

卒研ゼミ「深層学習」

5 機械学習の基礎

5.6 ベイズ統計

データ $D_m := \{x^{(1)}, \dots, x^{(m)}\}$ が与えられたとき、データによって決まるパラメータが θ という値をとりうる確率はベイズの定理より

$$p(\theta | D_m) = \frac{p(D_m | \theta)p(\theta)}{p(D_m)} \propto p(D_m | \theta)p(\theta) \quad (1)$$

となる。中辺の分母は

$$p(D_m) = \int p(D_m | \theta)p(\theta) d\theta$$

で計算でき、これは θ によらず、また $D_m = \{x^{(1)}, \dots, x^{(m)}\}$ の値は既知のものであるから、定数として扱ってよい。式 (1) の $p(\theta)$ を事前確率（分布）、 $p(\theta | D_m)$ を事後確率（分布）、 $p(D_m | \theta)$ を尤度という。事前確率はデータを手に入れる前に想定していた確率であるのに対して、事後確率はデータを得たあとに事前確率を修正（ベイズ修正）したものである。式 (1) からわかるように、ベイズ修正とは事前確率に尤度をかけて規格化し、よりもっともらしい分布に修正することである。

例：ベイズ線形回帰

m 個の訓練データ $(X^{(\text{train})}, \mathbf{y}^{(\text{train})}) = (X, \mathbf{y})$ が与えられたとする（右肩の (train) は省略する）。ここに $\mathbf{y} = [y_1, \dots, y_m]^T$ であり、基底関数を $\phi_k(x)$, $k = 0, \dots, n$ とすれば計画行列は $(X)_{ij} = \phi_j(x_i)$ である（推定する関数を n 次多項式と仮定したときは $\phi_k(x) = x^k$ であった）。 $\mathbf{x} = [\phi_0(x), \dots, \phi_n(x)]^T$, $\mathbf{w} = [w_0, w_1, \dots, w_n]^T$ とすれば、訓練データによる y の推定は $\hat{y} = \sum_{k=0}^n w_k \phi_k(x) = \mathbf{w}^T \mathbf{x}$ と表現できる。また各訓練データ x_i ($i = 1, \dots, m$) に対する y の推定を計算した $\hat{y}_i = \sum_{k=0}^n w_k \phi_k(x_i)$ を並べたベクトルを $\hat{\mathbf{y}} = [\hat{y}_1, \dots, \hat{y}_m]^T$ と書けば、 $\hat{\mathbf{y}} = X\mathbf{w}$ である。

さて X, \mathbf{y} が得られたもとで係数のパラメータが \mathbf{w} となる条件付き確率 $p(\mathbf{w} | X, \mathbf{y})$ を求めたい。

$$p(\mathbf{w} | X, \mathbf{y}) = \frac{p(\mathbf{y} | X, \mathbf{w})p(X | \mathbf{w})p(\mathbf{w})}{p(X, \mathbf{y})} \propto p(\mathbf{y} | X, \mathbf{w})p(\mathbf{w}). \quad (2)$$

尤度を $p(\mathbf{y} | X, \mathbf{w})$ としたのは理由がある。ベイズの定理を考えれば、尤度としての候補は $p(X, \mathbf{y} | \mathbf{w})$, $p(X | \mathbf{y}, \mathbf{w})$, $p(\mathbf{y} | X, \mathbf{w})$ の 3 つが挙げられる。しかし推定したい関数 \hat{y} は X と \mathbf{w} が決まれば決まるというもので、尤度としては $p(\mathbf{y} | X, \mathbf{w})$ が適当である。また \mathbf{w} の値は X, \mathbf{y} の値によって決められる不確かさを持っているが、 X の値は観測された既知のデータで固定されているため、 \mathbf{w} の値がわかっている条件のもと

でデータの値が X である条件付き確率 $p(X | \mathbf{w})$ は、 \mathbf{w} の値によらず常に 1 である。また式 (2) の分母は \mathbf{w} について積分をしたものなので、 \mathbf{w} によらない定数と考えることができる。

尤度 $p(\mathbf{y} | X, \mathbf{w})$ が正規分布に従うとき、式 (2) の事後確率を具体的に計算する例をみせる。

$$p(\mathbf{y} | X, \mathbf{w}) = \mathcal{N}(\mathbf{y}; X\mathbf{w}, I) \propto \exp\left(-\frac{1}{2}(\mathbf{y} - X\mathbf{w})^\top(\mathbf{y} - X\mathbf{w})\right). \quad (3)$$

