

機械学習基礎

赤穂 昭太郎

akaho@ism.ac.jp

2025 年 4 月 8 日

目次

第 4 章	ベイズモデリングと機械学習	1
4.1	ベイズ的アプローチの位置づけと有用性	1
4.2	ベイズモデリングの課題と展開	2
4.3	本章の構成	2
4.4	ベイズモデリングの一般的な手順	3
4.5	ベイズ線形回帰の導入	3
4.6	ガウス過程回帰 (Gaussian Process Regression)	5
4.7	ハイパーパラメータの推定と周辺尤度最大化	9
4.8	ナイーブベイズ法	11
4.9	確率的グラフィカルモデルとベイジアンネットワーク	13
4.10	因子グラフと推論アルゴリズムの一般化	16
4.11	変分ベイズ法 (Variational Bayes, VB) の導入	20
4.12	乱数生成法の基礎	27
4.13	MCMC 法の基本原理と代表的アルゴリズム	28
	参考文献	43

第4章

ベイズモデリングと機械学習

4.1 ベイズ的アプローチの位置づけと有用性

機械学習におけるベイズモデリングとは、観測データに対して確率的な生成過程（生成モデル）を仮定し、未知のパラメータや潜在変数に対してベイズの定理を用いて事後分布を推論する枠組みである。

$$\text{事後分布 } p(\theta | \mathcal{D}) = \frac{p(\mathcal{D} | \theta) p(\theta)}{p(\mathcal{D})}$$

ここで、 \mathcal{D} は観測データ、 θ はモデルパラメータ、 $p(\mathcal{D} | \theta)$ は尤度、 $p(\theta)$ は事前分布、 $p(\mathcal{D})$ は周辺尤度（証拠）である。

ベイズ的手法は以下のような有用な性質を持つ：

- **不確実性の定量化**：パラメータ推定だけでなく、予測分布を通じて予測の不確実性を表現できる。
- **事前情報の活用**：既存の知識や仮定を事前分布として組み込むことができる。
- **小規模データへの対応**：データ数が少ない状況でも頑健な推定が可能。
- **過学習の抑制**：事前分布や積分的推論によってモデルの複雑さが自動的に調整される。

4.2 ベイズモデリングの課題と展開

一方で、ベイズモデリングには以下のような課題も存在する：

- **解析的計算の困難さ**：事後分布の積分や正規化定数の計算が困難な場合が多い。
- **スケーラビリティ**：大規模データに対する推論の計算コストが高い。
- **モデル設計の自由度の高さ**：柔軟性が高い分、適切なモデル・事前分布の設計が求められる。

これらの課題に対処するために、近年では以下のようなベイズモデリングの技術が発展している：

- **生成モデルの設計**：観測変数と潜在変数を統一的に扱う表現力の高いモデル（例：混合モデル、潜在変数モデル、ベイズ線形回帰、ガウス過程など）。
- **グラフィカルモデル**：複雑な確率モデルを視覚的・構造的に表現する枠組み。推論アルゴリズム（変分推論、ベイズ推論、メッセージパッシング）と密接に関係。
- **推論アルゴリズムの発展**：マルコフ連鎖モンテカルロ法（MCMC）、変分ベイズ法、確率的推論などの近似的な手法によって実用的な推論が可能に。

4.3 本章の構成

本章では、ベイズ的アプローチに基づく代表的な手法とその応用を順に解説する：

- **ベイズ線形回帰**：予測とパラメータの両方に対する事後分布の導出
- **ガウス過程回帰**：関数空間上でのベイズ推論と非線形回帰への拡張
- **ナイーブベイズ法**：生成モデルに基づいたシンプルな分類器
- **ベイズ的推論アルゴリズム**：変分推論や MCMC など
- **ベイズ的意思決定と最適化**：能動学習、ベイズ最適化など

ベイズモデリングは単なるパラメータ推定にとどまらず、モデリング、推論、意思

決定を統一的に扱う機械学習の強力な枠組みである。

4.4 ベイズモデリングの一般的な手順

ベイズモデリングは、観測データに対して確率的な生成過程を仮定し、その中で未知のパラメータや潜在変数に対してベイズの定理を適用することによって推論を行う枠組みである。

基本的な手順は以下の2段階で構成される：

1. **生成モデルの設計**：観測変数 \mathbf{y} と入力 \mathbf{X} 、および未知パラメータ $\boldsymbol{\theta}$ に関して、次のようなモデルを設計する：

$$p(\mathbf{y}, \boldsymbol{\theta} \mid \mathbf{X}) = p(\mathbf{y} \mid \mathbf{X}, \boldsymbol{\theta}) p(\boldsymbol{\theta})$$

ここで、 $p(\mathbf{y} \mid \mathbf{X}, \boldsymbol{\theta})$ は尤度関数、 $p(\boldsymbol{\theta})$ は事前分布である。

2. **事後分布の推定**：観測データ $\mathcal{D} = (\mathbf{X}, \mathbf{y})$ に基づいて、未知パラメータの事後分布をベイズの定理により求める：

$$p(\boldsymbol{\theta} \mid \mathbf{X}, \mathbf{y}) = \frac{p(\mathbf{y} \mid \mathbf{X}, \boldsymbol{\theta}) p(\boldsymbol{\theta})}{p(\mathbf{y} \mid \mathbf{X})}$$

ここで分母 $p(\mathbf{y} \mid \mathbf{X})$ は周辺尤度 (evidence) であり、

$$p(\mathbf{y} \mid \mathbf{X}) = \int p(\mathbf{y} \mid \mathbf{X}, \boldsymbol{\theta}) p(\boldsymbol{\theta}) d\boldsymbol{\theta}$$

で与えられる。

4.5 ベイズ線形回帰の導入

ベイズモデリングの簡単な例として、線形回帰におけるパラメータ \mathbf{w} を確率的に扱う **ベイズ線形回帰**を考える。

モデルの設定

- 入力 $\mathbf{x} \in \mathbb{R}^d$

- 出力 $y \in \mathbb{R}$
- パラメータ $\mathbf{w} \in \mathbb{R}^d$
- ノイズ $\epsilon \sim \mathcal{N}(0, \sigma^2)$

線形モデルの生成過程は次のように表現される：

$$y = \mathbf{w}^\top \mathbf{x} + \epsilon, \quad \epsilon \sim \mathcal{N}(0, \sigma^2)$$

これにより、出力の条件付き分布（尤度）は

$$p(y \mid \mathbf{x}, \mathbf{w}) = \mathcal{N}(y \mid \mathbf{w}^\top \mathbf{x}, \sigma^2)$$

データ全体に対する尤度

N 個の独立なデータ (\mathbf{x}_i, y_i) に対して、尤度関数は次のようになる：

$$p(\mathbf{y} \mid \mathbf{X}, \mathbf{w}) = \prod_{i=1}^N \mathcal{N}(y_i \mid \mathbf{w}^\top \mathbf{x}_i, \sigma^2) = \mathcal{N}(\mathbf{y} \mid \mathbf{X}\mathbf{w}, \sigma^2 I)$$

ここで、 $\mathbf{X} \in \mathbb{R}^{N \times d}$ はデザイン行列で、 i 行目が \mathbf{x}_i^\top である。

事前分布の設定

パラメータ \mathbf{w} に対してゼロ平均・等方性の正規事前分布を仮定する：

$$p(\mathbf{w}) = \mathcal{N}(\mathbf{w} \mid \mathbf{0}, \tau^2 I)$$

事後分布の導出

ベイズの定理により、事後分布は次のように与えられる：

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

$$\propto \exp \left(-\frac{1}{2\sigma^2} \|\mathbf{y} - \mathbf{X}\mathbf{w}\|^2 - \frac{1}{2\tau^2} \|\mathbf{w}\|^2 \right)$$

これはガウス分布に比例する指数型分布であり、次のように事後分布は明示的なガウス分布となる：

$$p(\mathbf{w} \mid \mathbf{X}, \mathbf{y}) = \mathcal{N}(\mathbf{w} \mid \mathbf{w}_N, \Sigma_N)$$

ただし、

$$\Sigma_N = \left(\frac{1}{\sigma^2} \mathbf{X}^\top \mathbf{X} + \frac{1}{\tau^2} \mathbf{I} \right)^{-1}, \quad \mathbf{w}_N = \frac{1}{\sigma^2} \Sigma_N \mathbf{X}^\top \mathbf{y}$$

リッジ回帰との関係

最尤推定では、 \mathbf{w} を固定パラメータと見なして最小二乗法により推定する。一方、ベイズ回帰では事後分布全体が得られる。

事後分布の平均 \mathbf{w}_N は、次の最適化問題の解にもなっている：

$$\mathbf{w}_N = \arg \min_{\mathbf{w}} \{ \|\mathbf{y} - \mathbf{X}\mathbf{w}\|^2 + \lambda \|\mathbf{w}\|^2 \}, \quad \lambda = \frac{\sigma^2}{\tau^2}$$

これはまさに **リッジ回帰 (L2 正則化付き最小二乗) の解**であり、ベイズ的枠組みにおいてリッジ回帰は事後分布の平均 (MAP 推定) に一致することが分かる。

