

# ベイジアンネットワーク

白坂貴規

平成 32 年 3 月 17 日

## 1 ベイズ推論の基礎

ベイズ推論では学習や予測, モデル選択などを全て確率分布上の計算問題として扱う。

### 1.1 確率推論

$M$  次元ベクトル  $\mathbf{x} = (x_1, \dots, x_M)^T \in \mathbb{R}^M$  の実数関数  $p(\mathbf{x})$  が次の条件を満たす時,  $p(\mathbf{x})$  を確率密度関数 (probability density function) と呼ぶ。

$$p(\mathbf{x}) \geq 0 \int p(\mathbf{x}) d\mathbf{x} = 1 \quad (1)$$

また, 各要素が離散値である時,  $p(\mathbf{x})$  を確率質量関数 (probability mass function) と呼ぶ。以降, 確率密度関数や確率質量関数で決められる  $\mathbf{x}$  の分布を確率分布 (probabilistic distribution) あるいは確率モデル (probabilistic model) と呼ぶ。ある 2 つの変数  $x$  と  $y$  に関する確率分布  $p(x, y)$  を同時分布 (joint distribution) とよび,

$$p(y) = \int p(x, y) dx \quad (2)$$

のように一方の変数  $x$  を積分により除去する操作を周辺化 (marginalization) とよび, 結果として得られる確率分布  $p(y)$  を  $y$  の周辺分布と呼ぶ。また, 同時分布  $p(x, y)$  において,  $y$  に対して特定の値が決められた時の  $x$  の確率分布を条件付き分布 (conditional distribution) とよび, 次のように定義する。

$$p(x|y) = \frac{p(x, y)}{p(y)} \quad (3)$$

条件付き分布  $p(x|y)$  は  $x$  の確率分布であり,  $y$  はこの分布の特性を決めるパラメータのようなものであると解釈できる。

$$\int p(x|y) dx = \frac{\int p(x, y) dx}{p(y)} = \frac{p(y)}{p(y)} = 1 \quad (4)$$

であり,  $p(x, y)$  と  $p(y)$  が共に非負であることも考慮すれば, 条件付き分布  $p(x|y)$  は確率分布の要件を満たす。さらに, 同時分布を考える際に重要となるのが

独立 (independence) という概念である。同時分布が

$$p(x, y) = p(x)p(y) \quad (5)$$

を満たす時,  $x$  と  $y$  は独立であるという。ある同時分布が与えられた時, そこから興味の対象となる条件付き分布や周辺分布を算出することをここではベイズ推論 (Bayesian inference), あるいは単に推論 (inference) と呼ぶ。期待値 (expectation) は確率分布の特徴を定量的に表すことに使われる。 $\mathbf{x}$  をベクトルとした時に, 確率分布  $p(\mathbf{x})$  に対して, ある関数  $f(\mathbf{x})$  の期待値  $\mathbb{E}_{p(\mathbf{x})}[f(\mathbf{x})]$  は次のように計算される。

$$\mathbb{E}_{p(\mathbf{x})}[f(\mathbf{x})] = \int f(\mathbf{x})p(\mathbf{x}) d\mathbf{x} \quad (6)$$

2 つの確率分布  $p(\mathbf{x})$  及び  $q(\mathbf{x})$  に対して次のような期待値を KL ダイバージェンス (Kullback-Leibler divergence) と呼ぶ。

$$\begin{aligned} D_{KL}[q(\mathbf{x})][p(\mathbf{x})] &= - \int q(\mathbf{x}) \ln \frac{p(\mathbf{x})}{q(\mathbf{x})} d\mathbf{x} \\ &= \mathbb{E}_{q(\mathbf{x})}[\ln q(\mathbf{x})] \\ &\quad - \mathbb{E}_{q(\mathbf{x})}[\ln p(\mathbf{x})] \end{aligned} \quad (7)$$

KL ダイバージェンスは任意の確率分布のくみに対して非負であり, 0 となるのは 2 つの分布が完全に一致する場合  $p(\mathbf{x}) = q(\mathbf{x})$  に限られる。また, KL ダイバージェンスは 2 つの確率分布の”距離”を表していると解釈されるが,  $p(\mathbf{x})$  と  $q(\mathbf{x})$  が対称ではない為, 数学的な距離の公理は満たしていない。

