

ML 輪講: 19 章 近似推論

杉浦 圭祐

慶應義塾大学理工学部情報工学科 松谷研究室

April 12, 2019

① K-Means 法

② 混合ガウス分布

③ EM アルゴリズム

- EM アルゴリズムの解釈
- 混合ガウス分布の再解釈
- K-Means 法との関連
- 一般の EM アルゴリズム
- MAP 推定に対する EM アルゴリズム
- EM アルゴリズムの拡張
- EM アルゴリズムに関する補足
- 逐次型の EM アルゴリズム

④ 変分の導入

- EM アルゴリズムが困難な場合
- 変分法
- 変分法で解ける問題の例

- 変分法のまとめ

5 近似推論法

- 変分推論
- 座標降下法
- MAP 推定および最尤推定と変分推論
- 変分近似による弊害
- 離散潜在変数の変分推論
- 連続潜在変数の変分推論
- 正則化との関わり
- 雑多な話題

6 変分自己符号化器

- 生成モデル
- 変分自己符号化器 (VAE) の概要
- 変分自己符号化器 (VAE) の理論

目次

- 1 K-Means 法
- 2 混合ガウス分布
- 3 EM アルゴリズム
- 4 変分の導入
- 5 近似推論法
- 6 変分自己符号化器

K-Means 法によるクラスタリング

- 扱う問題

- 多次元空間上のデータ点集合を考える
- 各データが属するクラスタを決定する問題を考える

- 問題設定

- D 次元ユークリッド空間における、確率変数 x を観測
- x の N 個の観測点で構成されるデータ集合 $\mathcal{D} = \{x_1, \dots, x_N\}$
($x_i \in \mathbb{R}^D$)
- データ集合 $\mathcal{D} = \{x_1, \dots, x_N\}$ を、 K 個のクラスタに分割
- K は、既知の定数であるとする

- クラスタとは

- クラスタとは、簡単に言えば近接するデータの集合である
- クラスタの内部のデータ点間の距離が、クラスタの外側のデータとの距離と比べて、小さいようなデータの集合

K-Means 法によるクラスタリング

- クラスタを代表するベクトルの表現

- 各クラスタを代表する K 個の D 次元ベクトル $\mathcal{M} = \{\mu_1, \dots, \mu_K\}$ ($\mu_k \in \mathbb{R}^D$) を導入する
 - これらのベクトル $\mathcal{M} = \{\mu_k\}$ を、**プロトタイプ**という
 - k 番目のクラスタ (μ_k が支配するクラスタ) を $\mathcal{M}(\mu_k)$ と記す
 - ベクトル μ_k は、 k 番目のクラスタに対応するプロトタイプである
 - μ_k は、 $\mathcal{M}(\mu_k)$ に属するデータ点の平均、即ち k 番目のクラスタの中心である
-
- 解くべき問題
 - N 個の全データ点を、うまくクラスタに割り振る
 - 各データ点から、対応する (そのデータ点が属するクラスタの) プロトタイプ μ_k への、二乗距離の総和を最小化する

K-Means 法によるクラスタリング

- データ点のクラスタへの割り当てを表す変数
 - 各データ x_i ($1 \leq i \leq N$) に対して、二値の指示変数 $r_{ik} \in \{0, 1\}$ ($k = 1, \dots, K$) を定める
 - r_{ik} は、データ x_i が、 K 個あるクラスタのうちの、どれに割り当てられるのかを表す
 - データ点 x_i がクラスタ k に割り当てられるときに、 $r_{ik} = 1$ となり、 $j \neq k$ について $r_{ij} = 0$ である (1-of-K 符号化法という)

$$r_{ik} = \begin{cases} 1 & (x_i \in \mathcal{M}(\mu_k) \text{ の場合}) \\ 0 & \text{それ以外} \end{cases} \quad (1)$$

K-Means 法によるクラスタリング

- 目的関数の定義

- 目的関数を以下のように定義する
- 各データ点 x_i と、 x_i に割り当てたベクトル μ_k (x_i が属するクラスターのプロトタイプ) との、二乗距離の総和

$$J = \sum_{i=1}^N \sum_{k=1}^K r_{ik} \|x_i - \mu_k\|^2 \quad (2)$$

- 上式の J を最小化するような、 $\{r_{ik}\}$ と $\{\mu_k\}$ を求めるのが目標

K-Means 法によるクラスタリング

- 目的関数 J の $\{r_{ik}\}$ に関する最小化
 - 目的関数の式は、各 i について独立である

$$J = \sum_i \sum_k r_{ik} \|\mathbf{x}_i - \boldsymbol{\mu}_k\|^2 \quad (3)$$

- $\{r_{ik}\}$ を決定する方法は簡単
- 各 i について、 $\|\mathbf{x}_i - \boldsymbol{\mu}_k\|^2$ が最小となるような k に対し、 $r_{ik} = 1$ とする
- それ以外のクラスタ $j \neq k$ については、 $r_{ik} = 0$ とする

$$r_{ik} = \begin{cases} 1 & (k = \arg \min_j \|\mathbf{x}_i - \boldsymbol{\mu}_j\|^2 \text{ のとき}) \\ 0 & \text{それ以外} \end{cases} \quad (4)$$

