \frac{|V|}{2} \ln(2\pi\sigma_{t+1}^2) - \ln Z(\alpha_{t+1})$$

$$\mathbf{E}_{\text{pos}} [\dots]_t = \sum_{\mathbf{v}} (\dots) P(\mathbf{v} | \mathbf{y}, \boldsymbol{\theta}_t) \leftarrow \text{事後確率による期待値}$$

$|V| = (V \text{の要素の数})$   
 $= (\text{ピクセル数})$

最大化のステップ(Mステップ)のためパラメータに関する勾配を計算する

$$\left\{ \begin{array}{l} \frac{\partial Q(\boldsymbol{\theta}_{t+1} | \boldsymbol{\theta}_t, \mathbf{y})}{\partial \sigma_{t+1}^2} = \frac{1}{2\sigma_{t+1}^4} \sum_{i \in V} \mathbf{E}_{\text{pos}} \left[ (y_i - v_i)^2 \right]_t - \frac{|V|}{2\sigma_{t+1}^2} \\ \frac{\partial Q(\boldsymbol{\theta}_{t+1} | \boldsymbol{\theta}_t, \mathbf{y})}{\partial \alpha_{t+1}} = - \sum_{\{i,j\} \in E} \mathbf{E}_{\text{pos}} \left[ |v_i - v_j| \right]_t + \sum_{\{i,j\} \in E} \mathbf{E}_{\text{pri}} \left[ |v_i - v_j| \right]_{t+1} \end{array} \right.$$

$$\mathbf{E}_{\text{pri}} [\dots]_{t+1} = \sum_{\mathbf{v}} (\dots) P(\mathbf{v} | \alpha_{t+1}) \leftarrow \text{事前確率による期待値}$$

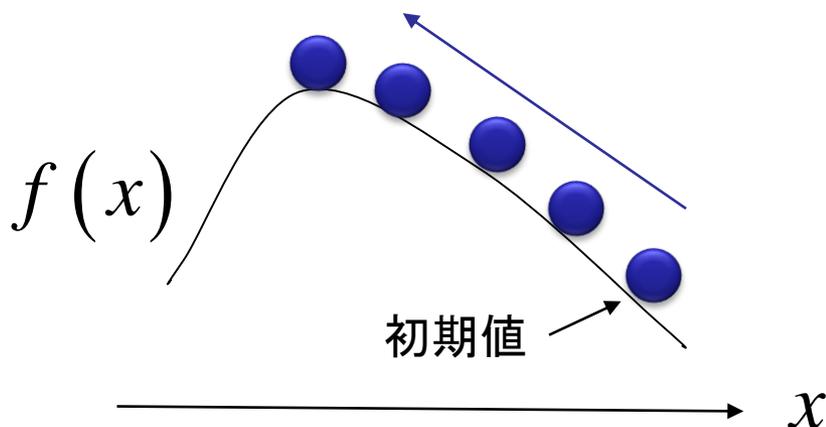
# ノイズ除去の場合のMステップ

82 / 141

$$\begin{cases} \frac{\partial Q(\boldsymbol{\theta}_{t+1} | \boldsymbol{\theta}_t, \mathbf{y})}{\partial \sigma_{t+1}^2} = \frac{1}{2\sigma_{t+1}^4} \sum_{i \in V} \mathbf{E}_{\text{pos}} \left[ (y_i - v_i)^2 \right]_t - \frac{|V|}{2\sigma_{t+1}^2} \\ \frac{\partial Q(\boldsymbol{\theta}_{t+1} | \boldsymbol{\theta}_t, \mathbf{y})}{\partial \alpha_{t+1}} = - \sum_{\{i,j\} \in E} \mathbf{E}_{\text{pos}} \left[ |v_i - v_j| \right]_t + \sum_{\{i,j\} \in E} \mathbf{E}_{\text{pri}} \left[ |v_i - v_j| \right]_{t+1} \end{cases}$$

各勾配は確率伝搬法 (sum-productアルゴリズム) で計算できる

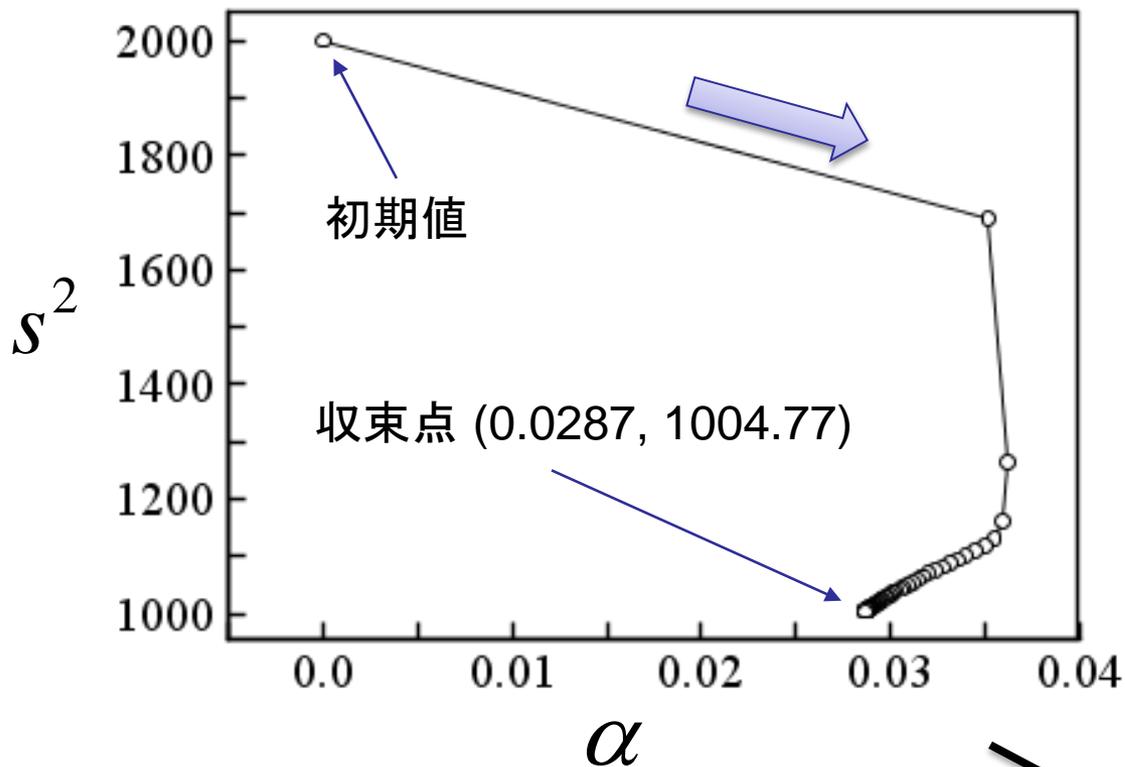
最大化はこれらの勾配を利用した**勾配上昇法**で実装可能！



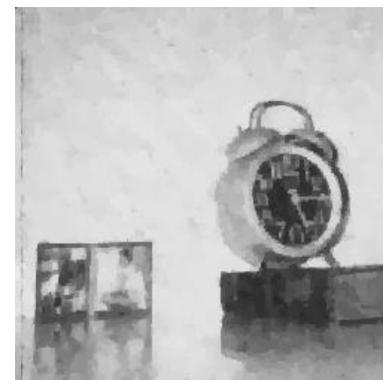
$$x^{\text{new}} \leftarrow x^{\text{old}} + \varepsilon \frac{\partial f(x^{\text{old}})}{\partial x}$$

$\varepsilon$ : 小さな正の数

# EMアルゴリズムの実行結果



それなりのパラメータが  
EMアルゴリズムにより  
推定された



PSNR 26.2 [dB]

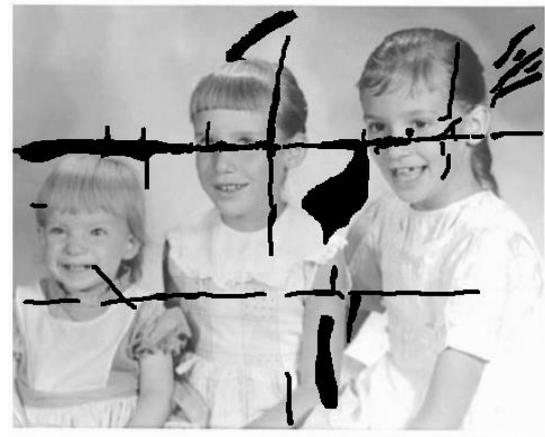
## 画像補修 (Image Inpainting) フィルタ

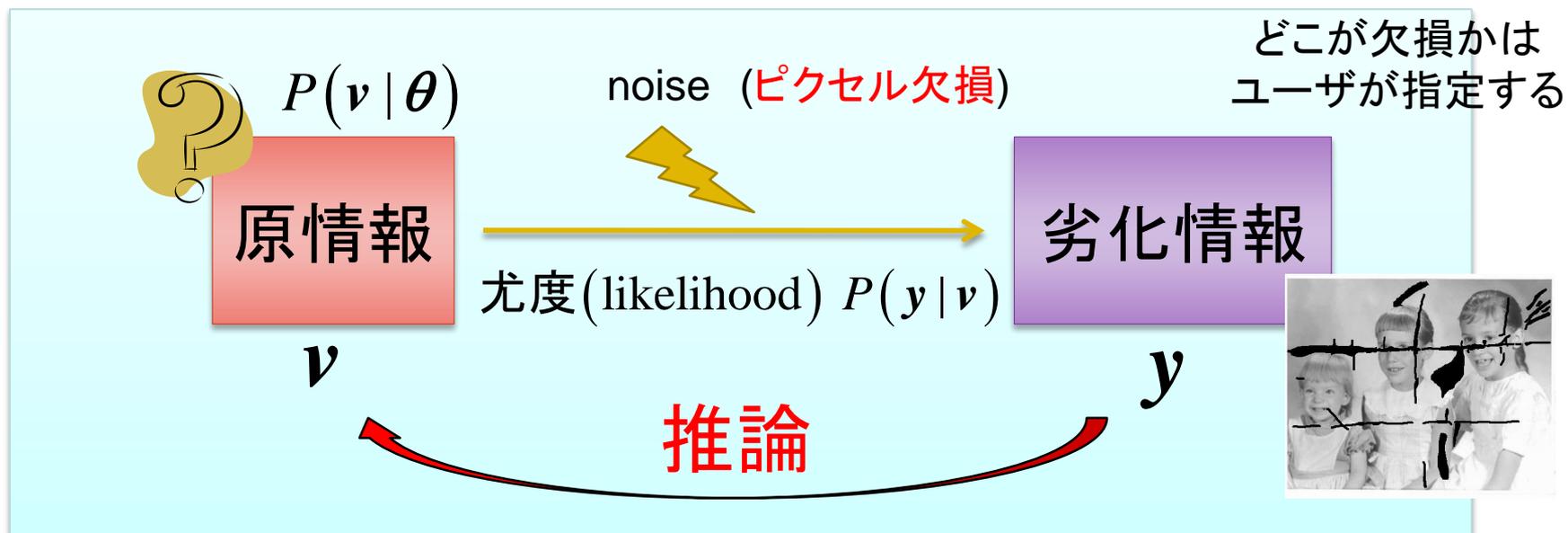
- 画像中の落書き傷を自動補修するための画像処理フィルタ
- 2000年初頭の提案から現在ではいくつかの商用ソフトに実装されてきている

大切な写真が傷付いてしまった



黒く塗りつぶした傷の部分を治したい





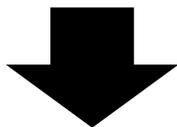
欠損した箇所情報はまったくない  
欠損していない箇所は原画像と同じ

デルタ関数

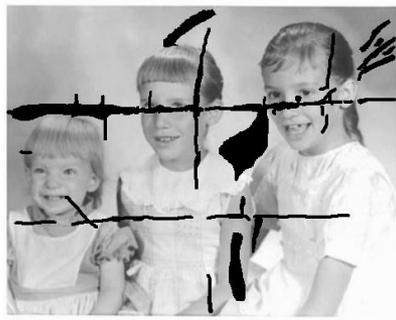
$$(\text{尤度}) = P(\mathbf{y} | \mathbf{v}) = \prod_{i \in (\text{非欠損部})} \delta(y_i, v_i)$$

# MRFを用いた画像補修フィルタ

大切な写真が傷付いてしまった



黒く塗りつぶした傷の部分を知りたい



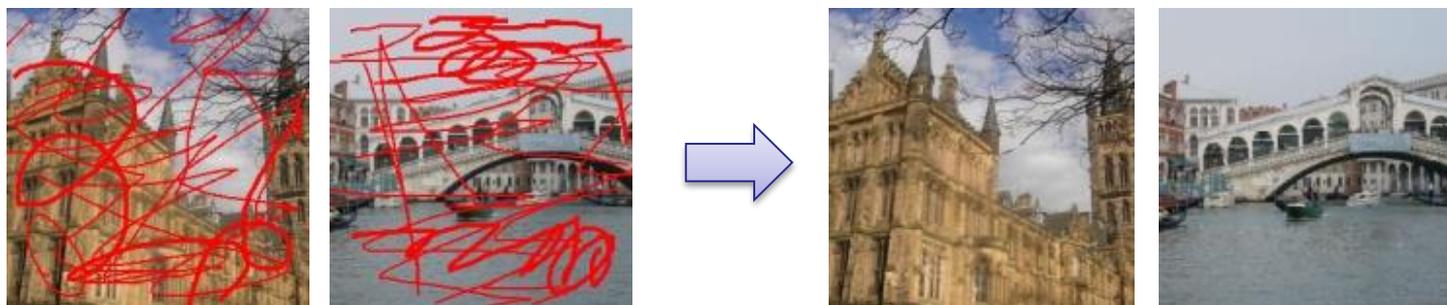
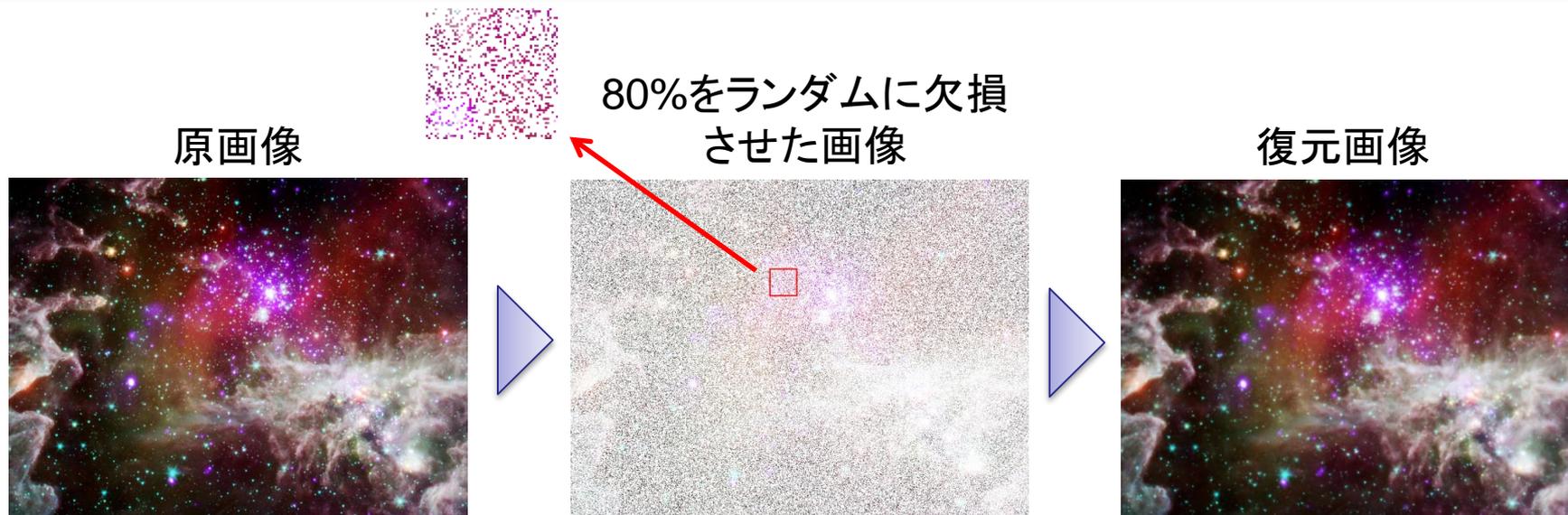
治った！！



[Yasuda et. al., 05, 13]

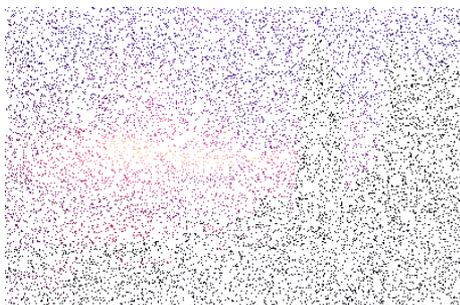
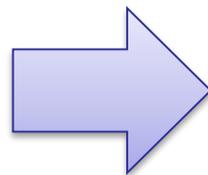
# MRFを用いた画像補修フィルタ

87 / 141



多くの自然画像を MRF で学習して  
それを事前分布の近似としたベイズ推定の方法で補修している

# MRFを用いた画像補修フィルタ



熊も城も断崖絶壁もなんでもござれ！！



仙台市

## 交通予測

車両感知器やプローブカーデータ等から得られる各道路の情報から  
現在(ナウ・キャスト)または未来(フォア・キャスト)の交通状態を予測すること

## 問題点

車両感知器・・・主要道路のみにしか設置されていない

プローブ情報・・・すべての車、道路から情報が得られるわけではない

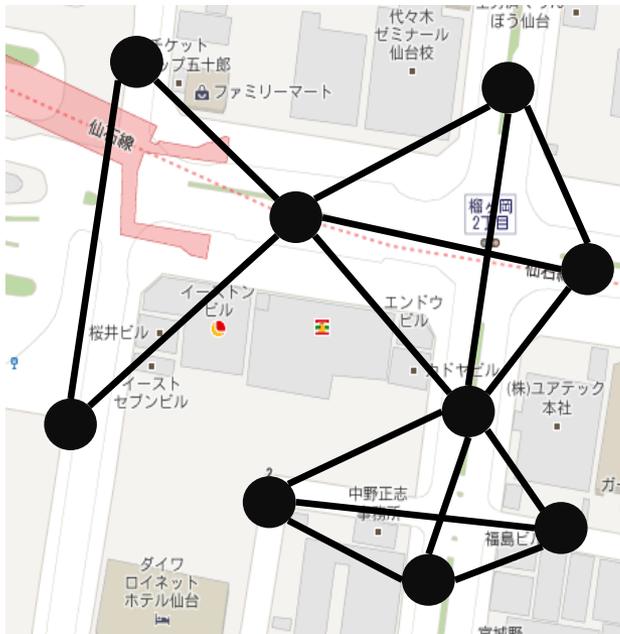


部分的な道路の情報から全体の交通状態を予測したい

# 道路交通に対するMRF

ノード ... 各道路上に配置 ← 予測対象が道路のため、道路がノードに対応

リンク ... 直接移動可能な道路同士を結ぶ ← 隣接道路間の平滑化性  
(一本前の道路が混雑していたら次の道路も混雑しやすいという仮定)



グラフの構造

各道路上に確率変数  $v$  を配置する

確率変数は配置された道路の交通密度を表わす

$$(\text{交通密度}) = (\text{車の台数}) \div (\text{道路面積})$$

確率変数は任意の実数値をとるとする

$$v_i \in (-\infty, +\infty)$$

## 道路交通に対する事前確率

$$P(\mathbf{v} | \mathbf{b}, \lambda, \alpha) \propto \exp \left( -\frac{\lambda}{2} \sum_{i \in V} v_i^2 + \sum_{i \in V} b_i v_i - \alpha \sum_{\{i, j\} \in E} (v_i - v_j)^2 \right)$$

バイアス項

平滑化項

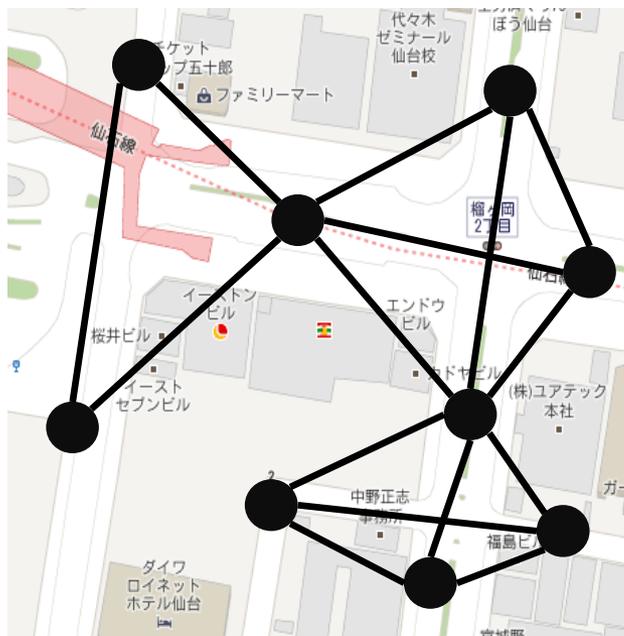
- 道路毎の個性を表わす
- 潜在的に混みやすい道路とそうでない道路があるであろう

つながった道路同士は似たような混雑状況であろう

以上のような仮定の下で MRF を定義する

# MRFを学習する

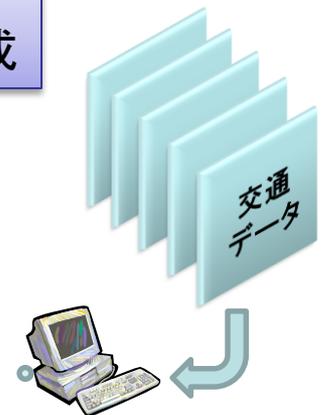
- ノード ... 各道路上に配置 ← 予測対象が道路のため、道路がノードに対応
- リンク ... 直接移動可能な道路同士を結ぶ ← 隣接道路間の平滑化性  
(一本前の道路が混雑していたら次の道路も混雑しやすいという仮定)