■**まとめ** ベイズ線形回帰は、線形モデルの係数ベクトルに確率的な不確実性を導入し、予測とパラメータの両方についての事後分布を提供する。事前分布にガウス分布を仮定した場合、解析的に事後分布を求めることが可能であり、その平均はリッジ回帰の解と一致する。

4.6 ガウス過程回帰 (Gaussian Process Regression)

ガウス過程の定義と直感

ガウス過程 (Gaussian Process, GP) とは、任意の有限個の点の集合に対して、その関数値が多変量正規分布に従うような確率過程である。

$$f(\cdot) \sim \mathcal{GP}(m(\cdot), k(\cdot, \cdot))$$

ここで、

- $m(\mathbf{x}) = \mathbb{E}[f(\mathbf{x})]$ は平均関数
- $k(\mathbf{x}, \mathbf{x}') = \text{Cov}(f(\mathbf{x}), f(\mathbf{x}'))$ は共分散関数（カーネル）

直感的には、ガウス過程は「関数空間上の確率分布」と捉えることができ、関数値を確率的に生成する無限次元の多変量正規分布と見なせる。

ベイズ線形回帰との関係

ベイズ線形回帰では、重みベクトル \mathbf{w} にガウス事前分布を置いた。基底関数 $\phi(\mathbf{x})$ を使って $f(\mathbf{x}) = \mathbf{w}^\top \phi(\mathbf{x})$ と表すと、 $f(\mathbf{x})$ 自体がガウス過程になる：

$$f(\mathbf{x}) \sim \mathcal{GP}(0, k(\mathbf{x}, \mathbf{x}')), \quad k(\mathbf{x}, \mathbf{x}') = \phi(\mathbf{x})^\top \phi(\mathbf{x}')$$

つまり、ベイズ線形回帰は特定のカーネル（有限次元内積）に対応するガウス過程の特殊例である。

観測データと予測分布

N 個の観測点 $\mathbf{X} = [\mathbf{x}_1, \dots, \mathbf{x}_N]^\top$ に対して、対応する出力ベクトルを $\mathbf{y} = [y_1, \dots, y_N]^\top$ とし、ノイズ付きモデルを考える：

$$y_i = f(\mathbf{x}_i) + \epsilon_i, \quad \epsilon_i \sim \mathcal{N}(0, \sigma^2)$$

このとき、観測点と新たな入力 \mathbf{x}_* における関数値の同時分布は以下のような多変量正規分布になる：

$$\begin{bmatrix} \mathbf{y} \\ f_* \end{bmatrix} \sim \mathcal{N} \left(\mathbf{0}, \begin{bmatrix} K(\mathbf{X}, \mathbf{X}) + \sigma^2 I & K(\mathbf{X}, \mathbf{x}_*) \\ K(\mathbf{x}_*, \mathbf{X}) & K(\mathbf{x}_*, \mathbf{x}_*) \end{bmatrix} \right)$$

ここで $K(\mathbf{X}, \mathbf{X})$ は観測点間のカーネル行列、 $K(\mathbf{X}, \mathbf{x}_*)$ は観測点と新入力間のカーネルベクトルである。

この共分散構造に基づき、条件付き分布の公式により、 f_* の事後分布は次の正規分布として与えられる：

$$p(f_* \mid \mathbf{x}_*, \mathbf{X}, \mathbf{y}) = \mathcal{N}(\mu_*, \sigma_*^2)$$

$$\mu_* = K(\mathbf{x}_*, X) [K(X, X) + \sigma^2 I]^{-1} \mathbf{y}$$

$$\sigma_*^2 = K(\mathbf{x}_*, \mathbf{x}_*) - K(\mathbf{x}_*, X) [K(X, X) + \sigma^2 I]^{-1} K(X, \mathbf{x}_*)$$

このようにして、新たな入力点における予測平均と不確実性（分散）を解析的に計算することができる。

クリギングとの関係

ガウス過程回帰は、空間統計学で用いられる **クリギング (Kriging)** と同一の数理的枠組みである。特に空間補間問題では、測定点の周囲の値を確率的に補間・予測するためにガウス過程が自然に使われる。

ベイズ最適化との関係

ガウス過程回帰は、関数評価のコストが高い最適化問題において、**獲得関数 (acquisition function)** を用いたベイズ最適化に応用される。予測分布から得られる平均・分散に基づいて「どこを評価すべきか」を決定できることが、ガウス過程の不確実性モデリングの大きな利点となっている。

補足：ガウス過程と再生核ヒルベルト空間 (RKHS) の関係

ガウス過程回帰と RKHS は密接に関連している。特に、ガウス過程の共分散関数 $k(\mathbf{x}, \mathbf{x}')$ が RKHS の再生核となると、ガウス過程の事後平均は RKHS 上の最適化問題の解と一致する。

再生核ヒルベルト空間の定義

RKHS \mathcal{H}_k は、ある正定値対称カーネル関数 $k(\mathbf{x}, \mathbf{x}')$ に対して定まるヒルベルト空間であり、以下の性質を満たす：

- 各点 \mathbf{x} に対して $k(\cdot, \mathbf{x}) \in \mathcal{H}_k$
- 任意の $f \in \mathcal{H}_k$ に対して「再生性」が成り立つ：

$$f(\mathbf{x}) = \langle f, k(\cdot, \mathbf{x}) \rangle_{\mathcal{H}_k}$$

この再生性により、関数値が内積で表現できるため、計算上の取り扱いが容易になる。

ガウス過程の平均関数と RKHS の最適化問題

ガウス過程回帰において、MAP 推定値すなわち予測分布の平均 $\mu_*(\mathbf{x})$ は、次の変分問題の解と一致することが知られている：

$$f^* = \arg \min_{f \in \mathcal{H}_k} \left\{ \sum_{i=1}^n (y_i - f(\mathbf{x}_i))^2 + \lambda \|f\|_{\mathcal{H}_k}^2 \right\}$$

この最適化問題は、RKHS ノルムを正則化項とする Tikhonov 正則化であり、**カーネルリッジ回帰**の変分形式と一致する。

このとき、Representer Theorem により、解 f^* は観測点におけるカーネルの線形結合の形を取る：

$$f^*(\mathbf{x}) = \sum_{i=1}^n \alpha_i k(\mathbf{x}, \mathbf{x}_i)$$

これはガウス過程回帰における事後平均と同一の形である。

共通点のまとめ

- カーネル k は、ガウス過程の共分散関数でもあり、RKHS の再生核でもある。
- ガウス過程回帰の事後平均（MAP）は、RKHS 上の正則化最小二乗問題の解と一致する。
- ガウス過程回帰は「確率的なカーネル回帰」とも言え、RKHS を背景に持つノンパラメトリックベイズ推論の一例である。

4.7 ハイパーパラメータの推定と周辺尤度最大化

ベイズモデリングにおいて、事前分布やカーネル関数にはハイパーパラメータが含まれる。これらの値はモデルの性能に大きな影響を与えるため、適切に設定する必要がある。

ベイズ的な立場では、ハイパーパラメータ θ を固定値とせず、観測データ $\mathcal{D} = (\mathbf{X}, \mathbf{y})$ に対する **周辺尤度 (marginal likelihood)** :

$$p(\mathbf{y} | \mathbf{X}, \theta) = \int p(\mathbf{y} | \mathbf{X}, \mathbf{w}, \theta) p(\mathbf{w} | \theta) d\mathbf{w}$$

を最大化することで、 θ をデータ駆動的に選択することができる。これは **evidence maximization** や **empirical Bayes** と呼ばれる手法である。

ベイズ線形回帰における周辺尤度

ベイズ線形回帰において、次の仮定を置く：

- 事前分布： $\mathbf{w} \sim \mathcal{N}(\mathbf{0}, \tau^2 I)$
- 尤度： $\mathbf{y} | \mathbf{w} \sim \mathcal{N}(\mathbf{X}\mathbf{w}, \sigma^2 I)$

このとき、 \mathbf{w} を積分消去することで、 \mathbf{y} の周辺分布はガウス分布：

$$p(\mathbf{y} | \mathbf{X}) = \mathcal{N}(\mathbf{y} | \mathbf{0}, \mathbf{X}\mathbf{X}^\top \tau^2 + \sigma^2 I)$$

よって、対数周辺尤度は次のように明示的に書ける：

$$\log p(\mathbf{y} | \mathbf{X}) = -\frac{1}{2} \mathbf{y}^\top \Sigma^{-1} \mathbf{y} - \frac{1}{2} \log \det \Sigma - \frac{n}{2} \log 2\pi, \quad \Sigma = \tau^2 \mathbf{X}\mathbf{X}^\top + \sigma^2 I$$

ここで、 τ^2 (事前分布の分散) や σ^2 (ノイズ分散) をハイパーパラメータとして最大化する。

ガウス過程回帰における周辺尤度

ガウス過程回帰では、観測モデルとして

$$y_i = f(\mathbf{x}_i) + \epsilon_i, \quad f \sim \mathcal{GP}(0, k), \quad \epsilon_i \sim \mathcal{N}(0, \sigma^2)$$

を仮定する。カーネル関数 k にもハイパーパラメータ（例：RBF カーネルの幅 ℓ やスケール α ）が含まれる。

このとき、観測データ \mathbf{y} の周辺分布は：

$$p(\mathbf{y} \mid \mathbf{X}) = \mathcal{N}(\mathbf{y} \mid \mathbf{0}, K + \sigma^2 I)$$

