

目次

第 1 章	生成モデルとベイズ脳仮説	3
1.1	確率的生成モデル	4
1.1.1	確率的生成モデルとベイズ推論	4
1.1.2	ベイズ線形回帰	6
1.1.3	最尤推定と MAP 推定	10
1.1.4	潜在変数モデル	12
1.1.5	階層ベイズモデル	12
1.2	スパース符号化と予測符号化	12
1.2.1	スパース符号化	12
1.2.2	予測符号化	13
1.3	近似ベイズ推論	13
1.3.1	変分推論	13
1.3.2	自由エネルギー原理	13
1.3.3	マルコフ連鎖モンテカルロ法	13
1.4	神経回路による不確実性の表現	13
1.4.1	神経サンプリング	13
1.4.2	確率的集団符号化	13
	参考文献	13

第 1 章

生成モデルとベイズ脳仮説

これまでの章では、知覚 (perception) のモデル、すなわち外界からの入力に対して、どのようにして神経回路網が意味のある出力を生成するのか、という問題を主に扱ってきた。ここで改めて知覚の基本的な定義を確認しておこう。知覚とは、外界からの刺激を感覚受容器によって受容し、それに意味を与える過程である。この「刺激に意味を与える」という個所を、より体系的に理解するために、「順問題」と「逆問題」という概念を導入しよう。

一般に、ある原因から結果を予測する問題は順問題 (forward problem) と呼ばれる。逆に、観測された結果からその原因を推定する問題は逆問題 (inverse problem) と呼ばれる。視覚を例にとって、順問題と逆問題について考えてみよう。たとえば、三次元の物体が光を反射し、それが二次元の網膜上にどのような像を結ぶか、という問いは順問題に分類される。これに対して、網膜上に投影された二次元像から、元の物体の三次元的な構造や大きさ、位置などを推定する課題が逆問題である^{*1}。光学の分野では、それぞれの問題は順光学 (forward optics)、逆光学 (inverse optics) と呼ばれている。逆問題は多くの場合、不良設定問題 (ill-posed problem) となる。すなわち、解が存在しない、解が一意に定まらない、あるいはわずかな誤差に対して解が大きく変化するという性質をもつ^{*2}。例えば、先ほどの例であれば同じ 2 次元像を示す 3 次元物体は複数 (あるいは無数に) 存在する。そのため、逆問題を解くには、事前知識や

^{*1} 他にも逆問題は数多く存在する。逆問題は様々な分野に現れるが、ここでは医学や神経科学に関連した例として、外部から脳の構造や機能を推定する問題を取り上げる。たとえば、医用画像解析では、コンピュータ断層撮影 (computed tomography; CT)、磁気共鳴画像法 (magnetic resonance imaging; MRI)、陽電子放射断層撮影 (positron emission tomography; PET) などにおいて、観測データから画像を再構成する必要がある。この再構成処理には、CT や PET では逆ラドン変換、MRI では逆フーリエ変換が用いられる。また、神経活動を非侵襲的に計測する手法として、脳波 (electroencephalography; EEG) や脳磁図 (magnetoencephalography; MEG) がある。これらにおける電流源推定 (source localization) も典型的な逆問題である。EEG や MEG における順問題は、脳内の神経電流源の位置・方向・強度から、頭皮上の電極 (EEG) や磁場センサ (MEG) によって観測される電位や磁場分布を予測することである。一方、逆問題は、実際に観測された電位や磁場データから、神経電流源の空間的位置と活動を推定することである。この逆問題は不良設定 (ill-posed) であるため、安定的に解くには、MRI から得られた頭部の構造データに基づいて構築された順モデル (forward model) が必要となる。

^{*2} これに対して、良設定問題 (well-posed problem) とは、解が存在し、一意であり、かつ入力の変動に対して連続的に変化する (安定性をもつ) ような問題を指す。良設定問題では、入力データに小さなノイズや誤差が含まれていても、求められる解は大きく変わることなく、安定に計算することができる。

仮定 (制約条件, 正則化) の導入が必要となる。

こうした逆問題を踏まえ, 知覚とは単なる入力情報の受動的な処理ではなく, 感覚入力という結果から外界に存在する潜在的な原因を推定する逆推論 (abductive reasoning) の過程とみなす考えがある (Helmholtz, 1867; Mumford, 1992; Kawato et al., 1993; Friston, 2003)^{*3}. この枠組みを推論的知覚 (perception as inference) と呼ぶ. 推論的知覚は, 外界の潜在的な原因から感覚入力生成される過程を記述する確率的生成モデル (probabilistic generative model) に基づいて説明される. 確率的生成モデルについて説明する前に, 前提となるベイズ推論について次節で説明する.

本章で触れる内容についてまとめ直す.