### 1.2 変数変換

既知の確率密度関数に対して変数関数を行うことで新たな確率密度関数を導出することを考える。全単射の関数  $f: \mathbb{R}^M \rightarrow \mathbb{R}^M$  によって変数を  $\mathbf{y} = f(\mathbf{x})$  のように一対一に変換する操作は, 既知の確率密度関数を  $p_x(\mathbf{x})$  とすれば, 変換によって得られる  $\mathbf{y}$  の確率密度関数は

$$p_y(\mathbf{y}) = p_x(g(\mathbf{y})) |\det(J_g)| \quad (8)$$

とかける。ただし、 $J$  は  $f$  の逆関数  $g$  のヤコビ行列 (Jacobian matrix)

$$J_g = \begin{pmatrix} \frac{\partial x_1}{\partial y_1} & \cdots & \frac{\partial x_1}{\partial y_M} \\ \vdots & \ddots & \vdots \\ \frac{\partial x_M}{\partial y_1} & \cdots & \frac{\partial x_M}{\partial y_M} \end{pmatrix} \quad (9)$$

であり、 $\det(J_g)$  は  $J_g$  の行列式である。

### 1.3 グラフィカルモデル

グラフィカルモデル (graphical model) は、確率モデルに存在する複数の変換の関係性をノードと矢印を使って表現する記法である。例として、次のようにあるパラメータ  $\theta$  に依存して  $N$  個の変数  $\mathbf{X} = x_1, \dots, x_N$  が発生するような確率モデルを考える。グラフィカルモデルは次図のようになる。

$$p(\mathbf{X}, \theta) = p(\theta)p(\mathbf{X}|\theta) = p(\theta) \prod_{n=1}^N p(x_n|\theta) \quad (10)$$

なお、上式における  $p(\mathbf{X}|\theta)$  を尤度関数 (likelihood function)、 $p(\theta)$  をパラメータ  $\theta$  の事前分布と呼ぶ。

### 1.4 指数分布族

ガウス分布やディクレ分布など、ベイズ推論で用いられる多くの実用的な確率分布は、指数分布族 (exponential family) と呼ばれるある形式をもつクラスに属する。その前に、代表的な確率分布をいくつか紹介する。一次元のガウス分布 (Gaussian distribution) または正規分布 (normal distribution) は次のような  $x \in \mathbb{R}$  の確率密度関数をもつ分布である。

$$\mathcal{N}(x|\mu, \sigma^2) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{(x-\mu)^2}{2\sigma^2}\right) \quad (11)$$

$\mu \in \mathbb{R}$  は平均パラメータで、 $\sigma^2 > 0$  は分散パラメータである。ガウス分布は次のように  $M$  次元の多変量に拡張される。

$$\mathcal{N}(\mathbf{x}|\boldsymbol{\mu}, \boldsymbol{\Sigma}) = \frac{1}{\sqrt{(2\pi)^D |\boldsymbol{\Sigma}|}} \exp\left(-\frac{1}{2}(\mathbf{x}-\boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1}(\mathbf{x}-\boldsymbol{\mu})\right) \quad (12)$$

ここで、 $\boldsymbol{\mu} \in \mathbb{R}^M$  は  $M$  次元の平均パラメータで、 $\boldsymbol{\Sigma}$  はサイズが  $M \times M$  の共分散行列 (covariance matrix) である。1次元ガウス分布の分散が正であったように、共分散行列  $\boldsymbol{\Sigma}$  は正定値行列 (positive definite matrix) である必要がある。ベルヌーイ分布 (Bernoulli distribution) はいわゆるコイン投げの

分布である。2値をとる変数  $x \in \{0, 1\}$  を生成するための確率分布で、単一のパラメータ  $\mu \in (0, 1)$  によって分布の性質が決まる。確率質量関数は

$$\text{Bern}(x|\mu) = \mu^x(1-\mu)^{1-x} \quad (13)$$

と定義される。カテゴリ分布 (categorical distribution) は、ベルヌーイ分布を任意の  $D$  値をとるよう拡張したもので、 $\mathbf{s} \in \{0, 1\}^D$  かつ、各要素  $s_d$  が  $\sum_{d=1}^D s_d = 1$  となるような確率変数  $\mathbf{s}$  を生成する分布である。

$$\text{Cat}(\mathbf{s}|\boldsymbol{\pi}) = \prod_{d=1}^D \pi_{s_d} \quad (14)$$

ここで、 $\boldsymbol{\pi} = (\pi_1, \dots, \pi_D)^T$  は分布を決める  $D$  次元のパラメータで、 $\pi_d \in (0, 1)$  かつ  $\sum_{d=1}^D \pi_d = 1$  を満たすように設定する必要がある。ガンマ分布 (gamma distribution) は正の実数  $\lambda > 0$  を生成してくれるような確率分布で、次のように定義される。