ただし $K = K(\mathbf{X}, \mathbf{X})$ はカーネル行列であり、ハイパーパラメータに依存する。

対数周辺尤度は以下の通り：

$$\log p(\mathbf{y} \mid \mathbf{X}) = -\frac{1}{2} \mathbf{y}^\top (K + \sigma^2 I)^{-1} \mathbf{y} - \frac{1}{2} \log \det(K + \sigma^2 I) - \frac{n}{2} \log 2\pi$$

これをハイパーパラメータ（カーネルのパラメータ、ノイズ分散など）で最大化することで、最適なモデルを選択できる。

解釈と利点

このような周辺尤度最大化によるハイパーパラメータ推定には以下の特徴がある：

- **モデル選択基準**：データに適したカーネルや正則化の強さを自動的に選択できる。
- **オーバーフィッティングの抑制**：尤度項とモデルの複雑さ（行列式項）のトレードオフが反映される。
- **ベイズ的解釈**：完全ベイズの代替として、事前・事後の積分の一部を解析的に行う方法として整合的。

特にガウス過程回帰では、カーネルの選択が性能に大きく影響するため、周辺尤度最大化は重要なモデル選択手法である。

4.8 ナイーブベイズ法

ナイーブベイズ法(Naive Bayes classifier)は、生成モデルに基づく単純な確率的分類器であり、クラスラベル $C \in \{1, \dots, K\}$ と観測された特徴ベクトル $\mathbf{x} = (x_1, \dots, x_d)$ の同時分布 $p(C, \mathbf{x})$ を仮定して推論を行う。

モデルの定義

ベイズの定理に基づく分類では、事後確率を最大化するクラスを選択する：

$$\hat{C} = \arg \max_{c \in \{1, \dots, K\}} p(C = c | \mathbf{x}) = \arg \max_c \frac{p(\mathbf{x} | C = c)p(C = c)}{p(\mathbf{x})}$$

ここでナイーブベイズ法では、条件付き確率 $p(\mathbf{x} | C = c)$ に対して、次のような条件付き独立性の仮定を置く：

$$p(\mathbf{x} | C = c) = \prod_{j=1}^d p(x_j | C = c)$$

これにより、クラスごとの特徴分布を次元ごとに分解し、それぞれ独立にモデル化することができる。

パラメータ推定と分類

学習データから以下を推定する：

- クラスの事前分布 $p(C = c)$ (相対頻度など)
- 各特徴のクラス条件付き分布 $p(x_j | C = c)$ (ベルヌーイ分布、ガウス分布など)

新しいデータ点 \mathbf{x} に対しては、以下のルールで分類する：

$$\hat{C} = \arg \max_c p(C = c) \prod_{j=1}^d p(x_j | C = c)$$

長所と短所

■長所

- **学習と推論が非常に高速**：特徴ごとの分布を独立に推定すればよく、計算量が小さい。
- **少ない学習データでも安定**：確率モデルに基づくため、小標本にも比較的強い。
- **理論的に明快**：ベイズ推論の基本形をシンプルに応用。

■短所

- **条件付き独立性の仮定が現実合わないことが多い**：多くの実データでは特徴量間に相関が存在する。
- **スコアは尤度に比例しており、確率の較正が困難**：確率値は過信すべきでない。
- **連続値の扱いには仮定が必要**：多くの場合、ガウス分布などを仮定する必要がある。

条件付き独立性仮定の破れとその影響

ナイーブベイズ法では、クラス条件付きで特徴量が独立であるという仮定を置くが、これは実データでは通常成立しない。にもかかわらず、ナイーブベイズ法が実践的に有効であることがしばしば観察されている。

この現象に対して、以下のような理論的・経験的知見がある：

- **分類には尤度比が関与するため、相関の影響がスコアの順位に現れにくい**ことがある (Domingos and Pazzani, 1997)。
- **特徴量の冗長性が高い場合、相関があっても分類境界にはそれほど影響を与えない**。
- **過剰な表現力を避けることで過学習を防いでいる**という観点から、むしろ仮定の単純さが汎化性能に寄与する場合もある。

そのため、ナイーブベイズ法は「モデルは間違っているが予測は正しい」ケースの典型例とされることが多い。より表現力の高いモデル（例：ベイジアンネットワーク、ツリーベースのモデルなど）との比較や、関連構造を部分的に取り入れる拡張も存在する。

4.9 確率的グラフィカルモデルとベイジアンネットワーク

複数の確率変数間の依存構造を明示的に記述し、効率的な表現と推論を可能にする枠組みが **確率的グラフィカルモデル (probabilistic graphical model)** である。

その中でも、変数間の因果的または生成的構造を **有向非巡回グラフ (DAG: Directed Acyclic Graph)** により記述するものは、**ベイジアンネットワーク (Bayesian network)** や、より広く **有向グラフィカルモデル**、**階層ベイズモデル**とも呼ばれる。

定義と構造

ベイジアンネットワークは、確率変数の集合 $\mathbf{X} = (X_1, \dots, X_n)$ に対して、それぞれの変数の親ノード $\text{Pa}(X_i)$ をグラフ構造で定義し、全体の結合確率分布を以下のように因数分解する：

$$p(\mathbf{X}) = \prod_{i=1}^n p(X_i \mid \text{Pa}(X_i))$$

この表現は、次のような利点を持つ：

- 条件付き独立性が構造的に組み込まれており、表現の冗長性を削減できる。
- 複雑なドメイン知識（例：因果関係や階層構造）をモデルとして明示的に記述可能。
- 一般の連続変数・離散変数・混合型変数にも対応できる（例：階層ベイズモデル）。

高次元の問題への対応と注意点

グラフィカルモデルは、条件付き独立性の構造を利用することで、**高次元の結合確率分布を効率的に表現・推定**できるという大きなメリットがある。

例えば、各変数が少数の親ノードにのみ依存する場合、必要なパラメータ数は指数的から線形または多項式オーダーに削減できる。

■次元の呪いの回避

- 高次元空間での直接的なモデル化は困難（次元の呪い）。
- 条件付き独立性を利用すれば、変数集合を小さな条件付き分布に分解でき、計算や推定が容易になる。
- これにより、少ないデータでも意味のあるモデル化が可能になる。

■モデルの複雑さと限界 ただし、以下のような問題も存在する：

- 親ノードの数が多くなると条件付き分布の次元が上がるため、利点が薄れる。
- 状態数が多い変数が親ノードにあると、指数的な条件付き確率表が必要になることがある。
- モデルの構造が固定でない場合、構造学習の問題も加わり計算が難しくなる。

推論アルゴリズムの概略

ベイジアンネットワークにおける推論（周辺化・条件付き分布の計算）には、構造に依存して様々なアルゴリズムが用いられる：

- **ループのない構造（木やポリツリー）**では、変数消去（variable elimination）やメッセージパッシング（belief propagation）により効率的に推論可能。
- **ループを含む構造（一般の DAG）**では、厳密な推論は指数時間を要するが、変分推論や MCMC（後述）といった近似推論法が活用される。

このように、ベイジアンネットワークは構造的な知識を活用しながら、柔軟かつス

ケーラブルに複雑な確率分布を扱うための基本的な枠組みを提供する。

補足：ベイジアンネットワークにおける「ループ」の異なる意味

ベイジアンネットワークにおいて「ループ (loop)」という用語は、文脈によって異なる意味を持つため、明確な区別が必要である。

■1. 有向閉路 (directed cycle) としてのループ これはグラフ理論的な意味でのループであり、次のような有向辺の列が存在する場合を指す：

$$X_1 \rightarrow X_2 \rightarrow \cdots \rightarrow X_k \rightarrow X_1$$

このような 有向閉路 (サイクル) が存在する場合、グラフは 有向非巡回グラフ (DAG) ではなく、ベイジアンネットワークとしては **well-defined** な結合分布を構成できない。よって、ベイジアンネットワークの構造は常に DAG である必要がある。

■2. 木でないグラフにおける「情報のループ」 一方で、推論アルゴリズムやネットワーク構造の議論において「ループがある」と言った場合、それは以下のような構造を指すことが多い：

- グラフが 木 (tree) 構造や ポリツリー (polytree) 構造でない。
- すなわち、あるノード間に 複数の無向経路 (undirected paths) が存在する。

このような構造では、有向閉路は存在しないものの、無向的な意味で「ループ (閉路)」があるため、推論アルゴリズムにおいては次のような影響が生じる：

- 厳密な推論 (例：変数消去法、Belief Propagation) は指数的な計算量となる可能性がある。
- Belief Propagation はポリツリーでは厳密解だが、ループがある場合は近似になる (Loopy Belief Propagation)。

■まとめ

- ベイジアンネットワークは構造として **有向非巡回グラフ (DAG)** である必要がある。
- 推論において「ループがある」と言うときは、無向的に閉路が存在する木構造でないグラフを指す。
- この意味でのループは、厳密推論の困難さや近似アルゴリズムの必要性をもたらす。

この区別を明確にすることで、グラフ構造・推論アルゴリズム・理論的性質の理解が混乱なく行えるようになる。

4.10 因子グラフと推論アルゴリズムの一般化

Belief Propagation は、ベイジアンネットワークやマルコフランダムフィールドのようなグラフィカルモデルにおいて、木構造であれば効率的に周辺分布を計算できる。