1.1 確率的生成モデル

1.1.1 確率的生成モデルとベイズ推論

外界から感覚入力などを通じて観測データ x を得る状況を考えよう^{*4}. 観測データが存在するという事は, それを生成する確率分布^{*5} $p_{\text{data}}(x)$ が存在する (すなわち $x \sim p_{\text{data}}(x)$ である) と仮定できる. この $p_{\text{data}}(\cdot)$ はしばしば真の確率分布と呼ばれるが, 実際にそのような分布が存在する保証はなく, 多くの場合は未知である. もし $p_{\text{data}}(\cdot)$ が既知であれば, 任意のサンプル x をそこから直接生成 (サンプリング) できるが, 現実にはこれを直接知ることはいできない.

このため, 観測データがある確率的な生成過程に従って生じたと仮定し, その過程を記述する生成モデルを構築する. 生成モデルは分布を明示的に表現するため, 新たなデータの生成や欠損値の補完, 潜在構造の抽出, 外界の状態推定など, 多様な推論を可能にする. ここではパラメータ θ をもつ条件付き確率密度関数 $p(x | \theta)$ を導入し, 観測の背後にある生成過程を近似的に表す. このような確率分布 $p(x | \theta)$ を定めるモデルを確率的生成モデル (probabilistic generative model) と呼ぶ. また, このように有限個のパラメータ θ で分布形状を規定するモデルをパラメトリックモデル (parametric model) と呼ぶ^{*6}.

^{*3} Helmholtz は, 知覚を単なる感覚の受容ではなく, 感覚入力に意味を与え, 対象として構成する過程であると捉えた. この過程には, 観念の連合 (*Vorstellungsverbindungen*) が関与している. 観念の連合とは, 過去の経験によって形成された (必ずしも言語化を伴わない) 観念や知識が, 現在の感覚入力と結び付けられる過程を指す. 通常, 推論とは意識的に行われるものと考えられているが, Helmholtz はこのような観念の連合を, 意識されることなく行われる推論として捉え, 無意識的推論 (*unbewusster Schluss*, unconscious inference) と表現した. なお, この脚注ではドイツ語を斜体で表記した.

^{*4} 厳密には, 確率変数は大文字 X , その実現値は小文字 x で表記して区別するのが原則である. しかし, 応用的な文脈では両者を混同しても支障をきたすことは少ないため, 本書では明確に区別しない方針をとる. 特に, 確率変数がスカラーの場合には大文字・小文字で容易に区別できるが, ベクトルや行列を扱う際には表記が煩雑となり, 可読性を損ねる恐れがある. このため, 本書では, 確率変数とその実現値の区別が必要となる場合にはその旨を明示し, それ以外では変数の次元に基づく記号表記を基本とする.

^{*5} 本書では確率分布を確率密度関数の意味で用いる.

^{*6} 有限個のパラメータで分布形状をあらかじめ規定せず, データ量に応じて表現可能な複雑さが変化するものをノンパラメトリックモデル (non-parametric model) と呼ぶ. 代表例にはヒストグラムやカーネル法による密度推定, ガウス過程, 分位

例えば, 正規分布 $p(x | \theta) = \mathcal{N}(x | \mu, \sigma^2)$ ($\theta = \{\mu, \sigma\}$) はパラメトリックな確率的生成モデルの一例である. この場合, 分布の形状 (正規分布) はあらかじめ固定され, 未知なのはパラメータ θ である. ここで改めて強調しておくと, 目標は真の分布 $p_{\text{data}}(\cdot)$ を近似できる生成モデルを構築することである. パラメトリックモデルの場合, この目標はモデルのパラメータを適切に推定することによって達成される.

パラメータ推定には, 大きく分けて二つの方法がある. 一つは, パラメータの最適な一点の値を求める点推定であり, もう一つはパラメータを確率変数として扱い, その不確実性を含めて推定する分布推定である. パラメータ推定には大きく二つの方法がある. 一つは, パラメータの最適な一点の値を求める点推定であり, もう一つはパラメータを確率変数として扱い, その不確実性を含めて推定する分布推定である. 分布推定には様々な方法があるが, ここではその代表例としてベイズ推論 (Bayesian inference) を取り上げる. ベイズ推論では, 観測前のパラメータ分布を事前分布 (prior) $p(\theta)$, 観測後の分布を事後分布 (posterior) $p(\theta | x)$ と呼び, 事前分布を事後分布へと更新する. この更新は, 尤度 (likelihood) とベイズの定理 (Bayes' theorem) に基づいて行われる.

尤度

先ほど導入した $p(x | \theta)$ は, 観測データ x とパラメータ θ のどちらを変数とみなすかによって確率モデルと尤度という 2 通りの解釈がある.

