

4/15 相当. 轟

個人的な期待: 表現学習

SVMより数%精度向上といっても応用としてはBS。
背後にあるデータの構造を自動的に抽出できないか?

職人芝 vs. ML

1章

NN: 多層化と強固な表現力を持つが、学習が難しい。90年代に下火
back propagation (S4) except CNN (S6)

① DBN (Deep belief Network) S8

Hinton et al.

NNに似た多層構造



RBM (restricted Boltzmann machine) に分解

+ greedy 法
で学習していき
DBNが。

} pre training

② autoencoder S5

教師なし。 $x_n \rightarrow y_n = x_n$ とする非有監督学習。
多層 NN の学習がてきと。

S 1.1.3 特徴量学習。

このあたりが面白い。

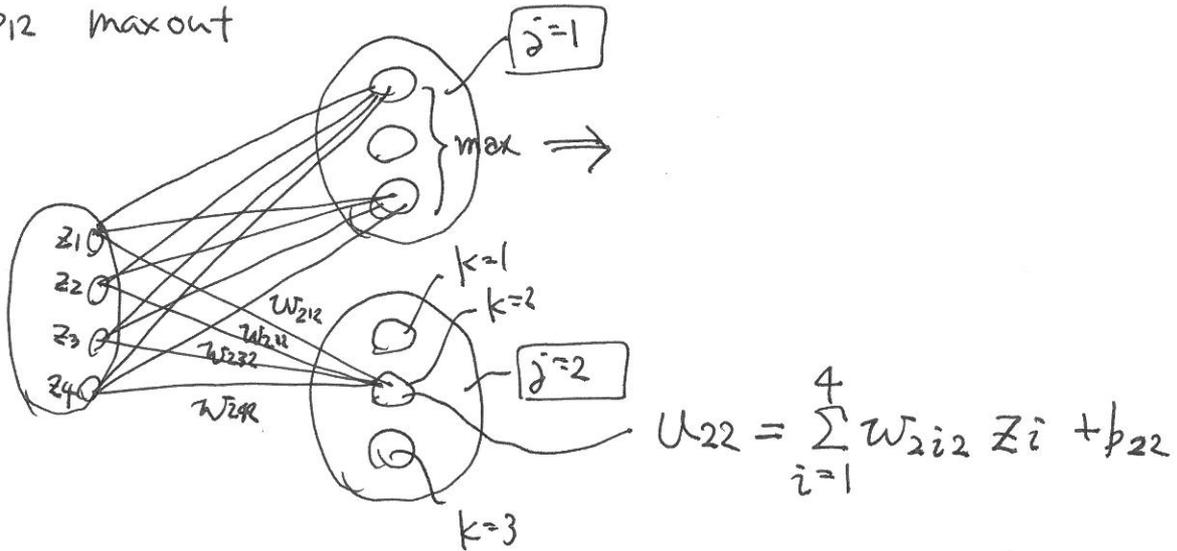
S 1.1.4

問題において多層の何層が最適か。
我々の場合はどれが良いか?

NLP, 音声認識 \rightarrow RNN S7

§ 2

p12 maxout



(2.5a) 式:

$$u^{(l+1)} = W^{(l+1)} z^{(l)} + b^{(l+1)}$$

$$\begin{bmatrix} u_1^{(l+1)} \\ \vdots \\ u_{J^{(l+1)}}^{(l+1)} \end{bmatrix} = \begin{bmatrix} W_{11}^{(l+1)} & \dots & W_{1I^{(l)}}^{(l+1)} \\ \vdots & \ddots & \vdots \\ W_{J^{(l+1)}1}^{(l+1)} & \dots & W_{J^{(l+1)}I^{(l)}}^{(l+1)} \end{bmatrix} \begin{bmatrix} z_1^{(l)} \\ \vdots \\ z_{I^{(l)}}^{(l)} \end{bmatrix} + \begin{bmatrix} b_1^{(l+1)} \\ \vdots \\ b_{J^{(l+1)}}^{(l+1)} \end{bmatrix}$$

$J^{(l+1)} \times 1$ $J^{(l+1)} \times I^{(l)}$ $I^{(l)} \times 1$

$$z^{(l+1)} = f(u^{(l+1)})$$

$$\begin{bmatrix} z_1^{(l+1)} \\ \vdots \\ z_{J^{(l+1)}}^{(l+1)} \end{bmatrix} = \begin{bmatrix} f(u_1^{(l+1)}) \\ \vdots \\ f(u_{J^{(l+1)}}^{(l+1)}) \end{bmatrix}$$

$$y \equiv z^{(L)} \quad \text{最終出力}$$

§2.4

学習 { 回帰 : 二乗誤差の最小化
 分類 : MLE : 最尤推定. = 尤度最大化
 同時確率.

