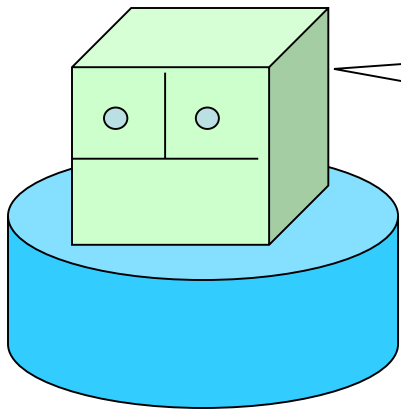
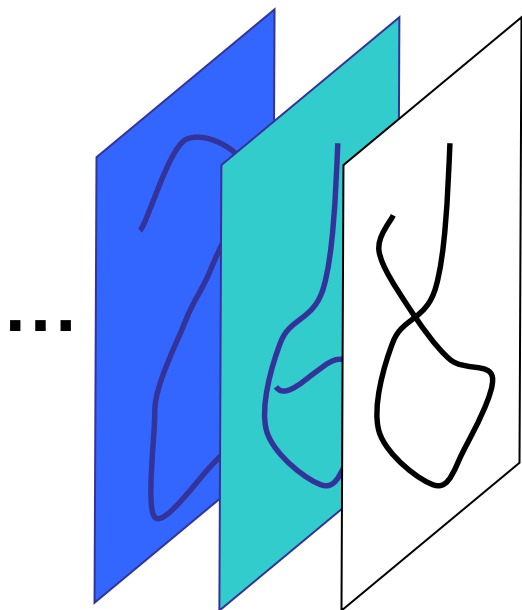


学習理論の練習 4

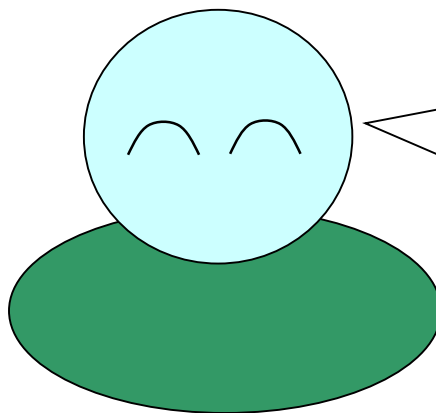
1 漫画でわかる学習理論

教師あり学習

例

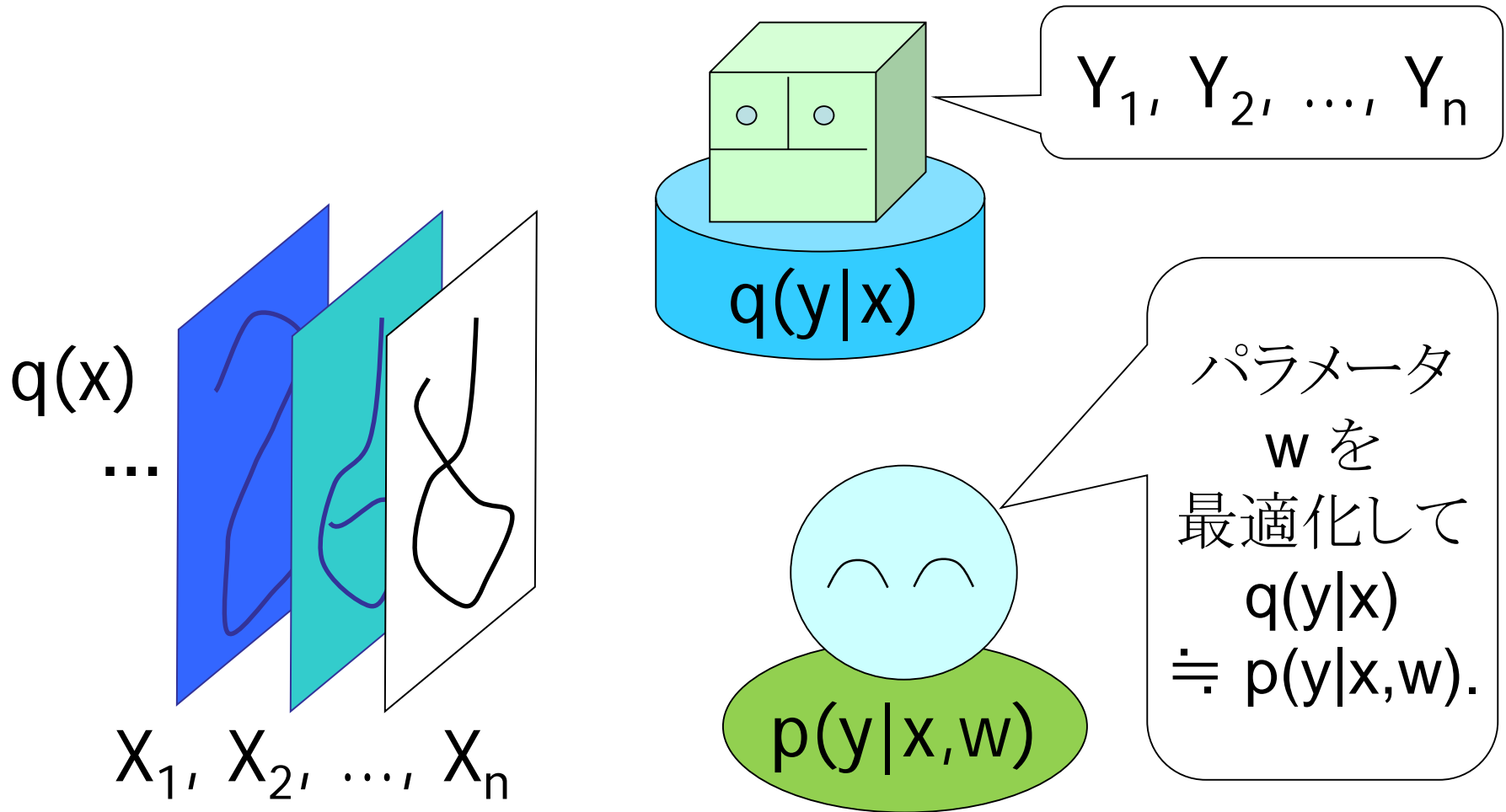


先生
8, 6, 2...

