

5.7 ベイズニューラルネットワーク

5.7.1 100X-タの事後分布

NNによる1次元の回帰をベイズ化

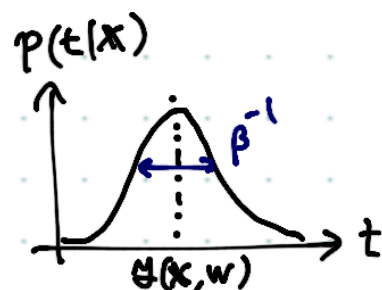
1次元・連続な目標変数: $t \xleftarrow{\text{予測}} \text{スカラベクトル } x$

t の条件付き確率 $p(t|x)$ は

$$p(t|x, w, \beta) = \mathcal{N}(t | y(x, w), \beta^{-1}) \quad (5.61)$$

$$\exp\left(-\frac{\beta}{2} \sum_n \{y_n - t_n\}^2\right)$$

$y(x, w)$: NNの出力

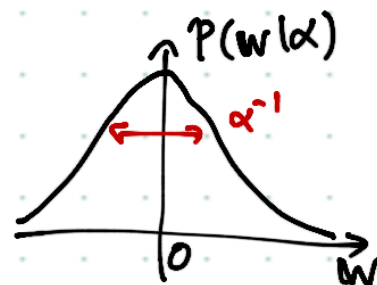


重みに関する事前分布を

$$p(w|\alpha) = \mathcal{N}(w | 0, \alpha^{-1} I) \quad (5.62)$$

とする.

$$\exp\left(-\frac{\alpha}{2} w^T w\right)$$



N 回の観測値 x_1, x_2, \dots, x_N | 目標集合 $\mathcal{D} = \{t_1, \dots, t_N\}$

尤度関数

$$p(\mathcal{D}|w, \beta) = \prod_{n=1}^N \mathcal{N}(t_n | y(x_n, w), \beta^{-1}) \quad (5.63)$$

事後分布は

$$p(w|\mathcal{D}, \alpha, \beta) \propto p(w|\alpha) p(\mathcal{D}|w, \beta) \quad (5.64)$$

$y(x, w)$ は w に非線形に依存 \Rightarrow (5.64) はガウス分布ではない.

なのでマクス近似でガウス分布に近似しよう.

事後分布 $P(w|D, \alpha, \beta)$ (5.64) を最大化しよう。

いつもの通り、事後分布の対数をとる

$$\ln P(w|D) = -\frac{\alpha}{2} w^T w - \frac{\beta}{2} \sum_{n=1}^N \{y(x_n, w) - t_n\}^2 + \text{定数} \quad (5.65)$$

(局所)

これを最大化する。これは、今まで行なった最適化アルゴリズムにより、事後分布の最大値 w_{MAP} を見つけることができる。

を見つけたら 事後分布の負の対数尤度の2階微分の行列を評価して、局所的にガウス分布で近似 する。(5.65) より、

$$A = -\nabla \nabla \ln P(w|D, \alpha, \beta) = \alpha I + \beta H \quad (5.66)$$

ここで H は二乗和誤差関数 $E = \frac{1}{2} \sum \{y_n - t_n\}^2$ の w の2階微分からなる行列である。←これを今までやった近似で求める。

この A を利用して事後分布を近似したガウス分布は (4.34) より

$$q(w|D) = \mathcal{N}(w|w_{MAP}, A^{-1}) \quad (5.67)$$

この事後分布について周辺化することによって予測分布が得られる

$$p(t|x, D) = \int p(t|x, w) q(w|D) dw \quad \dots (5.68)$$

となる。

(ただし $y(x, w)$ が w に対して非線形な場合、まだこの積分はあつかいづらい。

なので事後分布の分散が $y(x, w)$ が変化する w の特性スケールに比べて小さいと仮定する。

$$A = -\nabla \nabla \ln f(z) \big|_{z=z_0} \quad (P.215)$$

$$q(z) = \frac{|A|^{\frac{1}{2}}}{(2\pi)^{\frac{n}{2}}} \exp \left\{ -\frac{1}{2} (z - z_0)^T A (z - z_0) \right\} \\ = \mathcal{N}(z|z_0, A^{-1}) \quad \dots (4.34)$$