確率モデルとして解釈する場合, θ を固定し, x を確率変数として扱う. このとき, $p(x | \theta)$ は「 θ が与えられたときに, どのようなデータ x がどの確率で得られるか」を表す.

一方で, 尤度として解釈する場合は, 観測データを固定して, θ を変数として扱う. このとき,

$$L(\theta; x) := p(x | \theta) \quad (1.1)$$

を尤度関数 (likelihood function) と呼ぶ. 尤度は, 仮定した θ の下でデータが得られる「尤もらしさ」を定量化する指標であり, 値が大きいほどデータをよく説明すると解釈できる. なお, 尤度関数は θ に関する確率分布ではないため, θ について積分しても必ずしも 1 にはならない.

ベイズの定理

事前分布と尤度が設定されれば, 事後分布を次式で表されるベイズの定理によって求めることができる:

$$\underbrace{p(\theta | x)}_{\text{事後分布}} = \frac{\overbrace{p(x | \theta)}^{\text{尤度}} \overbrace{p(\theta)}^{\text{事前分布}}}{\underbrace{p(x)}_{\text{周辺尤度}}} \quad (1.2)$$

ここで

$$p(x) = \int p(x | \theta) p(\theta) d\theta \quad (1.3)$$

は周辺尤度 (marginal likelihood) あるいは 証拠 (evidence) , 正規化定数 (normalization constant) と呼ばれる。このベイズの定理は

$$p(\theta, x) = p(x | \theta) p(\theta) = p(\theta | x) p(x) \quad (1.4)$$

が成り立つことから導かれる。ここで $p(\theta, x)$ は θ と x の同時確率分布 (joint probability distribution) である。

予測分布

事後分布 $p(\theta | x)$ は、観測データを得た後のパラメータ θ の確率分布であり、ベイズ推論はこの分布を更新する過程とみなせる。十分に更新された事後分布が得られれば、新たなデータの予測も可能となる。予測においては事後分布全体を平均化した事後予測分布 (posterior predictive distribution) を用いる。観測されたデータの実現値を x , 予測対象のデータを \tilde{x} とすると、

$$p(\tilde{x} | x) = \int p(\tilde{x}, \theta | x) d\theta \quad (1.5)$$

$$= \int p(\tilde{x} | \theta, x) p(\theta | x) d\theta \quad (1.6)$$

$$= \int p(\tilde{x} | \theta) p(\theta | x) d\theta \quad (\because \tilde{x} \perp\!\!\!\perp x | \theta) \quad (1.7)$$

となる。ここで最後の式変形には、「 θ が与えられた条件下で \tilde{x} と x が独立」という条件付き独立性を用いた。ベイズ推論において、この事後予測分布 $p(\tilde{x} | x)$ が真の生成分布の推論結果となる。なお、データの観測をしていない場合の予測は周辺尤度と同様の形式

$$p(\tilde{x}) = \int p(\tilde{x} | \theta) p(\theta) d\theta \quad (1.8)$$

で与えられ、これを事前予測分布 (prior predictive distribution) と呼ぶ。

1.1.2 ベイズ線形回帰

ここでは線形回帰モデルをベイズ化した、すなわち予測の不確実性を表現できるようにしたベイズ線形回帰 (Bayesian linear regression) モデルを取り扱う。

多変量正規分布

まず、多変量正規分布 (ガウス分布) を導入する。1次元の場合、正規分布は次の確率密度関数で表される。

$$\mathcal{N}(x | \mu, \sigma^2) := \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{(x - \mu)^2}{2\sigma^2}\right) \quad (1.9)$$

ここで, $\mu \in \mathbb{R}$ は平均, $\sigma^2 > 0$ は分散を表し, σ は標準偏差である. この式を $x \in \mathbb{R}$ から d 次元のベクトル $\mathbf{x} \in \mathbb{R}^d$ に拡張すると, 分布は多変量正規分布 (multivariate normal distribution) となる.

$$\mathcal{N}(\mathbf{x} \mid \boldsymbol{\mu}, \boldsymbol{\Sigma}) := \frac{1}{\sqrt{(2\pi)^d |\boldsymbol{\Sigma}|}} \exp \left(-\frac{1}{2} (\mathbf{x} - \boldsymbol{\mu})^\top \boldsymbol{\Sigma}^{-1} (\mathbf{x} - \boldsymbol{\mu}) \right) \quad (1.10)$$