グラフの構造

交通シミュレーターSOUNDで  
交通データ(疑似現実データ)を生成

$$P(v | \theta)$$



統計的機械学習を利用して  
8時～9時の1時間分の交通データからモデルを設計

テストは同時間帯の新しいデータを用いておこなう

# MRFによる交通量予測の枠組み

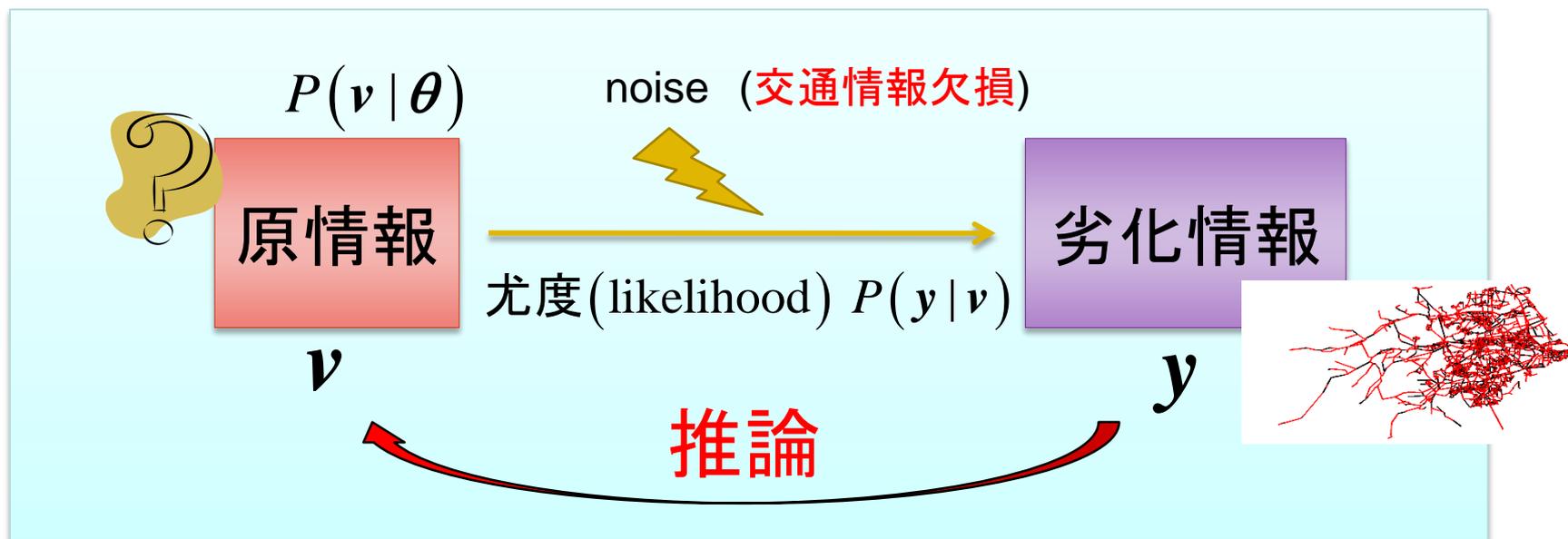
93 / 141



未観測道路(欠損道路)をランダムに70%選択

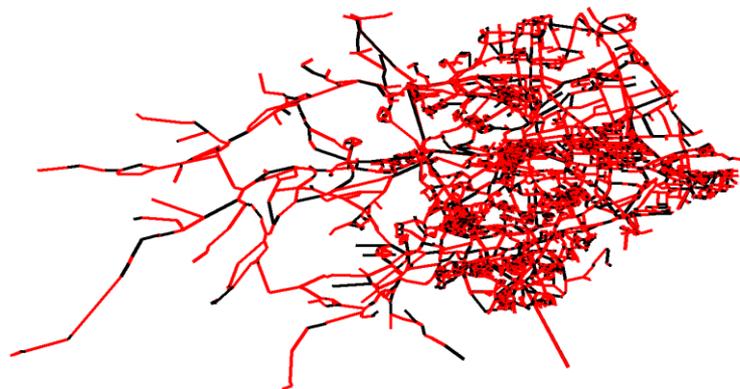
赤 ... 未観測道路

黒 ... 観測道路



枠組みは画像補修とまったく同じ！

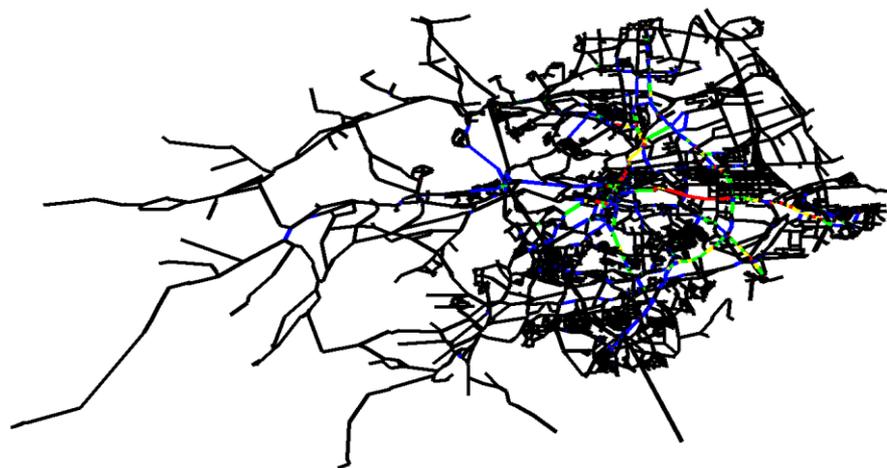
# MRFによる交通流予測



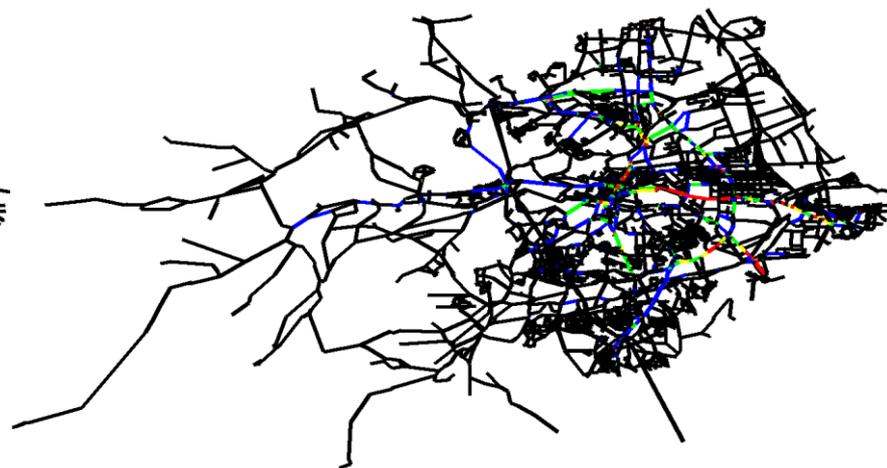
未観測道路(欠損道路)をランダムに70%選択

赤 ... 未観測道路

黒 ... 観測道路



正解のデータ



推定結果

黒 < 青 < 緑 < 黄 < 赤  
少 ← 交通密度 → 多

これも先の画像の例と同様に  
ベイズ推定の枠組みを用いて予測している

# MRFによる交通流予測

仙台中心部



未観測道路 (赤:未観測)



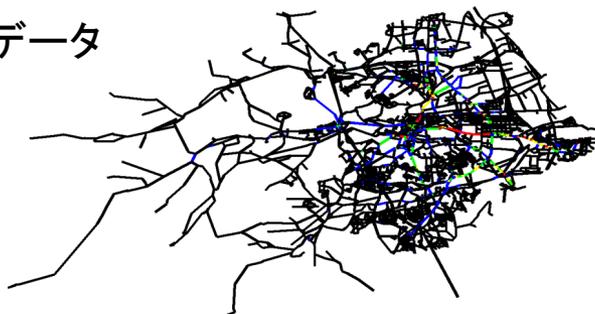
正解のデータ



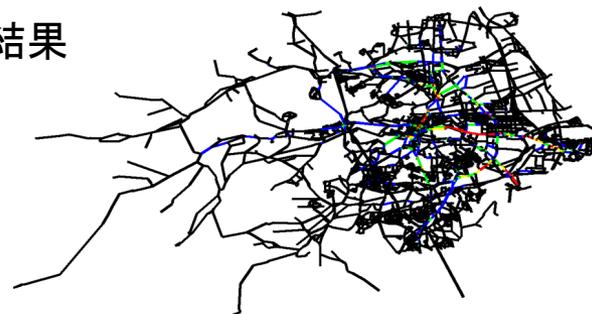
推定結果

# MRFによる交通流予測

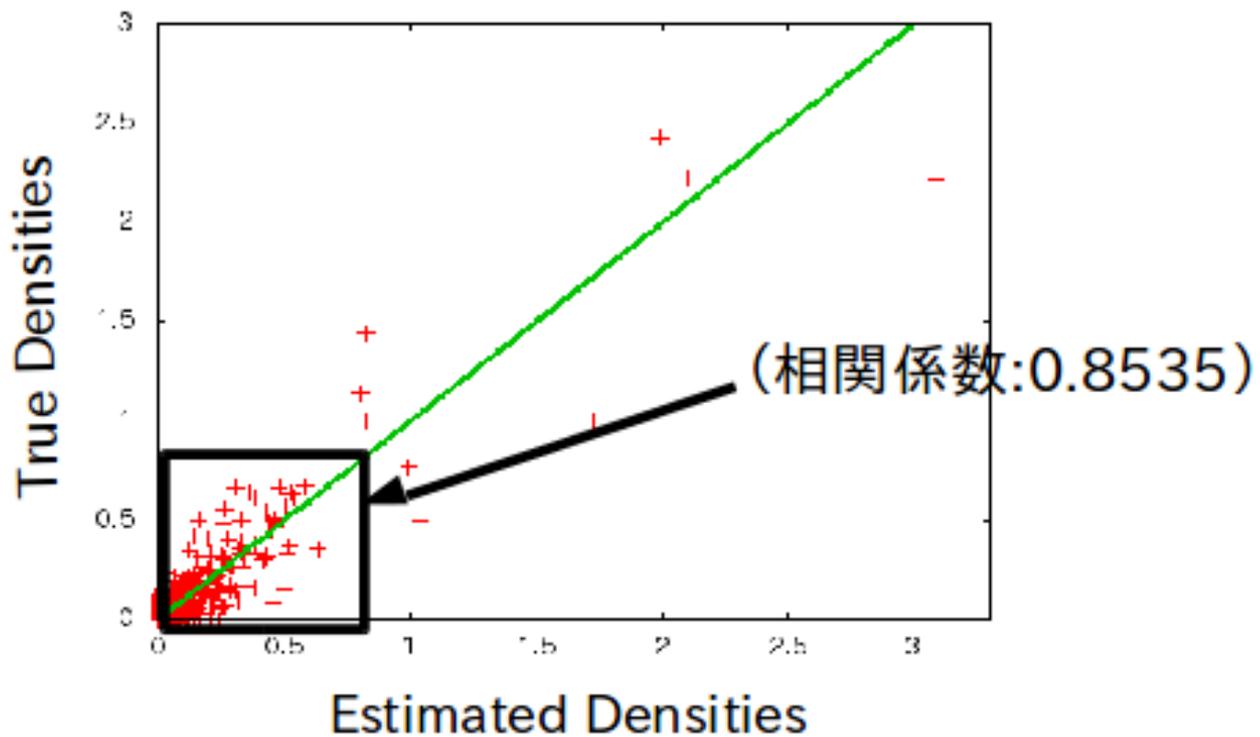
正解のデータ



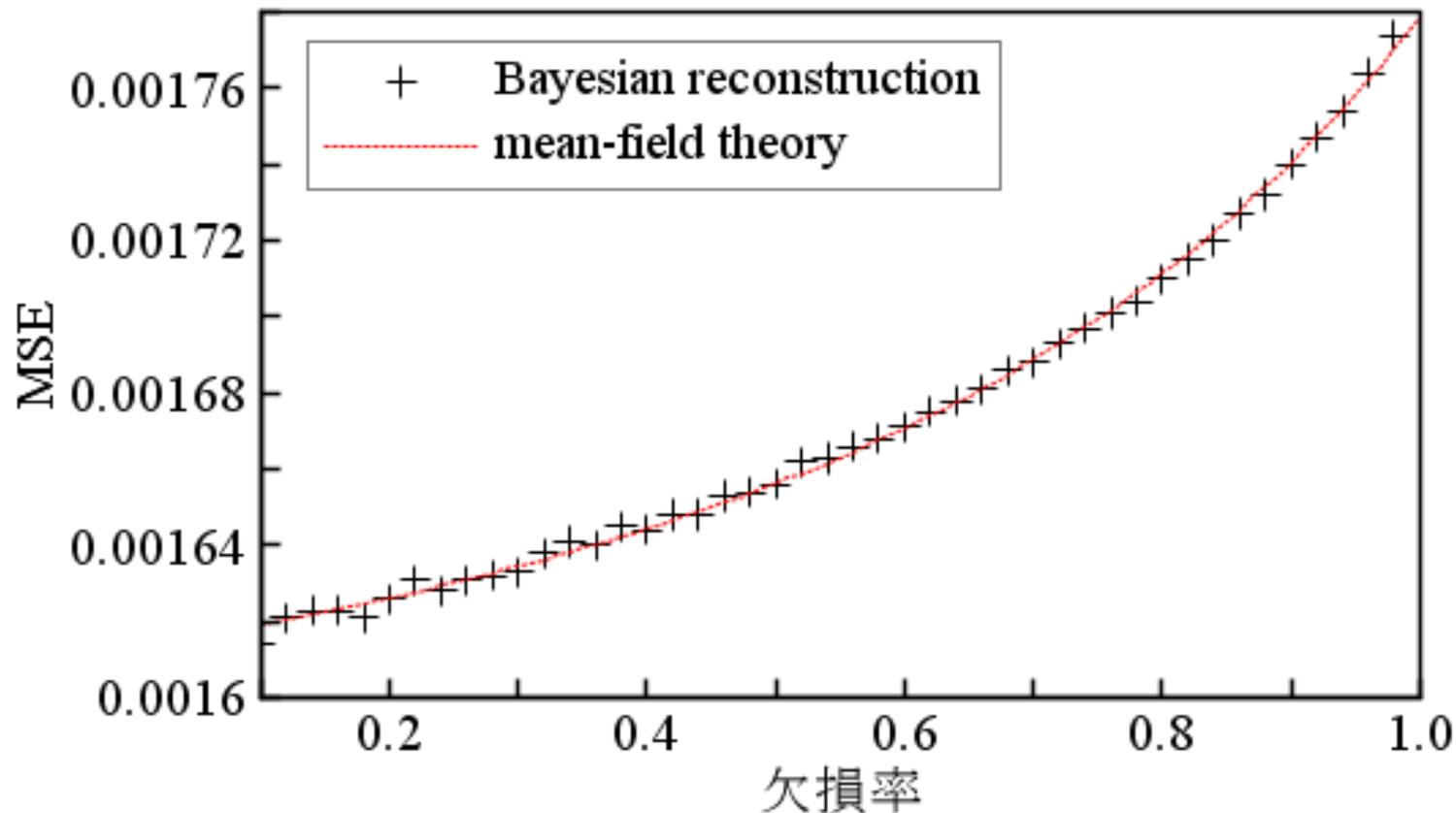
推定結果



真値と推定値の比較 (相関係数:0.9168)



## 欠損率と推定値のMSEの比較



赤の破線・・・平均場解析による理論予測曲線

[Kataoka, Yasuda and Tanaka, 15]

# ネットワークマイニング への応用とその他

$$P(\mathbf{v} | \mathbf{b}, \mathbf{w}) = \frac{1}{Z(\mathbf{b}, \mathbf{w})} \exp \left( \sum_{i \in V} b_i v_i + \sum_{\{i, j\} \in E} w_{ij} v_i v_j \right)$$

$$v_i \in \{-1, +1\}$$

確率変数が  $\{+1, -1\}$  の 2 値をとるボルツマンマシンを考えてみる

## パラメータが作る局所的な構造

- もし  $b_i > 0$  ならば局所的には  $v_i = +1$  で逆に  $b_i < 0$  なら  $v_i = -1$  が確率的に高い
- もし  $w_{ij} > 0$  ならば  $v_i$  と  $v_j$  が同じ値の方が局所的には確率が高い
- もし  $w_{ij} < 0$  ならば  $v_i$  と  $v_j$  が同じ値でない方が局所的には確率が高い

## 生成モデルのより深い構造理解 (データマイニング)

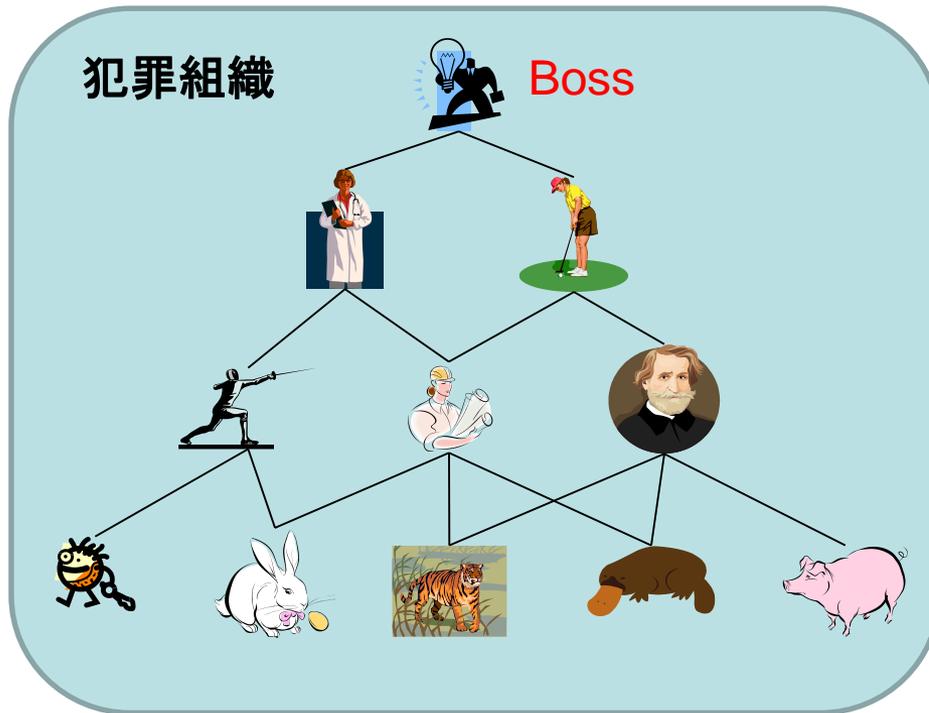
### リンク構造に注目する

リンクの重みはノードの局所的な関連性を表現している

$$w_{ij} \begin{cases} > 0: \text{ノード } i \text{ と } j \text{ は局所的に正の相関} \\ = 0: \text{ノード } i \text{ と } j \text{ は局所的には関連しない} \\ < 0: \text{ノード } i \text{ と } j \text{ は局所的に負の相関} \end{cases}$$

学習結果の重みを調べることで**データの詳細な構造**がわかる

## とある犯罪組織を調査する刑事さんがいる



ボスが誰かは判明しているが

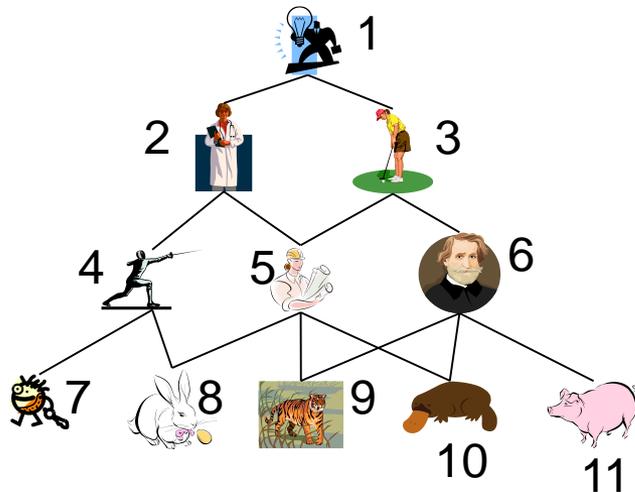
幹部と末端のつながりが分からないので

それを調査している

つまり刑事さんは左の組織図のリンク構造を知らない

つながりは命令系統をあらわしており  
リンクでつながった者同士は同時期に活動することが(確率的に)多いとする

# 犯罪組織の活動モデル



人(?)をノードとし, つながりをリンクとする

それぞれノードは活動中(+1)・休止中(-1)の2値をとるものとする

それぞれのノードは隣接するノードと同じ状態を取りやすいとする

## モデリング

$$\begin{cases} C_g(\mathbf{v}) = - \sum_{\{i,j\} \in E} v_i v_j \\ P_g(\mathbf{v}) = \frac{1}{Z_g} \exp(-C_g(\mathbf{v})) \end{cases}$$

この犯罪組織の活動状況をあらわす真のモデルとする

刑事さんはもちろんこの分布を知らない  
唯一知っているのは懸命に調べた**組織の活動記録**のみ

# 活動記録からつながりを推定

組織の活動モデル

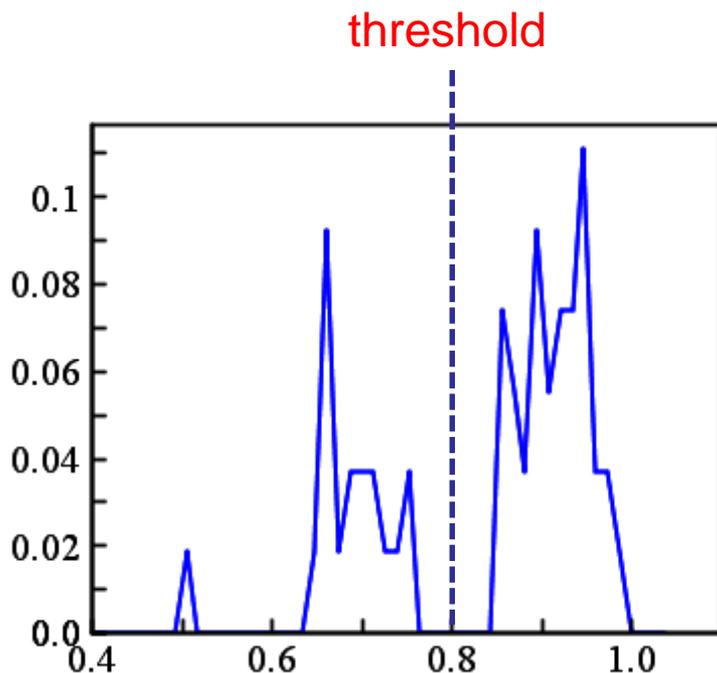
10000 データ点

$$P_g(\mathbf{v}) = \frac{1}{Z_g} \exp\left(\sum_{\{i,j\} \in E} v_i v_j\right)$$

MCMCで生成

$-1, 1, 1, 1, -1, -1, 1, \dots, -1, 1$   $\mathbf{d}^{(1)}$   
 $1, -1, -1, 1, 1, -1, 1, \dots, -1, -1$   $\mathbf{d}^{(2)}$   
 $\vdots$   
 $-1, 1, 1, 1, 1, -1, 1, \dots, -1, -1$   $\mathbf{d}^{(10000)}$

組織の活動記録



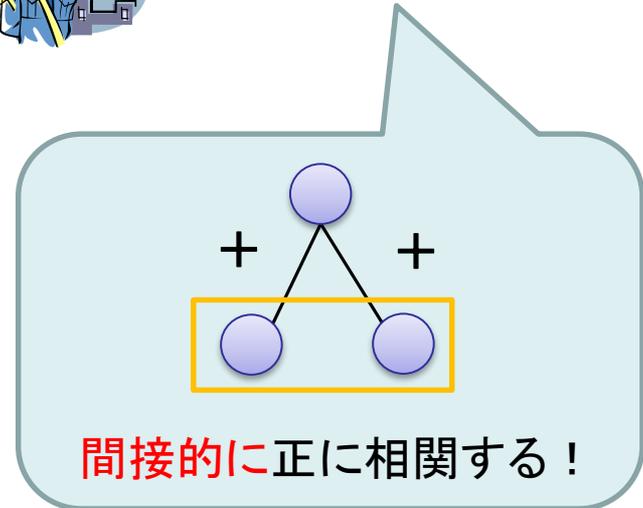
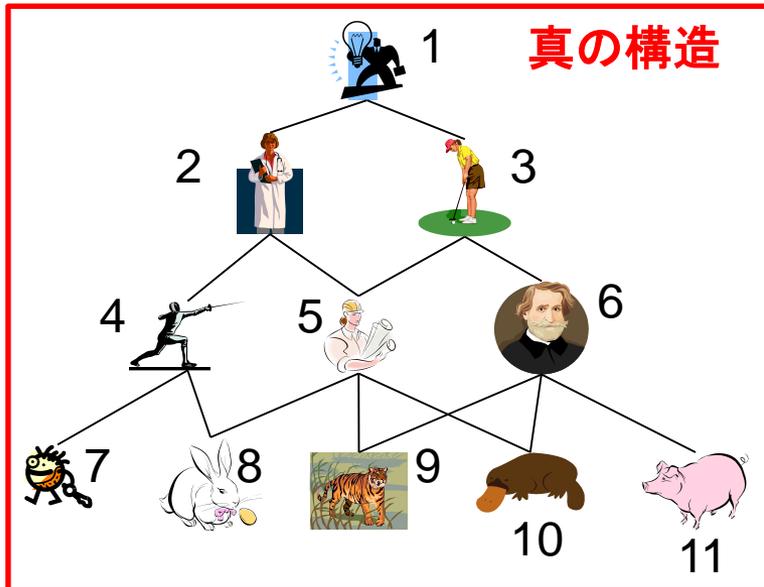
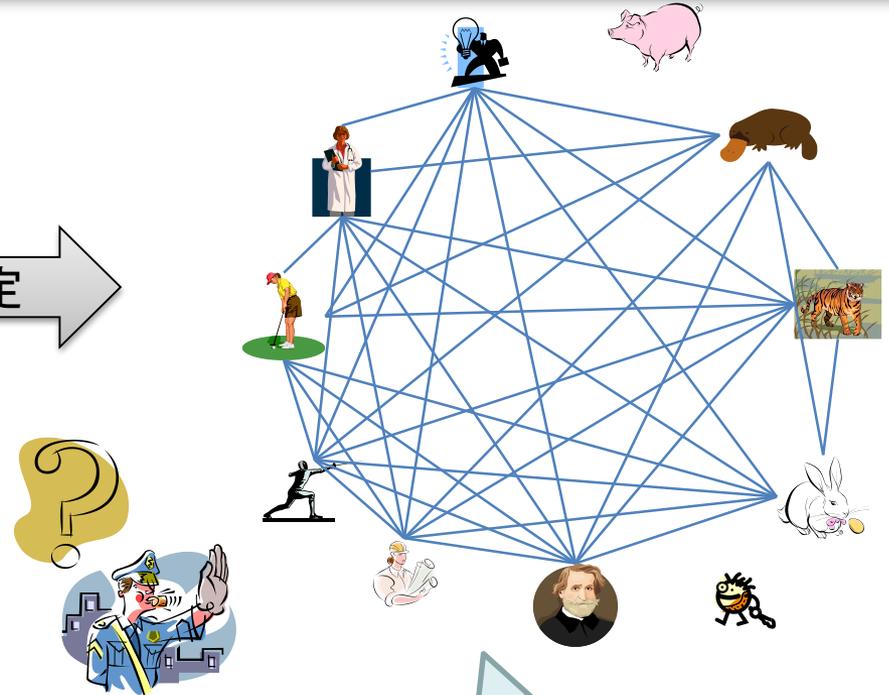
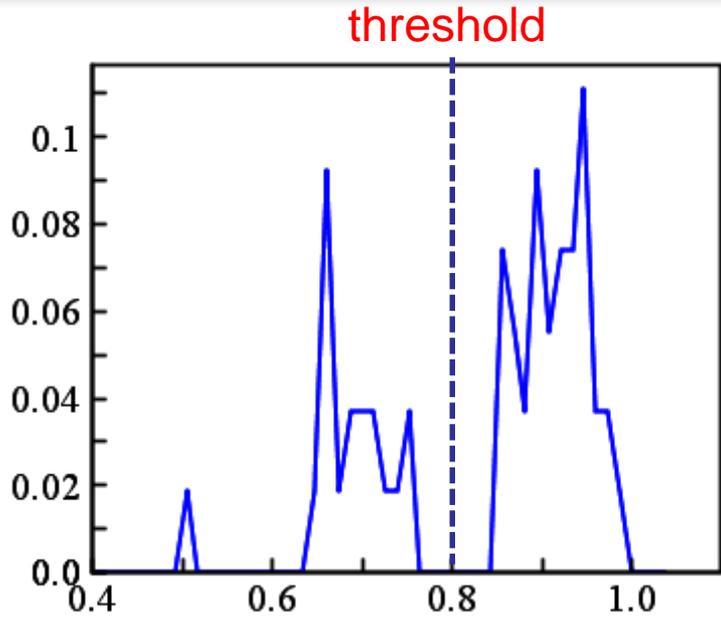
単純な例として