なお多変量正規分布は

$$\mathcal{N}(\mathbf{x}; \boldsymbol{\mu}, \Sigma) = \sqrt{\frac{1}{(2\pi)^n \det \Sigma}} \exp\left(-\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu})^\top \Sigma^{-1}(\mathbf{x} - \boldsymbol{\mu})\right) \quad (4)$$

であった (教科書の式 3.23)。この $\boldsymbol{\mu}$ は \mathbf{x} の各成分の平均を並べたベクトルで、 Σ は分散共分散行列である。ただし式 (3) では各 y_i の分散 s_i^2 は 1 で、 $i \neq j$ に対して y_i と y_j の共分散 s_{ij} はゼロ、すなわち相関はないと仮定している。式 (3) は、観測値 $\mathbf{y} = [y_1, \dots, y_m]^\top$ のずれ具合が正規分布に従うと仮定していることを意味している [1]。

事前確率はデータを得る前に想定していた確率のことであるから、過去の経験にもとづいた、想定していた分布を設定すればよい。つまり予想とか「直感的信頼度 (the degree of belief)」でよい。たとえば何も情報を得ておらず、どんな分布かわからない場合は、ただ単に定数とすることがある (無情報事前分布、non-informative prior という [2])。ここでは計算を容易にするため、事前分布は尤度と同じく正規分布と設定する。

$$p(\mathbf{w}) = \mathcal{N}(\mathbf{w}; \boldsymbol{\mu}_0, \Lambda_0) \propto \exp\left(-\frac{1}{2}(\mathbf{w} - \boldsymbol{\mu}_0)^\top \Lambda_0^{-1}(\mathbf{w} - \boldsymbol{\mu}_0)\right). \quad (5)$$

実用上、 $\Lambda_0 = \text{diag}(\boldsymbol{\lambda}_0)$ などとすれば計算はさらに簡単になる。

事後確率は尤度と事前確率の積に比例するので、式 (2)、式 (3)、式 (5) より

$$p(\mathbf{w} | X, \mathbf{y}) \propto \exp\left(-\frac{1}{2}(\mathbf{y} - X\mathbf{w})^\top(\mathbf{y} - X\mathbf{w})\right) \exp\left(-\frac{1}{2}(\mathbf{w} - \boldsymbol{\mu}_0)^\top \Lambda_0^{-1}(\mathbf{w} - \boldsymbol{\mu}_0)\right).$$

これは $e^x e^y = e^{x+y}$ のようにひとつの指数関数にまとめることができる。指数関数の肩の部分を展開して、 \mathbf{w} の二次形式 $-\frac{1}{2}(\mathbf{w} - \boldsymbol{\mu}_m)^\top \Lambda_m^{-1}(\mathbf{w} - \boldsymbol{\mu}_m)$ の形にまとめることを目指す。まず Λ_0^{-1} は対称行列であり、 $\mathbf{w}^\top X^\top \mathbf{y}$ と $\mathbf{w}^\top \Lambda_0^{-1} \boldsymbol{\mu}_0$ はスカラーであるから、 $\mathbf{w}^\top X^\top \mathbf{y} = (\mathbf{w}^\top X^\top \mathbf{y})^\top = \mathbf{y}^\top X \mathbf{w}$, $\mathbf{w}^\top \Lambda_0^{-1} \boldsymbol{\mu}_0 = (\mathbf{w}^\top \Lambda_0^{-1} \boldsymbol{\mu}_0)^\top = \boldsymbol{\mu}_0^\top \Lambda_0^{-1} \mathbf{w}$ のように、項の左側にある \mathbf{w}^\top を項の右側に移すことができる。

$$\begin{aligned} -\frac{1}{2}(\mathbf{y} - X\mathbf{w})^\top(\mathbf{y} - X\mathbf{w}) &= -\frac{1}{2}(\mathbf{y}^\top \mathbf{y} - \mathbf{y}^\top X \mathbf{w} - \mathbf{w}^\top X^\top \mathbf{y} + \mathbf{w}^\top X^\top X \mathbf{w}) \\ &= -\frac{1}{2}(\text{const.} - 2\mathbf{y}^\top X \mathbf{w} + \mathbf{w}^\top X^\top X \mathbf{w}), \\ -\frac{1}{2}(\mathbf{w} - \boldsymbol{\mu}_0)^\top \Lambda_0^{-1}(\mathbf{w} - \boldsymbol{\mu}_0) &= -\frac{1}{2}(\mathbf{w}^\top \Lambda_0^{-1} \mathbf{w} - \mathbf{w}^\top \Lambda_0^{-1} \boldsymbol{\mu}_0 - \boldsymbol{\mu}_0^\top \Lambda_0^{-1} \mathbf{w} + \boldsymbol{\mu}_0^\top \Lambda_0^{-1} \boldsymbol{\mu}_0) \\ &= -\frac{1}{2}(\mathbf{w}^\top \Lambda_0^{-1} \mathbf{w} - 2\boldsymbol{\mu}_0^\top \Lambda_0^{-1} \mathbf{w} + \text{const.}). \end{aligned}$$