ここで, $\boldsymbol{\mu} \in \mathbb{R}^d$ は各成分の平均を並べた平均ベクトル, $\boldsymbol{\Sigma} \in \mathbb{R}^{d \times d}$ は共分散行列 (covariance matrix) である. 共分散行列の対角成分は各次元の分散を表し, 非対角成分は共分散を表す. このため, 共分散行列は分散共分散行列 (variance-covariance matrix) とも呼ばれる. 多変量正規分布が定義可能であるためには, $\boldsymbol{\Sigma}$ が正定値行列 (positive definite matrix) であることが必要であり, これは任意の非ゼロベクトル $\mathbf{z} \in \mathbb{R}^d \setminus \{\mathbf{0}\}$ に対して^{*7}

$$\mathbf{z}^\top \boldsymbol{\Sigma} \mathbf{z} > 0 \quad (1.11)$$

が成り立つことを意味する. この条件を満たす場合, $\boldsymbol{\Sigma}^{-1}$ が存在してそれ自体も正定値となるため, 特に $\mathbf{x} = \boldsymbol{\mu}$ の場合も含めて

$$(\mathbf{x} - \boldsymbol{\mu})^\top \boldsymbol{\Sigma}^{-1} (\mathbf{x} - \boldsymbol{\mu}) \geq 0 \quad (1.12)$$

が常に成立する. なお, 共分散行列は $\boldsymbol{\Sigma} = \boldsymbol{\Sigma}^\top$ という対称性を持つが, これは非対角成分が共分散を表し, その定義から $\Sigma_{ij} = \Sigma_{ji}$ が必ず成り立つことによる. 正定値行列という概念は, 対称行列や, より一般にはエルミート行列に対して定義されるため, 多変量正規分布においても, 共分散行列はこのように対称性を持った上で正定値でなければならない. Julia 言語において, 行列が正定値行列か確認するには `LinearAlgebra` ライブラリの `isposdef` 関数を用いればよい. 行列 A が正定値ならば, `isposdef(A)` は `true` を返す.

また, 後ほど使用するため多変量正規分布の確率密度関数の対数を取った形式を確認しておこう. $p(\mathbf{x}) = \mathcal{N}(\mathbf{x} \mid \boldsymbol{\mu}, \boldsymbol{\Sigma})$ とすると,

$$\ln p(\mathbf{x}) = -\frac{d}{2} \ln(2\pi) - \frac{1}{2} \ln |\boldsymbol{\Sigma}| - \frac{1}{2} (\mathbf{x} - \boldsymbol{\mu})^\top \boldsymbol{\Sigma}^{-1} (\mathbf{x} - \boldsymbol{\mu}) \quad (1.13)$$

となる.

多変量正規分布のコード

ベイズ線形回帰のモデル定義

入力 $\mathbf{x} \in \mathbb{R}^d$ から実数値出力 $y \in \mathbb{R}$ を予測するモデルを考える. 基底関数 (basis function) $\phi: \mathbb{R}^d \rightarrow \mathbb{R}^{d'}$ を導入し, 入力を $\phi(\mathbf{x}) \in \mathbb{R}^{d'}$ に写像する. バイアス項を含める場合には, 基底関数の 1 つの成分を

^{*7} 集合 A, B があるとき, $A \setminus B$ は A から B を引いた差集合を意味する.

常に 1 を返す定数関数として組み込むことにする。訓練データを $\mathcal{D} = \{(\mathbf{x}_i, y_i)\}_{i=1}^n$ とすると、計画行列は

$$\Phi = \begin{bmatrix} \phi(\mathbf{x}_1)^\top \\ \vdots \\ \phi(\mathbf{x}_n)^\top \end{bmatrix} \in \mathbb{R}^{n \times d'} \quad (1.14)$$

で表される。重みパラメータ $\mathbf{w} \in \mathbb{R}^{d'}$ を導入すると、 $\mathbf{y} = [y_1, \dots, y_n]^\top \in \mathbb{R}^n$ の生成過程は

$$p(\mathbf{y} \mid \mathbf{w}, \mathbf{X}) = \mathcal{N}(\mathbf{y} \mid \Phi \mathbf{w}, \beta^{-1} \mathbf{I}_n) = \prod_{i=1}^n \mathcal{N}(y_i \mid \phi(\mathbf{x}_i)^\top \mathbf{w}, \beta^{-1}) \quad (1.15)$$

で与えられる。ここで、 $\beta (> 0)$ は尤度の精度 (precision) であり、分散の逆数を意味する。さらに、重みパラメータに対して正規分布の事前分布

$$p(\mathbf{w}) = \mathcal{N}(\mathbf{w} \mid \mathbf{0}, \alpha^{-1} \mathbf{I}_{d'}) \quad (1.16)$$

を仮定する。ここで、 $\alpha (> 0)$ はパラメータの事前分布の精度である。また、この事前分布は共役事前分布 (conjugate prior) となっている。共役であるとは、事前分布と尤度の組み合わせにより得られる事後分布が^{*}、同じ分布族で表現できることを意味する。