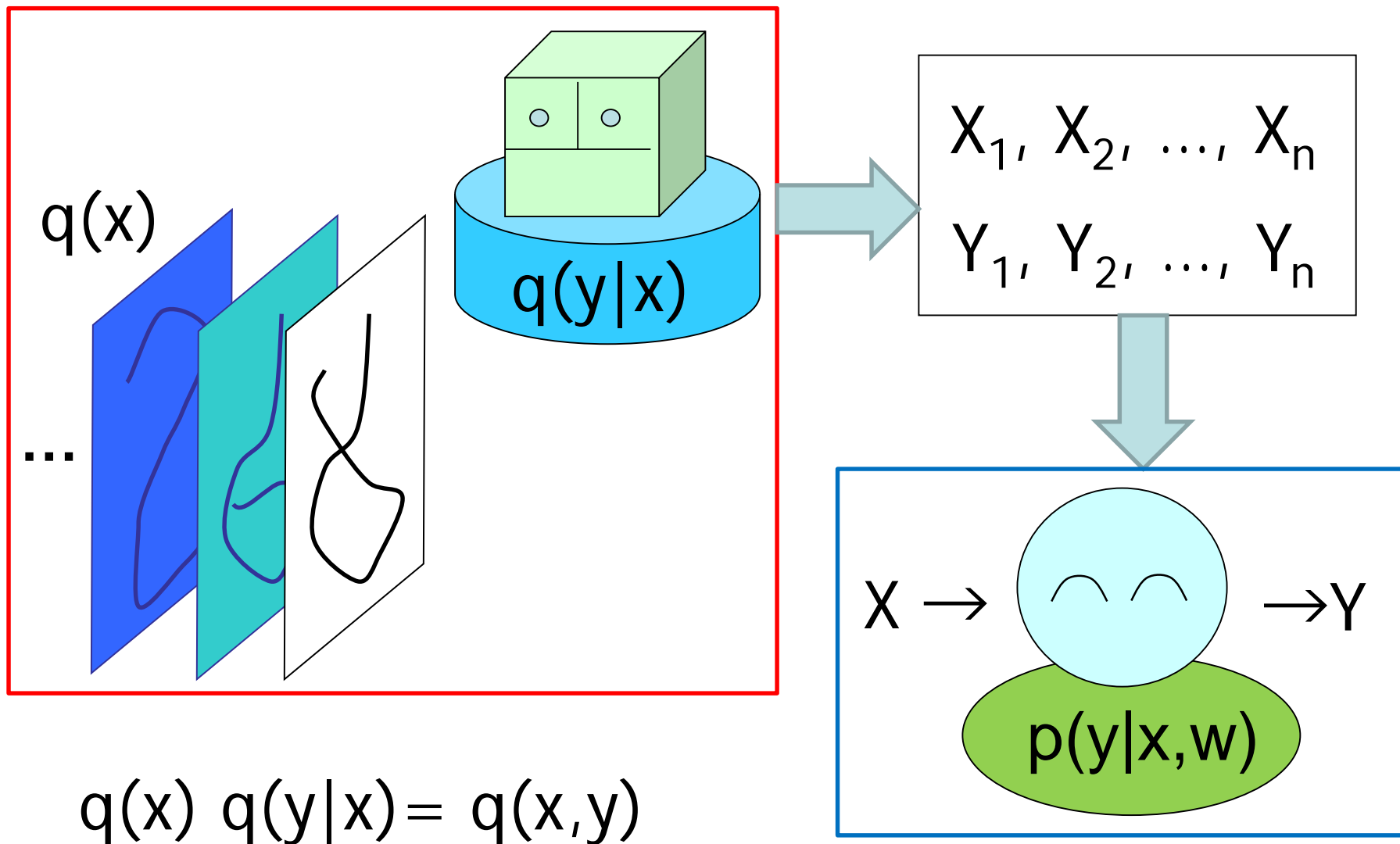


生徒
文字を
読みます

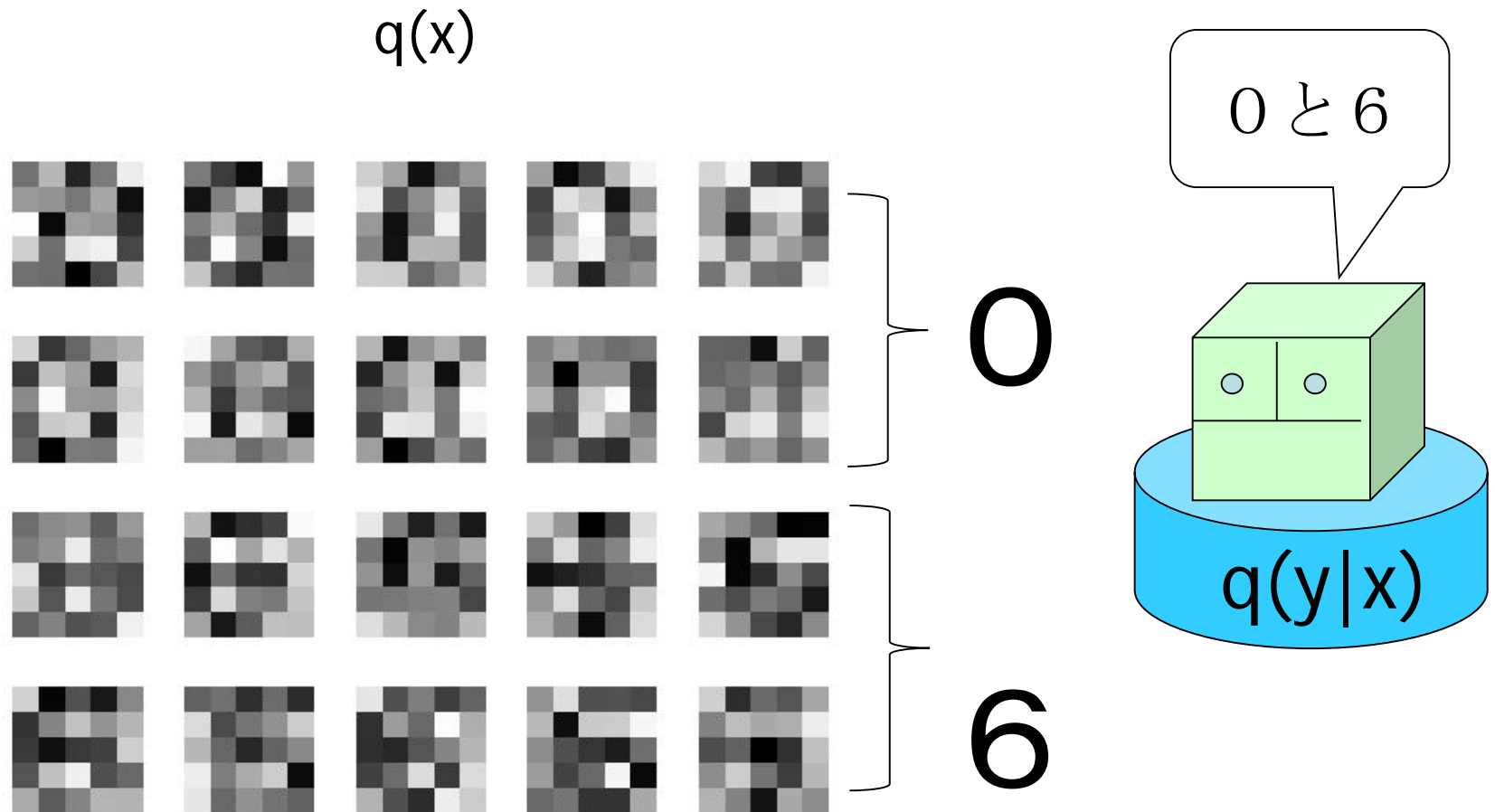
数学的に記述すると



情報源と学習モデル



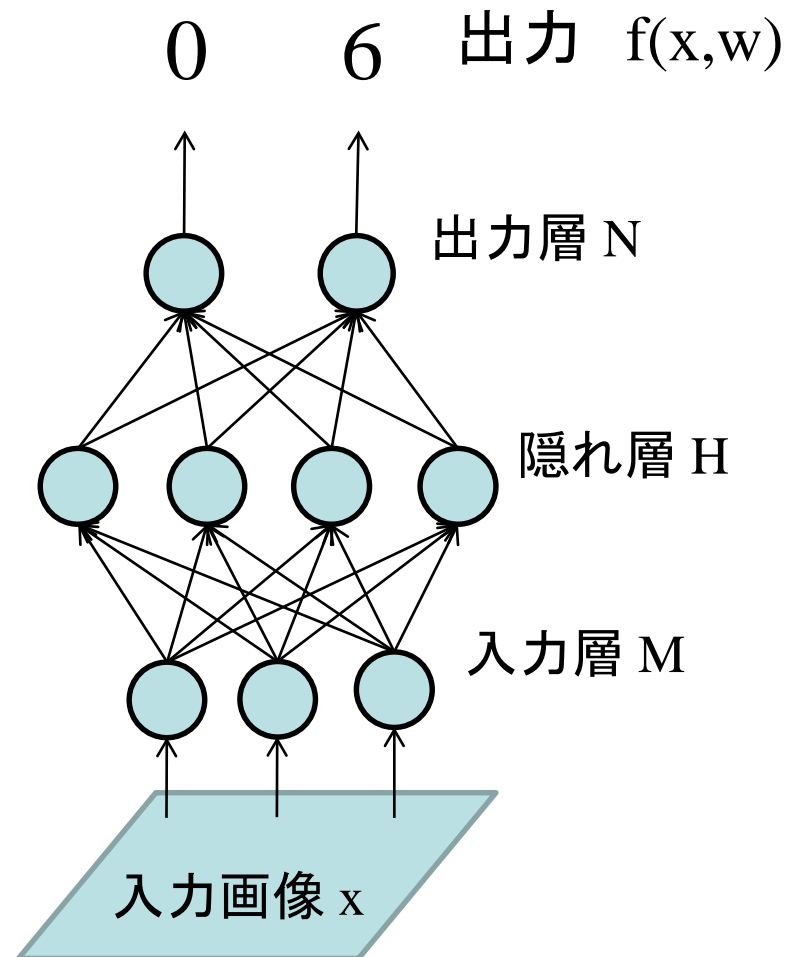