ネットワーク関数 $y(x, w)$ の w_{MAP} まわりの Taylor-展開を 1 次までやる

$$y(x, w) \simeq y(x, w_{MAP}) + g^T (w - w_{MAP}) \quad (5.196)$$

$$g = \nabla_w y(x, w) \big|_{w=w_{MAP}}$$

よって条件付き確率 $p(t|x, w, \beta)$ (5.161) は

$$P(t|x, w, \beta) \simeq \mathcal{N}(t | y(x, w_{MAP}) + g^T (w - w_{MAP}), \beta^{-1}) \quad (5.171)$$

で表せる。

したがって 2-変数に成り立つ (2.115) を周辺分布 $p(t)$ に利用して

$$p(w|D) = \mathcal{N}(w | \underbrace{w_{MAP}}_{\mu}, \underbrace{A^{-1}}_{\Sigma^{-1}}) \quad (5.167) \quad \leftarrow \begin{cases} p(x) = \mathcal{N}(x | \mu, \Sigma^{-1}) \\ p(y|x) = \mathcal{N}(Ax+b, L^{-1}) \end{cases}$$

$$P(t|x, w, \beta) \simeq \mathcal{N}(t | \underbrace{g^T w}_{A \mu} - \underbrace{g^T w_{MAP}}_{+b} + \underbrace{y(x, w_{MAP})}_{+b}, \underbrace{\beta^{-1}}_{L^{-1}}) \quad \leftarrow \begin{cases} \text{周辺分布} \\ p(y) = \mathcal{N}(y | A\mu+b, L^{-1} + A\Sigma^{-1}A^T) \end{cases}$$

右の (2.115) から周辺分布の平均 $A\mu+b$ は

$$A\mu+b = g^T w_{MAP} - g^T w_{MAP} + y(x, w_{MAP})$$

分散 $L^{-1} + A\Sigma^{-1}A^T$ は

$$L^{-1} + A\Sigma^{-1}A^T = \beta^{-1} + g^T A^{-1} g =: \sigma^2(x) \quad (5.173)$$

よって周辺分布 (予測分布) $P(t|x, D)$ は

$$P(t|x, D, \alpha, \beta) = \mathcal{N}(t | y(x, w_{MAP}), \sigma^2(x)) \quad (5.172)$$

$$\text{分散 } \sigma^2(x) = \underbrace{\beta^{-1}}_{\text{目標変数の}} + \underbrace{g^T A^{-1} g}_{\text{内在的なノイズ}} \quad (5.173)$$

目標変数の 内在的なノイズ x に依存, $E[y(x)]$ が w による不確実性に起因

5.7.2 超パラメータ最適化

ハイパーパラメータ α, β は固定に已知であるとしていた。

$P(t|x)$ の精度 \nearrow
 \nwarrow 事前分布 $P(w|\alpha)$ の精度

今から(5.7.1)近似で3.5節のエビデンス理論を使う!!

ハイパーパラメータの周辺化度あるいはエビデンスは

$$P(D|\alpha, \beta) = \int \underbrace{P(D|w, \beta) P(w, \alpha)}_{\downarrow} dw \quad (5.194)$$

で得らる。

$$\underbrace{\sum \cdot \frac{P(w|D, \alpha, \beta)}{\text{正規化係数}}}_{\downarrow} \rightarrow (5.165) \text{を使う。}$$

P.216のラプラス近似の結果(4.135)を使う。

$$\frac{1}{Z} P(D|w, \beta) P(w, \alpha) = P(w|D, \alpha, \beta)$$

$$\Rightarrow Z = \int \underbrace{P(D|w, \beta) P(w, \alpha)}_{f(z) \text{ に対応}} dw$$

$$P(D|\alpha, \beta) = Z \underbrace{\int P(w|D, \alpha, \beta) dw}_1$$

$$= Z$$

$$= P(D|w_{\text{MAP}}, \beta) P(w_{\text{MAP}}, \alpha) \frac{(2\pi)^{\frac{W}{2}}}{|A|^{\frac{1}{2}}}$$