§2.4.3 二値分類

入力: x

出力: $d \in \{0, 1\}$



事後確率: $P(d=1|x) \in \text{NN} \in \text{finite}$.

2層αβ11

$$\begin{aligned}
 & \begin{matrix} (x_1) \\ (x_2) \\ \vdots \\ (x_m) \end{matrix} \rightarrow \text{circle } u \rightarrow y = \frac{1}{1 + e^{-u}} = \frac{1}{1 + e^{-(w_1 x_1 + \dots + w_m x_m + b)}} = P(d=1|x) \\
 & u = w_1 x_1 + w_2 x_2 + \dots + w_m x_m + b \\
 & y(x; w)
 \end{aligned}$$

訓練データ: $\{(x_n, d_n) \mid n=1, \dots, N\}$ 入力数 = M 個.

$$P(d|x) = P(d=1|x)^d \cdot P(d=0|x)^{1-d}$$

$$P(d=0|x) = 1 - y(x; w)$$

$$\begin{aligned}
 & \text{事後確率} \\
 & P(d_n | x_n) = \frac{P(x_n | d_n) \cdot P(d_n)}{P(x_n) - \text{証拠}}
 \end{aligned}$$

$P(d_n), P(x_n)$ は w によらないので.

尤度最大 = 事後確率最大

正しくは

$$\begin{aligned}
 L(w) & \sim \prod_{n=1}^N p(d_n | x_n; w) \\
 & = \prod_{n=1}^N \{ y(x_n; w) \}^{d_n} \cdot \{ 1 - y(x_n; w) \}^{1-d_n}
 \end{aligned}$$

$$\log L(w) \sim \sum_{n=1}^N \{ d_n \cdot \log y(x_n; w) + (1-d_n) \log [1 - y(x_n; w)] \}$$

尤度最大 = 負の尤度最小 → (2.8) 式.

$$P(d=1|x) = \frac{P(d=1, x)}{P(x)} = \frac{P(x, d=1)}{P(x, d=0) + P(x, d=1)}$$

$$u \equiv \log \frac{P(x, d=1)}{P(x, d=0)} \quad \text{対数比をとる}$$

$$\frac{1}{1+e^{-u}} = \frac{e^u}{e^u + 1} = \frac{\frac{P(x, d=1)}{P(x, d=0)}}{\frac{P(x, d=1)}{P(x, d=0)} + 1} = \frac{P(x, d=1)}{P(x, d=1) + P(x, d=0)} = P(d=1|x)$$

(e^{log A} = A)

§ 2.4.4 多クラス分類

$$y_k \equiv z_k^{(L)} = \frac{\exp(u_k^{(L)})}{\sum_{j=1}^K \exp(u_j^{(L)})} = P(c=k|x) \quad K=2 \text{ のとき } y_1 + y_2 = 1$$

尤度

$$\hat{c} = \operatorname{argmax}_k y_k$$

推定クラス

入力: $x_n = \begin{bmatrix} x_1^{(n)} \\ \vdots \\ x_m^{(n)} \end{bmatrix} (n=1, \dots, N)$

目標: $d_n = \begin{bmatrix} d_{n1} \\ d_{n2} \\ \vdots \\ d_{nk} \end{bmatrix} \quad (3.1)$

$$\begin{bmatrix} 0 \\ 0 \\ 1 \\ 0 \\ \vdots \\ 0 \end{bmatrix} \leftarrow \text{正解} = c_3$$

事後分布

$$P(d|x) = \prod_{k=1}^K P(c_k|x)^{d_{nk}}$$

$$L(w) \equiv \text{尤度} \sim \prod_{n=1}^N P(d_n | x_n; w)$$

w の関数として

$$= \prod_{n=1}^N \prod_{k=1}^K P(c_k|x)^{d_{nk}} = \prod_{n=1}^N \prod_{k=1}^K \{y_k(x; w)\}^{d_{nk}}$$