$$\text{Gam}(\lambda|a, b) = C_G(a, b) \lambda^{a-1} e^{-b\lambda} \quad (15)$$

$$C_G(a, b) = \frac{b^a}{\Gamma(a)} \quad (16)$$

パラメータ  $a, b$  共に正の実数値として与える必要がある。 $\Gamma(\cdot)$  はガンマ関数 (gamma function) で

$$\Gamma(x) = \int_0^\infty t^{x-1} e^{-t} dt \quad (17)$$

で定義される。また、

$$\Gamma(x+1) = x\Gamma(x) \quad (18)$$

という性質をもつ。

### 1.5 定義

指数分布族は次のような形式でかける確率分布の族である。

$$p(\mathbf{x}|\boldsymbol{\eta}) = h(\mathbf{x}) \exp(\boldsymbol{\eta}^T \mathbf{t}(\mathbf{x}) - a(\boldsymbol{\eta})) \quad (19)$$

それぞれ  $\boldsymbol{\eta}$  を自然パラメータ、 $\mathbf{t}(\mathbf{x})$  を十分統計量、 $h(\mathbf{x})$  を基底測度、 $a(\boldsymbol{\eta})$  を対数分配関数とよぶ。ただし、 $\boldsymbol{\mu}$  は平均パラメータ。

$$a(\boldsymbol{\eta}) = \ln \int h(\mathbf{x}) \exp(\boldsymbol{\eta}^T \mathbf{t}(\mathbf{x})) d\mathbf{x} \quad (20)$$

であり、上式が積分して1になることを保証する。また、ガウス分布、ポアソン分布、多項分布、ベルヌーイ分布など多くの分布が指数分布族として表せる。

## 1.6 分布の共役性

指数分布族に対して次のような共役事前分布と呼ばれる分布族が存在する.

$$p_{\lambda}(\eta) = h_c(\eta) \exp(\eta^T \lambda_1 - a(\eta) \lambda_2 - a_c(\lambda)) \quad (21)$$

共役事前分布の重要な性質は, 次のような指数型分布族による尤度関数に対して, 事後分布も事前分布と同じような形式になることである. いま,  $N$  個のデータ  $\mathbf{X} = \mathbf{x}_1, \dots, \mathbf{x}_N$  を観測したとすると事後分布は

$$\begin{aligned} p(\eta|\mathbf{X}) &\propto p_{\lambda}(\eta) \prod_{n=1}^N p(\mathbf{x}_n|\eta) \\ &= h_c(\eta) \exp(\eta^T \lambda_1 - a(\eta) \lambda_2 - a_c(\lambda)) \\ &\quad \left\{ \prod_{n=1}^N h(\mathbf{x}_n) \right\} \exp \left( \eta^T \sum_{n=1}^N \mathbf{t}(\mathbf{x}_n) - N a(\eta) \right) \\ &\propto h_c(\eta) \exp \left( \eta^T \left( \lambda_1 + \sum_{n=1}^N \mathbf{t}(\mathbf{x}_n) \right) - a(\eta) (N \lambda_2 + 1) \right) \end{aligned}$$

となる. つまり, 事後分布のパラメータを  $\hat{\lambda}_1, \hat{\lambda}_2$  とすれば,

$$\hat{\lambda}_1 = \lambda_1 + \sum_{n=1}^N \mathbf{t}(\mathbf{x}_n), \hat{\lambda}_2 = \lambda_2 + N \quad (23)$$

となっている. よって, 指数分布族による尤度関数に対して共役事前分布を用いると事前分布が解析的に求められることがわかる. また, 共役性を利用すると事後分布を使って未観測のデータ  $\mathbf{x}_*$  の予測分布を次のように求められる.

$$\begin{aligned} p(\mathbf{x}_*|\mathbf{X}) &= \int p(\mathbf{x}_*|\eta) p(\eta|\mathbf{X}) d\eta \\ &= \int h(\mathbf{x}_*) \exp(\eta^T \mathbf{x}_* - a(\eta)) h_c(\eta) \exp(\eta^T \hat{\lambda}_1 - a(\eta) \hat{\lambda}_2 - a_c(\hat{\lambda})) d\eta \\ &= h(\mathbf{x}_*) \frac{\exp(a_c(\hat{\lambda}_1 + \mathbf{x}_*, \hat{\lambda}_2 + 1))}{\exp(a_c(bm\lambda_1, lam\hat{b}da_2))} \end{aligned}$$

結果の確率分布は一般的には指数分布族にはならない.