ただし \mathbf{w} によらない項は定数 const. とした。ゆえに事後確率は

$$\begin{aligned} p(\mathbf{w} | X, \mathbf{y}) &\propto \exp\left(-\frac{1}{2}(-2\mathbf{y}^\top X \mathbf{w} + \mathbf{w}^\top X^\top X \mathbf{w} + \mathbf{w}^\top \Lambda_0^{-1} \mathbf{w} - 2\boldsymbol{\mu}_0^\top \Lambda_0^{-1} \mathbf{w})\right) \\ &= \exp\left(-\frac{1}{2}[-2(\mathbf{y}^\top X + \boldsymbol{\mu}_0^\top \Lambda_0^{-1})\mathbf{w} + \mathbf{w}^\top (X^\top X + \Lambda_0^{-1})\mathbf{w}]\right). \end{aligned} \quad (6)$$

ここで、目指している二次形式を展開すると

$$-\frac{1}{2}(\mathbf{w} - \boldsymbol{\mu}_m)^\top \Lambda_m^{-1}(\mathbf{w} - \boldsymbol{\mu}_m) = -\frac{1}{2}(-2\boldsymbol{\mu}_m^\top \Lambda_m^{-1}\mathbf{w} + \mathbf{w}^\top \Lambda_m^{-1}\mathbf{w}) + \text{const.}$$

であるから、これと式 (6) の引数部分を比較して、 $\Lambda_m^{-1} = X^\top X + \Lambda_0^{-1}$, $\boldsymbol{\mu}_m^\top \Lambda_m^{-1} = \mathbf{y}^\top X + \boldsymbol{\mu}_0^\top \Lambda_0^{-1}$ を得る。
すなわち

$$\Lambda_m = (X^\top X + \Lambda_0^{-1})^{-1}, \quad \boldsymbol{\mu}_m = \Lambda_m(X^\top \mathbf{y} + \Lambda_0^{-1}\boldsymbol{\mu}_0)$$

として、事前分布が正規分布に従うとき事後確率は

$$p(\mathbf{w} \mid X, \mathbf{y}) \propto \exp\left(-\frac{1}{2}(\mathbf{w} - \boldsymbol{\mu}_m)^\top \Lambda_m^{-1}(\mathbf{w} - \boldsymbol{\mu}_m)\right)$$

となり、これもまた正規分布に従っていることがわかる。規格化するには式 (4) を参照すればよい。この計算結果は、はじめ \mathbf{w} は平均 $\boldsymbol{\mu}_0$, 分散共分散行列 Λ_0 の正規分布に従うと想定していたものが、データ X, \mathbf{y} を得たことにより、平均 $\boldsymbol{\mu}_m$, 分散共分散行列 Λ_m の正規分布に従うように修正されたことを表している。ただし Λ_m が正則になるように Λ_0 を設定する必要がある。この例のように、事前分布と事後分布が同じ分布族になるとき、これを共役自然分布という [2]。

事前確率において $\boldsymbol{\mu}_0 = \mathbf{0}$, $\Lambda_0 = \frac{1}{\lambda}I$ と設定した場合、事後確率において $\Lambda_m^{-1} = X^\top X + \lambda I$, $\boldsymbol{\mu}_m = \Lambda_m X^\top \mathbf{y}$ となるので、 \mathbf{w} の真の値 $\boldsymbol{\mu}_m$ についての方程式

$$\Lambda_m^{-1} \boldsymbol{\mu}_m = (X^\top X + \lambda I) \boldsymbol{\mu}_m = X^\top \mathbf{y}.$$