情報源 $q(x)q(y|x)$ の例



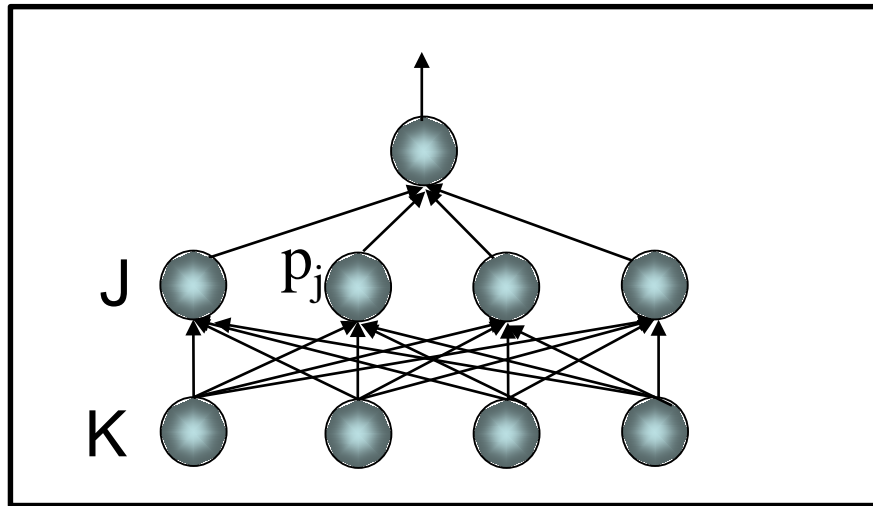
学習モデルの例



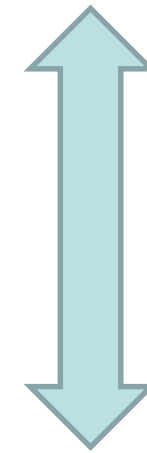
=



ここ数年の話題1: 敵対的データ生成