## 1.7 逐次学習

ベイズ推論によるモデルの学習では, 事後分布によって学習結果を保存することにより, 新規に入ってくる学習データに適応的に学習を進めることができる. これを逐次学習 (sequential learning) あるいはオンライン学習 (online learning) という. 特に, 共役事前分布を使った解析的な学習では, データ

の生成過程に順序の依存性を仮定しない場合, データを逐次的に与えた場合と一度にすべて与えた場合とで最終的に得られる事後分布が一致する.

## 2 マルコフ連鎖

確率変数の系列  $\mathbf{z}^{(1)}, \mathbf{z}^{(2)}, \dots$  に対して

$$p(\mathbf{z}^{(t)}|\mathbf{z}^{(1)}, \mathbf{z}^{(2)}, \dots, \mathbf{z}^{(t-1)}) = p(\mathbf{z}^{(t)}|\mathbf{z}^{(t-1)}) \quad (25)$$

が成り立つとき, 系列  $\mathbf{z}^{(1)}, \mathbf{z}^{(2)}, \dots$  を一次マルコフ連鎖 (first-order Markov chain) と呼ぶ. 遷移確率 (transition probability) を  $\mathcal{T}(\mathbf{z}^{(t-1)}, \mathbf{z}^{(t)})$  とおく. 次の式が成り立つとき, 分布  $p_*(\mathbf{z})$  は定常分布 (stationary distribution) であるという.

$$p_*(\mathbf{z}) = \int \mathcal{T}(\mathbf{z}', \mathbf{z}) p_*(\mathbf{z}') d\mathbf{z}' \quad (26)$$

(定常状態)  $p_*(\mathbf{z})$  がサンプルを取り出したい事後分布だとして,  $p_*(\mathbf{z})$  に分布収束するような遷移確率  $\mathcal{T}(\mathbf{z}^{(t-1)}, \mathbf{z}^{(t)})$  を設計するのがマルコフ連鎖モンテカルロ法のアイデアである.  $p_*(\mathbf{z})$  が定常分布となるための十分条件として詳細釣り合い条件 (detailed balance condition) がある.

$$p_*(\mathbf{z}) \mathcal{T}(\mathbf{z}^{(t-1)}, \mathbf{z}^{(t)}) = p_*(\mathbf{z}') \mathcal{T}(\mathbf{z}^{(t-1)}, \mathbf{z}^{(t)}) \quad (27)$$

詳細釣り合い条件に加え, サンプルサイズを無限大にしたとき, 遷移確率によって任意の初期状態  $p(\mathbf{z}_0)$  から定常分布  $p_*(\mathbf{z})$  に収束できなければならない. この特性をエルゴード性と呼ぶ. 具体的には, マルコフ連鎖において任意の状態から任意の状態へ有限回数で遷移できること (既約性), すべての状態が固定の周回性をもたないこと (非周期性), さらに同じ状態に有限回で戻ることができること (正再帰性) が求められる. (24)

### 2.1 最適化に基づく推論手法

マルコフ連鎖モンテカルロ法は, 無限に計算を続ければ得られるサンプルが真の分布から得られたものと同一視できるという理論面で優れた特性がある. しかし, 実用上は必要なサンプルサイズが明確に決めるににくいことや, 計算コストが膨大になるなどの欠点をもっています. 一方で, 機械学習の分野で主流となっている方法は勾配情報を用いた数値最適化に基づく手法で, 実験的には非常に速い収束速度を持つことが示されている. 最適化による近似推論アルゴリズム

ムの中で、現在最も広く用いられている手法が**変分推論法 (variational inference method)**である。この方法では、事前分布を計算する際に登場する解析不可能な積分を、最適化の問題に置き換えることによって近似的に数値計算する。周辺尤度  $p(\mathbf{X})$  の計算には潜在変数  $\mathbf{Z}$  の積分除去  $p(\mathbf{X}) = \int p(\mathbf{X}, \mathbf{Z}) d\mathbf{Z}$  が必要になるが、モデルが複雑になるとこの積分は解析的に解析的に実行できない。変分推論法では、エビデンス下界 (evidence lower bound, ELBO) と呼ばれる対数周辺尤度  $\ln p(\mathbf{X})$  の下界  $\mathcal{L}(\xi)$  を考える。