データから**相関係数**を計算し

0.8以上の値をもつペアのみに注目する

すべてのノードペアの相関係数のヒストグラム

# 単純な方法では大失敗



単純な相関係数ではこのような間接的な相関も検知してしまう

## ボルツマンマシンによる推定

$-1, 1, 1, 1, -1, -1, 1, \dots, -1, 1$   $\mathbf{d}^{(1)}$

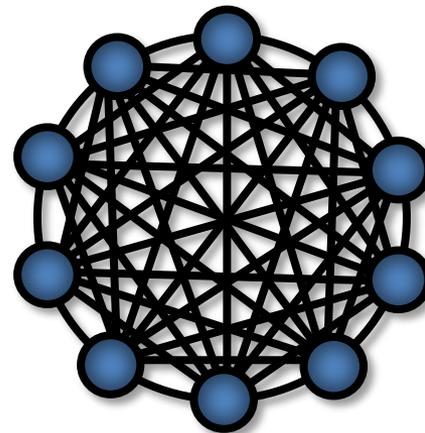
$1, -1, -1, 1, 1, -1, 1, \dots, -1, -1$   $\mathbf{d}^{(2)}$

$\vdots$

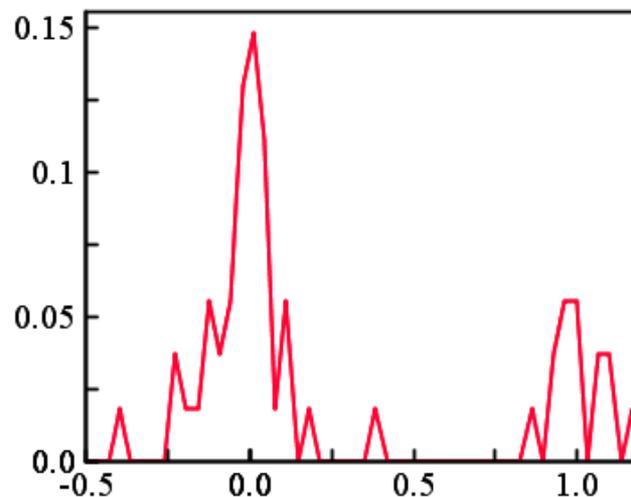
$-1, 1, 1, 1, 1, -1, 1, \dots, -1, -1$   $\mathbf{d}^{(10000)}$

学習

完全結合のボルツマンマシン

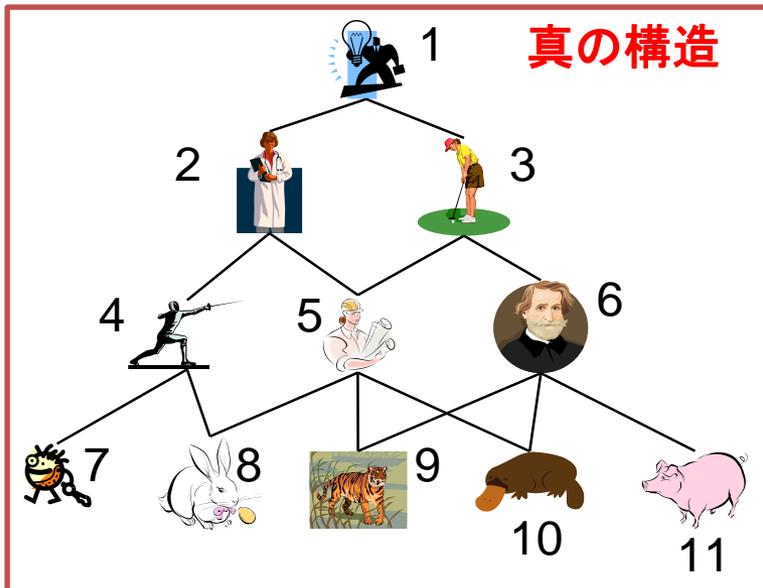
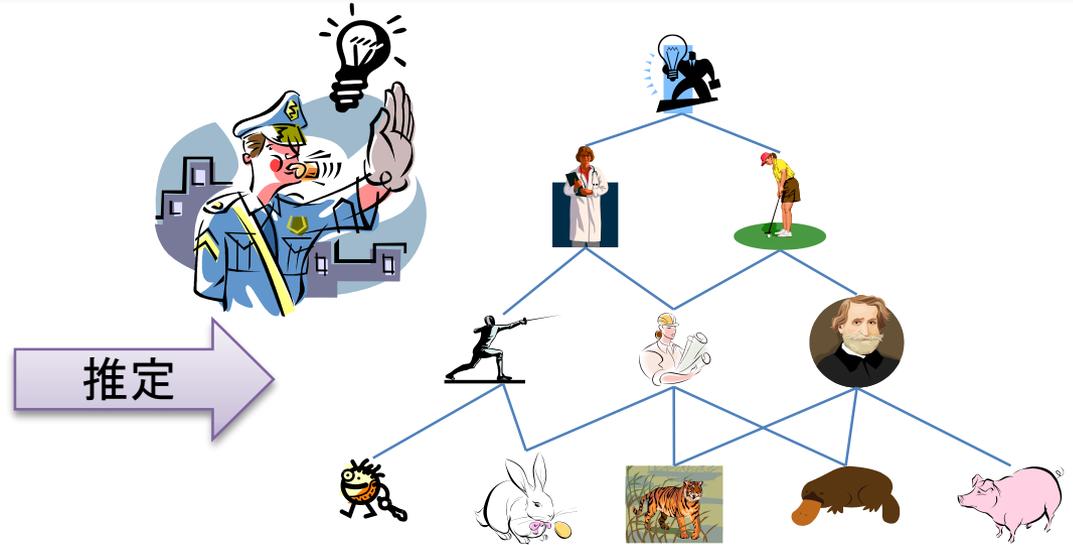
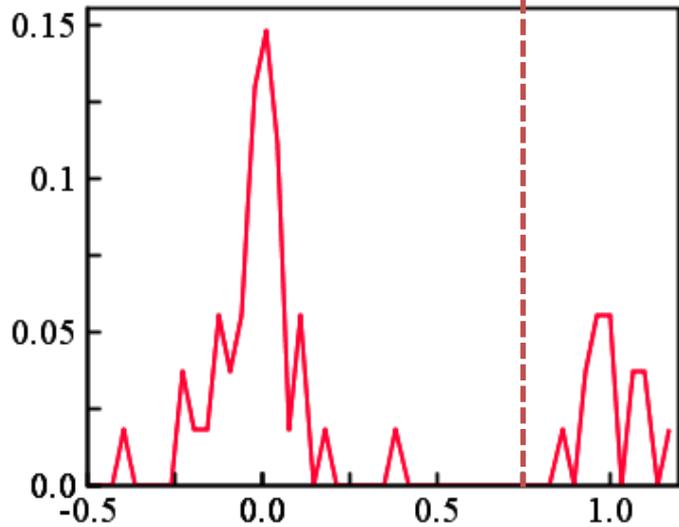


学習により得られた  
すべてのリンク重み  $w$  の  
ヒストグラム



# ネットワーク構造の抽出に成功

Threshold = 0.75



- 組織の**詳細な構造の抽出に成功**
- モデル自体が**間接的な関係性**も考慮にいらしているため相関係数を用いたときのような問題は起こりづらい
- 多少恣意的な閾値処理が入っているがここをある程度自動化する理論(**正則化理論**)が存在する

Threshold: 0.75

0.15

0

もちろんこの例は生成モデルと  
学習モデルが等しいので  
ある意味チートである

しかし重要なことは  
MRF的にデータが生成されている  
場合は  
MRFでないと正確に構造抽出が難しい  
ということである

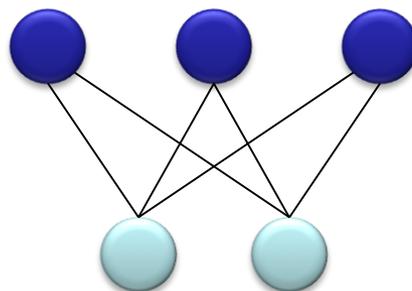


例えばニューラルネットワークは

入力データと出力データのセット(訓練セット)

を学習して入力から出力をだすシステムを構成できる

例えば2部グラフ上の MRF を考えてみる

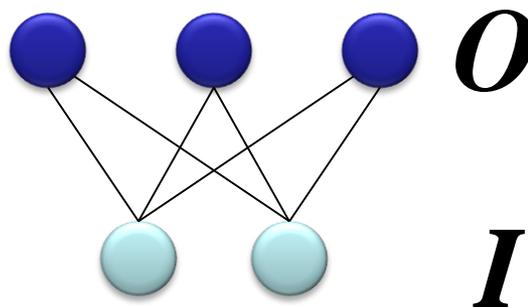


下層のノードを入力ノードとし上層を出力ノードと決めてしまう

入出力関係のデータがそろっていれば  
それを使ってこの MRF を学習することができる  
(MRF 自身はどれが入力でどれが出力なのかは区別していない  
人間が勝手に入出力を分けているだけ)

MRFが学習するのは入出力関係の確率分布となる

例えば2部グラフ上の MRF を考えてみる



便宜的に下層(入力)の変数を  $I$  で上層(出力)の変数を  $O$  で表す

学習される MRF はもちろん  $I$  と  $O$  の結合確率

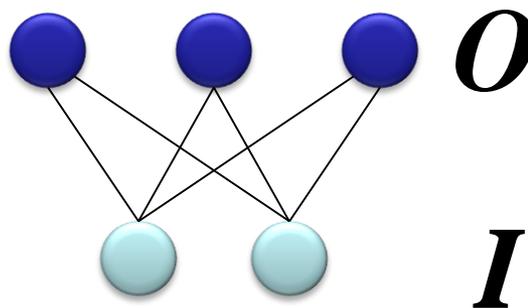
$$P(\mathbf{I}, \mathbf{O}) = \frac{1}{Z} \exp(-C(\mathbf{I}, \mathbf{O}))$$

# MRFでも入出力関係を作ることができる

111 / 141

例えば2部グラフ上の MRF を考えてみる

$$P(\mathbf{I}, \mathbf{O}) = \frac{1}{Z} \exp(-C(\mathbf{I}, \mathbf{O}))$$



2部グラフだと  
結構簡単に計算できる！

条件付き確率を計算する

$$P(\text{出力} | \text{入力}) = P(\mathbf{O} | \mathbf{I})$$

これはある入力  $\mathbf{I}$  がやって来たときの

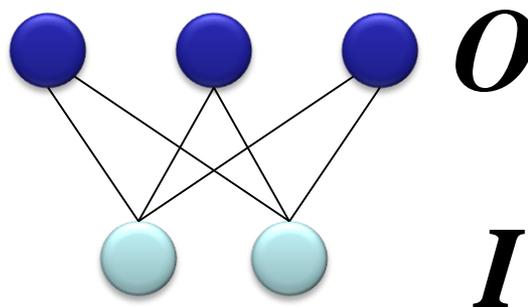
**出力  $\mathbf{O}$  の確率**として解釈できる

# MRFでも入出力関係を作ることができる

112 / 141

例えば2部グラフ上の MRF を考えてみる

$$P(\mathbf{I}, \mathbf{O}) = \frac{1}{Z} \exp(-C(\mathbf{I}, \mathbf{O}))$$



条件付き確率を計算する

$$P(\mathbf{O} | \mathbf{I})$$

2部グラフだと  
結構簡単に計算できる！

この条件付き確率を**最大にする O が出力**となるべきであろう

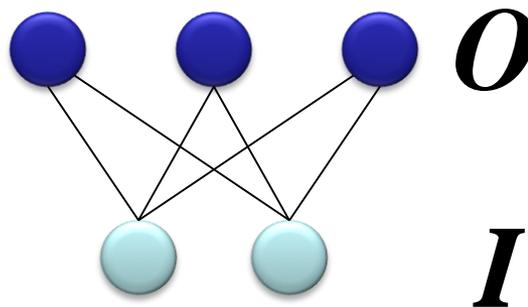
もちろん確率なので尤もあり得そうな出力だけでなくそこからの分散等も評価できる

# MRFでも入出力関係を作ることができる

113 / 141

例えば2部グラフ上の MRF を考えてみる

$$P(\mathbf{I}, \mathbf{O}) = \frac{1}{Z} \exp(-C(\mathbf{I}, \mathbf{O}))$$



2部グラフだと  
結構簡単に計算できる！

逆の条件付き確率を計算する

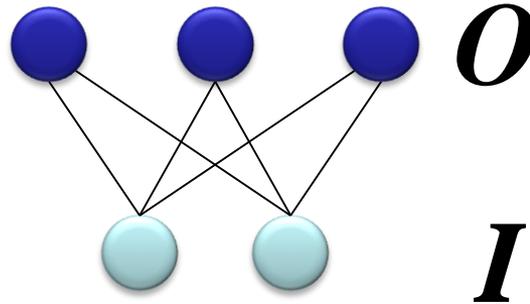
$$P(\text{入力} | \text{出力}) = P(\mathbf{I} | \mathbf{O})$$

MRF にはそもそも方向性がないので  
逆に出力から入力を推定することも可能である

通常のニューラルネットワークではこれは単一では出来ない

[Asari et. al., Master Thesis,14]

$$P(\mathbf{I}, \mathbf{O}) = \frac{1}{Z} \exp(-C(\mathbf{I}, \mathbf{O}))$$



入力に問診項目（体重・運動・飲酒・喫煙・歯磨きの頻度 等）

出力に検査結果（BMI・血圧・コレステロール・メタボリック 等）

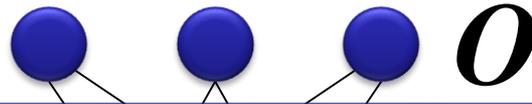
問診項目（生活習慣）から体内異常を推定するシステムができる

宮城県（2008年）の13209人分の人間ドックのデータで学習

ロジスティック回帰システムを上回る推定性能を実現

特に外れ値に対する推定に優れる

[Asari *et. al.*, Master Thesis, 14]



体内異常から逆に普段の生活習慣を  
推定することも可能である

これは単一のニューラルネットワーク  
では容易に実現できない

ロジスティック回帰システムを上回る推定性能を実現

特に外れ値に対する推定に優れる

# 制限ボルツマンマシン

(restricted Boltzmann machine; RBM)

これまで**すべての変数に対応した観測データ**がそろっていることが統計的機械学習において前提となっていた

データ数よりも多い数の変数をもったボルツマンマシンを考える

全ての変数のうち  
対応する観測データが得られた変数を  $v$   
得られない変数を  $h$  と記述する

$$P(v | \theta) \rightarrow P(v, h | \theta)$$

データに対応した変数  $v$   
→ **可視変数** (観測変数)

対応するデータがない変数  $h$   
→ **隠れ変数** (潜在変数)

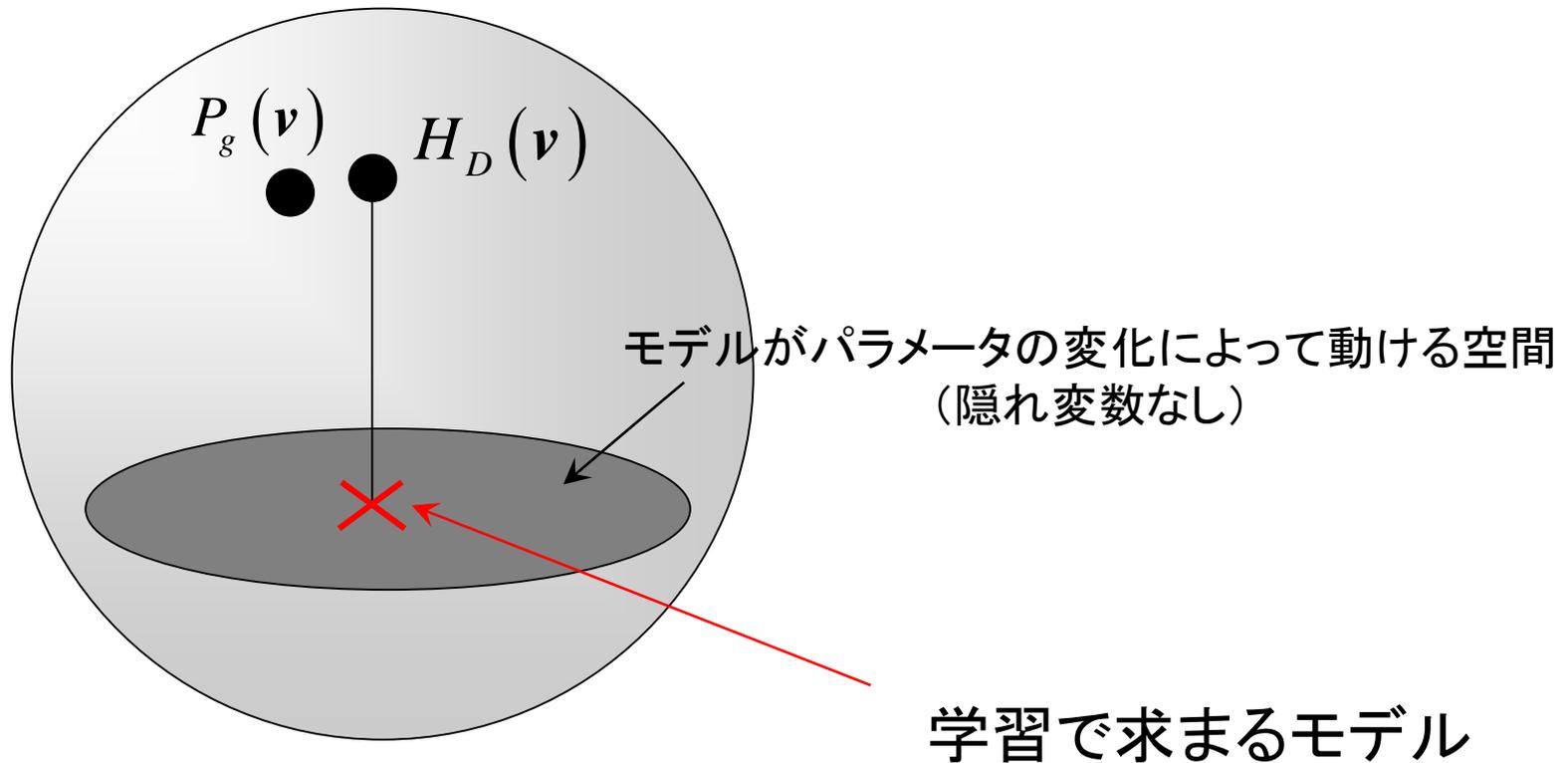
隠れ変数はシステムの内部変数となる

$$P(v|O) \rightarrow P(v, h|O)$$

# 隠れ変数を導入する意味

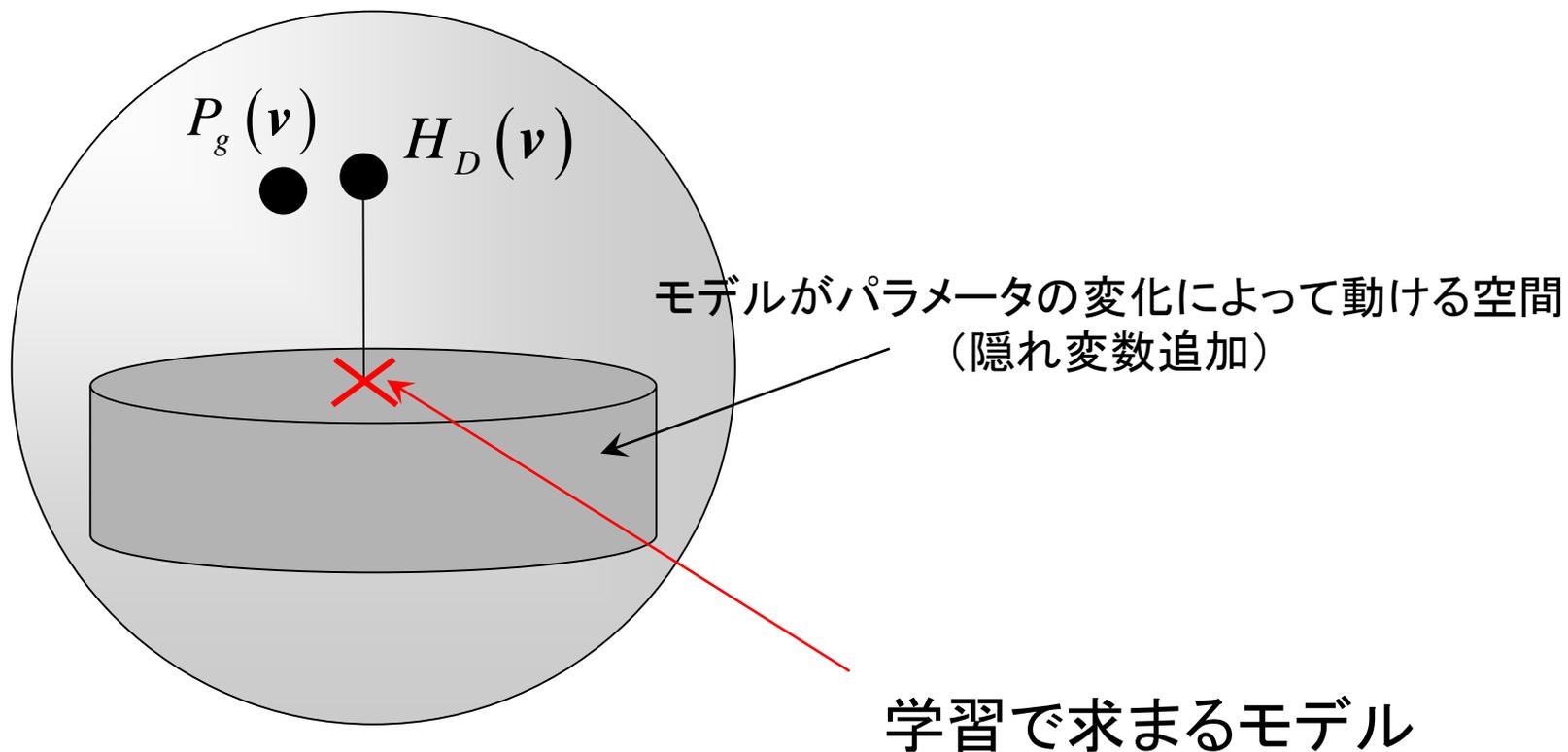
119 / 141

隠れ変数を追加することによりモデルの表現能力が増加する



## 隠れ変数を追加することによりモデルの表現能力が増加する

- ※ 隠れ変数の数がより多くなると一般により表現能力が増す  
現に後述の RBM は十分数の隠れ変数により任意の確率分布が表現できることが知られている



動ける空間が広がるので真の確率により近いモデルが得られる

隠れ変数を追加することによりモデルの表現能力が増加する

隠れ変数を導入すると確かに分布の表現力が上昇するが

実は良いことばかりではない

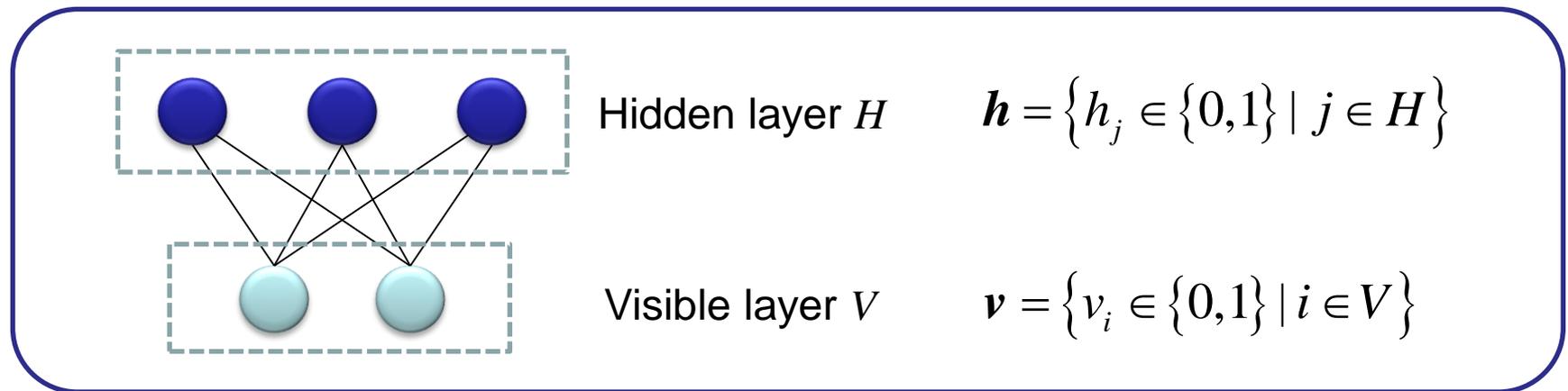
後述のように隠れ変数がある場合も  
基本的に最尤法で学習することになる



隠れ変数があると学習の**目的関数(対数尤度関数)**が  
**一般に凹関数でなくなる**ため  
最適解を見つけることが隠れ変数がない場合に比べ  
格段に困難になる