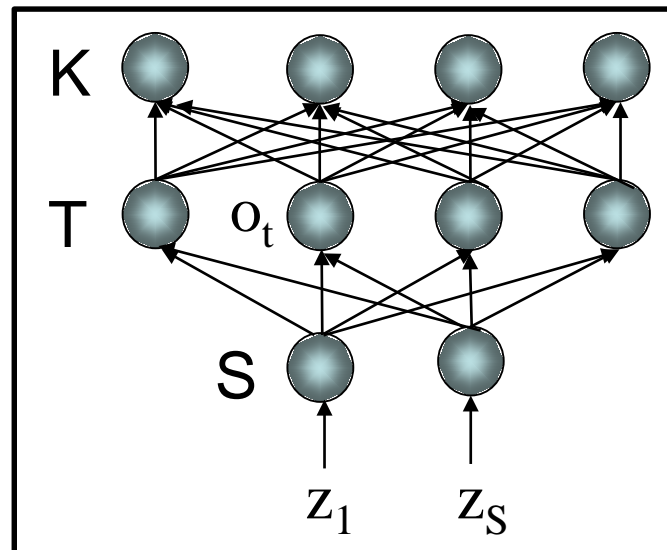
データと偽データを
識別するネットの
学習1



学習1と
学習2を
交互に
繰り返す

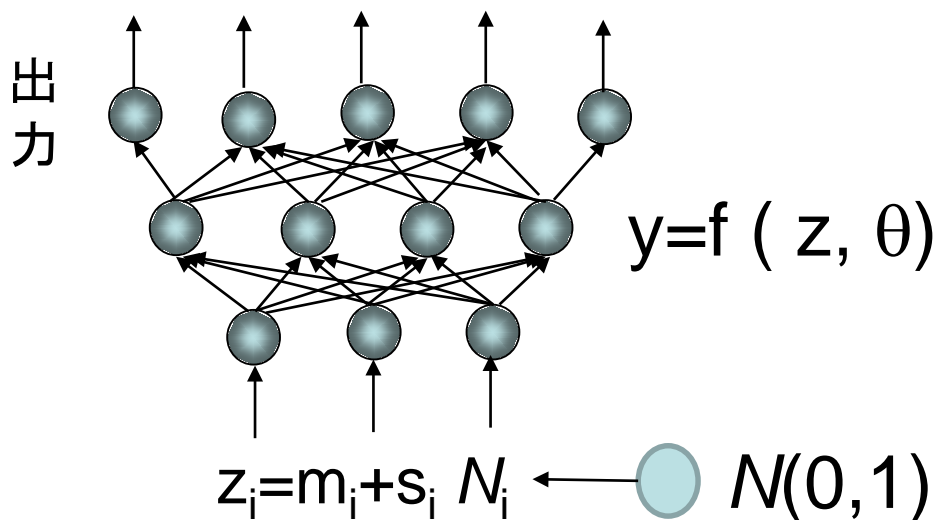
学習データ

$$X = \{x_k\}$$