事後分布

次に、事後分布を導出する。ベイズの定理より

$$p(\mathbf{w} \mid \mathbf{y}, \mathbf{X}) \propto p(\mathbf{y} \mid \mathbf{w}, \mathbf{X}) p(\mathbf{w}) \quad (1.17)$$

が成り立つ^{*8}。次に、事後分布を解析的に計算し、多変量正規分布の係数を見捨てた形状に式をまとめなおす。両辺の対数を取り、 \mathbf{w} に関する二次形式をまとめると次のようになる：

$$\ln p(\mathbf{w} \mid \mathbf{y}, \mathbf{X}) \propto \ln p(\mathbf{y} \mid \mathbf{w}, \mathbf{X}) + \ln p(\mathbf{w}) \quad (1.18)$$

$$= -\frac{\beta}{2} \|\mathbf{y} - \Phi \mathbf{w}\|^2 - \frac{\alpha}{2} \|\mathbf{w}\|^2 + \text{const.} \quad (1.19)$$

$$= -\frac{1}{2} [\mathbf{w}^\top (\beta \Phi^\top \Phi + \alpha \mathbf{I}_{d'}) \mathbf{w} - 2\beta \mathbf{y}^\top \Phi \mathbf{w}] + \text{const.} \quad (1.20)$$

が成り立つ (\mathbf{w} を含まない項は const. に吸収した)。ここで、事後分布を \mathbf{w} についての多変量正規分布として

$$p(\mathbf{w} \mid \mathbf{y}, \mathbf{X}) := \mathcal{N}(\mathbf{w} \mid \hat{\boldsymbol{\mu}}, \hat{\boldsymbol{\Sigma}}) \quad (1.21)$$

とおく。この確率密度関数の指数部 ($\exp(\cdot)$ の中身) は

$$-\frac{1}{2} (\mathbf{w} - \hat{\boldsymbol{\mu}})^\top \hat{\boldsymbol{\Sigma}}^{-1} (\mathbf{w} - \hat{\boldsymbol{\mu}}) = -\frac{1}{2} \left[\mathbf{w}^\top \hat{\boldsymbol{\Sigma}}^{-1} \mathbf{w} - 2\hat{\boldsymbol{\mu}}^\top \hat{\boldsymbol{\Sigma}}^{-1} \mathbf{w} + \hat{\boldsymbol{\mu}}^\top \hat{\boldsymbol{\Sigma}}^{-1} \hat{\boldsymbol{\mu}} \right] \quad (1.22)$$

^{*8} ここで、ベイズの定理における分母の周辺尤度 $p(\mathbf{y} \mid \mathbf{X})$ は事後分布 (\mathbf{w} の関数) の形状に影響しないため、無視した。また、 \mathbf{X} は常に実現値で与えられるので、 $p(\mathbf{X})$ を考える必要はない。

であるため, \mathbf{w} を含む項のみに関して両者を比べると,

$$\hat{\Sigma}^{-1} = \beta \Phi^{\top} \Phi + \alpha \mathbf{I}_{d'}, \quad \hat{\mu} = \hat{\Sigma} \beta \mathbf{y}^{\top} \Phi \quad (1.23)$$

とすれば良いことが分かる.

事後予測分布

最後に, 事後予測分布の導出を行う. 新しい入力を $\tilde{\mathbf{x}}$, 対応する出力を \tilde{y} で表す. 事後予測分布は前節での導出に基づくと,

$$p(\tilde{y} | \tilde{\mathbf{x}}, \mathbf{y}, \mathbf{X}) = \int p(\tilde{y} | \mathbf{w}, \tilde{\mathbf{x}}) p(\mathbf{w} | \mathbf{y}, \mathbf{X}) d\mathbf{w} \quad (1.24)$$

により計算することができる. ここで, 生成過程の仮定より

$$p(\tilde{y} | \mathbf{w}, \tilde{\mathbf{x}}) = \mathcal{N}(\tilde{y} | \phi(\tilde{\mathbf{x}})^{\top} \mathbf{w}, \beta^{-1}) \quad (1.25)$$

である. 積分を直接実行すると煩雑な計算を要するため, ここでは正規分布の閉じた性質を利用する. すなわち, 正規分布を正規分布で畳み込むと再び正規分布となる. このため, 事後予測分布は平均と分散を計算するだけで決定できる.