が得られる。ここで X と \mathbf{y} は有限個の観測されたデータであるため、この方程式を解いて得られる $\boldsymbol{\mu}_m$ の値は実際には \mathbf{w} の推定値 $\hat{\boldsymbol{\mu}}_m$ である。この方程式は、正則化項を重み減衰 $\lambda \|\mathbf{w}\|_2^2$ として正則化最小二乗法により得られた方程式 $(X^\top X + \lambda I)\mathbf{w} = X^\top \mathbf{y}$ に一致している。ただし $\lambda = 0$ としてしまうと、これは \mathbf{w} の分散がはじめ無限大であったことを表しているが、 $\Lambda_0, \Lambda_0^{-1}$ を定義できないため、ベイズ推定ではこの場合を取り扱うことができない。さらに重要な違いとしては、正則化最小二乗法では \mathbf{w} の推定値のみ得られていたが、ベイズ推定ではそれに加えて \mathbf{w} の分散共分散行列 Λ_m も得ることができる。

5.6.1 MAP 推定

事後確率 $p(\boldsymbol{\theta} \mid D_m)$ が最大になるときの $\boldsymbol{\theta}$ の値をその推定値とする方法を、MAP 推定（最大事後確率推定、maximum a posteriori estimation）という。すなわちその推定値を $\hat{\boldsymbol{\theta}}_{\text{MAP}}$ とかけば、

$$\hat{\boldsymbol{\theta}}_{\text{MAP}} = \underset{\boldsymbol{\theta}}{\operatorname{argmax}} p(\boldsymbol{\theta} \mid D_m) = \underset{\boldsymbol{\theta}}{\operatorname{argmax}} [\log p(\boldsymbol{\theta} \mid D_m)] = \underset{\boldsymbol{\theta}}{\operatorname{argmax}} [\log p(D_m \mid \boldsymbol{\theta}) + \log p(\boldsymbol{\theta})]$$

である。ただしここでは関数 $f(\boldsymbol{\theta})$ の最大点と $\log f(\boldsymbol{\theta})$ の最大点が一致することと、 $p(\boldsymbol{\theta} \mid D_m) \propto p(D_m \mid \boldsymbol{\theta})p(\boldsymbol{\theta})$ であることを用いている。

例として、線形回帰モデルにおいて事前分布 $p(\mathbf{w})$ が正規分布 $\mathcal{N}(\mathbf{w}; \mathbf{0}, \frac{1}{\lambda}I)$ に従うと設定する。すると

$$p(\mathbf{w}) \propto \exp\left(-\frac{1}{2}\lambda \mathbf{w}^\top \mathbf{w}\right), \quad \log p(\mathbf{w}) = -\frac{\lambda}{2}\mathbf{w}^\top \mathbf{w} + \text{const.}$$

であるから、

$$\log p(\mathbf{w} \mid X, \mathbf{y}) = \log p(\mathbf{y} \mid X, \mathbf{w}) + \log p(\mathbf{w}) + \text{const.} = -\text{MSE}(\mathbf{w}) - \frac{\lambda}{2}\mathbf{w}^\top \mathbf{w} + \text{const.}$$

である。これは最尤推定法で最大化する関数（式 5.65）に重み減衰を付加したものになっていることがわかる。符号がマイナスとなっているから、これを最大化することは、正則化最小二乗法における $\text{MSE}(\mathbf{w}) + \frac{\lambda}{2} \mathbf{w}^T \mathbf{w}$ を最小化することと等価である。なお式に定数が足されているが、最大・最小となる点はそれによって影響されないで、やっていることは正則化最小二乗法と同じである。

As with full Bayesian inference, MAP Bayesian inference has the advantage of leveraging information that is brought by the prior and cannot be found in the training data. This additional information helps to reduce the variance in the MAP point estimate (in comparison to the ML estimate). However, it does so at the price of increased bias. なぜ？

データからのみでは得られないような情報を事前確率から得ることができるので、MAP 推定は有用である。それによって推定量のバリエーションを最尤推定に比べて小さくすることができるが、そのときバイアスの増大が伴う。

正則化を含む推定方法の多くは、裏で MAP 推定を行なっていると考えることができ、その正則化項は $\log p(\boldsymbol{\theta})$ に対応している。ただしすべての正則化項が $\log p(\boldsymbol{\theta})$ に対応しているわけではない。例えば正則化項がデータを含むとき、 $p(\boldsymbol{\theta})$ は $\boldsymbol{\theta}$ のみによる確率密度関数であるため、「正則化項は $\log p(\boldsymbol{\theta})$ に対応している」と考えることはできない。

逆に MAP 推定によって正則化項を設計することができる。

参考文献

- [1] 中谷秀洋. 「第 12 回 ベイズ線形回帰 [前編]」. gihyo.jp. <https://gihyo.jp/dev/serial/01/machine-learning/0012>. 最終閲覧 2020 年 1 月 17 日.
- [2] 安道知寛, 「ベイズ統計モデリング」, 株式会社朝倉書店, 2010 年. pp. 28–41.