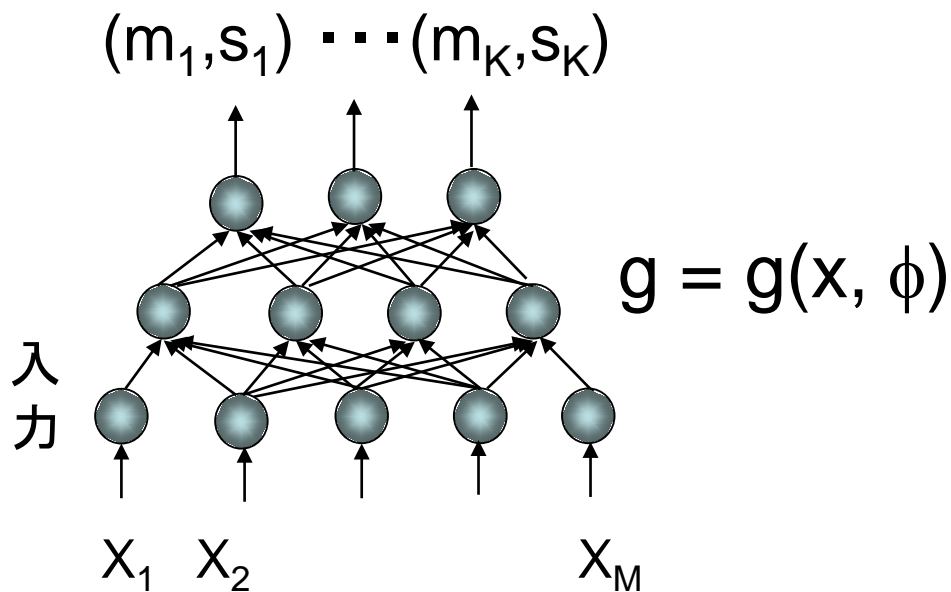


偽データを
生成するネットの
学習2

ここ数年の話題2: 変分オートエンコーダー



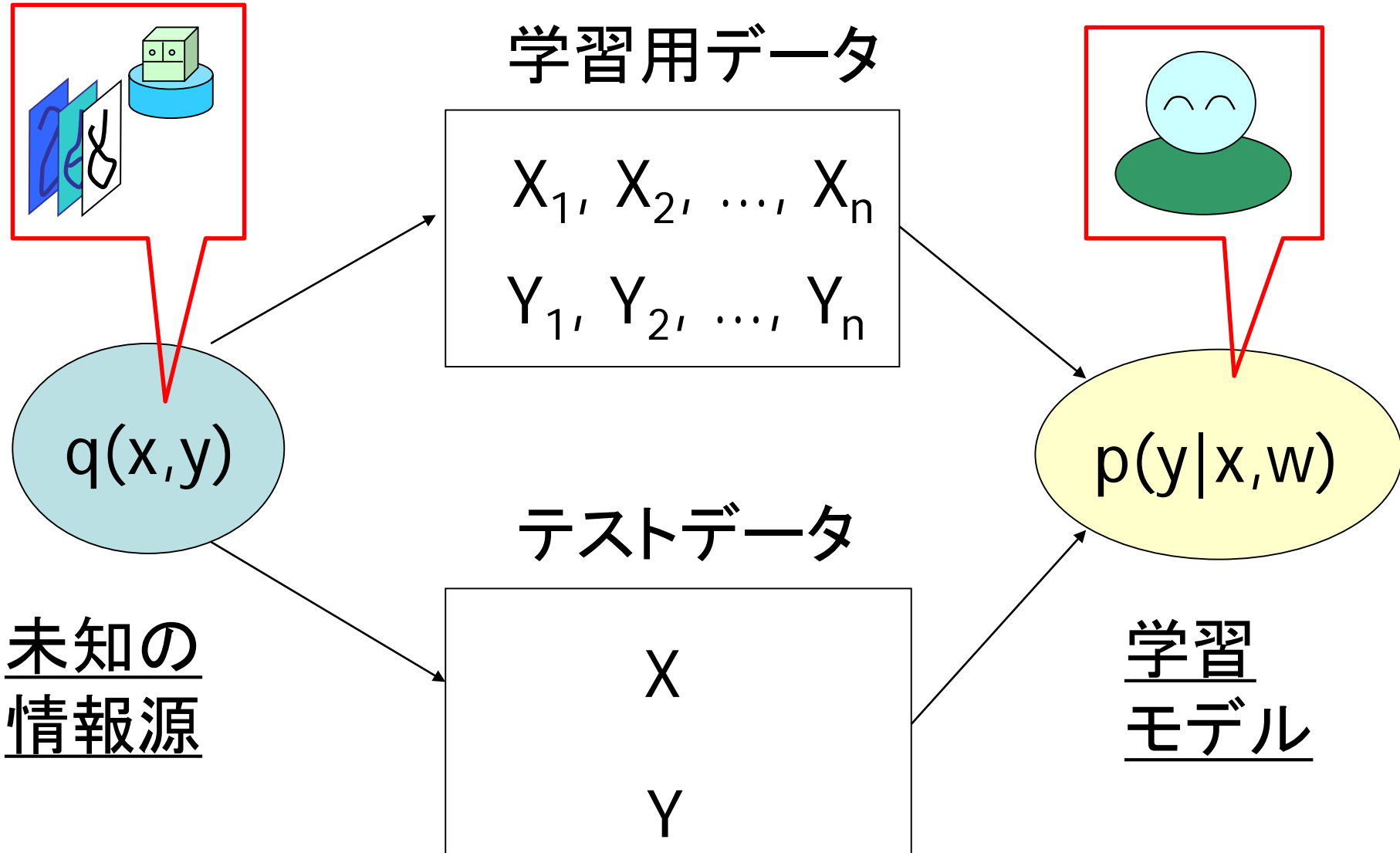
オートエンコーダー
の中間部分に
正規雑音を入れて
その平均と分散を
学習すると