対数尤度

→ (2.11) 式

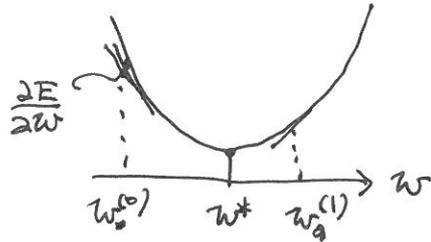
§3

学習: $W = \underset{W}{\operatorname{argmin}} E(W)$

↓
gradient descent method

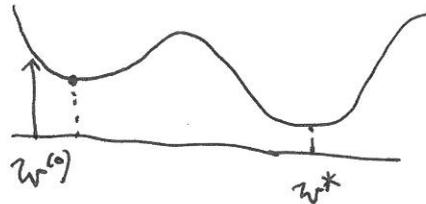
$$\nabla E \equiv \frac{\partial E}{\partial W} = \begin{bmatrix} \frac{\partial E}{\partial W_1} \\ \vdots \\ \frac{\partial E}{\partial W_M} \end{bmatrix}$$

$$W^{(t+1)} = W^{(t)} - \epsilon \nabla E$$



$$W_0^{(1)} = W_0^{(0)} - \epsilon \frac{\partial E}{\partial W} \Big|_{W=W_0^{(0)}}$$

local minimum: 局所的な可能性



§3.2

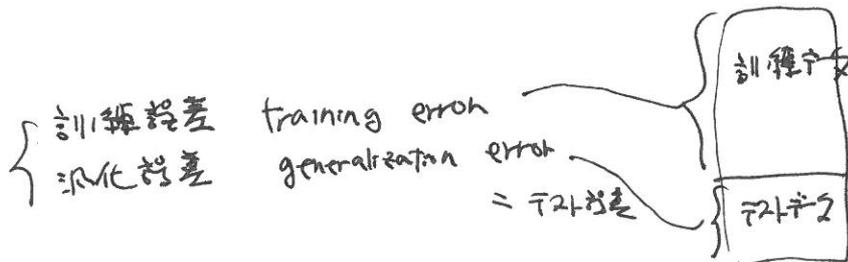
SGD: 非常に使われる

- ・ 計算が簡単
- ・ local minimum の 1 つに収束する問題がある可能性
- ・ 本-ライ=学習

§3.3 $\epsilon = 10^{-4} \dots$ 逐次処理

$$\left. \begin{array}{l} \bigcirc E_1(W) \\ \bigcirc E_2(W) \\ \vdots \\ \bigcirc E_{N_t}(W) \end{array} \right\} E_t(W) = \frac{1}{N_t} \sum_{i=1}^{N_t} E_i(W)$$

§3.4



NN 学習の正則化: 早期終了

overfitting, overlearning を止める

層数, 2=1 数, 1, 1K-1 の 2 の調整
正則化. 学習率 etc. を作る

$$e^{-u} = \frac{P(x, d=1)}{P(x, d=0)}$$

$$\frac{1}{1 + e^{-u}} = \frac{1}{1 + \frac{P(x, d=0)}{P(x, d=1)}} = \frac{P(x, d=1)}{P(x, d=1) + P(x, d=0)} = \frac{P(x, d=1)}{P(x)}$$

$$e^{\log x} = x$$

$$= \frac{1}{e^{-\log x}} = \frac{1}{\log x}$$

$$u = \log \frac{P(x, d=1)}{P(x, d=0)}$$

$$\frac{1}{1 + e^{-u}} = \frac{e^u}{e^u + 1} = \frac{\frac{P(x, d=1)}{P(x, d=0)}}{\frac{P(x, d=1)}{P(x, d=0)} + 1} = \frac{P(x, d=1)}{P(x, d=1) + P(x, d=0)} = \frac{P(x, d=1)}{P(x)}$$

$$= \frac{P(x|d=1)P(d=1)}{P(x|d=1)P(d=1) + P(x|d=0)P(d=0)}$$

$$= \frac{1}{p}$$

§ 3.5.1 正規化

PRML本 の § 1-1 参照. 重み減衰.

§ 3.5.3 ドロップアウト (67) 2014年 Hinton グループ

2.2.1と重み更新の度に確率的に選出

補11. 学習時にネットワークの自由度を統計的に落とす → 過適合を避けた
複数ネットワークの平均に相当する結果

§ 3.6 トリプ

§ 3.6.1 正規化

§ 3.6.2 処理

§ 3.6.3 複数ネットワーク

: textbook

$$\text{Var}[XY] = (E[X])^2 \text{Var}[Y] + [E[Y]]^2 \text{Var}[X] + \text{Var}[X] \text{Var}[Y]$$