動ける空間が広がるので真の確率により近いモデルが得られる

- ◆ RBM は**完全2部グラフ**上に定義された2層構造のボルツマンマシン
- ◆ 片方の層は**可視層**であり, もう片方の層は**隠れ層**である
- ◆ RBM は深層学習の基本単位



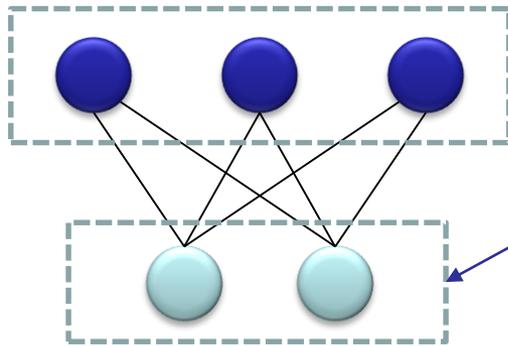
可視層は可視変数のみで構成され

隠れ層は隠れ変数のみで構成される

層内結合は存在しない

これまで  $V$  はすべてのノードの集合として使用してきたが

今後は可視変数に対応したノードの集合として用いる



データに対応した変数は可視変数なので  
データはこの層の変数にのみ対応する  
(可視変数の個数はデータの要素数と同数)

$$v_i, h_j \in \{0, 1\}$$

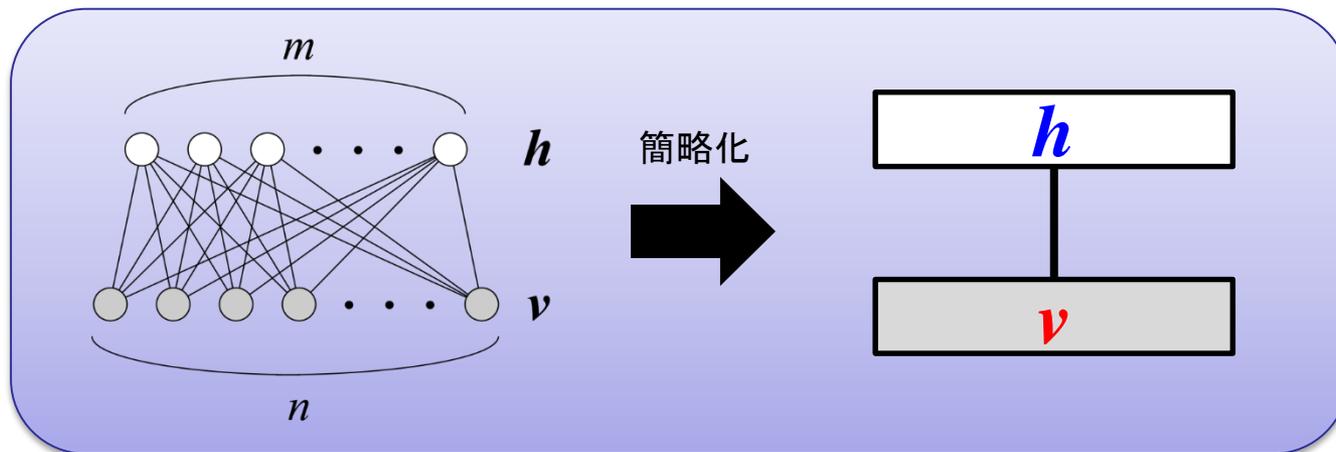
## RBMのエネルギー関数

$$C(\mathbf{v}, \mathbf{h}; \boldsymbol{\theta}) = -\sum_{i \in V} b_i v_i - \sum_{j \in H} c_j h_j - \sum_{i \in V} \sum_{j \in H} w_{ij} v_i h_j$$

バイアスパラメータ

結合パラメータ

$$P(\mathbf{v}, \mathbf{h} | \boldsymbol{\theta}) = \frac{1}{Z(\boldsymbol{\theta})} \exp(-C(\mathbf{v}, \mathbf{h}; \boldsymbol{\theta}))$$



片方の層の変数を条件としたもう片方の層の条件付き確率

$$P(\mathbf{v} | \mathbf{h}) = \frac{P(\mathbf{v}, \mathbf{h} | \boldsymbol{\theta})}{\sum_{\mathbf{v}} P(\mathbf{v}, \mathbf{h} | \boldsymbol{\theta})} = \prod_{i \in V} \frac{\exp(\lambda_i v_i)}{1 + \exp(\lambda_i)}$$

$$P(\mathbf{h} | \mathbf{v}) = \frac{P(\mathbf{v}, \mathbf{h} | \boldsymbol{\theta})}{\sum_{\mathbf{h}} P(\mathbf{v}, \mathbf{h} | \boldsymbol{\theta})} = \prod_{j \in H} \frac{\exp(\lambda_j h_j)}{1 + \exp(\lambda_j)}$$

$$\lambda_i = b_i + \sum_{j \in H} w_{ij} h_j, \quad \lambda_j = c_j + \sum_{i \in V} w_{ij} v_i$$

条件付き確率が各変数ごとに因数分解されえた形になっている  
→ 各変数が**統計的に独立**

## 条件付き独立の性質

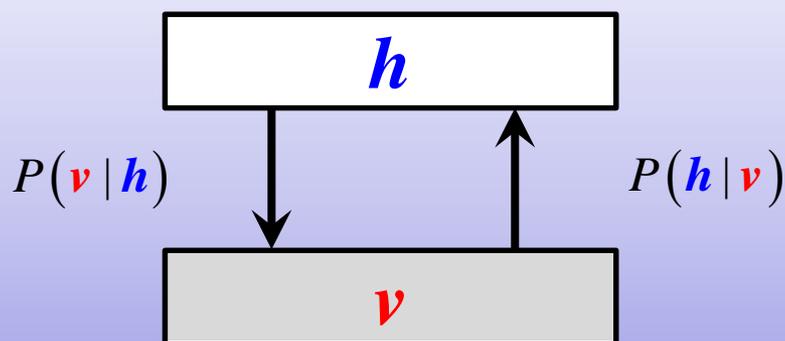
片方の層の変数を条件としたもう片方の層の条件付き確率

$$P(\mathbf{v} | \mathbf{h}) = \frac{P(\mathbf{v}, \mathbf{h} | \boldsymbol{\theta})}{\sum_{\mathbf{v}} P(\mathbf{v}, \mathbf{h} | \boldsymbol{\theta})} = \prod_{i \in V} \frac{\exp(\lambda_i v_i)}{1 + \exp(\lambda_i)}$$

$$P(\mathbf{h} | \mathbf{v}) = \frac{P(\mathbf{v}, \mathbf{h} | \boldsymbol{\theta})}{\sum_{\mathbf{h}} P(\mathbf{v}, \mathbf{h} | \boldsymbol{\theta})} = \prod_{j \in H} \frac{\exp(\lambda_j h_j)}{1 + \exp(\lambda_j)}$$

$$\lambda_i = b_i + \sum_{j \in H} w_{ij} h_j, \quad \lambda_j = c_j + \sum_{i \in V} w_{ij} v_i$$

各変数が統計的に独立なので  
(他の変数の値を気にせず) 各変数ごとに独立にサンプリングが可能



条件付き独立の性質より

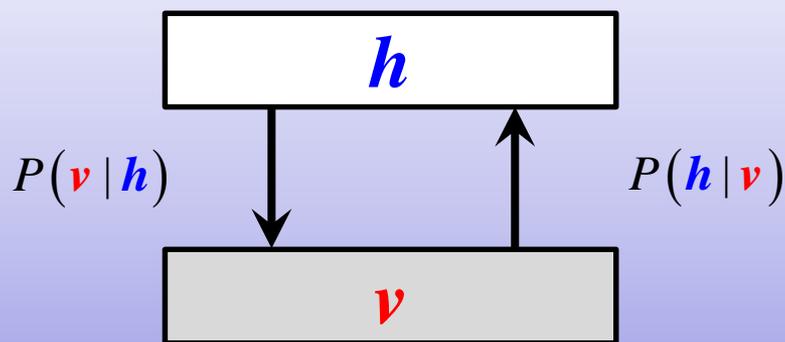
片方の層からもう片方の層のサンプリング  
がとても簡単

※統計的に独立でない则他の変数の値が確率に影響を与えるためサンプリングが面倒

RBMのギブスサンプリングは  
通常の BM のギブスサンプリングと違って  
層の変数の値をいっぺんに更新できる

ブロック化ギブスサンプリング

各変数が統計的に独立なので  
(他の変数の値を気にせず) 各変数ごとに独立にサンプリングが可能



条件付き独立の性質より

片方の層からもう片方の層のサンプリング  
がとても簡単

※統計的に独立でない则他の変数の値が確率に影響を与えるためサンプリングが面倒

層間のギブスサンプリングが容易であるというこの性質は  
RBM の学習法である**コントラストティブ・ダイバージェンス**  
と**深層学習**に大きく影響している

# RBMの学習 と コントラストティブ・ ダイバージェンス

## 最尤推定

尤度関数 (= 学習モデルがデータセットを実際に生成する確率)

$$L_D(\theta) = \prod_{\mu=1}^N P(\mathbf{v} = \mathbf{d}^{(\mu)} | \theta)$$

を最大化

データを可視変数に代入した確率

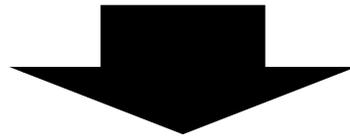
対応するデータをもたない隠れ変数が

存在する場合はどうしたらいい??



周辺化して可視変数のみの確率分布に変換する

$$P(\mathbf{v} | \theta) = \sum_h P(\mathbf{v}, \mathbf{h} | \theta)$$



周辺化により得られた可視変数のみの確率分布を使って  
尤度関数を構成する

$$L_D(\theta) = \prod_{\mu=1}^N P(\mathbf{v} = \mathbf{d}^{(\mu)} | \theta)$$

隠れ変数が存在しても周辺化を経由することで最尤推定が使える！

対数尤度関数を定義してそれぞれのパラメータに関する勾配を計算する

$$\frac{\partial l_D(\boldsymbol{\theta})}{\partial b_i} = \frac{1}{N} \sum_{\mu=1}^N d_i^{(\mu)} - \mathbb{E}_{\text{RBM}} [v_i | \boldsymbol{\theta}]$$

$$\frac{\partial l_D(\boldsymbol{\theta})}{\partial c_j} = \frac{1}{N} \sum_{\mu=1}^N \sigma_1 \left( c_j + \sum_{k \in V} w_{kj} d_k^{(\mu)} \right) - \mathbb{E}_{\text{RBM}} [h_j | \boldsymbol{\theta}]$$

$$\frac{\partial l_D(\boldsymbol{\theta})}{\partial w_{ij}} = \frac{1}{N} \sum_{\mu=1}^N d_i^{(\mu)} \sigma_1 \left( c_j + \sum_{k \in V} w_{kj} d_k^{(\mu)} \right) - \mathbb{E}_{\text{RBM}} [v_i h_j | \boldsymbol{\theta}]$$



$$\sigma_1(x) = \frac{1}{1 + e^{-x}}$$

$$\mathbb{E}_{\text{RBM}} [(\dots) | \boldsymbol{\theta}] = \sum_{\mathbf{v}} \sum_{\mathbf{h}} (\dots) P(\mathbf{v}, \mathbf{h} | \boldsymbol{\theta})$$

RBM の期待値

RBM の学習は下記の勾配を用いた勾配上昇法によって実行されることとなる

$$\frac{\partial l_D(\boldsymbol{\theta})}{\partial b_i} = \frac{1}{N} \sum_{\mu=1}^N d_i^{(\mu)} - \mathbb{E}_{\text{RBM}} [v_i | \boldsymbol{\theta}]$$
$$\frac{\partial l_D(\boldsymbol{\theta})}{\partial c_j} = \frac{1}{N} \sum_{\mu=1}^N \sigma_1 \left( c_j + \sum_{k \in V} w_{kj} d_k^{(\mu)} \right) - \mathbb{E}_{\text{RBM}} [h_j | \boldsymbol{\theta}]$$
$$\frac{\partial l_D(\boldsymbol{\theta})}{\partial w_{ij}} = \frac{1}{N} \sum_{\mu=1}^N d_i^{(\mu)} \sigma_1 \left( c_j + \sum_{k \in V} w_{kj} d_k^{(\mu)} \right) - \mathbb{E}_{\text{RBM}} [v_i h_j | \boldsymbol{\theta}]$$



各勾配の第1項はデータとパラメータから直ちに計算できるが

**第2項は RBM の期待値であるため組み合わせ爆発の問題が起こる**

※ 第1項が計算できるのはRBMの2部グラフ構造のおかげである

一般の構造の隠れ変数ありBMの場合は第1項の計算においても組み合わせ爆発の問題をもつ

**RBM の学習もやはり近似学習法を必要とする**

## ギブスサンプリングを基礎とした確率的近似学習法

コントラストティブダイバージェンス(contrastive divergence; CD)

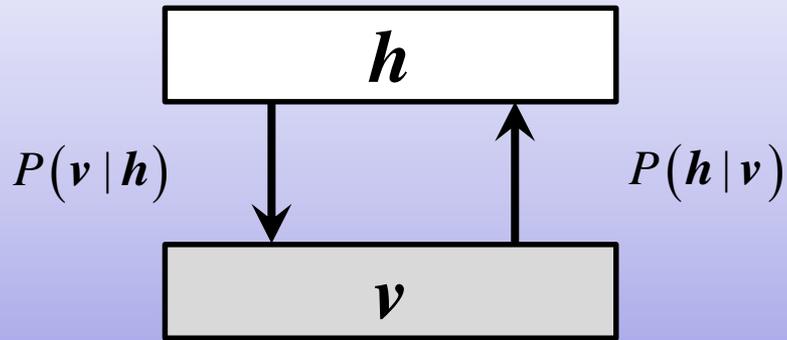
は実装が非常に容易で性能もそれなりに良いので広く普及している

RBM から深層学習へつながるための **立役者的アルゴリズム**

より性能の良いアルゴリズムも多数提案されてきているが

計算コスト等の兼ね合いも含めて総合的に見ると

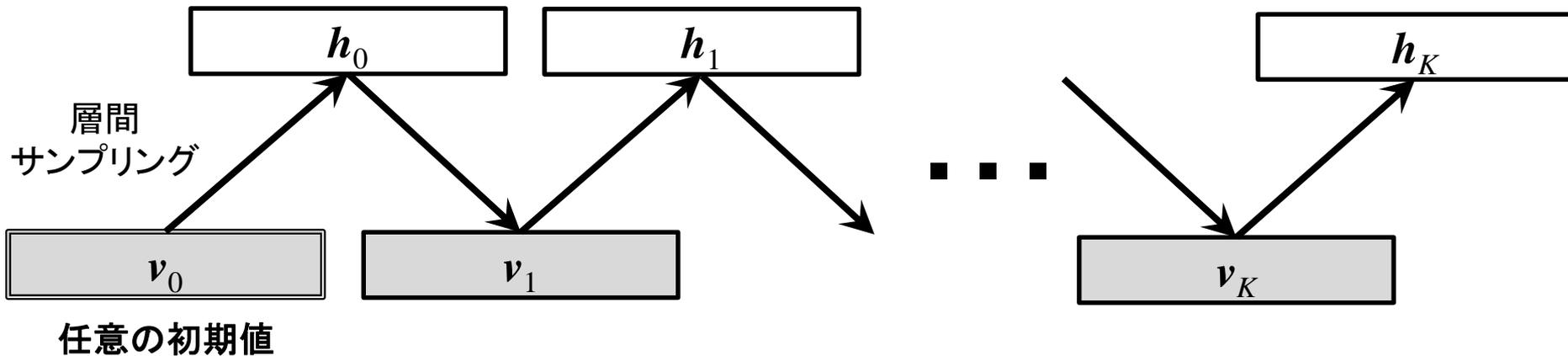
やはりまだ CD がバランス良いと思われる(私見)

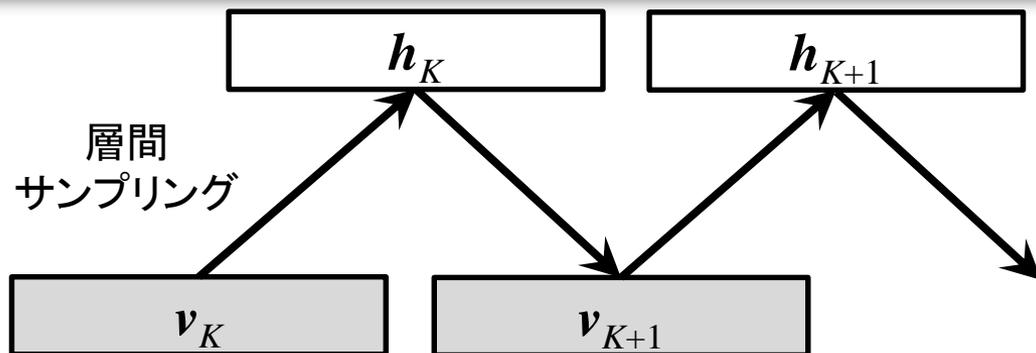


条件付き独立の性質より

片方の層からもう片方の層のサンプリング  
がとても簡単

可視変数の任意の初期値から出発して  
層間サンプリングを十分回繰り返す(緩和期間)





十分回層間サンプリング  
を行った後に  
一定間隔でサンプル点  
を取得していく

可視変数と隠れ変数に対応するサンプル点をそれぞれ  $S$  点取得したとする

$$\left\{ \mathbf{v}^{(1)}, \mathbf{h}^{(1)} \right\}, \left\{ \mathbf{v}^{(2)}, \mathbf{h}^{(2)} \right\}, \dots, \left\{ \mathbf{v}^{(S)}, \mathbf{h}^{(S)} \right\}$$

得られたサンプル点の標本平均で RBM の期待値は近似される  
(モンテカルロ積分)

$$\begin{aligned} E_{\text{RBM}} [v_i | \boldsymbol{\theta}] &\approx \frac{1}{S} \sum_{r=1}^S v_i^{(r)}, & E_{\text{RBM}} [h_j | \boldsymbol{\theta}] &\approx \frac{1}{S} \sum_{r=1}^S h_j^{(r)}, \\ E_{\text{RBM}} [v_i h_j | \boldsymbol{\theta}] &\approx \frac{1}{S} \sum_{r=1}^S v_i^{(r)} h_j^{(r)} \end{aligned}$$

$S \rightarrow \infty$  のときは大数の法則で厳密

$h_K$  $h_{K+1}$ 

ギブスサンプリング

原理的にはモンテカルロ積分でRBMの期待値を  
近似して勾配を計算すればよいのだが

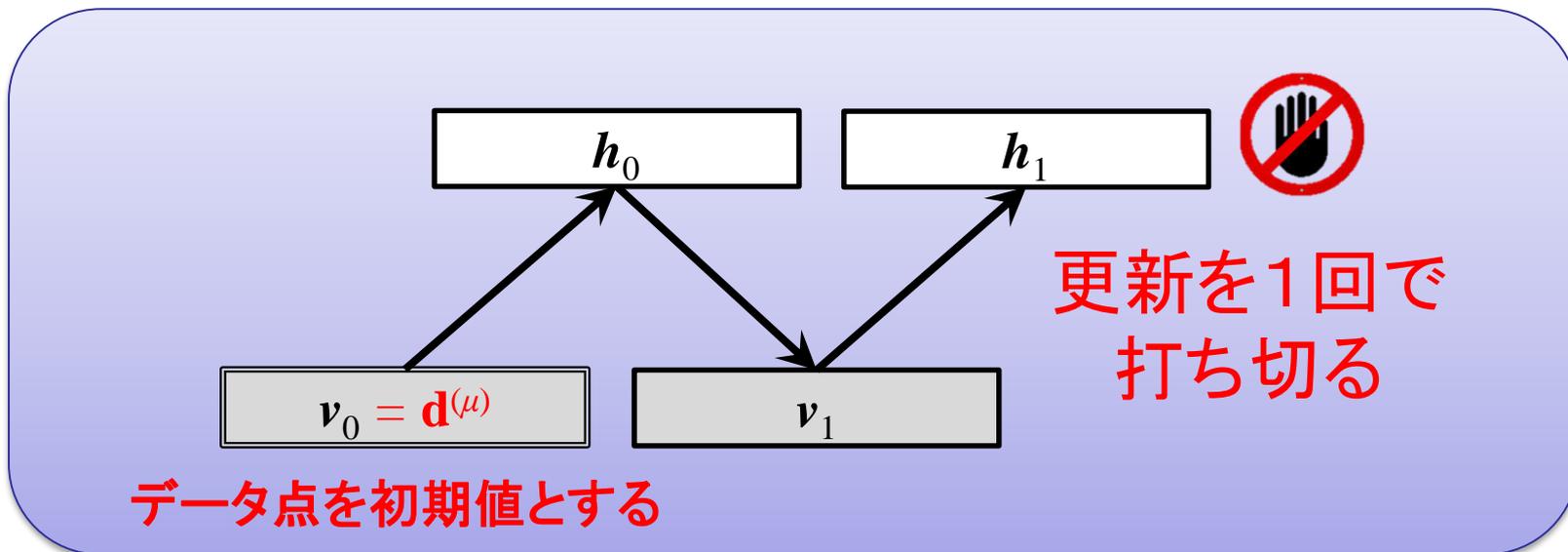
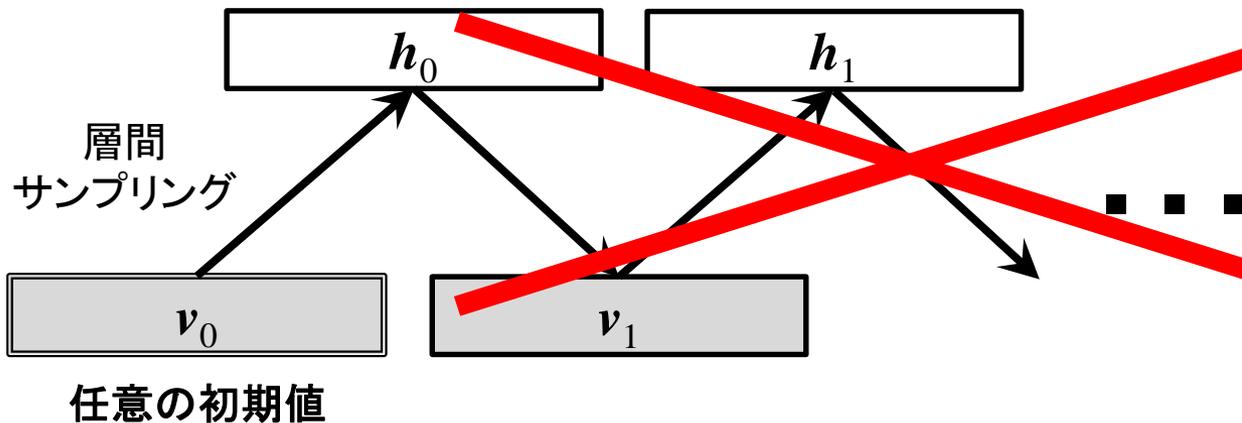
**勾配法の毎ステップ**において**緩和期間**を  
必要としてしまい