データがあるところ
だけでなくその周り
を確率的に補完
することができる

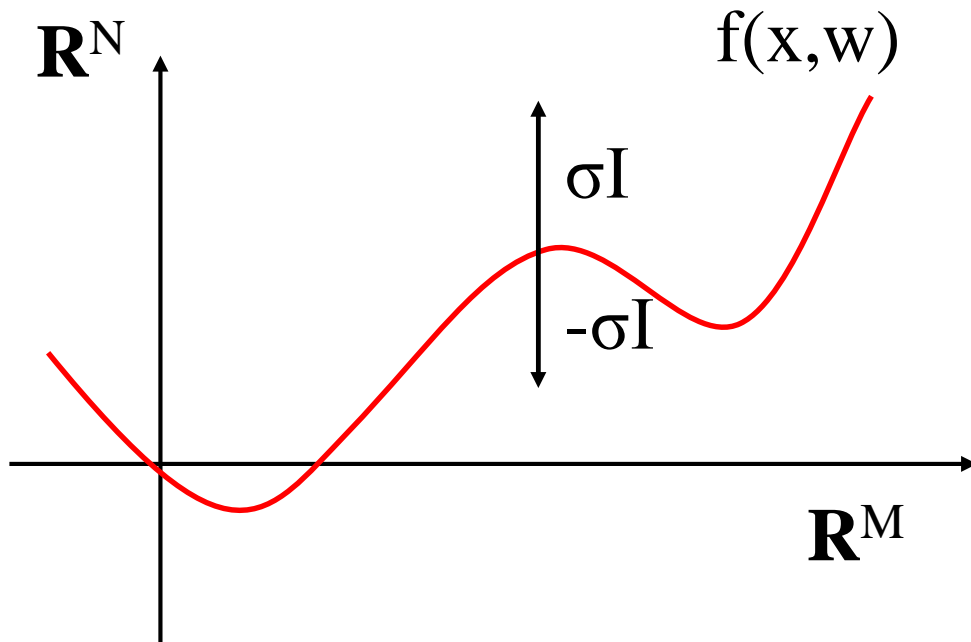
2 漫画から数式へ

教師あり学習の枠組み



学習モデルの例1

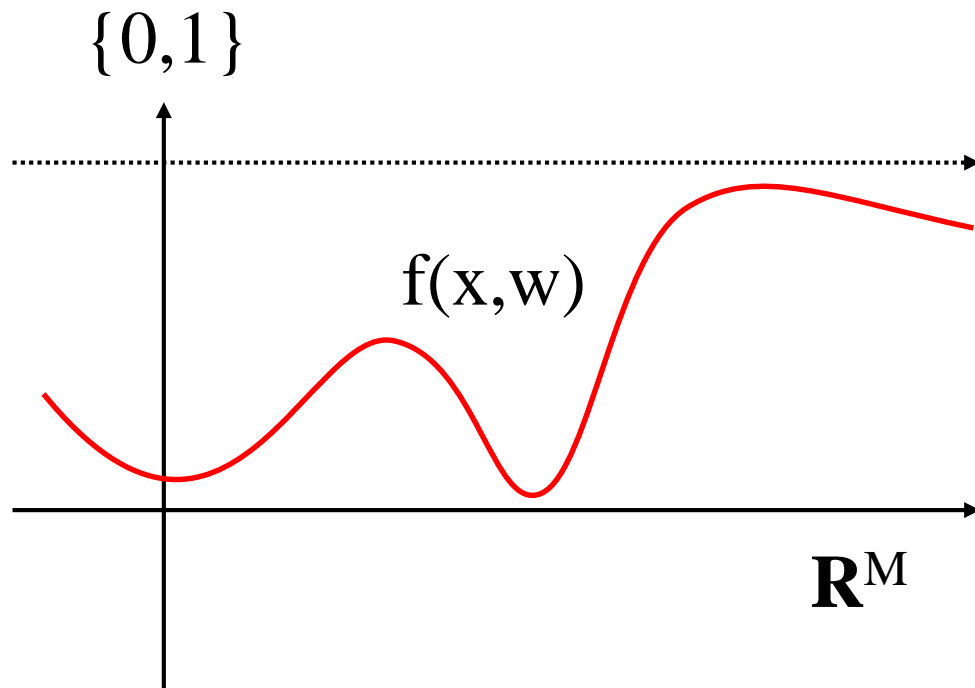
$$p(y|x, w, \sigma) = \frac{1}{(2\pi\sigma^2)^{N/2}} \exp \left\{ -\frac{1}{2\sigma^2} \|y - f(x, w)\|^2 \right\}$$



- $p(y|x, w, \sigma)$ は X が与えられたときの Y の条件つき確率密度
- Y は実数全体をとりうる
- Y は平均 $f(x, w)$ 分散 $\sigma^2 I$ の正規分布
- w と σ がパラメータ

学習モデルの例2

$$p(y|x,w) = f(x,w)^y (1-f(x,w))^{1-y}$$



- $p(y|x,w)$ は X が与えられたときの Y の条件つき確率
- Y は $\{0,1\}$ をとる
- Y は 確率 $f(x,w)$ で1 確率 $(1-f(x,w))$ で0
- w がパラメータ

条件つき確率の尤度関数

データ $\{(X_i, Y_i) ; i=1, 2, \dots, n\}$, モデル $p(y|x, w)$ が与えられたとき

$$L(w) = \prod_{i=1}^n p(Y_i | X_i, w)$$

を(条件つき確率の) **尤度関数**という。

- もしも「モデル=真」であれば、 $L(w)$ は $\{X_i\}$ が与えられたときの $\{Y_i\}$ の確率密度関数
- 尤度関数を最大にする w を **最尤推定量**という。
- 最尤推定量は良い推定量ではないがよく利用される。深層学習では悪すぎて使えない。学習理論の問題は一般には最適化問題には帰着しない。

対数尤度関数

尤度関数の対数をとって

$$\log L(w) = \sum_{i=1}^n \log p(Y_i|X_i, w)$$

を(条件つき確率の)対数尤度関数という。

- \log は単調増大関数だから対数尤度関数を最大にするパラメータが最尤推定量
- $H(w) = -(1/n) \log L(w)$ と書けば $L(w) = \exp(-nH(w))$
- 尤度関数と対数尤度関数はどちらを考えても同じだが計算上対数尤度関数のほうが便利なが多い。

対数尤度関数の例1

学習モデル

$$p(y|x, w, \sigma) = \frac{1}{(2\pi\sigma^2)^N} \exp \left\{ -\frac{1}{2\sigma^2} \|y - f(x, w)\|^2 \right\}$$

の対数尤度関数は

$$\log L(w) = -\frac{1}{2\sigma^2} \sum_{i=1}^n \|Y_i - f(X_i, w)\|^2 - Nn \log \sigma - (N/2) \log(2\pi)$$

最尤推定量(w^*, s^*) は

$$w^* = \operatorname{argmin}_w \sum_{i=1}^n \|Y_i - f(X_i, w)\|^2$$
$$s^* = \operatorname{sqrt} \left\{ \frac{1}{nN} \sum_{i=1}^n \|Y_i - f(X_i, w^*)\|^2 \right\}$$

対数尤度関数の例2

学習モデル $p(y|x,w) = f(x,w)^y (1-f(x,w))^{1-y}$

の対数尤度関数は

$$\log L(w) = \sum_{i=1}^n \{ Y_i \log f(X_i, w) + (1-Y_i) \log (1-f(X_i, w)) \}$$

最尤推定量 w^* はこれを最大にするパラメータ

- X を0と1に識別する問題ではどちらの学習モデルを用いることもできる。同じ神経回路網を使っても学習モデルとしては異なる。
- X を $\{(0,0), (1,0), (0,1), (1,1)\}$ に識別したい場合は多次元に拡張すればよい。個々に問題を解いてもよい。