$$Z \simeq f(z_0) \frac{(2\pi)^{\frac{W}{2}}}{|A|^{\frac{1}{2}}} \quad (4.135)$$

$P(w|D, \alpha, \beta)$
 はラプラス近似で

$$q(w|D) = \mathcal{N}(w|w_{\text{MAP}}, A^{-1}) \dots (5.67)$$

なので平均は w_{MAP} を使う。

$$P(D|W_{\text{MAP}}, \beta) = \prod_{n=1}^N \mathcal{N}(t_n | y(t_n, W_{\text{MAP}}), \beta^{-1}) \quad (5.163) \text{ 式}$$

$$= \prod_{n=1}^N \frac{1}{(2\pi)^{\frac{1}{2}}} \frac{1}{(\beta^{-1})^{\frac{1}{2}}} \exp(\sim)$$

$$P(W_{\text{MAP}}, \alpha) = \mathcal{N}(W|0, \alpha^{-1}I) \quad (5.162) \text{ 式}$$

$$= \frac{1}{(2\pi)^{\frac{N}{2}}} \frac{1}{(\alpha^{-1}I)^{\frac{1}{2}}} \exp(\sim)$$

\downarrow
 $\alpha^{-\frac{N}{2}}$

よ、2

$$\ln P(D|\alpha, \beta) = -E(W_{\text{MAP}}) - \frac{1}{2} \ln |A| + \frac{W}{2} \ln \alpha + \frac{N}{2} \ln \beta - \frac{N}{2} \ln (2\pi) \quad (5.175)$$

$$E(W_{\text{MAP}}) = \frac{\beta}{2} \sum_{n=1}^N \{y(x_n, W_{\text{MAP}}) - t_n\}^2 + \frac{\alpha}{2} W_{\text{MAP}}^T W_{\text{MAP}} \quad (5.176)$$

エビデンス理論では、 $\ln P(D|\alpha, \beta)$ を最大化して α と β の点推定を行う。

α については 3.5.2 節の 線形回帰 モデルの場合と同様にしる。

$$A = -\nabla \ln P(W|D, \alpha, \beta) = \alpha I + \beta H \quad \dots (5.166)$$

の固有値 ($W = W_{\text{MAP}}$) を考えるために固有値方程式

$$\beta H|_{W=W_{\text{MAP}}} u_i = \lambda_i u_i \quad \dots (5.177)$$

を考える。

これを $\ln P(D|\alpha, \beta)$ の $\ln |A|$ の評価で使う。 $|A| = \prod_i (\lambda_i + \alpha)$

すると $\ln P(D|\alpha, \beta)$ を最大化する α は (3.92) と同様に

$$\alpha = \frac{\gamma}{W_{MAP}^T W_{MAP}} \quad (5.178)$$

有効な X - Y 数
(3.5.3 節) $\rightarrow \gamma = \sum_{i=1}^N \frac{\lambda_i}{\alpha + \lambda_i} \quad (5.179)$

ただし今回の場合はこれは厳密でない。(線形回帰では厳密.)

[H が α によって変わる $\rightarrow \lambda_i$ が α に依存するので]

計算の途中で

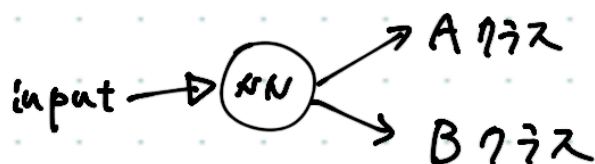
λ_i の α 依存分をムシしている。

同様に (5.177) から $\ln |A|$ を評価し、(3.95) より β の再推定の式は

$$\frac{1}{\beta} = \frac{1}{N-\gamma} \sum_{n=1}^N \left[y(x_n, W_{MAP}) - t_n \right]^2 \quad \dots \quad (5.180)$$

5.7.3 クラス分類のためのベイズ NN

1つのロジスティックモデル出力を持つ2クラス分類を考える。



このモデルの対数尤度関数は

$$\ln p(D|W) = \sum_{n=1}^N \{t_n \ln y_n + (1-t_n) \ln (1-y_n)\} \quad (5.181)$$

$$t_n = \begin{cases} 0 \\ 1 \end{cases}, \quad y_n = y(x_n, W) \leftarrow \text{NNが出力する分類の確率}$$

$$\left(P(D|W) = \prod_{n=1}^N y_n^{t_n} \{1-y_n\}^{1-t_n} \dots (4.89) \right)$$