この手法をさらに一般化したのが、**因子グラフ (factor graph)** に基づく **Sum-Product アルゴリズム** である。

因子グラフと Sum-Product アルゴリズム

因子グラフとは、変数ノードと因子ノードの2種類のノードから構成される **双部グラフ (bipartite graph)** である。

■**定義** 与えられた変数集合 $\mathbf{X} = \{X_1, \dots, X_n\}$ に対して、確率分布が以下のよう因子に分解されているとする：

$$p(\mathbf{X}) = \frac{1}{Z} \prod_a f_a(\mathbf{X}_a)$$

ここで f_a は因子 (局所的な関数)、 \mathbf{X}_a は因子 f_a が依存する変数の集合、 Z は正規化定数である。

この構造を因子グラフでは次のように表現する：

- 各変数 X_i を変数ノードとして表現
- 各因子 f_a を因子ノードとして表現

- f_a が変数 X_i に依存する場合、 X_i と f_a の間にエッジを引く

■Sum-Product アルゴリズム (SP アルゴリズム) 因子グラフにおける推論 (周辺分布の計算) は、**Sum-Product アルゴリズム**によって効率的に行える。これは Belief Propagation の因子グラフ版とも言える。

メッセージの送信規則は次の通り：

- 変数ノード i から因子ノード a へのメッセージ：

$$m_{i \rightarrow a}(x_i) = \prod_{b \in \text{nb}(i) \setminus \{a\}} m_{b \rightarrow i}(x_i)$$

- 因子ノード a から変数ノード i へのメッセージ：

$$m_{a \rightarrow i}(x_i) = \sum_{\mathbf{x}_a \setminus x_i} f_a(\mathbf{x}_a) \prod_{j \in \text{nb}(a) \setminus \{i\}} m_{j \rightarrow a}(x_j)$$

ここで $\text{nb}(i)$ はノード i の隣接因子ノードの集合を、 $\text{nb}(a)$ は因子ノード a に隣接する変数ノードの集合を表す。

すべてのメッセージが更新された後、各変数の周辺分布は次のように計算される：

$$b_i(x_i) \propto \prod_{a \in \text{nb}(i)} m_{a \rightarrow i}(x_i)$$

このアルゴリズムは木構造またはループのない因子グラフであれば、厳密な周辺分布を効率的に計算できる。

Junction Tree アルゴリズム：ループを含む場合の厳密推論

ループがある場合には、Belief Propagation や Sum-Product アルゴリズムは近似になる。ループを含むグラフに対しても厳密推論を可能にするのが、**Junction Tree アルゴリズム (結合木アルゴリズム)**である。

■基本的なアイデア 元のグラフから次の変換を行う：

1. グラフの **モラル化 (moralization)**：有向グラフの場合、すべての親ノード間に無向辺を追加し、有向辺をすべて無向に変換。

2. トリアングレーション (chordal 化) : ループが長い部分を分解し、グラフを「三角化」してクリーク構造に変換。
3. クリーク (団) から Junction Tree を構成 : 最大クリーク (完全連結部分グラフ) をノードとする木構造を作成し、クリーク間の共有変数 (separator) を通じてメッセージを送る。

■アルゴリズムの流れ

- 各クリークに確率因子を割り当てて初期化。
- Junction Tree 上で、Sum-Product のようなメッセージパッシングを行い、すべてのクリークが整合的な周辺分布を持つようにする。
- クリーク内での周辺分布を用いて、任意の変数の周辺分布や条件付き分布を計算。

■計算コスト Junction Tree における計算量は、最大クリークサイズに指数的に依存する。したがって、実用上は treewidth が小さいグラフでのみ現実的な方法である。

■まとめ

- Sum-Product アルゴリズムは因子グラフ上の BP として統一的に扱える。
- ループのある場合でも、Junction Tree に変換することで厳密推論が可能。
- 計算量は treewidth に依存し、大規模なループ構造では近似推論が現実的。

補足：Turbo 符号・LDPC 符号の実用通信システムへの応用

Belief Propagation や Sum-Product アルゴリズムは、誤り訂正符号の設計・復号アルゴリズムとして現実の通信システムにおいて広く用いられている。特に、**Turbo 符号**および **LDPC 符号**は、無線通信や衛星通信、Wi-Fi などにおける標準的な誤り訂正方式である。

Turbo 符号の応用

Turbo 符号は 1993 年に提案され、Shannon 限界に迫る性能を持つことが実証された。主に次のような通信規格で採用されている：

- 第 3 世代移動体通信 (3G, UMTS / W-CDMA)
- CDMA2000
- 一部の 4G LTE 制御チャネル
- 宇宙通信 (NASA の深宇宙探査ミッションなど)

Turbo 符号の復号アルゴリズムは、2 つの畳み込み復号器が交互にメッセージを交換する構造を持ち、これはループを持つ因子グラフに対する **Belief Propagation** と見なすことができる。

LDPC 符号の応用

LDPC 符号 (Low-Density Parity-Check code) は、1962 年に Gallager によって提案されたが、1990 年代以降の計算機性能の向上により実用化が進んだ。以下のような最新の通信規格で広く採用されている：

- 第 5 世代移動体通信 (5G NR)
- IEEE 802.11 (Wi-Fi 6 / Wi-Fi 7)
- デジタル衛星放送 (DVB-S2 / DVB-S2X)
- 光ファイバ通信 (10GBASE-T, 100GBASE-R など)
- ストレージインターフェース (NVMe over Fabrics など)

LDPC 符号の復号は、因子グラフ上の変数ノードと検査ノード間でのメッセージ更新により行われ、これは **Sum-Product アルゴリズム** の典型例である。

グラフィカルモデルとの関係と意義

Turbo 符号や LDPC 符号の復号は、いずれも **ループのある因子グラフに対する Loopy Belief Propagation (近似推論)** として解釈できる。

- 実際にはループを含むため厳密な推論ではないが、適切な構造設計により非常に高い性能を実現。
- 推論アルゴリズムの工夫が、シャノン限界に近い誤り訂正性能を支えている。
- 機械学習における確率的推論技術との数学的共通点が多く、理論的・実用的に興味深い接点となっている。

4.11 変分ベイズ法 (Variational Bayes, VB) の導入

複雑なベイズモデルでは、事後分布の厳密な計算は難しい。そのため、**近似推論**の枠組みが用いられる。変分ベイズ法は、事後分布 $p(\mathbf{Z} | \mathbf{X})$ を近似分布 $q(\mathbf{Z})$ によって近似し、その分布の最適化により推論を行う手法である。

一般的な定式化

ベイズ推論における中心的な対象は、観測データ \mathbf{X} に対する事後分布：

$$p(\mathbf{Z} | \mathbf{X}) = \frac{p(\mathbf{X}, \mathbf{Z})}{p(\mathbf{X})}$$

であるが、周辺尤度 $p(\mathbf{X}) = \int p(\mathbf{X}, \mathbf{Z}) d\mathbf{Z}$ の計算が困難な場合が多い。

変分推論では、近似分布 $q(\mathbf{Z})$ を導入し、**変分下限 (Evidence Lower Bound: ELBO)** を最大化する：

$$\log p(\mathbf{X}) \geq \mathbb{E}_{q(\mathbf{Z})} \left[\log \frac{p(\mathbf{X}, \mathbf{Z})}{q(\mathbf{Z})} \right] =: \mathcal{L}(q)$$

この下限は次のように Kullback-Leibler 距離を用いて表現できる：

$$\log p(\mathbf{X}) = \mathcal{L}(q) + \text{KL}(q(\mathbf{Z}) \parallel p(\mathbf{Z} | \mathbf{X}))$$

よって、 $\mathcal{L}(q)$ を最大化することは、 q が $p(\mathbf{Z} | \mathbf{X})$ に近づくことを意味する。

■**変分分布の仮定** 通常、近似分布 $q(\mathbf{Z})$ は計算可能性のために、以下のような独立性仮定 (mean-field 分解) を置く：

$$q(\mathbf{Z}) = \prod_i q_i(Z_i)$$

この仮定により、各因子 q_i の更新式が解析的に導出可能な場合が多い。

Loopy BP との関係

Belief Propagation (BP) は木構造上の厳密推論アルゴリズムであるが、ループのあるグラフに拡張した **Loopy Belief Propagation** は近似推論となる。

興味深いことに、Loopy BP もまた ELBO の最適化問題における **特定の変分分布族 (Bethe 近似)** の最適化とみなすことができ、**変分推論の一種**と解釈される (Yedidia et al., 2001)。これは後の Free Energy 形式でも重要な視点となる。

混合分布モデルへの適用と EM との比較

変分ベイズ法は、混合分布モデル (例：混合ガウス分布) のように、**潜在変数モデル** に対して特に有効である。たとえば混合ガウス分布では、以下のような構造になる：

- 潜在変数 \mathbf{Z} (どの成分に属するかのラベル)
- モデルパラメータ $\boldsymbol{\theta} = \{\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k, \pi_k\}$