3 数式から数理へ

対数尤度関数の符号反転

対数尤度関数を(1/n)倍して符号反転したものを $H(w)$ と書く

$$H(w) = - (1/n) \sum_{i=1}^n \log p(Y_i|X_i, w)$$

データを発生している確率密度が $q(x) q(y|x)$ であるとする。
データの数 n が大きいとき 大数の法則から

$$H(w) \doteq - \iint q(x) q(y|x) \log p(y|x, w) dx dy$$

条件つきエントロピーを下記で定義する

$$S = - \iint q(x) q(y|x) \log q(y|x) dx dy$$

尤度とKL情報量

$$\begin{aligned} H(w) &\doteq - \iint q(x) q(y|x) \log q(y|x) \, dx dy \\ &\quad + \iint q(x) q(y|x) \log \{ q(y|x) / p(y|x, w) \} \, dx dy \\ &= S + D(q(y|x) \parallel p(y|x, w)) \end{aligned}$$

ここで $D(q(y|x) \parallel p(y|x, w))$ は条件つき確率のKL情報量であり、常に非負で、零になるのは $q(y|x) = p(y|x, w)$ のときだけである。

データの数 n が十分大きければ

尤度が大きい = 対数尤度が大きい = $H(w)$ が小さい

⇔ 真とモデルの条件つき確率のKL情報量が小さい

学習理論は何をしたいのか

- 真の分布がわからなくても、尤度関数はデータとモデルだけで記述できる。
- 尤度関数を大きくすると「真とモデルのKL情報量」がほぼ小さくなる。つまり未知の真がほぼ推定できる。
- 「尤度関数最大 \Leftrightarrow 真とモデルのKL情報量が小さい」が等価であるためには $n=\infty$ が必要。
- n が有限なら「対数尤度関数 \div 真とモデルのKL情報量」のずれは、関数が複雑であるほど大きくなる。深層学習では、そのずれはベイズ法以外では解明されていない。
- 対数尤度関数とKL情報量のずれがわかると、データに対して最適なモデルやハイパーパラメータを決めることや統計的検定を作ることができるようになる。

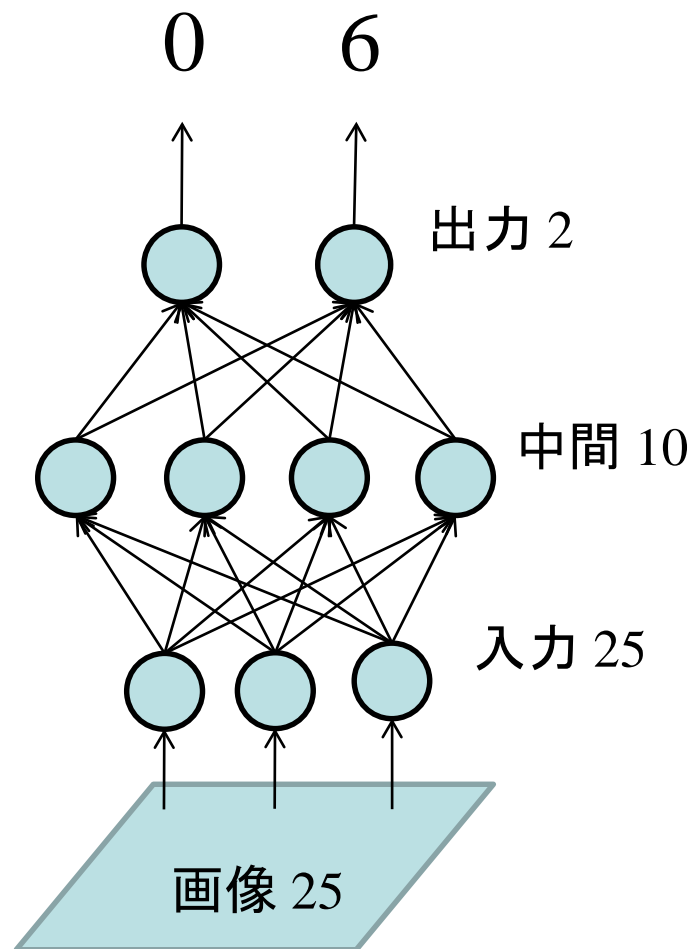
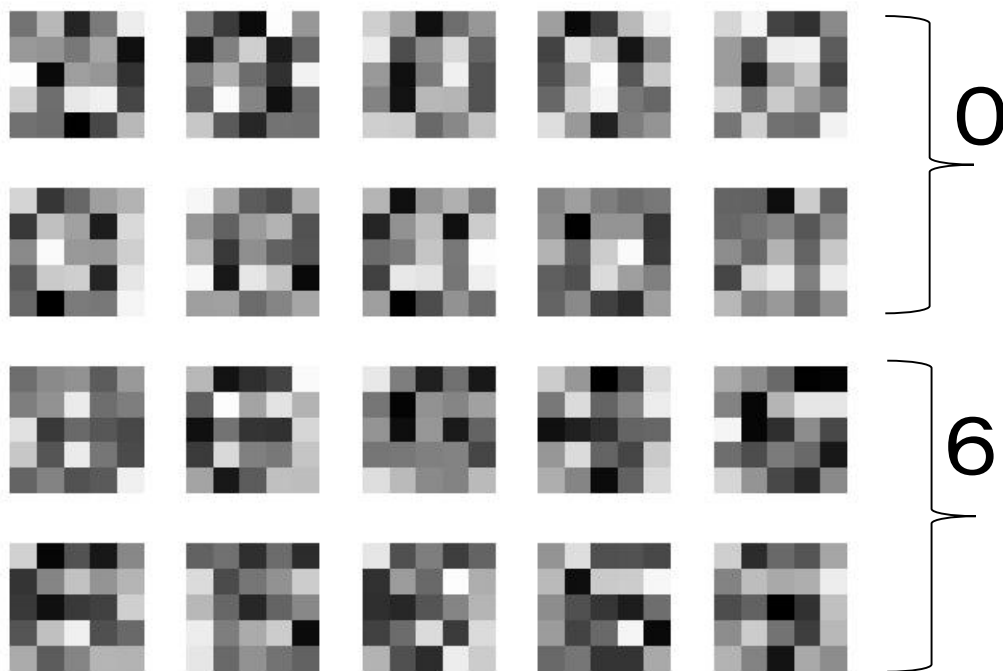
4 実験してみる

画像認識の例

学習用データ, $n=20, 40, \dots, 2000$

テスト用データ, $m=2000$.