$$\mathbb{E}[\tilde{y}] = \mathbb{E}_{\mathbf{w}}[\mathbb{E}[\tilde{y} | \mathbf{w}]] = \mathbb{E}_{\mathbf{w}}[\phi(\tilde{\mathbf{x}})^{\top} \mathbf{w}] = \phi(\tilde{\mathbf{x}})^{\top} \mathbb{E}[\mathbf{w}] = \phi(\tilde{\mathbf{x}})^{\top} \hat{\mu} \quad (1.26)$$

$$\text{Var}[\tilde{y}] = \mathbb{E}_{\mathbf{w}}[\text{Var}[\tilde{y} | \mathbf{w}]] + \text{Var}_{\mathbf{w}}[\mathbb{E}[\tilde{y} | \mathbf{w}]] = \mathbb{E}_{\mathbf{w}}[\beta^{-1}] + \text{Var}_{\mathbf{w}}[\phi(\tilde{\mathbf{x}})^{\top} \mathbf{w}] \quad (1.27)$$

$$= \beta^{-1} + \phi(\tilde{\mathbf{x}}) \text{Var}[\mathbf{w}] \phi(\tilde{\mathbf{x}})^{\top} = \beta^{-1} + \phi(\tilde{\mathbf{x}})^{\top} \hat{\Sigma} \phi(\tilde{\mathbf{x}}) \quad (1.28)$$

ここで, 期待値の計算では, 期待値の線形性の法則を用いた. 分散の計算では, 全分散の法則 (law of total variance) および分散の線形変換の法則を用いた^{*9}. よって, 事後予測分布は

$$p(\tilde{y} | \tilde{\mathbf{x}}, \mathbf{y}, \mathbf{X}) := \mathcal{N}\left(\tilde{y} \mid \phi(\tilde{\mathbf{x}})^{\top} \hat{\mu}, \beta^{-1} + \phi(\tilde{\mathbf{x}})^{\top} \hat{\Sigma} \phi(\tilde{\mathbf{x}})\right) \quad (1.29)$$

となる.

ベイズ線形回帰のコード

^{*9} ここで用いた 3 つの法則を補足しておく. 確率変数 X, Y と行列 \mathbf{A} に関して, 以下の基本的な性質が成り立つ:

- | | |
|-----------------------|---|
| (1) 期待値の線形性: | $\mathbb{E}[\mathbf{A}X] = \mathbf{A} \mathbb{E}[X]$ |
| (2) 全分散の法則 (分散分解の法則): | $\text{Var}[Y] = \mathbb{E}[\text{Var}[Y X]] + \text{Var}[\mathbb{E}[Y X]]$ |
| (3) 分散の線形変換: | $\text{Var}[\mathbf{A}X] = \mathbf{A} \text{Var}[X] \mathbf{A}^{\top}$ |

1.1.3 最尤推定と MAP 推定

ベイズ推論の枠組みから派生し、パラメータ θ の分布を推定するのではなく、その点推定値を与える手法について説明する。ここでは、パラメータの不確実性は扱わず、単一の推定値を得ること、すなわちパラメータの点推定について取り扱う。代表的な点推定の方法には、最尤推定 (maximum likelihood estimation; MLE) と 最大事後確率推定 (maximum a posteriori estimation; MAP 推定) がある。

まず、一般的な生成モデルを考える。観測データを \mathbf{x} , 未知パラメータを θ , 真のデータ分布を $p_{\text{data}}(\mathbf{x})$, 確率モデルあるいは尤度を $p(\mathbf{x} | \theta)$, 事前分布を $p(\theta)$ とする。このとき、事後分布はベイズの定理により $p(\theta | \mathbf{x}) \propto p(\mathbf{x} | \theta) p(\theta)$ と表される。この設定の下で、各推定法におけるパラメータの推定値は次のように定式化できる：

$$\text{最尤推定: } \hat{\theta}_{\text{ML}} = \arg \max_{\theta} p(\mathbf{x} | \theta) = \arg \max_{\theta} \ln p(\mathbf{x} | \theta) \quad (1.30)$$

$$\text{MAP 推定: } \hat{\theta}_{\text{MAP}} = \arg \max_{\theta} p(\theta | \mathbf{x}) = \arg \max_{\theta} [\ln p(\mathbf{x} | \theta) + \ln p(\theta)] \quad (1.31)$$

両者はよく似た形をしているが、最大化する対象 (目的関数) が異なる。最尤推定は尤度そのものを最大化するのに対し、MAP 推定は事前分布を考慮した事後分布を最大化する。なお、最後の式変形では、 θ に依存しない周辺尤度 $p(\mathbf{x})$ が最適化に影響しないことを利用している。

最尤推定に関しては、次のような最小化問題からの導出も可能である。生成モデルの学習における目的は、パラメータ θ を調整して、生成モデルが定める確率密度関数 $p(\mathbf{x} | \theta)$ を、学習データが従う真の分布 $p_{\text{data}}(\mathbf{x})$ に近づけることである。この「近づける」という操作には、両分布間の差異を定量化する指標、すなわち確率分布間の距離 (あるいは不一致度) を定義する必要がある。ここではその尺度として、KL ダイバージェンス (Kullback – Leibler divergence) を用いる：

$$D_{\text{KL}}(p_{\text{data}}(\mathbf{x}) \| p(\mathbf{x} | \theta)) := \int p_{\text{data}}(\mathbf{x}) \ln \frac{p_{\text{data}}(\mathbf{x})}{p(\mathbf{x} | \theta)} d\mathbf{x} \quad (1.32)$$