変分分布は次のように分解して仮定する：

$$q(\mathbf{Z}, \boldsymbol{\theta}) = q(\mathbf{Z})q(\boldsymbol{\theta})$$

そして、変分下限 $\mathcal{L}(q)$ を最大化する反復更新法 (coordinate ascent) を行う。これは次のように構造的に **EM アルゴリズムと類似**している：

- **E-step (VB 版)** : $q(\mathbf{Z})$ を $q(\boldsymbol{\theta})$ に対して最適化
- **M-step (VB 版)** : $q(\boldsymbol{\theta})$ を $q(\mathbf{Z})$ に対して最適化

この構造は、標準的な EM アルゴリズム：

$$\boldsymbol{\theta}^{(t+1)} = \arg \max_{\boldsymbol{\theta}} \mathbb{E}_{p(\mathbf{Z}|\mathbf{X}, \boldsymbol{\theta}^{(t)})} [\log p(\mathbf{X}, \mathbf{Z} | \boldsymbol{\theta})]$$

と極めて類似している。ただし、VB ではパラメータ θ を点推定するのではなく、分布 $q(\theta)$ を最適化する点が異なる。

■比較のまとめ

- EM は点推定（最尤または MAP）、VB は分布推定（近似ベイズ）。
- 両者とも潜在変数を持つ生成モデルに対して反復最適化を行う。
- VB の方が不確実性を表現でき、過学習を抑制する傾向がある。

変分ベイズによるガウス混合分布の推論（具体的な更新式）

観測データ $\mathbf{X} = \{\mathbf{x}_1, \dots, \mathbf{x}_N\} \subset \mathbb{R}^D$ に対して、 K 成分のガウス混合モデル (GMM) を考える：

$$p(\mathbf{x}_n \mid \boldsymbol{\pi}, \boldsymbol{\mu}_{1:K}, \boldsymbol{\Lambda}_{1:K}) = \sum_{k=1}^K \pi_k \mathcal{N}(\mathbf{x}_n \mid \boldsymbol{\mu}_k, \boldsymbol{\Lambda}_k^{-1})$$

ここで：- π_k ：混合比 ($\sum_k \pi_k = 1$) - $\boldsymbol{\mu}_k$ ：平均ベクトル- $\boldsymbol{\Lambda}_k$ ：精度行列（共分散の逆）

潜在変数 $z_n \in \{1, \dots, K\}$ を導入し、 $z_n = k$ のとき \mathbf{x}_n は第 k 成分から生成されるとする。

事前分布の仮定（共役事前分布）

- $\boldsymbol{\pi} \sim \text{Dir}(\boldsymbol{\alpha}_0)$
- $\boldsymbol{\mu}_k \mid \boldsymbol{\Lambda}_k \sim \mathcal{N}(\mathbf{m}_0, (\beta_0 \boldsymbol{\Lambda}_k)^{-1})$
- $\boldsymbol{\Lambda}_k \sim \text{Wishart}(W_0, \nu_0)$

変分分布の形式（mean-field 分解）

$$q(\mathbf{Z}, \boldsymbol{\pi}, \boldsymbol{\mu}_{1:K}, \boldsymbol{\Lambda}_{1:K}) = \left[\prod_{n=1}^N q(z_n) \right] q(\boldsymbol{\pi}) \prod_{k=1}^K q(\boldsymbol{\mu}_k, \boldsymbol{\Lambda}_k)$$

VB-E ステップ: $q(z_n)$ の更新

各データ点 \mathbf{x}_n に対して、所属確率 (責任度):

$$r_{nk} := q(z_n = k) \propto \exp(\mathbb{E}_q[\log \pi_k] + \mathbb{E}_q[\log \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_k, \boldsymbol{\Lambda}_k^{-1})])$$

具体的には:

$$\log r_{nk} = \mathbb{E}_q[\log \pi_k] - \frac{D}{2} \log(2\pi) + \frac{1}{2} \mathbb{E}_q[\log \det \boldsymbol{\Lambda}_k] - \frac{1}{2} \mathbb{E}_q[(\mathbf{x}_n - \boldsymbol{\mu}_k)^\top \boldsymbol{\Lambda}_k (\mathbf{x}_n - \boldsymbol{\mu}_k)]$$

すべての k に対して $\{r_{nk}\}_k$ を計算し、正規化する。

VB-M ステップ: $q(\boldsymbol{\pi}), q(\boldsymbol{\mu}_k, \boldsymbol{\Lambda}_k)$ の更新

まず、責任度の合計:

$$N_k := \sum_{n=1}^N r_{nk} \quad \bar{\mathbf{x}}_k := \frac{1}{N_k} \sum_{n=1}^N r_{nk} \mathbf{x}_n \quad S_k := \sum_{n=1}^N r_{nk} (\mathbf{x}_n - \bar{\mathbf{x}}_k)(\mathbf{x}_n - \bar{\mathbf{x}}_k)^\top$$

更新式は以下ようになる:

- **混合比の分布 (ディリクレ分布): **

$$q(\boldsymbol{\pi}) = \text{Dir}(\boldsymbol{\alpha}), \quad \alpha_k = \alpha_{0k} + N_k$$

- **平均と精度の結合分布: **

$$q(\boldsymbol{\mu}_k, \boldsymbol{\Lambda}_k) = \mathcal{N}(\boldsymbol{\mu}_k | \mathbf{m}_k, (\beta_k \boldsymbol{\Lambda}_k)^{-1}) \cdot \text{Wishart}(\boldsymbol{\Lambda}_k | W_k, \nu_k)$$

ここで:

$$\beta_k = \beta_0 + N_k \quad \nu_k = \nu_0 + N_k$$

$$\mathbf{m}_k = \frac{\beta_0 \mathbf{m}_0 + N_k \bar{\mathbf{x}}_k}{\beta_k} \quad W_k^{-1} = W_0^{-1} + S_k + \frac{\beta_0 N_k}{\beta_0 + N_k} (\bar{\mathbf{x}}_k - \mathbf{m}_0)(\bar{\mathbf{x}}_k - \mathbf{m}_0)^\top$$

反復更新と収束判定

以上の更新を E/M ステップで交互に繰り返し、変分下限 (ELBO) やパラメータの収束を基準に停止する。

4.11.1 ラプラス近似 (Laplace Approximation)

ラプラス近似は、解析的に扱うのが困難な確率分布を、**多変量ガウス分布によって局所的に近似する方法**である。特にベイズ推論においては、**事後分布や周辺尤度の近似評価**によく用いられる。

1. 基本的な定式化

D 次元の変数 $\boldsymbol{\theta} \in \mathbb{R}^D$ に対する非正規化分布 $p(\boldsymbol{\theta}) \propto \tilde{p}(\boldsymbol{\theta})$ が与えられており、その近似を考える。

このとき、ラプラス近似では、 $\tilde{p}(\boldsymbol{\theta})$ のモード $\hat{\boldsymbol{\theta}}$ のまわりで 2 次のテイラー展開を行い、正規分布で近似する：

$$\begin{aligned} \log \tilde{p}(\boldsymbol{\theta}) &\approx \log \tilde{p}(\hat{\boldsymbol{\theta}}) - \frac{1}{2}(\boldsymbol{\theta} - \hat{\boldsymbol{\theta}})^\top \mathbf{H}(\boldsymbol{\theta} - \hat{\boldsymbol{\theta}}) \\ \Rightarrow \tilde{p}(\boldsymbol{\theta}) &\approx \tilde{p}(\hat{\boldsymbol{\theta}}) \cdot \exp \left\{ -\frac{1}{2}(\boldsymbol{\theta} - \hat{\boldsymbol{\theta}})^\top \mathbf{H}(\boldsymbol{\theta} - \hat{\boldsymbol{\theta}}) \right\} \end{aligned}$$

ここで \mathbf{H} は $\hat{\boldsymbol{\theta}}$ におけるヘッセ行列 (負の 2 階微分)：

$$\mathbf{H} = -\nabla^2 \log \tilde{p}(\boldsymbol{\theta})|_{\boldsymbol{\theta}=\hat{\boldsymbol{\theta}}}$$

結果として、 $p(\boldsymbol{\theta})$ は次のような正規分布で近似される：

$$p(\boldsymbol{\theta}) \approx \mathcal{N}(\boldsymbol{\theta} \mid \hat{\boldsymbol{\theta}}, \mathbf{H}^{-1})$$

2. 周辺尤度 (モデル証拠) の近似

ベイズモデル選択では、事後分布だけでなく、周辺尤度 (evidence)：

$$p(\mathbf{x}) = \int p(\mathbf{x} \mid \boldsymbol{\theta}) p(\boldsymbol{\theta}) d\boldsymbol{\theta}$$

を評価する必要がある。ラプラス近似では、積分の主寄与を与えるモード $\hat{\theta}$ のまわりで近似することで、次の近似式が得られる：

$$\log p(\mathbf{x}) \approx \log p(\mathbf{x} | \hat{\theta}) + \log p(\hat{\theta}) + \frac{D}{2} \log(2\pi) - \frac{1}{2} \log \det \mathbf{H}$$

ここでも \mathbf{H} は $\log p(\mathbf{x}, \theta)$ の $\hat{\theta}$ におけるヘッセ行列の負値。

この式は、**BIC (ベイズ情報量規準)** の導出においても基礎となる。

3. ラプラス近似の特徴と限界

- 局所的な 2 次近似のため、**モードのまわりでは良好な近似**を与える。
- 非対称な分布、多峰性のある分布では近似精度が悪くなることがある。
- 解析的な手計算や、確率変数の周辺化が困難な場面での**簡便な近似法**として有用。