Lasso 学習



学習誤差と汎化誤差

各データについて500回Lasso 学習を繰り返してして止める。
パラメータ w^* が得られる。

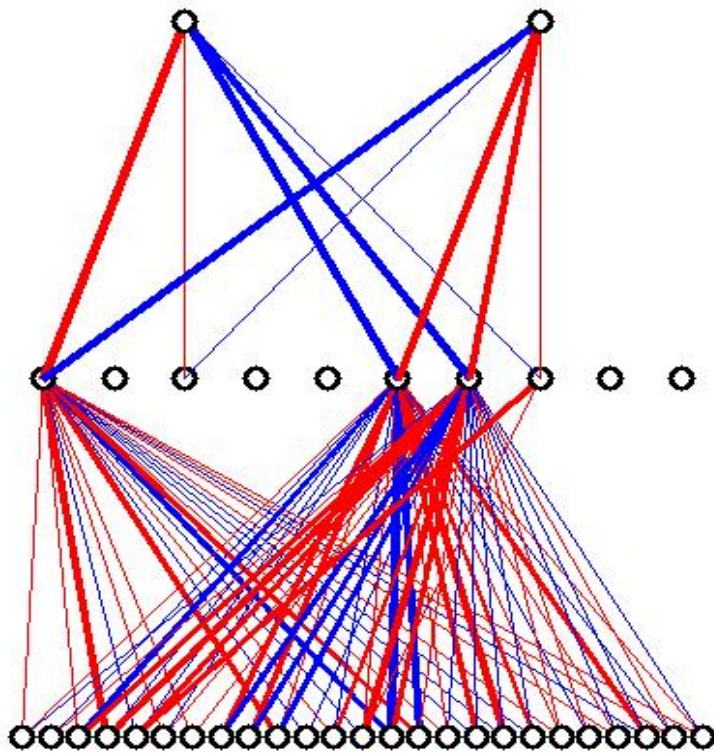
学習誤差関数 $E(w) = (1/n) \sum_{i=1}^n ||Y_i - f(X_i, w)||^2$

汎化誤差関数 $G(w) = (1/m) \sum_{i=1}^m ||Y_i - f(X_i, w)||^2$

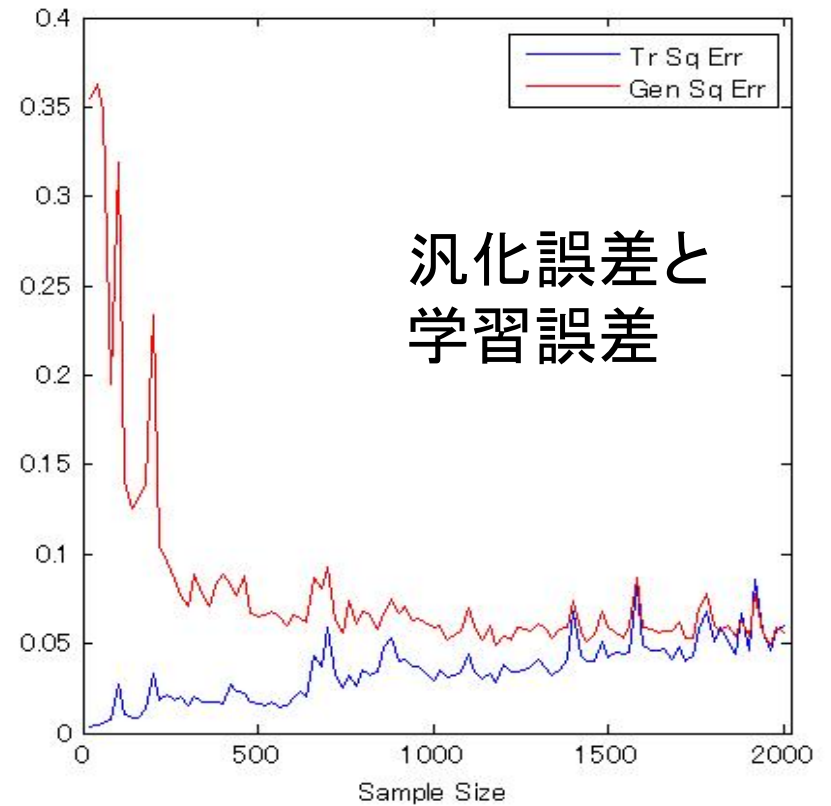
学習誤差 $E(w^*)$, 汎化誤差 $G(w^*)$ をプロットする。

- ☆ 繰り返し数500は十分多いとはいえない。
学習を途中で止めたことになる。

実験の例



Lasso 学習の結果($n=2000$)



学習データの数(n)

学習誤差と汎化誤差

- 平均的には 学習誤差 < 汎化誤差 であるが、個々の例ではそうなるとは限らない。データの数 n が大きくなると学習誤差と汎化誤差は近づいていく。
- 学習誤差も汎化誤差も Lasso ハイパーパラメータの影響が大きい。今回は人間力でハイパーパラメータを決めている。
- まったく同じデータを使ってもパラメータの初期値に依存して学習結果は異なる。