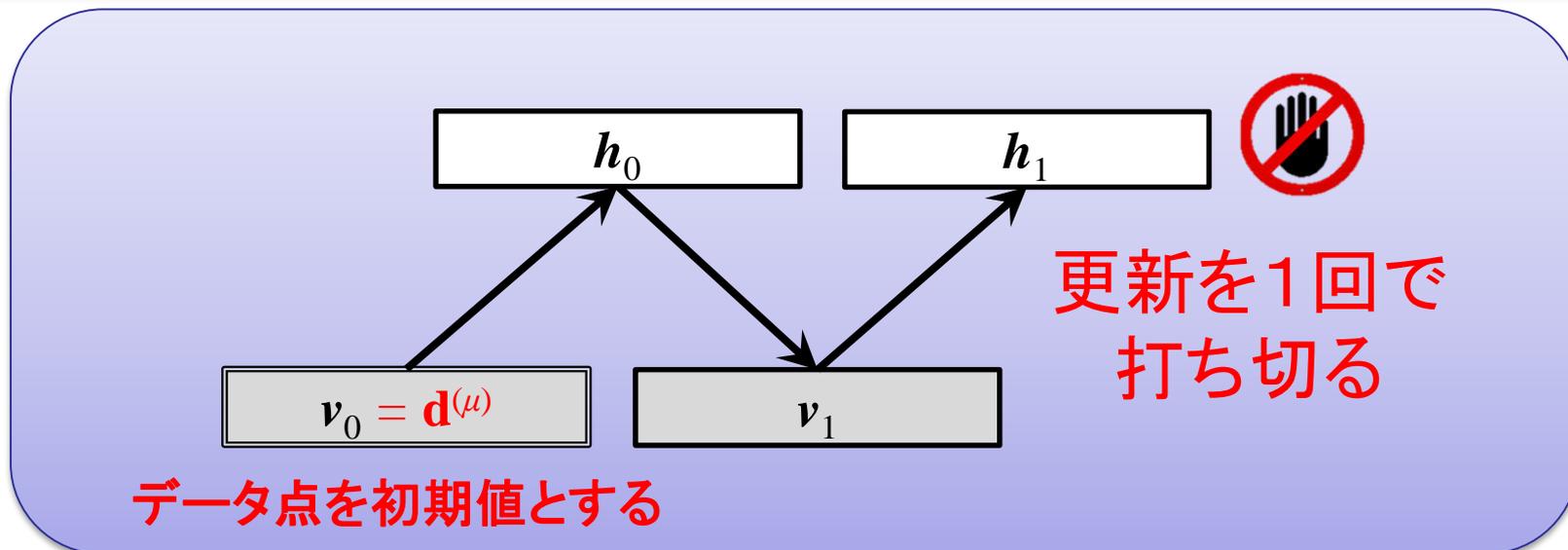
また**サンプル間隔**もある程度空ける必要があるため

計算コストが非常に高くなってしまふ

$S \rightarrow \infty$  のときは大数の法則で厳密

# コントラストティブ・ダイバージェンス



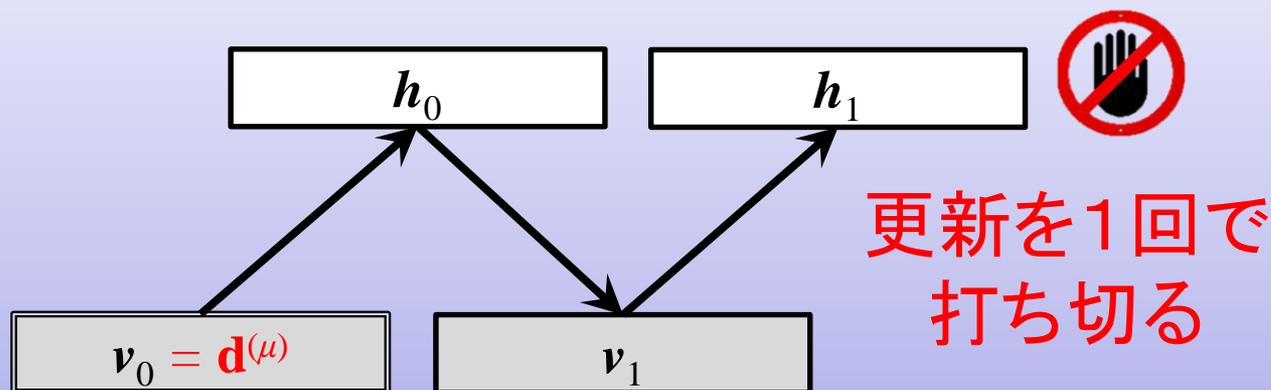


データ点を初期値とした1回の更新でサンプル点を取得する

$\{ \mathbf{v}^{(\mu)}, \mathbf{h}^{(\mu)} \}$   $\mu$ 番目のデータ点に対するサンプル点

同様な方法によりデータの個数分手に入る

$\{ \mathbf{v}^{(1)}, \mathbf{h}^{(1)} \}, \{ \mathbf{v}^{(2)}, \mathbf{h}^{(2)} \}, \dots, \{ \mathbf{v}^{(N)}, \mathbf{h}^{(N)} \}$



データ点を初期値とする

各データ点に対するサンプル点の標本平均で  
RBM の期待値を近似する

$$\mathbb{E}_{\text{RBM}} [v_i | \theta] \approx \frac{1}{N} \sum_{\mu=1}^N v_i^{(\mu)}, \quad \mathbb{E}_{\text{RBM}} [h_j | \theta] \approx \frac{1}{N} \sum_{\mu=1}^N h_j^{(\mu)},$$

$$\mathbb{E}_{\text{RBM}} [v_i h_j | \theta] \approx \frac{1}{N} \sum_{\mu=1}^N v_i^{(\mu)} h_j^{(\mu)}$$

データ点を初期値とすることで  
コストの高い緩和期間を省略する

また

サンプル点は各データ点に対して  
1つだけ取得するため  
一定間隔空けて複数サンプルする  
必要がない

CD は以下の**近似勾配**を用いた勾配上昇法によって実行されることとなる

$$\frac{\partial l_D(\boldsymbol{\theta})}{\partial b_i} \approx \frac{1}{N} \sum_{\mu=1}^N d_i^{(\mu)} - \frac{1}{N} \sum_{\mu=1}^N v_i^{(\mu)}$$

$$\frac{\partial l_D(\boldsymbol{\theta})}{\partial c_j} \approx \frac{1}{N} \sum_{\mu=1}^N \sigma_1 \left( c_j + \sum_{k \in V} w_{kj} d_k^{(\mu)} \right) - \frac{1}{N} \sum_{\mu=1}^N h_j^{(\mu)}$$

$$\frac{\partial l_D(\boldsymbol{\theta})}{\partial w_{ij}} \approx \frac{1}{N} \sum_{\mu=1}^N d_i^{(\mu)} \sigma_1 \left( c_j + \sum_{k \in V} w_{kj} d_k^{(\mu)} \right) - \frac{1}{N} \sum_{\mu=1}^N v_i^{(\mu)} h_j^{(\mu)}$$

$$\mathbb{E}_{\text{RBM}} [v_i | \boldsymbol{\theta}] \approx \frac{1}{N} \sum_{\mu=1}^N v_i^{(\mu)}, \quad \mathbb{E}_{\text{RBM}} [h_j | \boldsymbol{\theta}] \approx \frac{1}{N} \sum_{\mu=1}^N h_j^{(\mu)},$$

$$\mathbb{E}_{\text{RBM}} [v_i h_j | \boldsymbol{\theta}] \approx \frac{1}{N} \sum_{\mu=1}^N v_i^{(\mu)} h_j^{(\mu)}$$

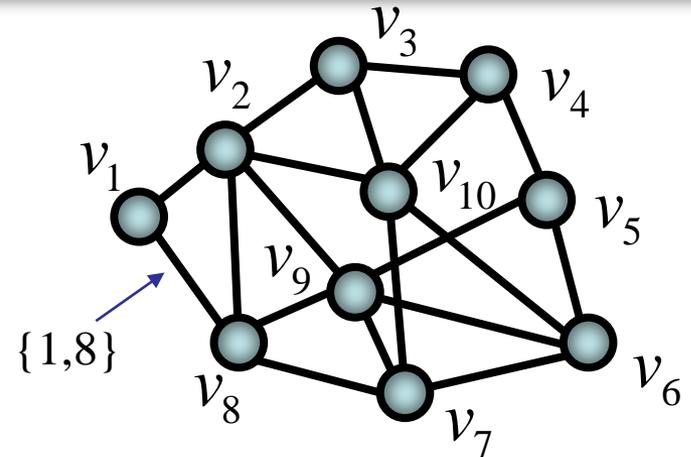
# 最新の研究結果

**An Effective Algorithm  
for  
Boltzmann Machine Learning**

# Boltzmann Machine

$G(V, E) \leftarrow$  undirected graph

$V = \{1, 2, \dots, 10\} \leftarrow$  set of nodes



$E = \{\{1, 2\}, \{1, 8\}, \{2, 3\}, \{2, 8\}, \{2, 9\}, \dots\} \leftarrow$  set of links

energy function of BM

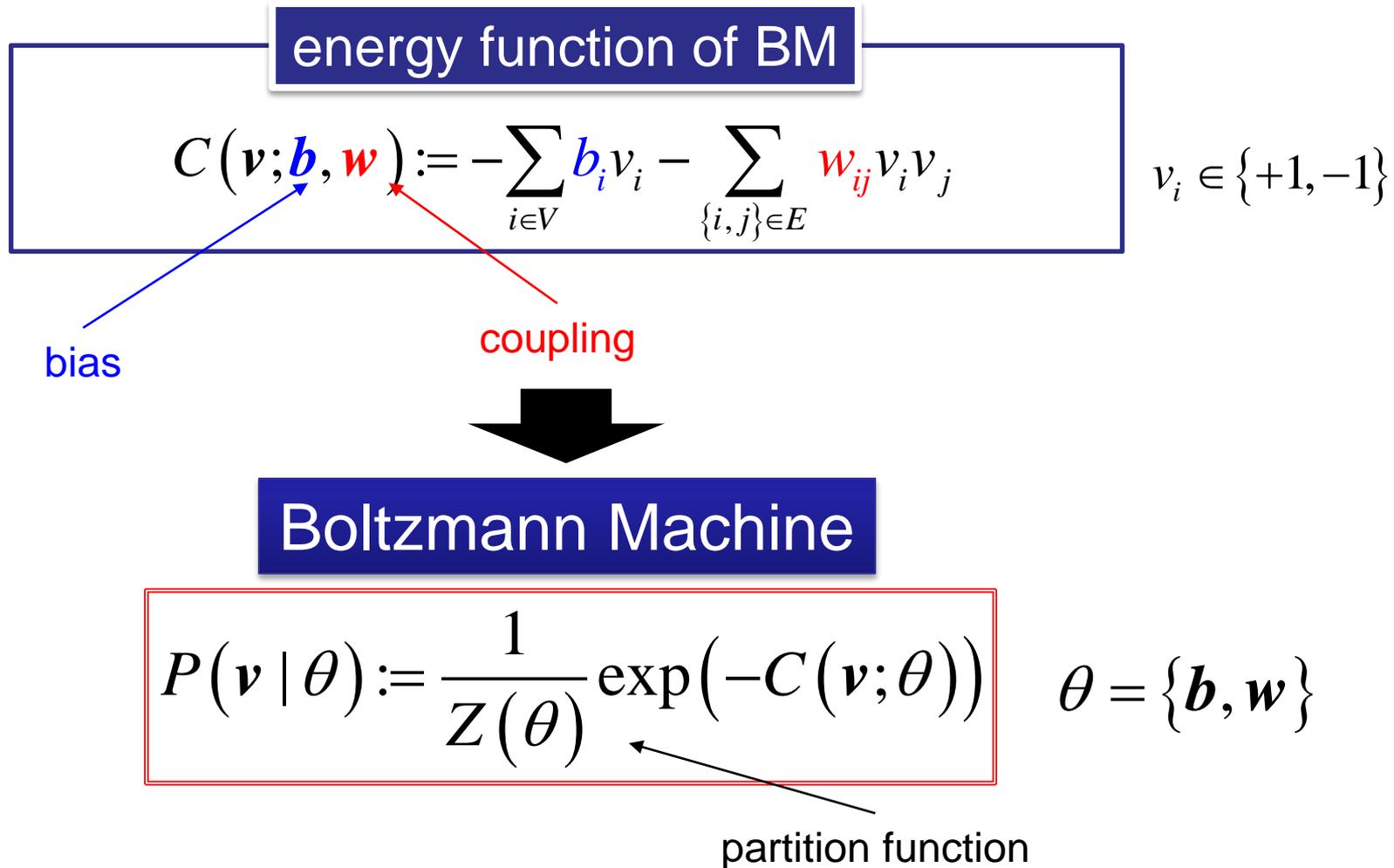
$$C(\mathbf{v}; \mathbf{b}, \mathbf{w}) := - \sum_{i \in V} b_i v_i - \sum_{\{i, j\} \in E} w_{ij} v_i v_j$$

$$v_i \in \{+1, -1\}$$

bias

coupling

$$w_{ij} = w_{ji} \text{ (symmetric couplings)}$$

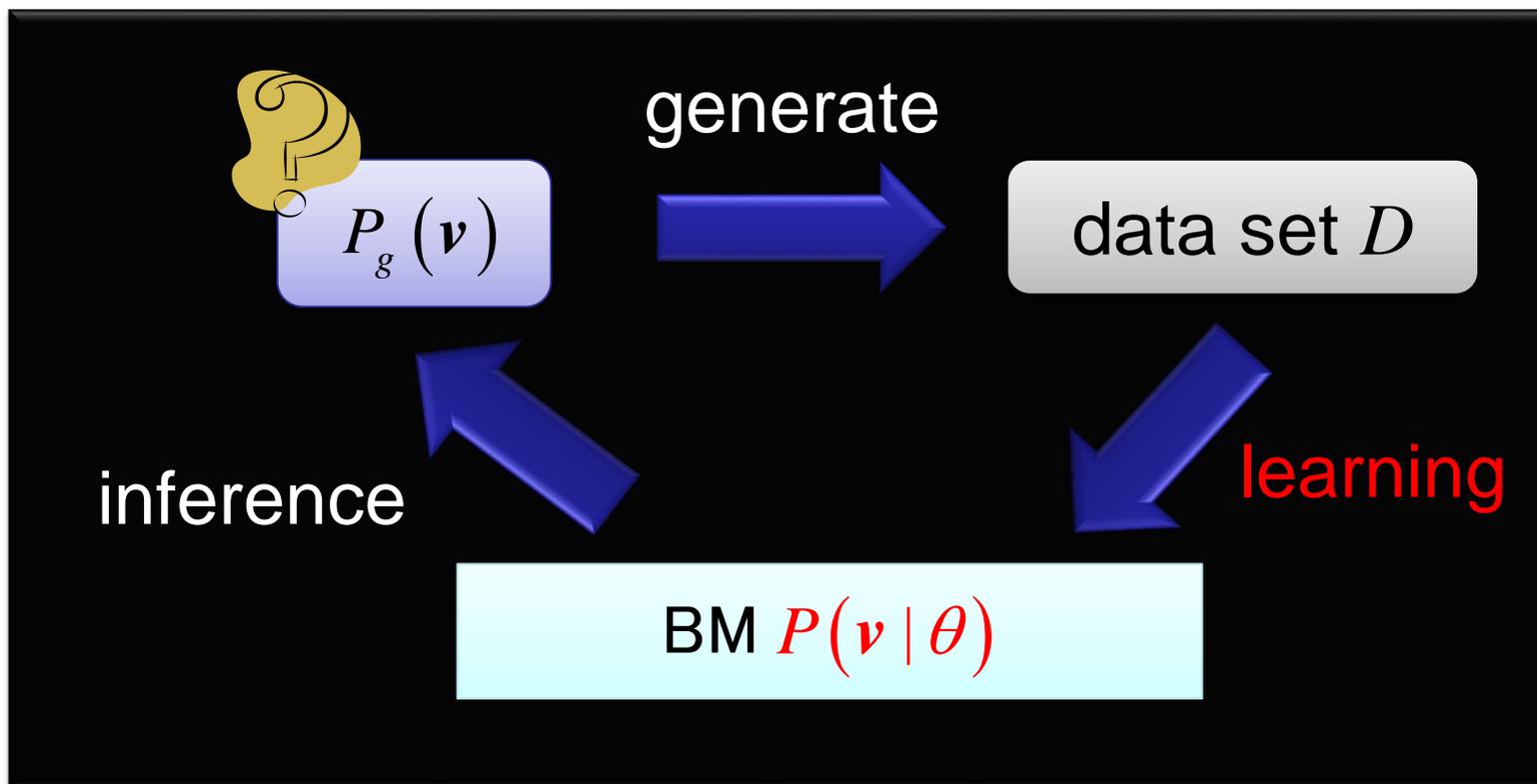


BM is the simplest probabilistic model in Markov random fields

# Statistical Machine Learning

146 / 141

The aim of the statistical machine learning is to find the unknown **generative model**.



Using data set  $D$  to tune the values of the parameters,  
we approximate the generative model by our model (BM).

set of  $N$  data points (each data point contains  $n$  elements):

$$D = \underbrace{\{\mathbf{d}^{(1)}, \mathbf{d}^{(2)}, \dots, \mathbf{d}^{(N)}\}}_{N \text{ data points}}, \quad \mathbf{d}^{(\mu)} = \underbrace{\{d_1^{(\mu)}, d_2^{(\mu)}, \dots, d_n^{(\mu)}\}}_{n \text{ elements}}$$

$$d_i^{(\mu)} \in \{+1, -1\}$$

Each data point is generated from a certain generative model

$$\begin{array}{ccc} \text{?} P_g(\mathbf{v}) & \xrightarrow{\text{generate (i.i.d.)}} & \mathbf{d}^{(1)}, \mathbf{d}^{(2)}, \dots, \mathbf{d}^{(N)} \end{array}$$

log-likelihood function  $l_D(\theta) := \frac{1}{N} \sum_{\mu=1}^N \ln P(\mathbf{v} = \mathbf{d}^{(\mu)} \mid \theta)$

**Maximum likelihood estimation**  
(MLE)

$$\theta^* = \arg \max_{\theta} l_D(\theta)$$

Maximum likelihood estimation  
(MLE)

$$\theta^* = \arg \max_{\theta} l_D(\theta)$$

gradients

$$\frac{\partial l_D(\theta)}{\partial b_i} = \frac{1}{N} \sum_{\mu=1}^N d_i^{(\mu)} - \mathbf{E}[v_i | \theta]$$

$$\frac{\partial l_D(\theta)}{\partial w_{ij}} = \frac{1}{N} \sum_{\mu=1}^N d_i^{(\mu)} d_j^{(\mu)} - \mathbf{E}[v_i v_j | \theta]$$

$$\mathbf{E}[(\dots) | \theta] := \sum_{\mathbf{v}} (\dots) P(\mathbf{v} | \theta) = \sum_{v_1} \sum_{v_2} \dots \sum_{v_n} (\dots) P(\mathbf{v} | \theta)$$

The second terms are the expectations of BM and they include **the intractable multiple-summations**.

we need to evaluate the expectations with an **approximate method**

# Monte Carlo Integration

149 / 141

**Monte Carlo Integration** (MCI) is one of the most fundamental method to approximate the expectations.

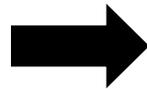
sampling phase

averaging phase

MC sampling

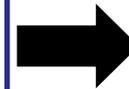
$$P(\mathbf{v} | \theta)$$

model



$$\mathbf{s}^{(1)}, \mathbf{s}^{(2)}, \dots, \mathbf{s}^{(K)}$$

sampled points



$$E[f(\mathbf{v}) | \theta] \approx \frac{1}{K} \sum_{k=1}^K f(\mathbf{s}^{(k)})$$

sample average

MCI is the approximate method on the basis of **the random sampling**.

In the sampling phase,  
we **use** information of the model (parameters, structures, and so on)

In the averaging phase,  
we **do not use** any information of the model

# Spatial Monte Carlo Integration

150 / 141

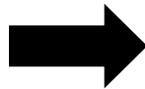
**Monte Carlo Integration** (MCI) is one of the most fundamental method to approximate the expectations.

sampling phase

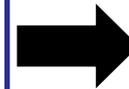
averaging phase

MC sampling

$$P(\mathbf{v} | \theta)$$



$$\mathbf{s}^{(1)}, \mathbf{s}^{(2)}, \dots, \mathbf{s}^{(K)}$$



~~$$E[f(\mathbf{v}) | \theta] \approx \frac{1}{K} \sum_{k=1}^K f(\mathbf{s}^{(k)})$$~~

model

sampled points

sample average

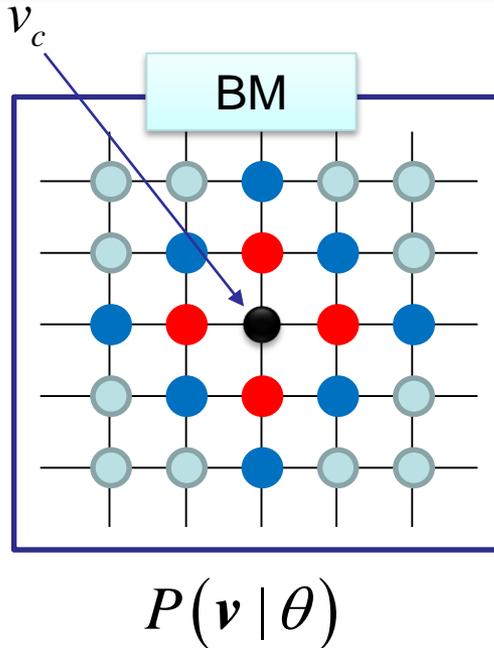
Change this part



In the new method, **spatial Monte Carlo Integration** (SMCI), we **use available information** of the model in the averaging phase.