この量は、真の分布 $p_{\text{data}}(\mathbf{x})$ を基準としたときに、モデル分布 $p(\mathbf{x} | \theta)$ がどれだけ情報的に乖離しているかを測る指標である。すなわち、モデルが生成する分布が、実際のデータ分布からどの程度逸脱しているかを定量化するものである。この KL ダイバージェンスを展開すると、

$$D_{\text{KL}}(p_{\text{data}}(\mathbf{x}) \| p(\mathbf{x} | \theta)) = \int p_{\text{data}}(\mathbf{x}) \ln p_{\text{data}}(\mathbf{x}) d\mathbf{x} - \int p_{\text{data}}(\mathbf{x}) \ln p(\mathbf{x} | \theta) d\mathbf{x} \quad (1.33)$$

$$= \text{const.} - \mathbb{E}_{\mathbf{x} \sim p_{\text{data}}} [\ln p(\mathbf{x} | \theta)] \quad (1.34)$$

となる。ここで第 1 項は θ に依存しない定数であるため、パラメータ θ を最適化する際には、第 2 項の期待値 (すなわち対数尤度の期待値) を最大化することに等しい。したがって、最適なパラメータ θ^* は、

$$\theta^* = \arg \min_{\theta} D_{\text{KL}}(p_{\text{data}}(\mathbf{x}) \| p(\mathbf{x} | \theta)) = \arg \max_{\theta} \mathbb{E}_{\mathbf{x} \sim p_{\text{data}}} [\ln p(\mathbf{x} | \theta)] \quad (1.35)$$

として求められる。しかし実際には、真の分布 $p_{\text{data}}(\mathbf{x})$ の形は不明であり、観測されるのは有限個のデータ点 $\{\mathbf{x}_i\}_{i=1}^n$ のみである。そこで、真の分布の代替として、以下のような経験分布 (empirical distribution) $\hat{p}_{\text{data}}(\mathbf{x})$ を用いる：

$$\hat{p}_{\text{data}}(\mathbf{x}) := \frac{1}{n} \sum_{i=1}^n \delta(\mathbf{x} - \mathbf{x}_i) \quad (1.36)$$

ここで、 $\delta(\cdot)$ は Dirac のデルタ関数であり、この経験分布 $\hat{p}_{\text{data}}(\mathbf{x})$ は、観測された各データ点の位置にのみ確率を集中させるような離散的な点分布として解釈できる。すなわち、サンプル $\{\mathbf{x}_i\}_{i=1}^n$ 以外の点では確率密度がゼロであり、各 \mathbf{x}_i に等しい重み $1/n$ を割り当てているとみなせる。この近似を用いることで、最適化問題は次のように書き換えられる：

$$\theta^* \approx \arg \max_{\theta} \mathbb{E}_{\mathbf{x} \sim \hat{p}_{\text{data}}} [\ln p(\mathbf{x} | \theta)] = \arg \max_{\theta} \sum_{i=1}^n \ln p(\mathbf{x}_i | \theta) \quad (1.37)$$

この最適化は最尤推定と一致する。

一方で、MAP 推定はこの枠組みに事前分布に基づく項を加えたものであり、「正則化付き最尤推定」とみなすことができる。さらに後に述べるように、MAP 推定は変分推論 (variational inference) の特殊な場合 (近似分布をデルタ分布に制限した場合) として位置づけることもできる。本節では扱わず、変分推論の議論の後に改めて説明する。

前項で説明したベイズ線形回帰モデルに対し、パラメータの最尤推定と MAP 推定に基づいた推定を行ってみよう。それぞれの推定における最適化問題は、以下のように表すことができる。

$$\hat{\mathbf{w}}_{\text{ML}} = \arg \max_{\mathbf{w}} p(\mathbf{y} | \mathbf{w}, \mathbf{X}) = \arg \max_{\mathbf{w}} \ln p(\mathbf{y} | \mathbf{w}, \mathbf{X}) \quad (1.38)$$

$$\hat{\mathbf{w}}_{\text{MAP}} = \arg \max_{\mathbf{w}} p(\mathbf{w} | \mathbf{y}, \mathbf{X}) = \arg \max_{\mathbf{w}} [\ln p(\mathbf{y} | \mathbf{w}, \mathbf{X}) + \ln p(\mathbf{w})] \quad (1.39)$$

まず、最尤推定の場合、対数尤度は

$$\ln p(\mathbf{y} | \mathbf{w}, \mathbf{X}) = -\frac{\beta}{2} \|\mathbf{y} - \Phi \mathbf{w}\|^2 + \text{const.} \quad (1.40)$$