補足：ラプラス近似と他の近似推論法の比較

ラプラス近似は、事後分布や周辺尤度の簡便な近似法として有効であるが、他にもベイズ推論のための代表的手法として次のような方法がある：

- **HMC (Hamiltonian Monte Carlo)**：マルコフ連鎖モンテカルロ法の一つで、サンプリングによって事後分布の情報を直接的に得る。
- **変分ベイズ法 (Variational Bayes, VB)**：近似分布の族を設定し、KL ダイバージェンスを最小化することで事後分布を近似する。

これらの手法はラプラス近似と同様に解析的に扱えないベイズ事後分布を近似・評価する目的で用いられるが、それぞれに特有の利点と限界がある。

■ラプラス近似の特徴

- モード近傍での 2 次近似に基づく。
- 局所的な構造を利用するため、**計算が簡便**で解析的に扱いやすい。
- 事後分布が**単峰で対称に近い**場合には有効。
- 多峰性や強い非対称性がある場合には精度が劣る。

■HMC との比較

- HMC は事後分布の全体構造をサンプリングによって探索するため、多峰性や非対称性にも対応できる。
- 高次元でも効率的にサンプルを生成できる（物理モデルに基づく探索）。
- 反復的なサンプリングに時間がかかり、計算コストは高め。
- 精度の高い信頼区間や周辺化が可能。

■変分推論との比較

- 変分ベイズ法では、近似分布 $q(\theta)$ の族（通常は独立成分の積など）を仮定して、KL ダイバージェンスの最小化により近似する。
- 分布全体を明示的に表現できるため、不確実性の評価がしやすい。
- 局所的モードに依存しないが、近似分布族の選び方に依存。
- 下限（ELBO）の最大化による近似なので、近似精度は下限のタイトさに依存する。

■比較のまとめ

特徴	ラプラス近似	HMC	変分ベイズ法
近似精度	中（単峰）	高（正確）	中～高（族依存）
計算コスト	低	高	中
分布の形状への対応	弱い（対称・単峰）	強い（任意）	中程度（族次第）
周辺尤度の計算	可（解析）	難しい（サンプル必要）	可（ELBO）
事後分布の表現	正規近似	サンプル集合	近似分布明示

■選択指針（状況に応じた使い分け）

- 高速な解析的近似が必要：ラプラス近似
- 精度が重要で複雑な分布に対応したい：HMC
- 中庸な精度と効率を両立したい：変分推論

4.12 乱数生成法の基礎

ベイズ推論や MCMC 法をはじめとする多くの確率的手法では、**特定の分布に従う乱数の生成**が基本的な操作となる。本節では、基本的な乱数生成の方法として、**逆関数法**と**棄却法（受容棄却法）**を紹介する。

1. 逆関数法（inverse transform sampling）

確率密度関数（PDF） $f(x)$ を持つ連続分布において、その累積分布関数（CDF）を $F(x) := \int_{-\infty}^x f(t) dt$ とする。

■基本原理 一様分布 $U \sim \text{Uniform}(0, 1)$ の乱数 u に対して、

$$x = F^{-1}(u)$$

とすれば、 x は分布 $f(x)$ に従う乱数となる。

■例：指数分布の乱数生成 指数分布（パラメータ λ ）の累積分布関数は $F(x) = 1 - e^{-\lambda x}$ 。よって、

$$x = -\frac{1}{\lambda} \log(1 - u) \quad (\text{または } -\frac{1}{\lambda} \log u)$$

は指数分布に従う乱数を与える。

■注意点 逆関数 F^{-1} が解析的に求まらない分布ではこの方法は使えない。

2. 棄却法（acceptance-rejection method）

逆関数法が使えない場合、**提案分布**（proposal distribution） $q(x)$ とその上限を使って目的分布 $f(x)$ から乱数を生成する方法である。

■基本アルゴリズム 以下の手順で乱数 x を得る：

1. $x \sim q(x)$ からサンプルを生成
2. $u \sim \text{Uniform}(0, 1)$ を生成

3. $u < \frac{f(x)}{Mq(x)}$ なら x を受容、そうでなければ棄却（再試行）

ここで $M \geq \sup_x \frac{f(x)}{q(x)}$ は十分に大きな定数。

■性質と注意点

- 提案分布 $q(x)$ は $f(x)$ に近い方が効率が高くなる。
- M が大きすぎると多くのサンプルが棄却され、効率が悪くなる。
- $f(x)$ の正規化定数が不明でも比が計算できれば使用可能。

■例：標準正規分布の近似 標準正規分布に対して、提案分布として一様分布や二項分布などが試されるが、通常は Box-Muller 法や Ziggurat 法などの特殊な手法が使われる（詳細は割愛）。

まとめ

- 逆関数法は累積分布関数の逆写像が解析的に求まる場合に有効。
- 棄却法は柔軟であり、複雑な分布にも適用できるが効率に注意が必要。
- いずれも MCMC 法の構成要素（提案分布や受容ステップ）の基本を理解するために重要。

4.13 MCMC 法の基本原理と代表的アルゴリズム

複雑な確率分布（特に高次元のベイズ事後分布）からのサンプリングを行うために、マルコフ連鎖モンテカルロ法（Markov Chain Monte Carlo, MCMC）が広く用いられている。

MCMC 法は、目的の分布 $p(x)$ を定常分布とするマルコフ連鎖を構成し、その連鎖から得られるサンプル列を用いて、事後分布や期待値の近似を行う。

1. MCMC の基本的な枠組み

MCMC 法の本質は、以下の条件を満たすマルコフ連鎖を設計することにある：

- 定常分布 (stationary distribution) が $p(\mathbf{x})$ に一致する。
- 十分な反復の後、連鎖の分布が $p(\mathbf{x})$ に収束する (エルゴード性)。

これにより、サンプル $\{\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(T)}\}$ を用いて、任意の関数 $f(\mathbf{x})$ に対する期待値 $\mathbb{E}_p[f(\mathbf{x})]$ を次のように近似できる：

$$\mathbb{E}_p[f(\mathbf{x})] \approx \frac{1}{T} \sum_{t=1}^T f(\mathbf{x}^{(t)})$$

2. メトロポリス法 (Metropolis Algorithm)

提案分布 $q(\mathbf{x}' | \mathbf{x})$ に従って候補点 \mathbf{x}' を生成し、以下の確率で現在の状態 \mathbf{x} を \mathbf{x}' に更新する：

$$\alpha(\mathbf{x}, \mathbf{x}') = \min \left(1, \frac{p(\mathbf{x}')}{p(\mathbf{x})} \right)$$

この方法は、**対称な提案分布** ($q(\mathbf{x}' | \mathbf{x}) = q(\mathbf{x} | \mathbf{x}')$) の場合に使える。

3. メトロポリス-ヘイスティングス法 (Metropolis - Hastings Algorithm)

メトロポリス法を非対称な提案分布に拡張したもの。提案分布 $q(\mathbf{x}' | \mathbf{x})$ に基づき、以下の受容確率で状態を更新する：

$$\alpha(\mathbf{x}, \mathbf{x}') = \min \left(1, \frac{p(\mathbf{x}')q(\mathbf{x} | \mathbf{x}')}{p(\mathbf{x})q(\mathbf{x}' | \mathbf{x})} \right)$$

このアルゴリズムは、**詳細釣り合い条件 (detailed balance)** を満たし、 $p(\mathbf{x})$ を定常分布とする。

4. ギブスサンプラー (Gibbs Sampling)

変数を成分ごとに更新する方法。多変量分布 $p(x_1, \dots, x_d)$ に対して、各成分の条件付き分布 $p(x_i | \text{rest})$ から順にサンプリングする：

$$x_1^{(t+1)} \sim p(x_1 | x_2^{(t)}, \dots, x_d^{(t)})$$

$$x_2^{(t+1)} \sim p(x_2 | x_1^{(t+1)}, x_3^{(t)}, \dots) \quad \text{など}$$

この方法は、すべての条件付き分布が簡単にサンプリングできるときに非常に効率的。

5. ハミルトニアンモンテカルロ法 (Hamiltonian Monte Carlo, HMC)

連続値空間で高次元な分布に適した MCMC 手法。状態 \mathbf{x} に対して「運動量変数」 \mathbf{p} を導入し、ハミルトン力学に基づく時間発展を行う：

$$H(\mathbf{x}, \mathbf{p}) = U(\mathbf{x}) + K(\mathbf{p})$$

$$U(\mathbf{x}) = -\log p(\mathbf{x}), \quad K(\mathbf{p}) = \frac{1}{2} \mathbf{p}^\top \mathbf{M}^{-1} \mathbf{p}$$

時間発展は数値的にシミュレーション (Leapfrog 法) し、その結果を Metropolis - Hastings により受容する。大きな空間でも高速かつ効率よく探索可能。