データは正しくラベル付けされているとして11110-1101x-1βは見えない。

事前分布を $p(W|\alpha) = \mathcal{N}(W|0, \alpha^{-1}I) \dots (5.162)$ とする。

1. まず α を初期化する。

2. 2.2 対数事後分布 $p(W|D, \alpha) \propto p(W|\alpha) p(D|W)$

を最大化して1101x-1ベクトル W を決める。

これは正則化誤差関数

$$E(W) = -\ln p(D|W) + \frac{\alpha}{2} W^T W \quad \dots (5.182)$$

を最小化することと等価であり Backpropagation などできる。

3. W_{MAP} をみつけたら 負の対数尤度関数の2階微分からなるヘッセ行列 H を評価する。

4. 事後分布のガウス近似は $q(W|D) = \mathcal{N}(W|W_{MAP}, A^{-1})$ (5.167) 与えられる。

5. $1110-1103X$ -タ α を最適化するには、やはり周辺尤度を最大化する。
(5.175) を導出したときと同様にして

$$\ln p(D|\alpha) \simeq -E(W_{MAP}) - \frac{1}{2} \ln |A| + \frac{W}{2} \ln \alpha \dots (5.183)$$

$$E(W_{MAP}) = -\sum_{n=1}^N \{t_n \ln y_n + (1-t_n) \ln (1-y_n)\} + \frac{\alpha}{2} W_{MAP}^T W_{MAP} \dots (5.184)$$

$$y_n \equiv y_n(x_n, W_{MAP})$$

$$(5.183) \text{ を } \alpha \text{ について最大化すると再び } \alpha = \frac{\gamma}{W_{MAP}^T W_{MAP}} \quad (5.198) \text{ 与えられる。}$$

与えられた再推定方程式が導かれる。

6. 最後に (5.168) で定義される予測分布が必要である。
この積分もネットワークの非線形性からムズかしい。

事後分布の幅がせまいとして

$$p(t|x, D) \simeq p(t|x, W_{MAP})$$

とするか、分散をさらに考慮するかで近似する。

(しかし出力がロジスティックモイド活性化関数であることから出力は $(0, 1)$

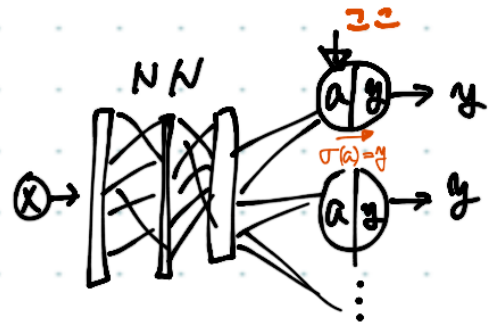
に制約されるため、回帰のときのように (5.169) みたいな線形近似はよくない。

なので $\sigma(z)$ の1歩手前の出力ユニット活性を線形近似する。

$$Q(X, W) \simeq a_{\text{MAP}}(X) + b^T (W - W_{\text{MAP}})$$

$$a_{\text{MAP}}(X) = a(X, W_{\text{MAP}}), \quad b \equiv \nabla a(X, W_{\text{MAP}})$$

は逆伝播により求めらる。



7. 今までの結果を 4.5.2 節の結果に当てはめることができる。

出力 $y = +1$ の活性化の値の分布は

$$p(a|X, D) = \int \delta(a - a_{\text{MAP}}(X) - b^T (W - W_{\text{MAP}})) \mathcal{P}(W|D) dW \quad \dots (5.189)$$

$$\mathcal{P}(W|D) \propto \exp \left[-\frac{1}{2} (W - W_{\text{MAP}})^T A (W - W_{\text{MAP}}) \right]$$

平均 a_{MAP} , 分散 $\sigma_a^2(X) = b^T A^{-1} b$
のガウス分布になる。

予測分布を得るために a に関して周辺化する。

$$p(t=1|X, D) = \int \sigma(a) p(a|X, D) da \quad (5.189)$$

けど σ をもつガウス分布のたたみ込みはムズかしい。

(4.153) の近似を (5.189) に使, 2

$$p(t=1|X, D) = \sigma(K(\sigma_a^2) a_{\text{MAP}}) \quad (5.190)$$

$$K(\sigma^2) = (1 + \pi \sigma^2 / 8)^{-\frac{1}{2}} \quad (4.153)$$