であったので、最適化問題は

$$\arg \max_{\mathbf{w}} \ln p(\mathbf{y} | \mathbf{w}, \mathbf{X}) = \arg \min_{\mathbf{w}} \|\mathbf{y} - \Phi \mathbf{w}\|^2 \quad (1.41)$$

となる。これは最小二乗法に一致し、通常の (ベイズでない) 線形回帰モデルと同じ解を与える。

次に、MAP 推定の場合、対数尤度に関しては上と同様であるため、対数事前分布 $\ln p(\mathbf{w})$ ののみを考える。前項では事前分布を正規分布 $p(\mathbf{w}) = \mathcal{N}(\mathbf{w} | \mathbf{0}, \alpha^{-1} \mathbf{I}_{d'})$ と設定したため、

$$\ln p(\mathbf{w}) = -\frac{\alpha}{2} \|\mathbf{w}\|^2 + \text{const.} \quad (1.42)$$

であった。このため、最適化問題は対数尤度と合わせて、

$$\arg \max_{\mathbf{w}} p(\mathbf{w} | \mathbf{y}, \mathbf{X}) = \arg \min_{\mathbf{w}} [\|\mathbf{y} - \Phi \mathbf{w}\|^2 + \lambda \|\mathbf{w}\|^2] \quad (1.43)$$

となる。ただし、 $\lambda = \alpha/\beta (> 0)$ とした。この λ は正則化の強度を表すハイパーパラメータである。よって、この最適化問題は最小二乗法の目的関数に、パラメータの二乗和 (L^2 ノルムの二乗) を正則化項 (罰則項) として追加したものであり、これは Ridge 回帰と一致する。

ここで、事前分布をラプラス分布にする場合を考える。ラプラス分布の確率密度関数は $x \in \mathbb{R}$ について、

$$\text{Lap}(x \mid \mu, b) := \frac{1}{2b} \exp\left(-\frac{|x - \mu|}{b}\right) \quad (1.44)$$

で与えられる。ここで、 $p(w_i) = \text{Lap}(w_i \mid 0, 2\alpha^{-1}) = \frac{\alpha}{4} \exp\left(-\frac{\alpha}{2}|w_i|\right)$ に設定した場合は

$$\ln p(\mathbf{w}) = -\frac{\alpha}{2} \sum_{i=1}^n |w_i| + \text{const.} \quad (1.45)$$

となるので、最適化問題は

$$\arg \max_{\mathbf{w}} p(\mathbf{w} \mid \mathbf{y}, \mathbf{X}) = \arg \min_{\mathbf{w}} \left[\|\mathbf{y} - \Phi \mathbf{w}\|^2 + \lambda \sum_{i=1}^n |w_i| \right] \quad (1.46)$$

と書ける。この最適化問題は最小二乗法の目的関数に、パラメータの絶対値 (L^1 ノルム) を正則化項として追加したものであり、これは Lasso 回帰と一致する。

以上より、最尤推定は最小二乗法に一致し、MAP 推定では事前分布の選択に応じて Ridge 回帰や Lasso 回帰と対応づけられる。すなわち、正則化はパラメータに対する事前分布の設定として解釈できる。

1.1.4 潜在変数モデル

ここで、 $p_{\theta}(\mathbf{x})$ は、観測変数 \mathbf{x} に対する条件付き分布 $p(\mathbf{x} \mid \theta)$ の略記である。
確率的主成分分析

エネルギーベースモデル

1.1.5 階層ベイズモデル

1.2 スパース符号化と予測符号化

1.2.1 スパース符号化

スパース符号化のコード

1.2.2 予測符号化

予測符号化のコード

1.3 近似ベイズ推論

1.3.1 変分推論

MAP 推定は変分推論から導出できる.

1.3.2 自由エネルギー原理

FEP

Active inference

1.3.3 マルコフ連鎖モンテカルロ法

ボルツマンマシン

1.4 神経回路による不確実性の表現

1.4.1 神経サンプリング

1.4.2 確率的集団符号化

参考文献

- Helmholtz, H. von (1867). *Handbuch der physiologischen Optik*. Allgemeine Encyklopädie der Physik, Band IX. Leipzig: Leopold Voss.
- Mumford, D. (1992). “On the computational architecture of the neocortex: II The role of cortico-cortical loops”. *Biological cybernetics* 66.3, pp. 241–251.
- Kawato, M., Hayakawa, H., and Inui, T. (1993). “A forward-inverse optics model of reciprocal connections between visual cortical areas”. *Network: computation in neural systems* 4.4, p. 415.
- Friston, K. (2003). “Learning and inference in the brain”. *Neural Networks* 16.9, pp. 1325–1352.