6. MCMC の収束性と定理

MCMC 法の理論的背景として、次のような定理がある：

- **定常分布の存在**：連鎖がエルゴード的であれば、定常分布 $p(\mathbf{x})$ に収束する。
- **エルゴード定理**：長期平均は期待値に収束する：

$$\frac{1}{T} \sum_{t=1}^T f(\mathbf{x}^{(t)}) \xrightarrow{T \rightarrow \infty} \mathbb{E}_p[f(\mathbf{x})]$$

- **詳細釣り合い条件 (detailed balance)** が成立すれば、 $p(\mathbf{x})$ は定常分布。

まとめ

- MCMC は、目的分布を定常分布とするマルコフ連鎖からのサンプルを用いた推論手法。

- メトロポリス法、MH 法、Gibbs 法、HMC など様々なアルゴリズムがある。
- いずれも詳細釣り合いとエルゴード性に基づく収束性が理論的に保証される。

補足：詳細釣り合い条件 (Detailed Balance Condition)

MCMC 法において、マルコフ連鎖の状態遷移確率を $T(\mathbf{x} \rightarrow \mathbf{x}')$ とし、目標とする定常分布を $p(\mathbf{x})$ とする。

このとき、**詳細釣り合い条件 (detailed balance condition)** とは次の式で表される：

$$p(\mathbf{x}) T(\mathbf{x} \rightarrow \mathbf{x}') = p(\mathbf{x}') T(\mathbf{x}' \rightarrow \mathbf{x}) \quad \text{for all } \mathbf{x}, \mathbf{x}'$$

これは、状態 \mathbf{x} から \mathbf{x}' への遷移と、 \mathbf{x}' から \mathbf{x} への遷移が、確率の重み付きでバランスしていることを意味する。

■意味と役割 この条件が成立すると、全体として $p(\mathbf{x})$ は不変となり、連鎖の定常分布となる。すなわち、

$$\sum_{\mathbf{x}} p(\mathbf{x}) T(\mathbf{x} \rightarrow \mathbf{x}') = p(\mathbf{x}')$$

が自動的に成り立つ（これを**定常性条件**という）。

■例：メトロポリス-ヘイスティングス法 MH 法では、提案分布 $q(\mathbf{x}' | \mathbf{x})$ と受容確率 $\alpha(\mathbf{x}, \mathbf{x}')$ により遷移確率が定義される：

$$T(\mathbf{x} \rightarrow \mathbf{x}') = q(\mathbf{x}' | \mathbf{x}) \alpha(\mathbf{x}, \mathbf{x}')$$

受容確率：

$$\alpha(\mathbf{x}, \mathbf{x}') = \min \left(1, \frac{p(\mathbf{x}') q(\mathbf{x} | \mathbf{x}')}{p(\mathbf{x}) q(\mathbf{x}' | \mathbf{x})} \right)$$

とすることで、詳細釣り合い条件を満たすように設計されている。

■まとめ

- 詳細つり合い条件は、定常分布の一 sufficient condition。
- 通常の MCMC 手法（MH 法など）はこの条件を満たすように設計される。
- 定常分布の保証に加え、エルゴード性があれば収束性が理論的に確保される。

4.13.1 能動学習とベイズ最適化：実験計画法の観点からの位置づけ

機械学習において、学習データは与えられるものではなく、自ら選択できる場合がある。このような状況では、限られたデータ取得コストの中で、どのように効率よく情報を集めるかという問題が生じる。

能動学習（Active Learning） および **ベイズ最適化（Bayesian Optimization）** は、こうした文脈において極めて有効な手法であり、ともに「**データを自ら選びながら学習を進める戦略**」であるという共通点をもつ。

この考え方は、統計学における **実験計画法（Design of Experiments, DOE）** の発展形として位置づけることができる。

1. 古典的な実験計画法とその目的

実験計画法は、限られた実験回数の中で、統計的に効率よく因果効果やパラメータ推定を行うための方法論であり、主に次のような目標をもつ：

- 推定量の分散を最小化する（D 最適設計、A 最適設計など）
- パラメータの識別性や非交絡性を確保する
- 実験コストや制約を考慮した設計を行う

これはあくまで**統計モデルの推定精度の向上**が主目的である。

2. 機械学習における能動的なデータ取得戦略

一方、機械学習では、モデルの予測精度や意思決定の最適化を目的として、以下のような戦略的データ選択が行われる：

- **能動学習**：識別問題において、予測が不確かなサンプルを優先的にラベル付けしてもらうことで、効率よく分類性能を向上させる。

- **ベイズ最適化**：高価な評価関数を最小化／最大化するために、少ない試行回数で最適値を探索する。

これらは、統計的最適設計の発想をベースにしつつ、**予測モデルの更新と次の観測点の選択を逐次的に繰り返す**という点で、より動的な手法である。

3. 共通構造とベイズ的枠組みの利点

能動学習やベイズ最適化は、どちらも次のような枠組みを共有する：

- 現時点でのモデル（多くは確率モデル）に基づき、将来の観測の「価値」を定量化する。
- その価値に基づいて、次に取得すべきデータ点（クエリ）を選択する。
- 新しいデータに基づき、モデルを更新し、再び次の選択へ。

このような繰り返しは、**逐次的ベイズ更新と意思決定の結合**という視点から、ベイズモデリングと非常に親和性が高い。

4. 概要のまとめ

- **実験計画法**は、限られた実験資源の中での最適な観測設計を目的とした古典的手法。
- **能動学習**は、分類モデルなどにおける不確実性を利用し、情報量の高いデータを優先的に取得する戦略。
- **ベイズ最適化**は、高価な評価関数の最適化において、獲得関数（acquisition function）を用いて次の探索点を選択する戦略。
- いずれも、**データ選択と学習のループ**を通じて、効率的に性能を向上させる方法であり、ベイズ的な考え方と密接に結びついている。

4.13.2 能動学習の基本とプールベースドな枠組み

1. 能動学習の基本概念

能動学習 (Active Learning) は、限られたラベル付けコストの下で、**効率的に学習データを取得してモデルの性能を高めるための戦略**である。

教師あり学習の前提として、大量のラベル付きデータが必要とされるが、実際にはラベル付けが高価である場面が多い (例：医学画像診断、自然言語アノテーションなど)。

能動学習では、ラベルのないデータから「今ラベル付けすべきサンプル」をモデルが自ら選択する。

2. プールベースドな能動学習

能動学習の代表的な枠組みのひとつが、**プールベースド (pool-based) 能動学習**である。

- 未ラベルのサンプル集合 (プール) $\mathcal{U} = \{x_1, \dots, x_n\}$
- 初期のラベル付き集合 \mathcal{L}
- モデルを \mathcal{L} で学習し、 \mathcal{U} からラベルを追加するクエリ点 x^* を選ぶ

このような選択と学習のサイクルを繰り返すことで、効率的に性能を向上させる。

3. 代表的なサンプル選択戦略

以下はプールベースド能動学習で用いられる代表的なクエリ戦略である。

■(1) 不確実性サンプリング (Uncertainty Sampling) モデルが最も予測に自信を持ってないサンプルを選ぶ戦略。

2 クラス分類の場合、予測確率 $P(y = 1 | x)$ が 0.5 に近いものを優先する：

$$x^* = \arg \min_{x \in \mathcal{U}} |P(y = 1 | x) - 0.5|$$

他の形式として：- 最小信頼度 (Least Confidence) - マージン (2 クラスの予測確率差) 最小- エントロピー最大化： $-\sum_y P(y | x) \log P(y | x)$

■(2) クエリ・バイ・コミッティ (Query by Committee, QBC) モデルの不確実性を、異なる学習器の意見の不一致として捉える戦略。

複数のモデルから構成されるコミッティ $\{h_1, \dots, h_k\}$ を用いて、予測のばらつきが最大のサンプルを選ぶ。

■(3) Expected Model Change / Error Reduction

- あるサンプルをラベル付けしたときに、モデルが大きく変化する、あるいは
- テスト誤差が大きく減少することが期待される点を選ぶ。

理論的には望ましいが、計算コストが高くなるため、近似手法が用いられる。

■(4) 最後に：代表点サンプリング (Diversity Sampling) 未ラベル集合内で「代表的なサンプル」や「密度の高い領域」を優先的に選ぶことで、全体の分布を広くカバーすることを狙う。

4. Exploration – Exploitation のトレードオフ

能動学習における根本的な問題は、次のようなトレードオフである：

- **Exploitation (活用)**：現在のモデルが不確実と判断する点に集中して精度を高める
- **Exploration (探索)**：未知の領域を積極的に探索し、モデルの偏りを是正する

多くの能動学習アルゴリズムはこのバランスを暗に、あるいは明示的にとっている。

- 不確実性サンプリング：exploitation 寄り
- QBC：exploration の要素も含む
- 組合せ戦略：diversity sampling + uncertainty

このトレードオフの制御は、学習の効率や汎化性能に大きく影響を与える重要な要素である。

5. ベイズモデルに基づく能動学習

能動学習においては、「どのデータ点がより情報を与えてくれるか」を判断する必要がある。そのためには、現在のモデルが持つ不確実性（予測分布の幅や分散）を定量化する必要がある、ベイズモデルはこの目的に非常に適している。