# Spatial Monte Carlo Integration

151 / 141



We focus on the expectation of a function of the target variable(s)  $v_c$

$$E[f(v_c) | \theta] = \sum_{\mathbf{v}} v_c P(\mathbf{v} | \theta).$$

$v_c$  : target variable(s)

$\mathbf{v}_c^{(1)}$  : set of first-nearest variables of the target (shown by red)

$\mathbf{v}_c^{(2)}$  : set of second-nearest variables of the target (shown by blue)

We can easily obtain the distribution conditioned by the  $k$ th-nearest variables:

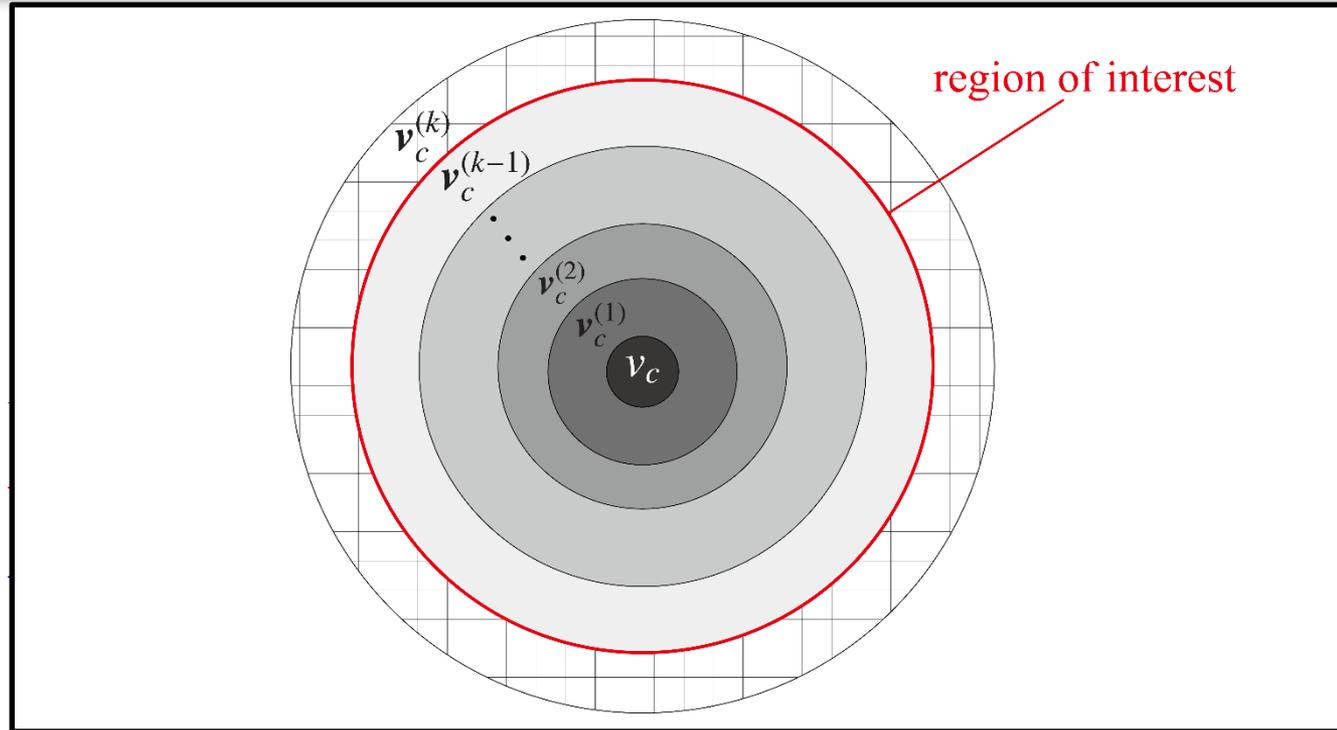
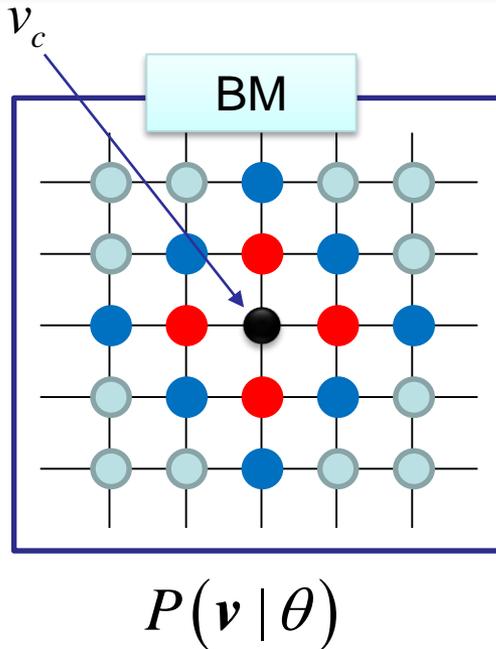
$$P(v_c, \mathbf{v}_c^{(1)}, \mathbf{v}_c^{(2)}, \dots, \mathbf{v}_c^{(k-1)} | \mathbf{v}_c^{(k)}, \theta)$$

This distribution is also BM on the **region of interest**.

$$\{v_c, \mathbf{v}_c^{(1)}, \mathbf{v}_c^{(2)}, \dots, \mathbf{v}_c^{(k-1)}\}$$

# Spatial Monte Carlo Integration

152 / 141



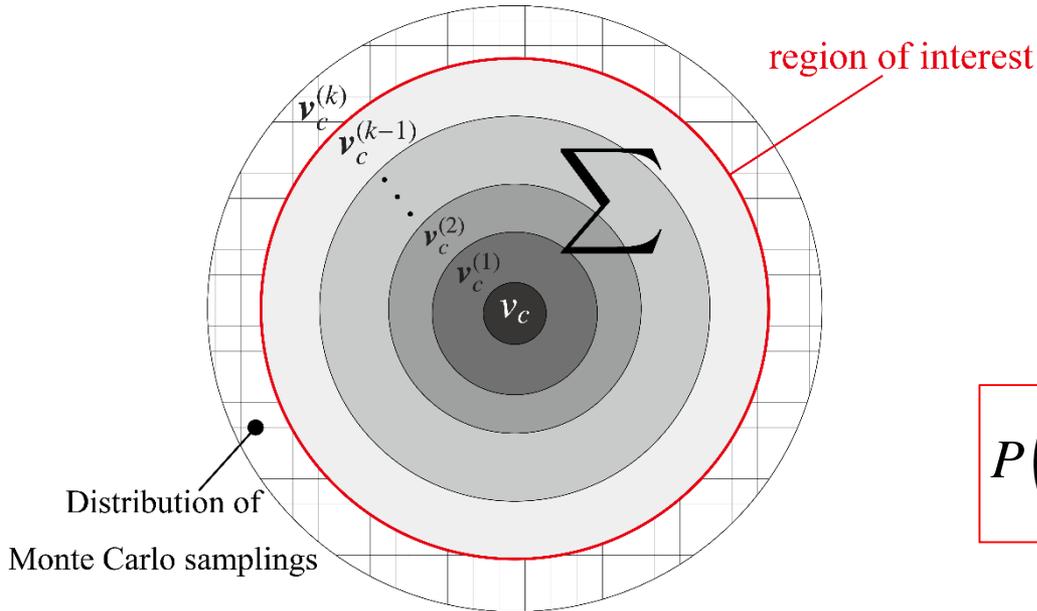
We can easily obtain the distribution conditioned with the  $k$ th-nearest variables:

$$P\left(v_c, \mathbf{v}_c^{(1)}, \mathbf{v}_c^{(2)}, \dots, \mathbf{v}_c^{(k-1)} \mid \mathbf{v}_c^{(k)}, \theta\right)$$

This distribution is also BM on the **region of interest**.

$$\left\{v_c, \mathbf{v}_c^{(1)}, \mathbf{v}_c^{(2)}, \dots, \mathbf{v}_c^{(k-1)}\right\}$$

# Spatial Monte Carlo Integration



The marginal distribution on the outside of the region

is approximated by the distribution of the sampled points

$$P(\mathbf{v}_c^{(k)} | \theta) \approx Q(\mathbf{v}_c^{(k)}) := \frac{1}{K} \sum_{k=1}^K \delta(\mathbf{v}_c^{(k)}, \mathbf{s}_c^{(k)})$$

$$\begin{aligned} E[f(v_c) | \theta] &= \sum_{v_c} \sum_{v_c^{(1)}} \sum_{v_c^{(2)}} \cdots \sum_{v_c^{(k-1)}} \sum_{v_c^{(k)}} f(v_c) P(v_c, v_c^{(1)}, v_c^{(2)}, \dots, v_c^{(k-1)} | v_c^{(k)}, \theta) P(v_c^{(k)} | \theta) \\ &\approx \sum_{v_c} \sum_{v_c^{(1)}} \sum_{v_c^{(2)}} \cdots \sum_{v_c^{(k-1)}} \sum_{v_c^{(k)}} f(v_c) P(v_c, v_c^{(1)}, v_c^{(2)}, \dots, v_c^{(k-1)} | v_c^{(k)}, \theta) Q(v_c^{(k)}) \\ &= \frac{1}{K} \sum_{k=1}^K \sum_{v_c} \sum_{v_c^{(1)}} \sum_{v_c^{(2)}} \cdots \sum_{v_c^{(k-1)}} f(v_c) P(v_c, v_c^{(1)}, v_c^{(2)}, \dots, v_c^{(k-1)} | \mathbf{s}_c^{(k)}, \theta) \end{aligned}$$

***k*th-order SMCI (*k*-SMCI)**

## $k$ th-order SMCI ( $k$ -SMCI)

$$E[f(v_c) | \theta] \approx \frac{1}{K} \sum_{k=1}^K \sum_{v_c} \sum_{v_c^{(1)}} \sum_{v_c^{(2)}} \cdots \sum_{v_c^{(k-1)}} f(v_c) P(v_c, v_c^{(1)}, v_c^{(2)}, \dots, v_c^{(k-1)} | \mathbf{s}_c^{(k)}, \theta)$$

The first- and second-order moments given by 1-SMCI

$$E[v_i | \theta] \approx \frac{1}{K} \sum_{k=1}^K \sum_{v_i} v_i P(v_i | \mathbf{s}_i^{(1)}, \theta) = \frac{1}{K} \sum_{k=1}^K \tanh U_i^{(k)}$$

$$E[v_i v_j | \theta] \approx \frac{1}{K} \sum_{k=1}^K \sum_{v_i} \sum_{v_j} v_i v_j P(v_i v_j | \mathbf{s}_{\{i,j\}}^{(1)}, \theta) = \frac{1}{K} \sum_{k=1}^K \tanh \left[ \tanh^{-1} \left[ \tanh U_{i,j}^{(k)} \tanh U_{j,i}^{(k)} \right] + w_{ij} \right]$$

$$U_i^{(k)} := b_i + \sum_{j \in \partial i} w_{ij} s_j^{(k)}, \quad U_{i,j}^{(k)} := U_i^{(k)} - w_{ij} s_j^{(k)}$$

$\partial i := \{j | \{i, j\} \in E\}$ : the set of nearest-neighbor nodes of node  $i$ .

The important two facts can be proven:

- SMCI is statistically **better** than the normal MCI
- a higher-order SMCI is statistically **better** than a lower-order one

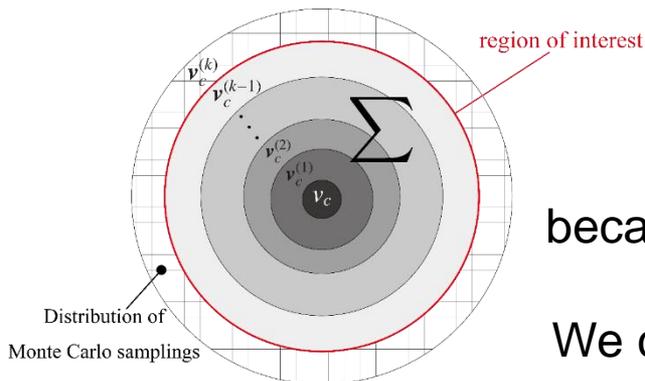
## Merit of SMCI

SMCI gives **high-accuracy** approximations

## Demerit of SMCI

Higher-order SMCI is slower, because we have to take the sum over the region of interest.

We cannot use SMCI in a dense system, except for 1-SMCI.



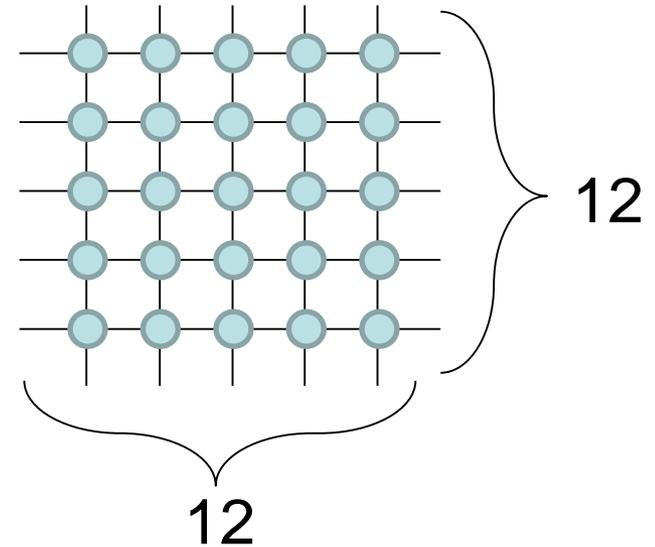
# Numerical Experiment

## Generative Model

Boltzmann machine

$$P_g(\mathbf{v}) = \frac{1}{Z} \exp\left(\sum_{\{i,j\} \in E} w_{ij} v_i v_j\right), \quad v_i \in \{+1, -1\}$$

$$w_{ij} \text{ i.i.d. } \sim U[-0.5, 0.5]$$

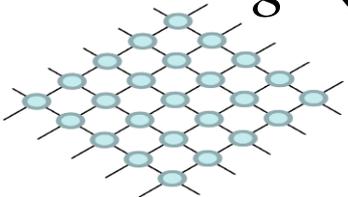


Gibbs sampling method

$P_g(\mathbf{v})$

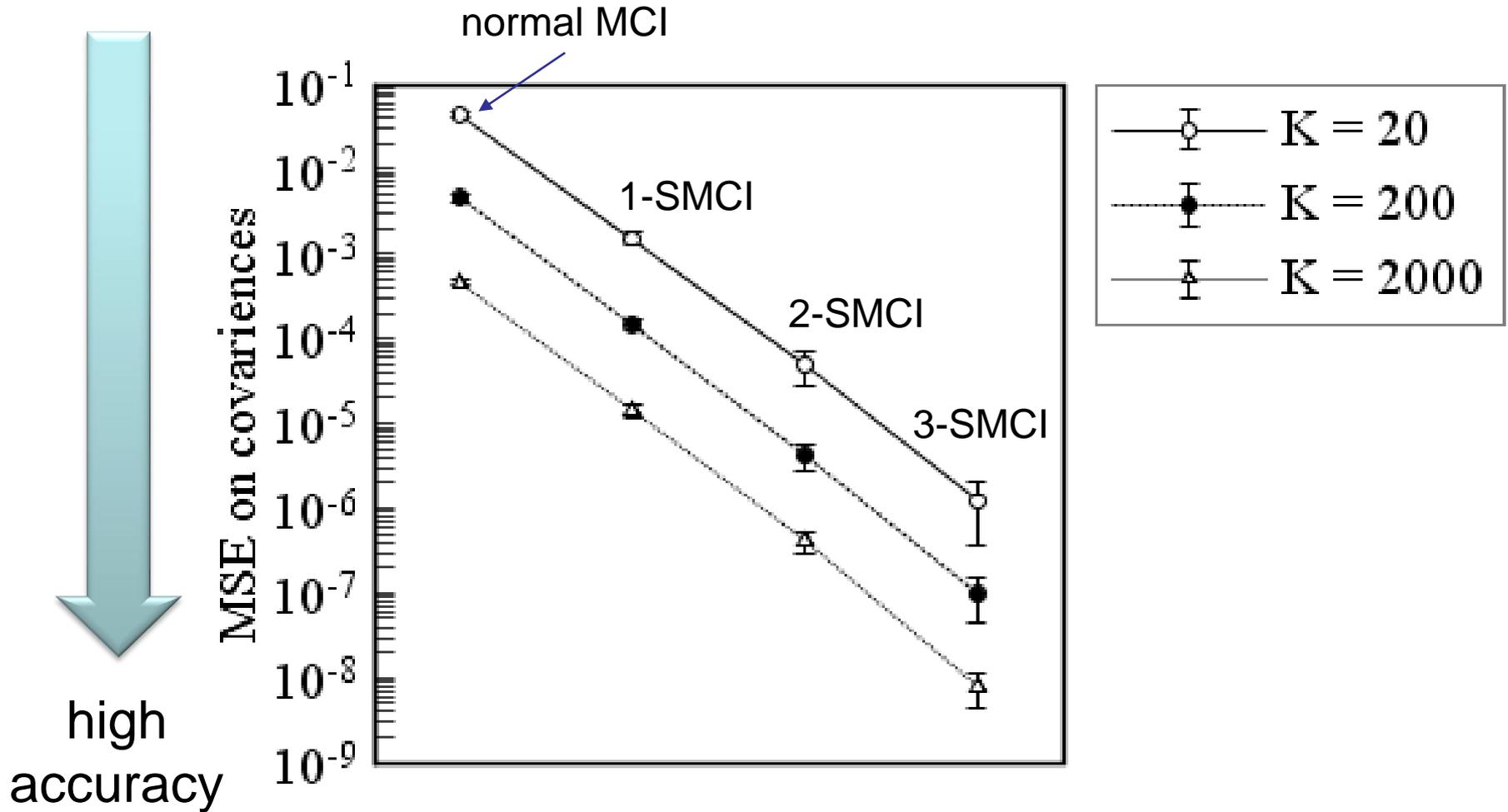
SAMPLING

$\{\mathbf{s}^{(k)} \mid k = 1, 2, \dots, K\}$



# Numerical Experiment

157 / 141



A higher-order SMCI is at least 10 times better than a lower-order one.

MSE : mean square error between the exact values and the approximated values

## Basic Idea

$$\frac{\partial l_D(\theta)}{\partial b_i} = \frac{1}{N} \sum_{\mu=1}^N d_i^{(\mu)} - \mathbf{E}[v_i | \theta], \quad \frac{\partial l_D(\theta)}{\partial w_{ij}} = \frac{1}{N} \sum_{\mu=1}^N d_i^{(\mu)} d_j^{(\mu)} - \mathbf{E}[v_i v_j | \theta]$$

Approximate these intractable expectations by SMCI

By regarded the **observed data set  $D$**  as the **MC samplings  $S$**  generated from the model and by **using  $D$  instead of  $S$** ,

we can skip the sampling phase in the MCI.

# Numerical Experiment (Learning)

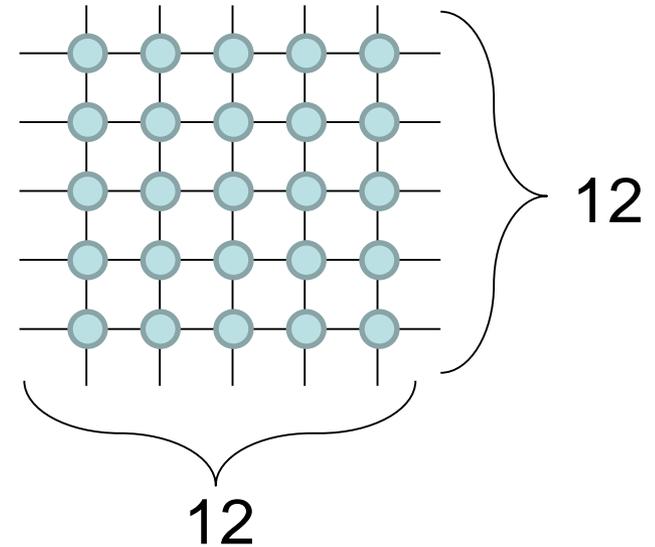
159 / 141

## Generative Model

Boltzmann machine

$$P_g(\mathbf{v}) = \frac{1}{Z} \exp \left( \sum_{\{i,j\} \in E} w_{ij} v_i v_j \right), \quad v_i \in \{+1, -1\}$$

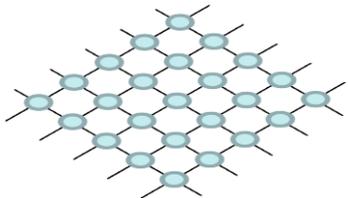
$$w_{ij} \text{ i.i.d. } \sim U[-0.5, 0.5]$$



Gibbs sampling  
method

SAMPLING

$P_g(\mathbf{v})$

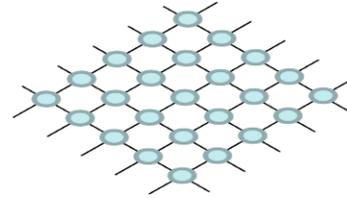


$D = \{\mathbf{d}^{(\mu)} \mid \mu = 1, 2, \dots, N\}$

synthetic data set

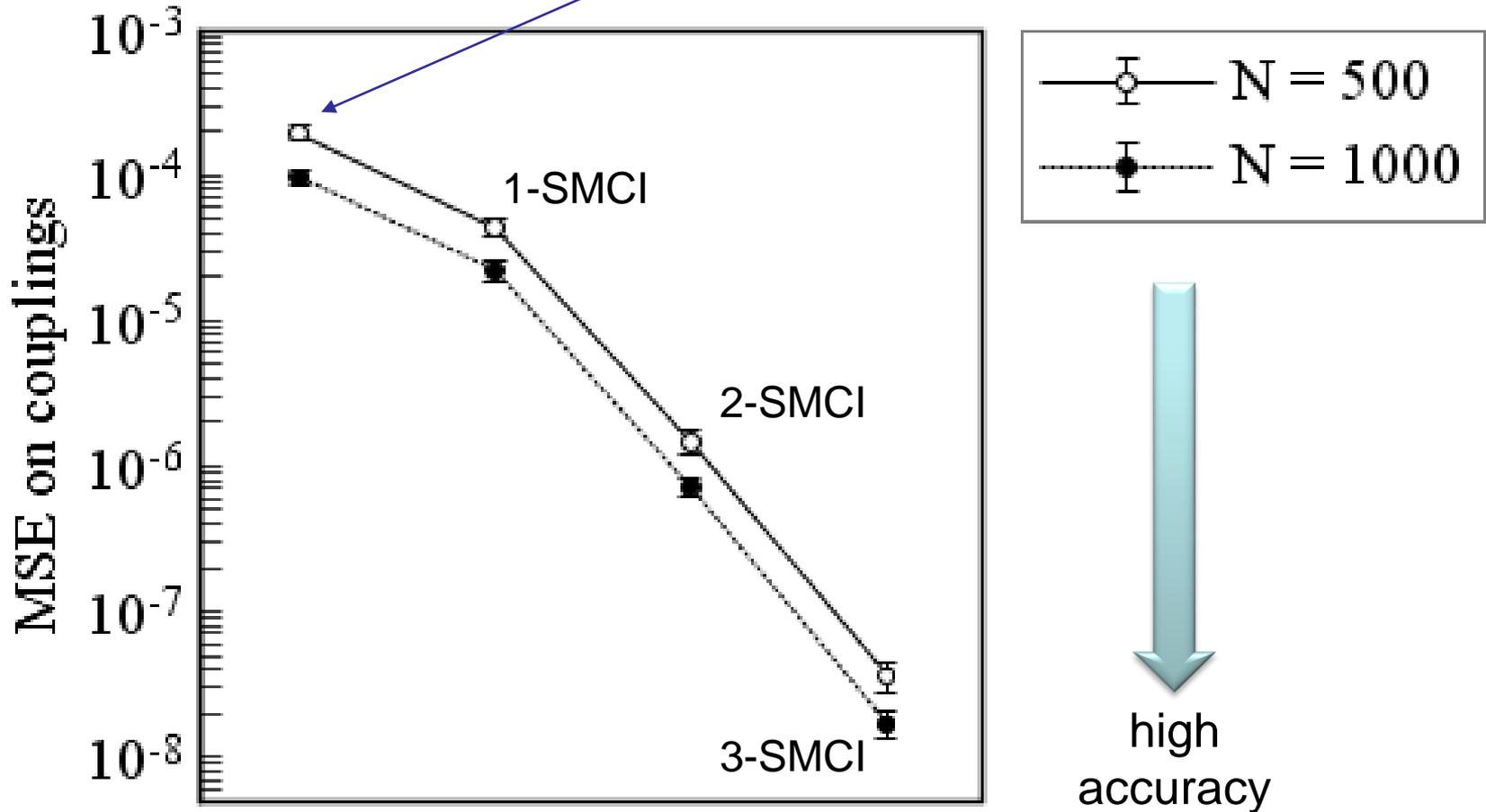
LEARNING

$P(\mathbf{v} \mid \theta)$



# Numerical Experiment (Learning)

maximum pseudo-likelihood estimation (MPLE) [Basag: *The Statistician*, 75]



MSE : mean square error between the ML estimators and the approximated values

# Numerical Experiment (Learning)

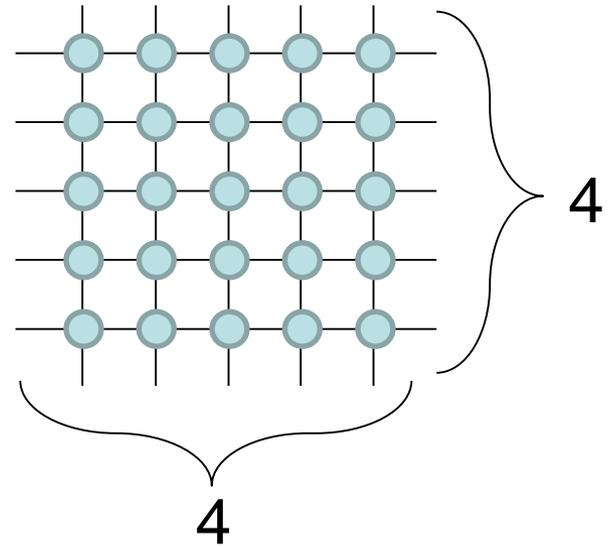
161 / 141

## Generative Model

Boltzmann machine

$$P_g(\mathbf{v}) = \frac{1}{Z} \exp \left( \sum_{\{i,j\} \in E} w_{ij} v_i v_j \right), \quad v_i \in \{+1, -1\}$$

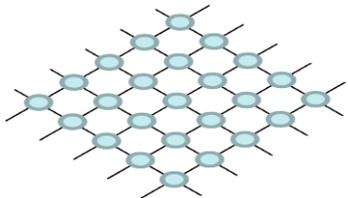
$$w_{ij} \text{ i.i.d. } \sim U[-0.5, 0.5]$$