$$\ln p(\mathbf{X}) \geq \mathcal{L}(\xi) \quad (28)$$

ここで、 $\xi$  **変分パラメータ (variational parameter)** と呼ばれてるもので、変分推論報における近似分布の平均や分散などを指す。ちなみに、ELBO を負にした  $\mathcal{F} = -\mathcal{L}(\xi)$  は**変分エネルギー (variational energy)** と呼ばれる。勾配降下法などの一般的な最適化手法を用いて  $\mathcal{L}(\xi)$  を  $\xi$  によって最大化すれば、対数周辺尤度  $\ln p(\mathbf{X})$  の近似解が得られることになる。ELBO の設計の仕方はいくつかあり、モデルや目的に応じて使い分ける。最もよく使われる手法は、事後分布  $p(\mathbf{Z}|\mathbf{X})$  をパラメータ  $\xi$  で定められるある分布  $q(\mathbf{Z}; \xi)$  によって近似することである。近似の良さを測る手法にもいくつかあるが、変分推論法では次のような KL ダイバージェンスを使い、変分パラメータ  $\xi$  に関して最小化することによって近似分布  $q(\mathbf{Z}; \xi)$  を得る。

$$q(\mathbf{Z}; \xi_{opt}) = \arg \min_{opt} D_{KL}[q(\mathbf{Z}; \xi) || p(\mathbf{Z}|\mathbf{X})] \quad (29)$$

また、対数周辺尤度  $\ln p(\mathbf{X})$  は次のように ELBO と上式の KL ダイバージェンスに分解できる。

$$\ln p(\mathbf{X}) = \mathcal{L}(\xi) + D_{KL}[q(\mathbf{Z}; \xi) || p(\mathbf{Z}|\mathbf{X})] \quad (30)$$

ただし、

$$\mathcal{L}(\xi) = \int q(\mathbf{Z}; \xi) \ln \frac{p(\mathbf{Z}|\mathbf{X})}{q(\mathbf{Z}; \xi)} d\mathbf{Z} \quad (31)$$

上式からわかるように、対数周辺尤度  $\ln p(\mathbf{X})$  自体は  $\xi$  の値にかかわらず一定なので、 $D_{KL}[q(\mathbf{Z}; \xi) || p(\mathbf{Z}|\mathbf{X})]$  を  $\xi$  に関して最小化する問題は、 $\mathcal{L}\xi$  を  $\xi$  に関して最小化する問題と等価になる。近似分布  $q$  の置き方には様々な選択があり、潜在変数の集合  $\mathbf{Z}$  が  $\mathbf{Z} = \{\mathbf{Z}_1, \dots, \mathbf{Z}_M\}$  のように  $M$  個に分割できるとする。複雑なモデルに対しては、事前分布に独立性を仮定して近似する方法がよく使われている。

$$q(\mathbf{Z}) = \prod_{i=1}^M q(\mathbf{Z}_i) \quad (32)$$

この手法は特に**平均場近似 (mean field approximation)** と呼ばれている。平均場近似では、各近似分布  $q(\mathbf{Z}_1, \dots, \mathbf{Z}_M)$  を交互に更新していく手続きを繰り返すため、アルゴリズムの特性はギブスサンプリングと非常に似たものになっている。

### 3 ニューラルネットワークのベイズ推論

バッチ学習によるニューラルネットワーク (以下, NN) の基本的な学習方法を示す。順伝播型 NN をベイズ化する方法は、ネットワークの挙動を支配するパラメータに事前分布を設定することによって確率的な学習や予測が行えるようにするものである。入力データ  $\mathbf{X} = \{x_1, \dots, x_N\}$  が与えられたもとでの、観測データ  $\mathbf{Y} = \{y_1, \dots, y_N\}$  及びパラメータの同時分布を

$$p(\mathbf{Y}, \mathbf{W}|\mathbf{X}) = p(\mathbf{W}) \prod_{n=1}^N p(y_n | x_n, \mathbf{W}) \quad (33)$$

とおく。ここでは、 $x_n \in \mathbb{R}^{H_0}$  から  $y_n \in \mathbb{R}^D$  を予測する回帰問題であるとし、観測モデルには次のようなガウス分布を使う。

$$p(y_n | x_n, \mathbf{W}) = \mathcal{N}(y_n | f(\mathbf{X}_n; \mathbf{W}), \sigma_y^2 \mathbf{I}) \quad (34)$$

ここで  $\sigma_y^2$  は固定のノイズパラメータで、 $f(\mathbf{X}_n; \mathbf{W})$  は出力次元が  $D$  である NN である。ベイズ推論の枠組みでは、学習データが与えられたあとのパラメータの事後分布を計算する。したがって、NN のパラメータには事前分布を明示的に設定する必要がある。ここでは簡単のため各重みパラメータを  $w \in \mathbf{W}$  とし、次のような独立なガウス分布を考える。

$$p(w) = \mathcal{N}(w | 0, \sigma_w^2) \quad (35)$$

次の図には活性化関数に双曲線正接関数、入力ベクトルを  $(x, 1)^T$  とし、隠れ層の数  $H_1$  及び重みパラメータに仮定するノイズ  $\sigma_w^2$  を変えた場合の関数のサンプル例を示す。  $H_1$  が大きくなる程事前分布によって生成される関数が複雑化し、また  $\sigma_w^2$  が大きくなるほど急激な変化を持つ関数が生成されていることがわかる。