■ベイズ的枠組みの特徴

- モデルのパラメータや予測に対して、確率分布を付与することができる
- 不確実性の大きさに応じた「価値に基づくサンプル選択」が可能になる

以下では、代表的な2つのベイズ的手法における能動学習の考え方を紹介する。

(1) ガウス過程回帰における能動学習

回帰問題において、ガウス過程（Gaussian Process, GP）モデルを用いると、ある入力 x に対して次のような予測分布が得られる：

$$f(x) \mid \mathcal{D} \sim \mathcal{N}(\mu(x), \sigma^2(x))$$

ここで：- $\mu(x)$ ：予測平均- $\sigma^2(x)$ ：予測分散（＝不確実性）

このとき、能動学習における選択戦略として次のようなものが考えられる：

- **最大分散サンプリング**： $\sigma^2(x)$ が最大となる点を選ぶ

$$x^* = \arg \max_{x \in \mathcal{U}} \sigma^2(x)$$

- **ベイズ最適化との関係**：目的関数の最小化を狙う場合、獲得関数（acquisition function）に基づいて選択する（例：UCB, EI など）

(2) 確率的分類器を用いた能動学習

分類問題において、ロジスティック回帰やナイーブベイズ分類器のような**確率的分類器**は、入力 x に対して各クラスに属する確率を出力する：

$$P(y \mid x, \mathcal{D})$$

これを活用して、不確実性サンプリングを以下のように定式化できる：

- 最小信頼度：

$$x^* = \arg \min_{x \in \mathcal{U}} \max_y P(y | x)$$

- 最大エントロピー：

$$x^* = \arg \max_{x \in \mathcal{U}} - \sum_y P(y | x) \log P(y | x)$$

■ベイズ分類器を用いた QBC ベイズ的な視点から、クエリ・バイ・コミッティ (QBC) を実装することも可能である。例えば：

- モデルの事後分布からサンプリングされた複数の分類器 (posterior samples) を構成員とする
- 予測のばらつき (分散や投票の不一致) を元にクエリ点を選ぶ

まとめ

- ベイズモデルは不確実性の定量化を自然に行えるため、能動学習と非常に相性が良い。
- ガウス過程回帰では予測分散が明示的に得られるため、分散最大化戦略がとれる。
- 確率的分類器では、予測確率に基づいてエントロピーや信頼度などの情報指標が利用できる。
- モデルの事後分布から多様な予測を得ることで、QBC 型の戦略もベイズ的に構築可能である。

4.13.3 ベイズ最適化の基本と能動学習との関係

ベイズ最適化 (Bayesian Optimization) は、評価コストが高く、導関数が得られないような関数 (ブラックボックス関数) を、できるだけ少ない評価回数で最適化するための戦略である。

$$x^* = \arg \min_{x \in \mathcal{X}} f(x)$$

ただし、 $f(x)$ は明示的な形式を持たず、1 回の評価が非常に高価である（実験、シミュレーション、チューニング等）。

1. ベイズ最適化の基本枠組み

ベイズ最適化は以下のステップで構成される：

1. 事前分布（通常はガウス過程）を用いて、 f に関する予測モデルを構築する。
2. 獲得関数（acquisition function）を定義し、それを最大化する点 x_t を選ぶ：

$$x_t = \arg \max_x a_t(x)$$

3. 選ばれた x_t で $f(x_t)$ を評価し、モデルを更新する。
4. 所望の予算に達するまで手順を繰り返す。

このように、モデルの予測と不確実性を利用して、**探索（exploration）**と**活用（exploitation）**のバランスを取りながら最適化を進める。

2. 能動学習との違いと共通点

- **共通点**：どちらもベイズモデルに基づき、次に観測すべき点を選ぶという逐次的な戦略。
- **主な違い**：
 - － 能動学習：予測モデル全体の精度向上が目的（分類性能など）
 - － ベイズ最適化：目的関数の最小値（または最大値）をいかに効率よく見つけるかが目的

3. 代表的な獲得関数（Acquisition Functions）

ベイズ最適化のコアは、**獲得関数**によって探索点を選ぶ部分である。以下に主要な例を示す。

■(1) 上限信頼境界 (Upper Confidence Bound, UCB)

$$a_{\text{UCB}}(x) = \mu(x) + \beta \cdot \sigma(x)$$

あるいは最大化問題では $\mu(x) - \beta \cdot \sigma(x)$ 。 $\beta > 0$ により探索-活用 のバランスを調整する。

備考： この形式はマルチアームバンディット問題における UCB 戦略と同型であり、後述の理論解析（後悔最小化など）と関連がある。

■(2) 期待改善量 (Expected Improvement, EI) 現在の最良値 f_{best} に対して、改善される期待値：

$$a_{\text{EI}}(x) = \mathbb{E}[\max(f_{\text{best}} - f(x), 0)]$$

ガウス過程予測下では閉形式で計算可能：

$$a_{\text{EI}}(x) = (f_{\text{best}} - \mu(x))\Phi(z) + \sigma(x)\phi(z), \quad z = \frac{f_{\text{best}} - \mu(x)}{\sigma(x)}$$

■(3) 確率的改善 (Probability of Improvement, PI) 最良値より良くなる確率：

$$a_{\text{PI}}(x) = \mathbb{P}(f(x) < f_{\text{best}}) = \Phi\left(\frac{f_{\text{best}} - \mu(x)}{\sigma(x)}\right)$$

探索の度合いを操作するためにスケーリングパラメータ ξ を加えることもある。

■(4) 情報利得 (Entropy Search, Knowledge Gradient など) 獲得関数を「最小値の位置に関する不確実性の減少量」として定式化する。計算は複雑だが、情報理論的に最も望ましい形式とされる。

4. バンディット問題との関連

UCB のような獲得関数は、バンディット問題 (multi-armed bandit) における後悔最小化の枠組みに基づいており、ベイズ最適化もその拡張として捉えることができる。

- バンディット問題：複数の選択肢のうち、報酬が最も大きいものを選ぶ戦略
- ベイズ最適化：連続空間上で同様の「後悔 (regret)」を最小にする選択戦略

理論的には、ベイズ最適化においても逐次選択による後悔の上界を示す研究が行われている（例：GP-UCB アルゴリズム）。

まとめ

- ベイズ最適化は、ベイズモデルを用いてブラックボックス関数を効率的に最適化する手法。
- モデルの予測平均と分散を活用し、獲得関数によって探索点を選択する。
- 代表的な獲得関数には、UCB, EI, PI, Entropy Search などがある。
- UCB は特にバンディット問題と関係が深く、後悔最小化の枠組みと接続される。

GP-UCB と後悔境界 (Regret Bound)

ベイズ最適化における GP-UCB (Srinivas et al., 2010) は、ガウス過程による予測の平均 $\mu_t(x)$ と分散 $\sigma_t(x)$ を用いて、次のような獲得関数を定義する：

$$x_t = \arg \max_{x \in \mathcal{X}} \left(\mu_{t-1}(x) + \beta_t^{1/2} \cdot \sigma_{t-1}(x) \right)$$

ここで： $\mu_{t-1}(x), \sigma_{t-1}(x)$ は時刻 $t-1$ におけるガウス過程の予測平均・標準偏差- β_t は時間依存の探索係数（理論的な制御パラメータ）

この選択戦略は、**上限信頼区間**を最大化する点を逐次的に評価していく。

■**累積後悔 (Cumulative Regret)** 最適値 $x^* = \arg \min f(x)$ に対して、時刻 t に選択した点の後悔は：

$$r_t := f(x_t) - f(x^*)$$

累積後悔は次のように定義される：

$$R_T := \sum_{t=1}^T r_t$$

■**理論的な上界** 適切な条件の下で、GP-UCB は高い確率で次のような後悔境界を満たすことが示されている：

$$R_T = O\left(\sqrt{T \cdot \gamma_T \cdot \beta_T}\right)$$

ここで：- γ_T は T 回の観測におけるガウス過程の情報ゲイン (information gain) であり、カーネル関数と探索空間に依存する。- β_T は探索-活用 のバランスを調整する係数 (通常は対数的に増加)

■情報ゲイン γ_T の例 カーネルが RBF (Gaussian) である場合、 $\gamma_T = O((\log T)^{d+1})$ のような増加をする。したがって：

$$R_T = O\left(\sqrt{T \cdot (\log T)^{d+1}}\right)$$

これは、ガウス過程を用いたベイズ最適化が理論的にも **サブ線形の後悔 (no-regret)** を達成できる ことを意味する。

まとめ

- GP-UCB は探索 (分散) と活用 (平均) のバランスをとる戦略。
- カーネルに応じた情報ゲイン γ_T により後悔の速度が決まる。
- 適切なパラメータ設定により、サブ線形後悔 $R_T = o(T)$ が保証される。

参考文献

- [1] G. James, D. Witten, T. Hastie, R. Tibshirani “An Introduction to Statistical Learning, Second Edition”, Springer 2023
- [2] C.M. Bishop, “Pattern recognition and machine learning”, Springer 2006 (ビショップ：パターン認識と機械学習（上下），丸善）
- [3] 萩原，入門 統計的回帰とモデル選択，共立出版 2022
- [4] 赤穂，カーネル多変量解析，岩波書店 2008
- [5] 青嶋，矢田，高次元の統計学，共立出版 2019
- [6] S. Watanabe, M. Opper, Asymptotic equivalence of Bayes cross validation and widely applicable information criterion in singular learning theory. Journal of machine learning research, 11(12) 2010