Gibbs sampling  
method

SAMPLING

$P_g(\mathbf{v})$

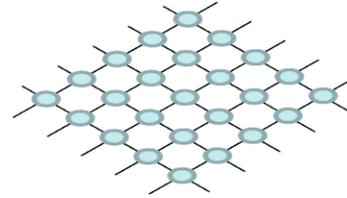


$D = \{ \mathbf{d}^{(\mu)} \mid \mu = 1, 2, \dots, 200 \}$

synthetic data set

LEARNING

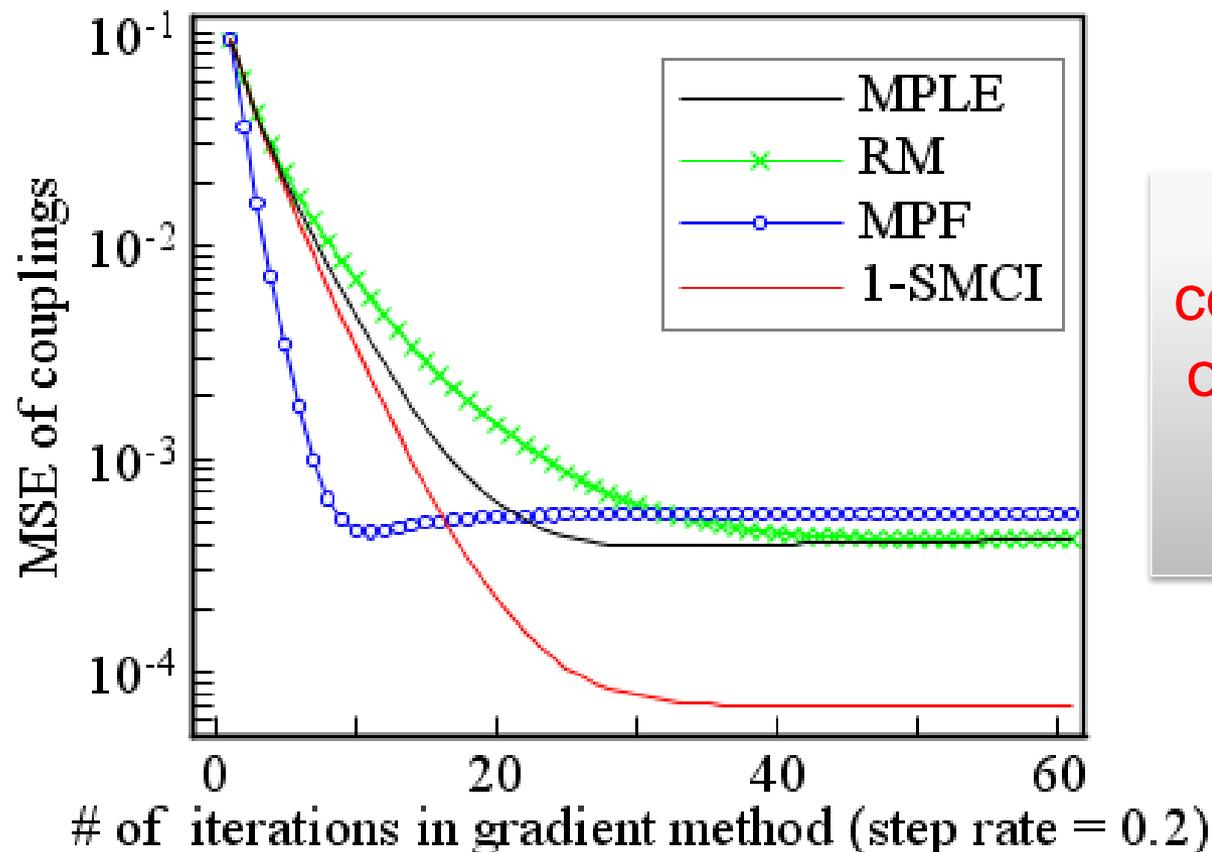
$P(\mathbf{v} \mid \theta)$



# Numerical Experiment (Learning)

162 / 141

comparison with the other recent learning algorithms



The orders of computational costs of these algorithms are the same

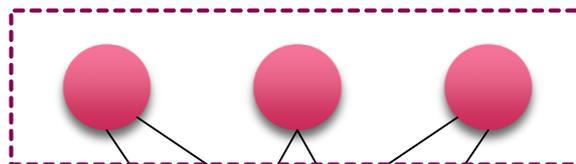
MPLE: maximum pseudo-likelihood estimation [Basag: *The Statistician*, 75]  
RM: ratio matching [Hyvärinen: *CSDA*, 07]  
MPF: minimum probability flow [Sohl-Dickstein: *PRL*, 11]

# 平均場アルゴリズムを用いた ガウス型制限ボルツマンマシン

# 背景

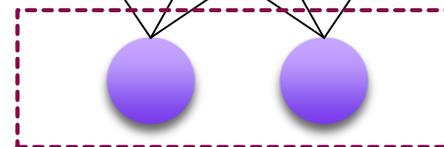
制限ボルツマンマシンは2部グラフ上に定義される確率モデル

隠れ層  $H$   
隠れ変数  $h$ : データと関連しない



$$h = \{h_j \mid j \in H\}$$

可視層  $V$   
可視変数  $v$ : データと関連する

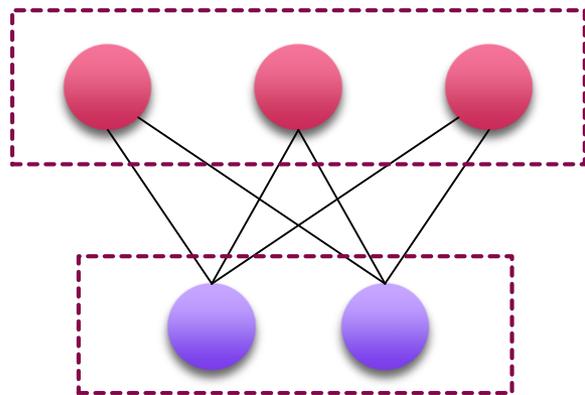


$$v = \{v_i \mid i \in V\}$$

制限ボルツマンマシン (RBM) → 可視変数: 離散 隠れ変数: 離散

ガウス型制限ボルツマンマシン (GRBM) → 可視変数: **連続** 隠れ変数: 離散

# 背景



## 学習フェーズ

得られたデータを用いて最適なパラメータの値を決定

e.g. コントラストティブ・ダイバージェンス

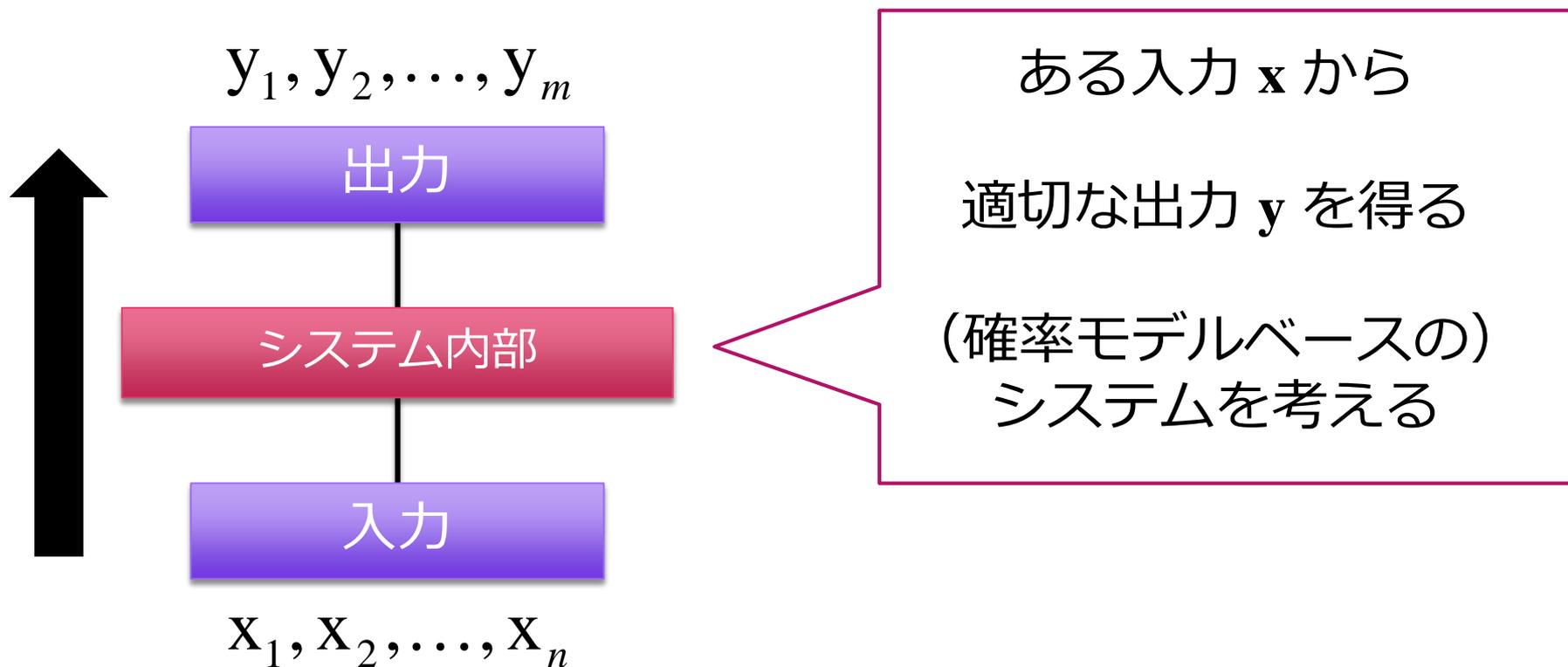
## 推定フェーズ

学習により得られたモデルを用いた入出力値推定

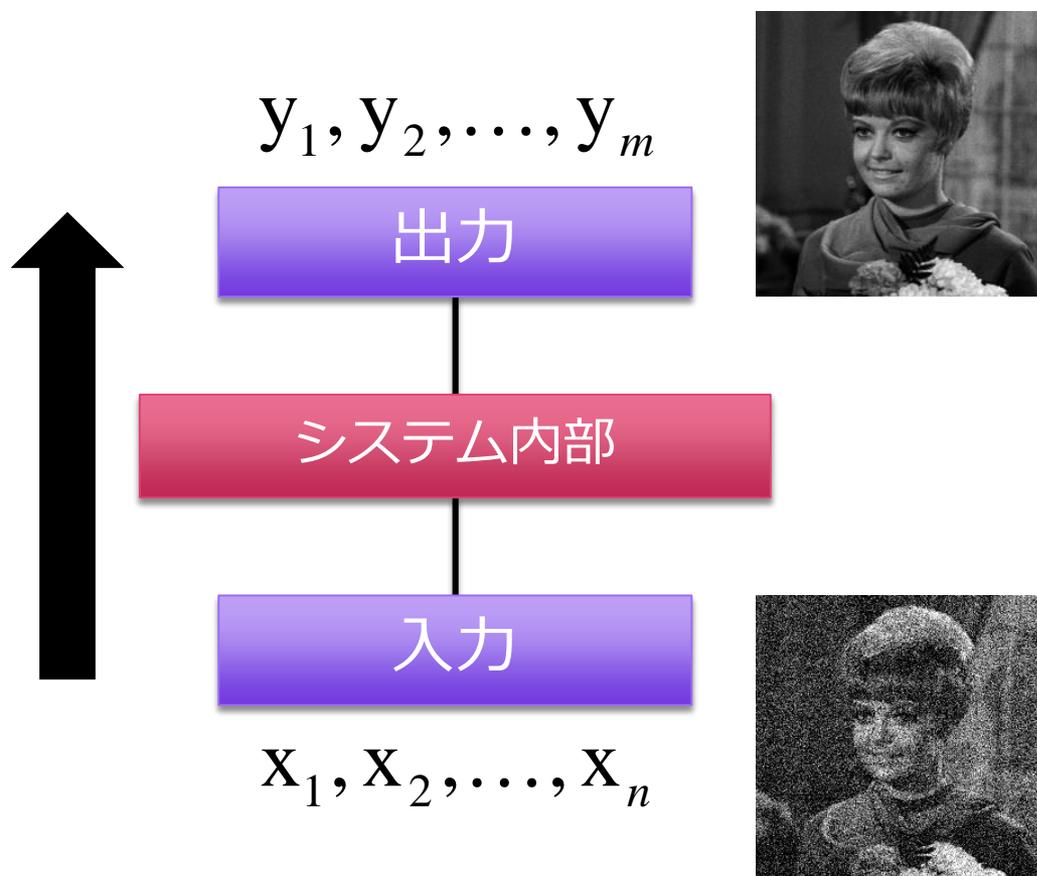
出力**変数の期待値**が出力値となる (MMSE推定法)

ギブスサンプリング  
や  
平均場近似法

## 学習？推定？ ～ 研究領域のイメージのために～



# 画像修復システム設計を例として



例えば  
画像修復システムなら

入力  $x$  が劣化画像  
で  
出力  $y$  が修復画像

という具合

# 機械学習

 $y_1, y_2, \dots, y_m$ 

出力



システム内部

入力

 $x_1, x_2, \dots, x_n$ 

入力から適切な出力を出すように  
システムの**内部パラメータを最適化**する必要がある

教師あり機械学習では

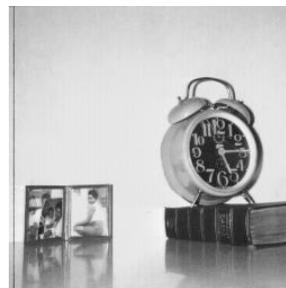
入力とそれに対応する出力（教師）  
の集合（**訓練データ集合**）  
をシステムに与えて

その与えた入出力データの関係が  
正しく成立するように  
内部パラメータを最適化する

# 機械学習

$y_1, y_2, \dots, y_m$

出力



...

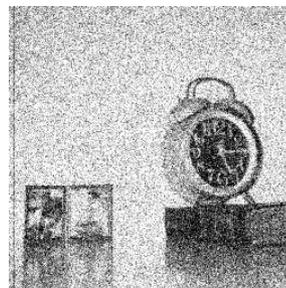
システム内部



訓練データ集合の入出力関係が  
(高確率で) 実現するようなパラメータ  
を探すことが機械学習の目的

入力

$x_1, x_2, \dots, x_n$



...



# 出力の推定

 $y_1, y_2, \dots, y_m$ 

出力

システム内部

入力

 $x_1, x_2, \dots, x_n$ 

機械学習により  
システムの内部パラメータ  
の最適化が完了したら

得られたシステムを用いて  
実際に画像修復を行うことが出来る

# 出力の推定

 $y_1, y_2, \dots, y_m$ 

出力

システム内部

入力

 $x_1, x_2, \dots, x_n$ 

確率モデルベースのシステムでは  
出力  $y$  は確率的にサンプリングされる  
ので  
サンプリング毎に違った値となる

我々が欲しいのは唯一の修復画像なので  
出力を何か一つに決めなければならない



機械学習により  
システムの内部パラメータ  
の最適化が完了したら  
  
得られたシステムを用いて  
実際に画像修復を行うことが出来る

# 出力の推定

 $y_1, y_2, \dots, y_m$ 

出力

システム内部

入力

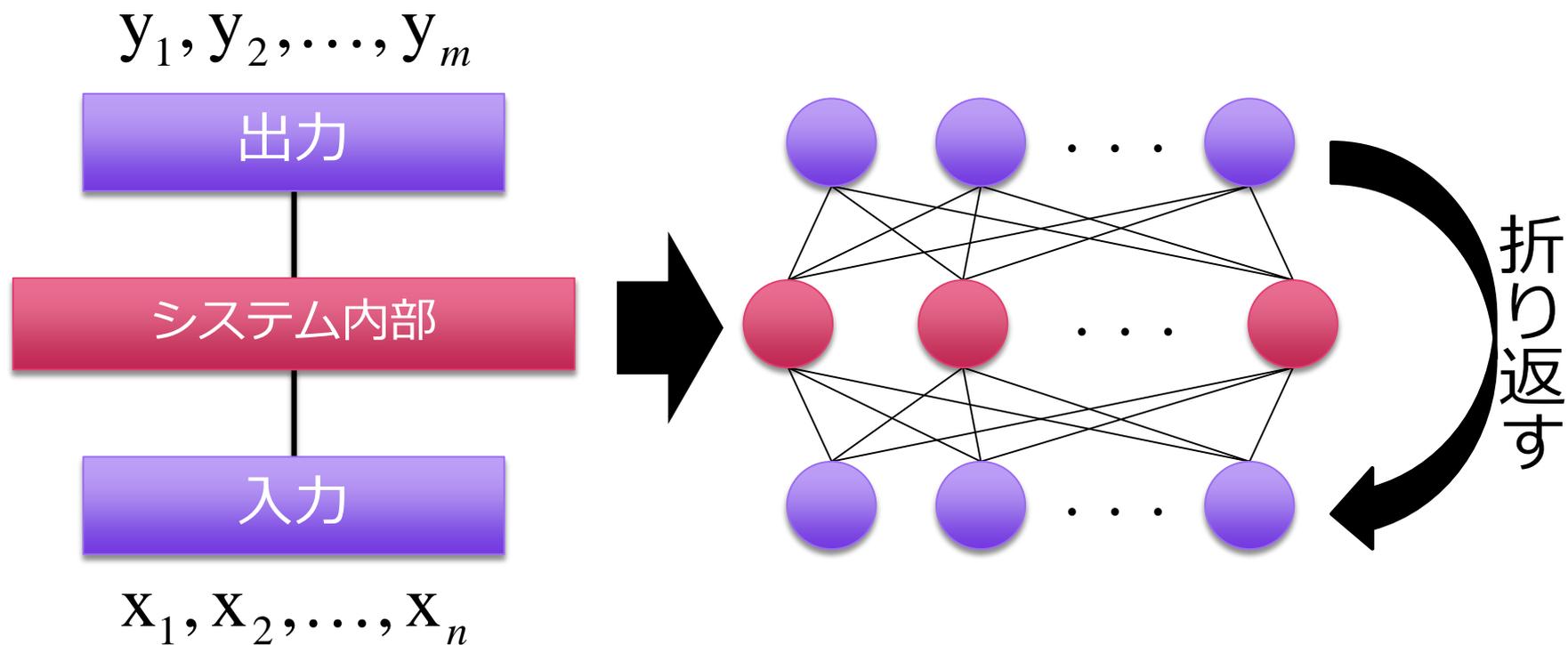
 $x_1, x_2, \dots, x_n$ 

確率モデルベースのシステムでは  
出力  $y$  は確率的にサンプリングされる  
ので  
サンプリング毎に違った値となる

我々が欲しいのは唯一の修復画像なので  
出力を何か一つに決めなければならない

推定法（出力を一つに決める方法）は様々だが  
MMSE（最小平均 2 乗誤差）推定では  
サンプリングの**期待値**を出力とする

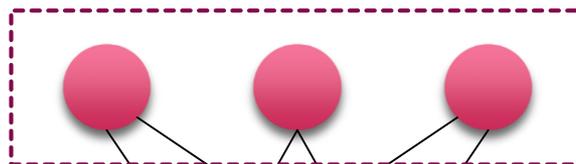
## RBM との架け橋



# 背景（再掲）

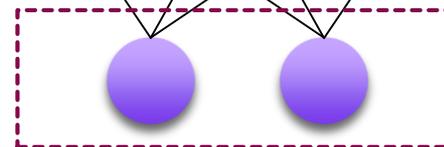
制限ボルツマンマシンは2部グラフ上に定義される確率モデル

隠れ層  $H$   
隠れ変数  $h$ : データと関連しない



$$h = \{h_j \mid j \in H\}$$

可視層  $V$   
可視変数  $v$ : データと関連する

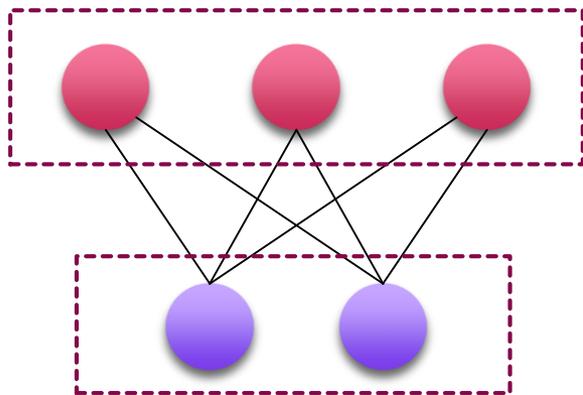


$$v = \{v_i \mid i \in V\}$$

制限ボルツマンマシン (RBM) → 可視変数：離散 隠れ変数：離散

ガウス型制限ボルツマンマシン (GRBM) → 可視変数：**連続** 隠れ変数：離散

# 背景（再掲）



## 学習フェーズ

得られたデータを用いて最適なパラメータの値を決定

e.g. コントラストティブ・ダイバージェンス

## 推定フェーズ

学習により得られたモデルを用いた入出力値推定

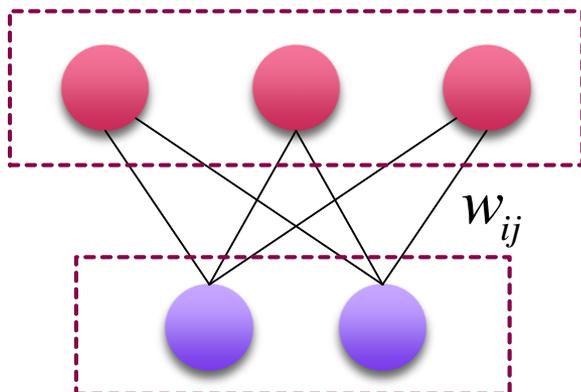
出力**変数の期待値**が出力値となる（MMSE推定法）

ギブスサンプリング  
や  
平均場近似法

# 目的

ガウス型制限ボルツマンマシン上での  
**平均場近似**を基礎とした  
推定アルゴリズム（期待値近似アルゴリズム）  
の性能を**情報理論的に調べる**

# ガウス型制限ボルツマンマシン



$$\mathbf{h} = \{h_j \in \mathcal{X} \mid j \in H\}$$

離散変数

$$\mathbf{v} = \{v_i \in (-\infty, \infty) \mid i \in V\}$$

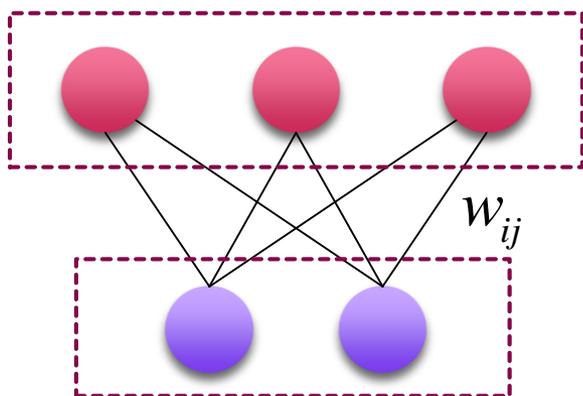
連続変数

エネルギー関数

$$E(\mathbf{v}, \mathbf{h}; \theta) := \sum_{i \in V} \frac{(v_i - b_i)^2}{2\sigma_i^2} - \sum_{j \in H} c_j h_j - \sum_{i \in V} \sum_{j \in H} \frac{w_{ij}}{\sigma_i^2} v_i h_j$$

$$P(\mathbf{v}, \mathbf{h} \mid \theta) := \frac{1}{Z(\theta)} \exp(-E(\mathbf{v}, \mathbf{h}; \theta))$$

# ガウス型制限ボルツマンマシン



$b_i, c_j$  : 可視変数, 隠れ変数の  
 バイアスパラメータ  
 $\sigma_i$  : 可視変数の分散パラメータ  
 $w_{ij}$  : 結合パラメータ

 $\theta$ 

エネルギー関数

$$E(\mathbf{v}, \mathbf{h}; \theta) := \sum_{i \in V} \frac{(v_i - b_i)^2}{2\sigma_i^2} - \sum_{j \in H} c_j h_j - \sum_{i \in V} \sum_{j \in H} \frac{w_{ij}}{\sigma_i} v_i h_j$$

$$P(\mathbf{v}, \mathbf{h} | \theta) := \frac{1}{Z(\theta)} \exp(-E(\mathbf{v}, \mathbf{h}; \theta))$$