#### 3.1 ラプラス近似による学習

ベイズニューラルネットワーク (以下, BNN) に対するラプラス近似による学習と予測を導出する。簡

単のため、出力次元は  $D = 1$  とする。はじめにモデルの事後分布の MAP 推定値を最適化により求めたあと、その周辺をガウス分布によって近似する。BNN においては重みパラメータ  $\mathbf{W}$  の事後分布の MAP 推定値を求めることが最初のステップになる。事後分布は

$$p(\mathbf{W}|\mathbf{X}, \mathbf{Y}) = \frac{p(\mathbf{W})p(\mathbf{Y}|\mathbf{X}, \mathbf{W})}{p(\mathbf{Y}|\mathbf{X})} \propto p(\mathbf{W})p(\mathbf{Y}|\mathbf{X}, \mathbf{W}) \quad (36)$$

と書ける。分母には  $\mathbf{W}$  がないので、 $\mathbf{W}$  の最適化の過程では無視できる。局所最適解  $\mathbf{W}_{MAP}$  は、対数事後分布の勾配を利用して

$$\mathbf{W}_{new} = \mathbf{W}_{old} + \alpha \nabla_{\mathbf{W}} \ln p(\mathbf{W}|\mathbf{Y}, \mathbf{X})|_{\mathbf{W}=\mathbf{W}_{old}} \quad (37)$$

のように繰り返し更新すれば求まる。ここで、 $\alpha > 0$  は学習率である。対数事後分布は、

$$\begin{aligned} \ln p(\mathbf{W}|\mathbf{Y}, \mathbf{X}) &= \ln p(\mathbf{W}|\mathbf{Y}, \mathbf{X}) + \ln p(\mathbf{W}) + c \\ &= \sum_{n=1}^N \ln p(y_n|\mathbf{x}_n, \mathbf{W}) + \sum_{w \in \mathbf{W}} \ln p(w) \end{aligned} \quad (38)$$

と書ける。あるパラメータ  $w \in \mathbf{W}$  の偏微分を計算すると

$$\frac{\partial}{\partial w} \ln p(\mathbf{W}|\mathbf{Y}, \mathbf{X}) = -\left\{ \frac{1}{\sigma_y^2} \frac{\partial}{\partial w} E(\mathbf{W}) + \frac{1}{\sigma_w^2} \frac{\partial}{\partial w} \Omega_{L2}(\mathbf{W}) \right\} \quad (39)$$

となる。上式において  $\Omega_{L2}(\mathbf{W})$  は各パラメータ  $w$  に関するガウス事前分布  $p(w)$  に由来するもびであるが、こちらの方が微分が容易である。また、 $E(\mathbf{W})$  は NN の誤差関数であるが、これは通常の誤差逆伝播を用いて微分を評価する。

$$q(\mathbf{W}) = \mathcal{N}(\mathbf{W}|\mathbf{W}_{MAP}, \{\Lambda(\mathbf{W}_{MAP})\}^{-1}) \quad (40)$$

ここで、 $\Lambda$  は精度行列であり、

$$\begin{aligned} \Lambda &= -\nabla_{\mathbf{W}}^2 \ln p(\mathbf{W}|\mathbf{Y}, \mathbf{X}) \\ &= \frac{1}{\sigma_w^2} \mathbf{I} + \frac{1}{\sigma_y^2} \mathbf{H} \end{aligned} \quad (41)$$

である。また、 $\mathbf{H}$  は NN の誤差関数に対するヘッセ行列である。

### 3.2 予測分布の近似

パラメータの事後分布の近似を得た後は、テストの入力  $\mathbf{x}_*$  に対する出力  $y_*$  の予測分布を

$$p(y_*|\mathbf{x}_*, \mathbf{Y}, \mathbf{X}) \approx \int p(y_*|\mathbf{x}_*, \mathbf{W})q(\mathbf{W})d\mathbf{W} \quad (42)$$

として近似する。パラメータの事後分布を簡単なガウス分布  $q(\mathbf{W})$  によって近似したものの、 $p(y_*|\mathbf{x}_*, \mathbf{W})$  の中には NN が含まれているため、依然として一般的に予測分布の計算は解析的に行えない。ここでは予測分布を計算するために、NN の関数の線形近似を行う。この近似では、パラメータの事後分布が MAP 推定値の周辺に集中しており、かつその小さな範囲においては NN の関数値  $f(\mathbf{x}_*; \mathbf{W})$  が  $\mathbf{W}$  の線形関数でよく近似できるという仮定をおく。テイラー展開で  $\mathbf{W}$  の関数  $f(\mathbf{x}_*; \mathbf{W})$  を  $\mathbf{W}_{MAP}$  まわりで一次近似すれば

$$f(\mathbf{x}_*; \mathbf{W}) \approx f(\mathbf{x}_*; \mathbf{W}_{MAP}) + \mathbf{g}^T(\mathbf{W} - \mathbf{W}_{MAP}) \quad (43)$$

ただし、 $\mathbf{g}$  は次のように関数の勾配を  $\mathbf{W}_{MAP}$  で評価したものである。

$$\mathbf{g} = \nabla_{\mathbf{W}} f(\mathbf{x}_*; \mathbf{W})|_{\mathbf{W}=\mathbf{W}_{MAP}} \quad (44)$$

この近似を用いれば、予測分布の計算に NN 特有の非線形関数がなくなるので、解析的に計算することが可能になる。したがって、求めたい予測分布の近似は、

$$\begin{aligned} p(y_*|\mathbf{x}_*, \mathbf{Y}, \mathbf{X}) &\approx \int p(y_*|\mathbf{x}_*, \mathbf{W})q(\mathbf{W})d\mathbf{W} \\ &\approx \int \mathcal{N}(y_*|f(\mathbf{x}_*; \mathbf{W}_{MAP})) + \mathbf{g}^T(\mathbf{W} - \mathbf{W}_{MAP}), \\ &\quad \mathcal{N}(\mathbf{W}|\mathbf{W}_{MAP}, \{\Lambda(\mathbf{W}_{MAP})\}^{-1})d\mathbf{W} \\ &= \mathcal{N}(y_*|f(\mathbf{x}_*; \mathbf{W}_{MAP}), \sigma^2(\mathbf{x}_*)) \end{aligned} \quad (45)$$

と計算できる。ただし、

$$\sigma^2(\mathbf{x}_*) = \sigma_y^2 + \mathbf{g}^T + \{\Lambda(\mathbf{W}_{MAP})^{-1}\}\mathbf{g} \quad (46)$$

である。

### 3.3 重みパラメータの推論

BNN への正規化されていない事後分布を利用すると対応するポテンシャルエネルギーは

$$\mathcal{U}(\mathbf{W}) = -\{\ln p(\mathbf{Y}|\mathbf{X}, \mathbf{W}) + \ln p(\mathbf{W})\} \quad (47)$$

となる。リープフロッグ法を使うためにはポテンシャルエネルギーの微分が必要になるが、正則化項をもつ基本的な NN のコスト関数の微分と等価であるから、誤差逆伝播による勾配計算が利用できる。

### 3.4 ハイパーパラメータの推論

重みパラメータ  $\mathbf{W}$  の事前分布を支配する  $\sigma_w^2$  や観測モデルのノイズパラメータ  $\sigma_y^2$  はハイパーパラメータとして扱われ、通常は学習を実行する前に適切なものを固定値として与えておく必要がある。これらのハイパーパラメータはデータの“大まかな”スケール感を反映していると考えられるが、もし直観的に値を指定するのが困難である場合や、学習データに対して当てはまりの良い学習結果を得たい場合は、ハイパーパラメータに対しても事前分布を設定することによって、サンプリングの枠組みの中で重みパラメータ  $\mathbf{W}$  と同時に推論することもできる。ハイパーパラメータも同時推論するために、これらに事前分布を与える。パラメータ  $\mathbf{W}$  は分散  $\sigma_w^2$  をもつガウス分布にしたがって決定される。表記を簡単にするため、精度パラメータ  $\gamma_w = \sigma_w^{-2}$  を導入し、事前分布として、

$$p(\gamma_w) = \text{Gam}(\gamma_w | a_w, b_w) \quad (48)$$

を考える。ここで、 $a_w > 0, b_w > 0$  は固定値として与える。観測ノイズに対する精度パラメータ  $\gamma_y = \sigma_y^{-2}$  の事前分布も同様にして

$$p(\gamma_y) = \text{Gam}(\gamma_y | a_y, b_y) \quad (49)$$

と設定する。ただし、 $a_y > 0, b_y > 0$  である。これらの精度パラメータの事前分布を導入した場合のモデルを改めて書き下すと

$$p(\mathbf{Y}, \mathbf{W}, \gamma_w, \gamma_y | \mathbf{X}) = p(\gamma_w) p(\gamma_y) p(\mathbf{W} | \gamma_w) \prod_{n=1}^N p(y_n | \mathbf{x}_n, \mathbf{W}, \gamma_y) \quad (50)$$