# 期待値計算の組み合わせ爆発

$$E[v_i | \theta] = \int_{-\infty}^{\infty} d\mathbf{v} \sum_{\mathbf{h} \in \mathcal{X}^{|\mathcal{H}|}} v_i P(\mathbf{v}, \mathbf{h} | \theta), \quad E[h_j | \theta] = \int_{-\infty}^{\infty} d\mathbf{v} \sum_{\mathbf{h} \in \mathcal{X}^{|\mathcal{H}|}} h_j P(\mathbf{v}, \mathbf{h} | \theta)$$

$$\int_{-\infty}^{\infty} d\mathbf{v} (\dots) = \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} \dots \int_{-\infty}^{\infty} (\dots) dv_1 dv_2 \dots dv_{|\mathcal{V}|} \quad \leftarrow \text{多重積分}$$

$$\sum_{\mathbf{h} \in \mathcal{X}^{|\mathcal{H}|}} (\dots) = \sum_{h_1 \in \mathcal{X}} \sum_{h_2 \in \mathcal{X}} \dots \sum_{h_{|\mathcal{H}|} \in \mathcal{X}} (\dots) \quad \leftarrow \text{多重和}$$

計算量は変数の数に対して指数的  $\rightarrow$  NP困難

近似が必須！  $\rightarrow$  平均場近似

# 平均場近似の枠組み

多次元結合分布  $P(\mathbf{x}) = P(x_1, x_2, \dots, x_n)$

に対して

$$\text{テスト分布 } T(\mathbf{x}) = \prod_{i=1}^n \rho_i(x_i)$$

各変数についての  
統計的独立性を  
仮定した分布  
(平均場制約)

を準備する

カルバック・ライブラー情報量

$$K[T] := \sum_{\mathbf{x}} T(\mathbf{x}) \ln \frac{T(\mathbf{x})}{P(\mathbf{x})} \geq 0$$

を最小にする  $T$  が元の  $P$  の平均場近似分布となる

## 平均場近似の枠組み

カルバック・ライブラー情報量は  
異なる分布間の距離的な尺度



元の分布  $P$  にもっとも近い近似分布  $T$  が得られる

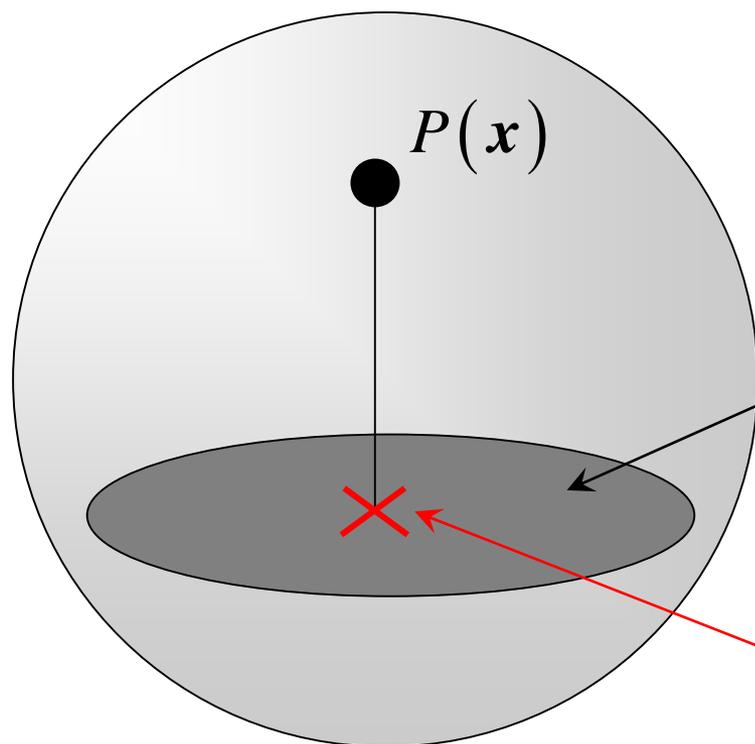
カルバック・ライブラー情報量

$$K[T] := \sum_x T(\mathbf{x}) \ln \frac{T(\mathbf{x})}{P(\mathbf{x})} \geq 0$$

を最小にする  $T$  が元の  $P$  の平均場近似分布となる

# カルバック・ライブラー情報量

カルバック・ライブラーダイバージェンス  
カルバック・ライブラー擬距離  
などとも呼ばれる



$T(x)$  が動ける空間

平均場近似で求まる近似分布  $T$

# GRBM に対する平均場近似

2 種類の平均場近似が定式化可能

Type I 平均場近似

$P(\mathbf{v}, \mathbf{h} | \theta)$  に対して平均場近似  
(全変数に対する平均場近似)

Type II 平均場近似

$P(\mathbf{h} | \theta) = \int_{-\infty}^{\infty} d\mathbf{v} P(\mathbf{v}, \mathbf{h} | \theta)$  に対して平均場近似  
(周辺分布に対する平均場近似)

# GRBM に対する平均場近似

2 種類の平均場近似が定式化可能

Type I 平均場近似

$$K_1[T_1] := \int_{-\infty}^{\infty} d\mathbf{v} \sum_{\mathbf{h} \in \mathcal{X}^{|\mathcal{H}|}} T_1(\mathbf{v}, \mathbf{h}) \ln \frac{T_1(\mathbf{v}, \mathbf{h})}{P(\mathbf{v}, \mathbf{h} | \theta)} \text{ の最小化}$$

Type II 平均場近似

$$K_2[T_2] := \int_{-\infty}^{\infty} d\mathbf{v} \sum_{\mathbf{h} \in \mathcal{X}^{|\mathcal{H}|}} P(\mathbf{v} | \mathbf{h}, \theta) T_2(\mathbf{h}) \ln \frac{P(\mathbf{v} | \mathbf{h}, \theta) T_2(\mathbf{h})}{P(\mathbf{v}, \mathbf{h} | \theta)}$$

の最小化

# GRBM に対する平均場近似

2 種類の平均場近似が定式化可能

Type I 平均場近似

Type II 平均場近似



計算量としては全く同等の近似計算アルゴリズムを導く

$$O(|V| |H|)$$

Which method is better ??

## 定理（情報理論的立場）

任意の GRBM に対して  
$$\min K_1 [T_1] \geq \min K_2 [T_2]$$
  
が常に成り立つ

カルバック・ライブラー情報量の観点で

**Type II 平均場近似**の方が**より元の分布に近い**近似分布を与える

※この定理はより強い定理の系として導かれる

※元の定理は任意の確率モデルで上記と類似の不等式が成り立つことを主張している

## 定理 (統計物理的立場)

$$\min K_1 [T_1] = \underbrace{F_1(\theta)}_{\text{平均場自由エネルギー}} - \underbrace{F(\theta)}_{\text{真の自由エネルギー}}$$

$$\min K_2 [T_2] = \underbrace{F_2(\theta)}_{\text{平均場自由エネルギー}} - \underbrace{F(\theta)}_{\text{真の自由エネルギー}}$$

任意の GRBM に対して

$$F_1(\theta) \geq F_2(\theta) \geq F(\theta)$$

が常に成り立つ

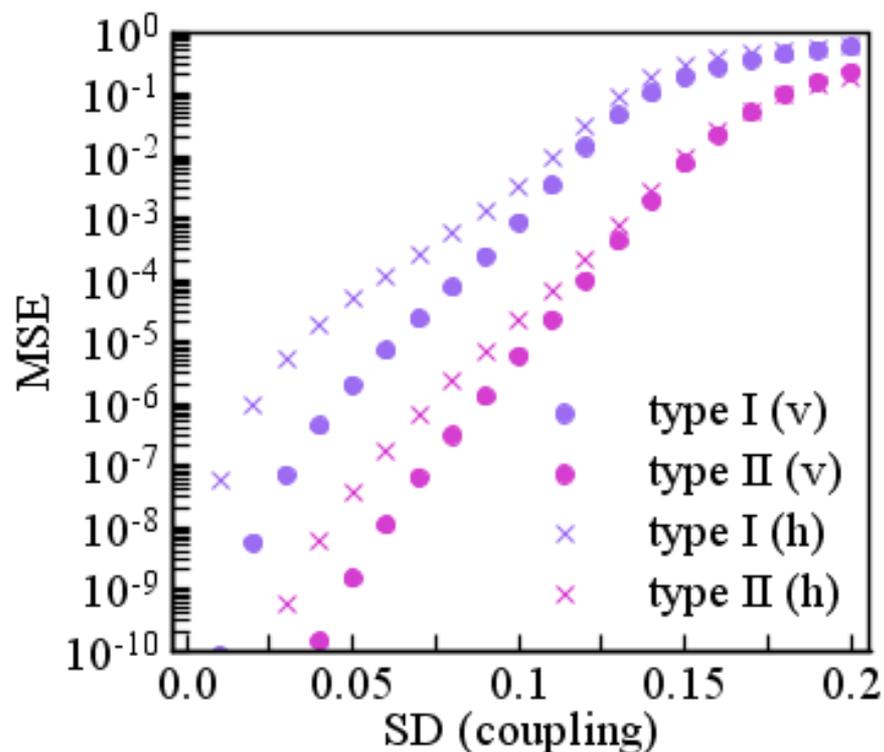
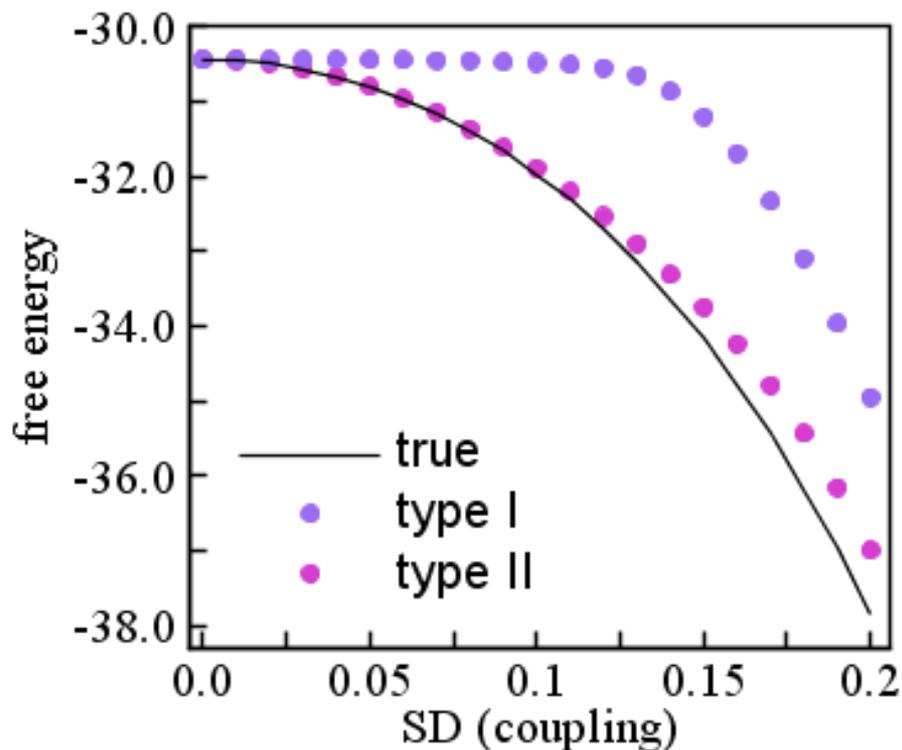
Type II 平均場自由エネルギーの方が真の自由エネルギーに近い

## 数値実験による定量的比較

$$|V|=12, |H|=24, \sigma_i^2=1$$

$$\mathcal{X} \in \{+1, -1\}$$

$b_i, c_j \leftarrow \mathcal{N}(0, 0.1^2)$  で生成,  $w_{ij} \leftarrow \mathcal{N}(0, \text{SD}^2)$  で生成



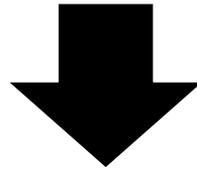
## 前半のまとめ

### GRBM 上での平均場近似による推定アルゴリズム

- 計算量的に同等な 2 種類のア​​ルゴリズム (Type I & II) を定式化
- 解析的・数値的比較から Type II の方が常に高性能であることが分かった
- 平均場法による推定アルゴリズムに対する重要な知見

# 確率伝搬法

計算量爆発の問題はどう頑張っても解決できない



真の値を得ることは諦めて**近似値**で良しとしましょうよ

高効率な近似計算法の発展

確率伝搬法 (belief propagation; BP)

近似計算法の発展が現在の MRF 実装の楚となっている

人工知能研究の中で生まれた計算技術 [Pearl, 82]

物性物理理論の中にも実は同様の手法がある [Bethe, 35]

Bethe 近似

情報科学分野と物理分野の狭間で発展してきた

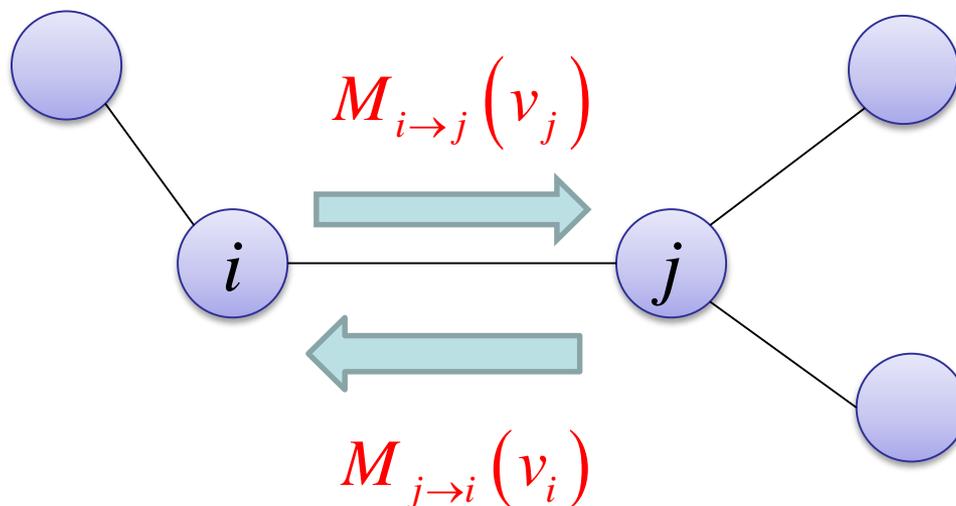
強力な近似計算技法

動的計画法的に再帰式を解く形式の計算アルゴリズム

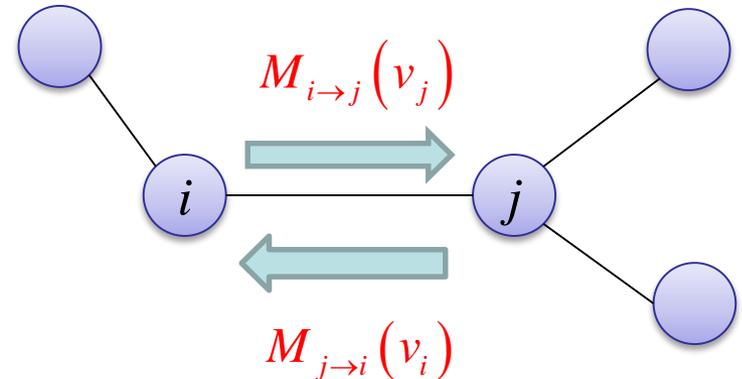
標記の簡便さのため MRF を次の形式に変形しておく

$$P(\mathbf{v}) = \frac{1}{Z} \exp \left( \sum_{i \in V} \Phi_i(v_i) + \sum_{\{i,j\} \in E} \Psi_{\{i,j\}}(v_i, v_j) \right) = \frac{1}{Z} \prod_{i \in V} \phi_i(v_i) \prod_{\{i,j\} \in E} \varphi_{\{i,j\}}(v_i, v_j)$$

ノードからノードへの**メッセージ**と呼ばれる量を定義する



- メッセージはすべてのリンク上に定義される
- メッセージは双方向に伝搬する
- メッセージは行き先側のノードの変数の関数



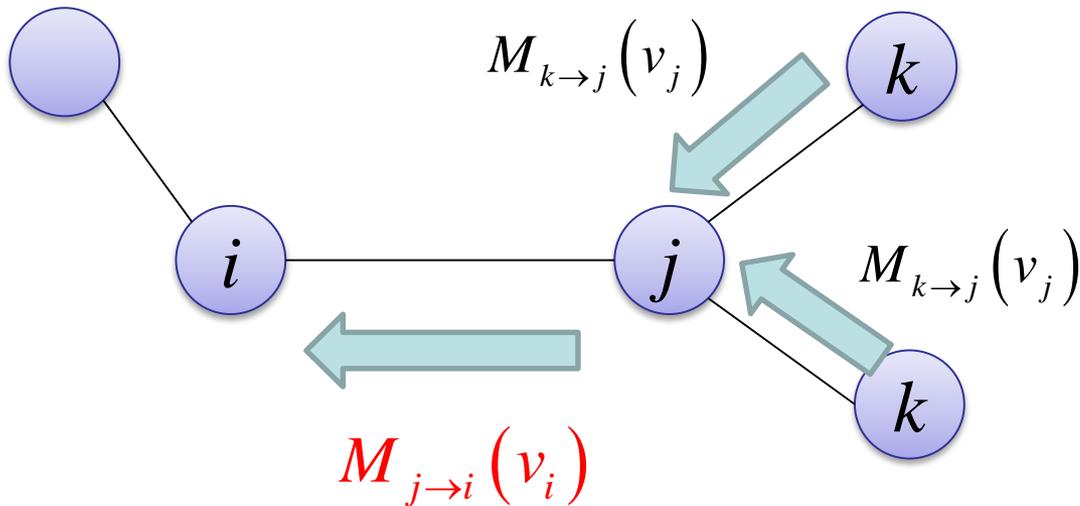
## メッセージ伝搬則 (message-passing rule)

$$M_{j \rightarrow i}(v_i) \propto \sum_{x_j} \phi_j(v_j) \varphi_{\{i, j\}}(v_i, v_j) \prod_{k \in \partial(j) \setminus \{i\}} M_{k \rightarrow j}(v_j)$$

メッセージはこの連立非線形方程式を満足する必要がある

## メッセージ伝搬則 (message-passing rule)

$$M_{j \rightarrow i}(v_i) \propto \sum_{x_j} \phi_j(v_j) \phi_{\{i,j\}}(v_i, v_j) \prod_{k \in \partial(j) \setminus \{i\}} M_{k \rightarrow j}(v_j)$$



行き先以外のノード  
からやってくる

メッセージを集めて

行き先に送る

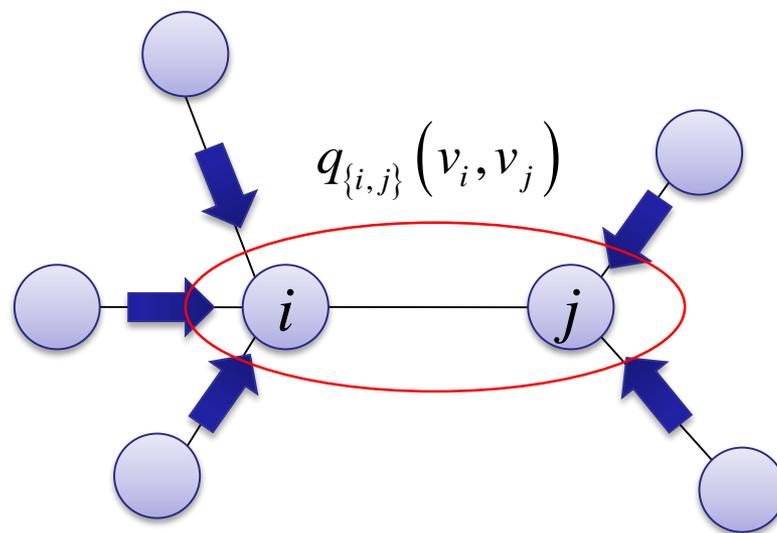
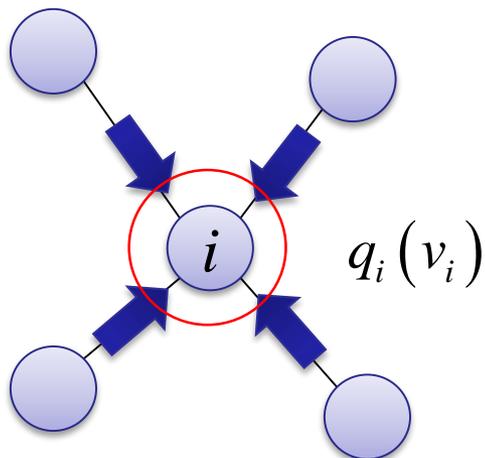
メッセージ伝搬則の式はリンクの本数の2倍の数存在する  
それを連立方程式と見立てて**数值的に解く**

求めたメッセージを用いて**近似周辺確率**を計算する

$$q_i(v_i) = \frac{1}{z_i} \phi_i(v_i) \left( \prod_{j \in \partial(i)} M_{j \rightarrow i}(v_i) \right)$$

$z_i, z_{\{i,j\}} \leftarrow$  規格化定数

$$q_{\{i,j\}}(v_i, v_j) = \frac{1}{z_{\{i,j\}}} \phi_i(v_i) \phi_j(v_j) \varphi_{\{i,j\}}(v_i, v_j) \left( \prod_{k \in \partial(i) \setminus \{j\}} M_{k \rightarrow i}(v_i) \right) \left( \prod_{l \in \partial(j) \setminus \{i\}} M_{l \rightarrow j}(v_j) \right)$$



求めた近似周辺確率を用いて**期待値**を次の要領で近似する

$$E[f(v_i)] = \sum_{\mathbf{v}} f(v_i) P(\mathbf{v}) \approx \sum_{v_i} f(v_i) q_i(v_i)$$

$$E[f(v_i, v_j)] = \sum_{\mathbf{v}} f(v_i, v_j) P(\mathbf{v}) \approx \sum_{v_i} \sum_{v_j} f(v_i, v_j) q_{\{i,j\}}(v_i, v_j)$$

1 変数の期待値のみでなく 2 変数の期待値も計算することができる！

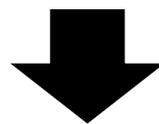
## 確率伝搬法の手続きまとめ

メッセージ

→ 近似周辺確率

→ 近似期待値

グラフ構造が木構造(閉路を持たない構造)なら真の期待値を与える



リンクがスパースなグラフ上で非常に高性能

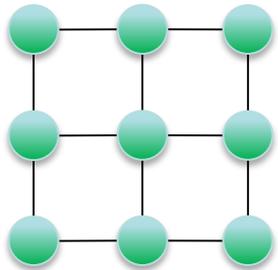
**sum-product アルゴリズム**という名でも知られている

確率伝搬法

が

MRF上での計算処理の基礎を成している

# 実際どんなもんか



$$P(\mathbf{v} | \mathbf{b}, \mathbf{w}) = \frac{1}{Z} \exp \left( \sum_{i \in V} b_i v_i + \sum_{\{i, j\} \in E} w_{ij} v_i v_j \right) \quad v_i \in \{+1, -1\}$$

パラメータは平均 0 で標準偏差 0.4 のガウス乱数によりランダムに生成

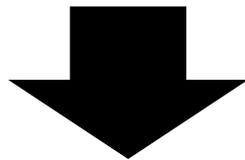
$E[v_i]$



sum-product アルゴリズムは MRF の期待値を計算する手法

確率伝搬法にはもう一つの手法がある！

MRF の確率を最大とする  $\nu$  の値を見つける手法



max-product アルゴリズム

$$M_{j \rightarrow i}(v_i) \propto \sum_{v_j} \phi_j(v_j) \phi_{\{i,j\}}(v_i, v_j) \prod_{k \in \partial(j) \setminus \{i\}} M_{k \rightarrow j}(v_j)$$



sum-product

max-product アルゴリズムの  
メッセージ伝搬則 (message-passing rule)

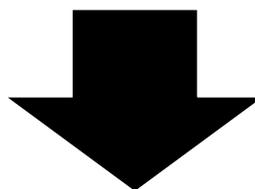
$$M_{j \rightarrow i}(v_i) \propto \max_{v_j} \left\{ \phi_j(v_j) \phi_{\{i,j\}}(v_i, v_j) \prod_{k \in \partial(j) \setminus \{i\}} M_{k \rightarrow j}(v_j) \right\}$$

sum-product アルゴリズムのメッセージ伝搬則の和を

マックス演算に変更すると max-product アルゴリズム

求めたメッセージを用いて**近似最大値**を計算する

$$v_i^{\text{BP}} = \max_{v_i} \phi_i(v_i) \left( \prod_{j \in \partial(i)} M_{j \rightarrow i}(v_i) \right)$$



$$\boldsymbol{v}^{\text{BP}} \approx \arg \max_{\boldsymbol{v}} P(\boldsymbol{v})$$

MRF を最大とする  
変数の値の  
近似値を得る

sum-product アルゴリズムと同様に  
木構造のグラフ上なら真の値が求まる

MRFの期待値を計算したい  
→ sum-product アルゴリズム

MRFの最大値を計算したい  
→ max-product アルゴリズム