となる。したがって事後分布全体は

$$p(\mathbf{W}, \gamma_w, \gamma_y | \mathbf{Y}, \mathbf{X}) \quad (51)$$

となる。ここではギブスサンプリングを用いて各確率変数のサンプルを得ることを考える。各確率変数  $\mathbf{W}, \gamma_w, \gamma_y$  を条件付き確率を用いて別々にサンプリングすることである。 $\gamma_w, \gamma_y$  がサンプルされた値で条件付けされた下での  $\mathbf{W}$  の分布は

$$p(\mathbf{W} | \mathbf{Y}, \mathbf{X}, \gamma_w, \gamma_y) \quad (52)$$

となるが、これはパラメータの事後分布そのものであるため、通常通りハミルトニアンモンテカルロ法を実行することによって  $\mathbf{W}$  のサンプルを得る。 $\mathbf{W}, \gamma_y$  が与えられた下での  $\gamma_w$  の分布は、 $\gamma_w$  に関わらない部分を無視すると

$$p(\gamma_w | \mathbf{Y}, \mathbf{X}, \mathbf{W}, \gamma_y) \propto p(\mathbf{W} | \gamma_w) p(\gamma_w) \quad (53)$$

とかける。 $p(\mathbf{W} | \gamma_w)$  がガウス分布であり、精度  $\gamma_w$  の事前分布  $p(\gamma_w)$  には共役事前分布であるガンマ分布を用いているので、この条件付き分布もガンマ分布として解析的に以下のようにもとまる。

$$\gamma_w \sim \text{Gam}(\hat{a}_w, \hat{b}_w) \quad (54)$$

$$\hat{a}_w = a_w + \frac{K_w}{2} \quad (55)$$

$$\hat{b}_w = b_w + \frac{1}{2} \sum_{w \in \mathbf{W}} w^2 \quad (56)$$

のようにして  $\gamma_w$  のサンプルが得られる。ただし、 $K_w$  は重みパラメータ  $\mathbf{W}$  の総数である。次に、 $\mathbf{W}, \gamma_w$  が与えられた下での  $\gamma_y$  の分布も

$$p(\gamma_y | \mathbf{Y}, \mathbf{X}, \mathbf{W}, \gamma_w) \propto p(\gamma_y) \prod_{n=1}^N p(y_n | \mathbf{x}_n, \mathbf{W}, \gamma_y) \quad (57)$$

と計算できる。観測モデル  $p(\mathbf{Y} | \mathbf{X}, \mathbf{W}, \gamma_y)$  がガウス分布であり、精度  $\gamma_y$  の事前分布  $p(\gamma_y)$  には共役なガンマ分布を用いたので、こちらも解析的に分布を計算でき、

$$\gamma_y = \text{Gam}(\hat{a}_y, \hat{b}_y) \quad (58)$$

$$\hat{a}_y = a_y + \frac{N}{2} \quad (59)$$

$$\hat{b}_y = b_y + \frac{1}{2} \sum_{n=1}^N \{y_n - f(\mathbf{x}_n; \mathbf{W})\}^2 \quad (60)$$

のようにして  $\gamma_y$  のサンプルが得られる。直観的に上式を解釈すると、 $\gamma_y$  は NN の関数  $f(\mathbf{x}; \mathbf{W})$  で表現できない誤差を学習していると言える。ガンマ分布の平均は  $\hat{a}_y / \hat{b}_y$  であるため、 $\hat{b}_y$  が大きいほど関数  $f(\mathbf{x}; \mathbf{W})$  による  $y$  の推定の精度が低く、観測に対する分散が大きくなるように学習される。なお、ここでは簡単のため共通なハイパーパラメータ  $\gamma_w$  が与えられていると仮定したが、ハイパーパラメータを複数のグループに分割して与えることもできる。

## 参考文献

- [1] 斎藤康毅. ゼロから作る Deep Learning 2. p57-129
- [2] 【機械学習】誤差逆伝播法による速度改善（その2）. <https://qiita.com/m-hayashi/items/fa4749f8080e542787d2>
- [3] 【機械学習 誤差逆伝播法】 word2vec メモ (1) <https://qiita.com/sand/items/85ea76f9c26aabb849e7>

- [4] Improving Distributional Similarity with  
Losses Learned from Word Embeddings  
<https://www.aclweb.org/anthology/Q15-1016/>