- 各データ点 \mathbf{x}_i を、 \mathbf{x}_i と最も近い $\boldsymbol{\mu}_k$ に割り当てることに相当
- 各データ点 \mathbf{x}_i を、クラスタを代表するベクトル (**クラスタの中心**) との二乗距離が最小になるような、クラスタに割り当てる

K-Means 法によるクラスタリング

- 目的関数 J の $\{\mu_k\}$ に関する最小化
 - J を μ_k について偏微分して 0 とおく

$$\begin{aligned}\frac{\partial}{\partial \mu_k} J &= \frac{\partial}{\partial \mu_k} \sum_i r_{ik} \|\mathbf{x}_i - \mu_k\|^2 \\&= \frac{\partial}{\partial \mu_k} \sum_i r_{ik} (\mathbf{x}_i - \mu_k)^T (\mathbf{x}_i - \mu_k) \\&= \sum_i r_{ik} \frac{\partial}{\partial \mu_k} (\mathbf{x}_i - \mu_k)^T (\mathbf{x}_i - \mu_k) \\&= \sum_i r_{ik} \frac{\partial}{\partial \mu_k} (\mathbf{x}_i^T \mathbf{x}_i - \mathbf{x}_i^T \mu_k - \mu_k^T \mathbf{x}_i + \mu_k^T \mu_k) \\&= \sum_i r_{ik} \frac{\partial}{\partial \mu_k} (\mathbf{x}_i^T \mathbf{x}_i - \mathbf{x}_i^T \mu_k - \mathbf{x}_i^T \mu_k + \mu_k^T \mu_k) \\&= \sum_i r_{ik} (2\mu_k - 2\mathbf{x}_i) = 0\end{aligned}\tag{5}$$

K-Means 法によるクラスタリング

- ここで、次の関係を用いている

$$\begin{aligned} & \boldsymbol{\mu}_k^T \mathbf{x}_i \\ = & (\boldsymbol{\mu}_k^T \mathbf{x}_i)^T \quad (\because \text{スカラーであるため転置してもよい}) \\ = & \mathbf{x}_i^T (\boldsymbol{\mu}_k^T)^T \quad (\because (\mathbf{a}\mathbf{b})^T = \mathbf{b}^T \mathbf{a}^T) \\ = & \mathbf{x}_i^T \boldsymbol{\mu}_k \end{aligned}$$

$$\frac{\partial}{\partial \mathbf{x}} \mathbf{x}^T \mathbf{a} = \frac{\partial}{\partial \mathbf{x}} \mathbf{a}^T \mathbf{x} = \mathbf{a} \quad (6)$$

$$\frac{\partial}{\partial \mathbf{x}} \mathbf{x}^T \mathbf{B} \mathbf{x} = (\mathbf{B} + \mathbf{B}^T) \mathbf{x} \quad (2 \text{ 次形式}) \quad (7)$$

$$\frac{\partial}{\partial \mathbf{x}} \mathbf{x}^T \mathbf{x} = 2\mathbf{x} \quad (8)$$

$$\because \frac{\partial}{\partial \mathbf{x}} \mathbf{x}^T \mathbf{x} = \frac{\partial}{\partial \mathbf{x}} \mathbf{x}^T \mathbf{I} \mathbf{x} = (\mathbf{I} + \mathbf{I}^T) \mathbf{x} = 2\mathbf{I} \mathbf{x} = 2\mathbf{x}$$

K-Means 法によるクラスタリング

- これより、 μ_k について解くと次のようになる

$$\begin{aligned}\sum_i 2r_{ik}(\mu_k - \mathbf{x}_i) &= 0 \\ \sum_i r_{ik}(\mu_k - \mathbf{x}_i) &= 0 \\ \sum_i r_{ik}\mu_k &= \sum_i r_{ik}\mathbf{x}_i \\ \mu_k \left(\sum_i r_{ik} \right) &= \sum_i r_{ik}\mathbf{x}_i \\ \mu_k &= \frac{1}{\sum_i r_{ik}} \sum_i r_{ik}\mathbf{x}_i\end{aligned}\tag{9}$$

- $\sum_i r_{ik}$ は、クラスタ k に属するデータの個数である
- $\sum_i r_{ik} = N_k$ と表すことがある

K-Means 法によるクラスタリング

- r_{ik} は、 i 番目のデータ x_i が、クラスタ k に割り当てられているときのみ、1 となる
- $\sum_i r_{ik} x_i$ は、クラスタ k に属しているデータのベクトル x_i の合計である
- 従って、 μ_k は、 k 番目のクラスタに割り当てられた、全てのデータ点 x_i の平均値である
- この意味で、 μ_k のことを、クラスタ k の平均ベクトル、重心、セントロイドということもある

K-Means 法によるクラスタリング

- $\{\mu_k\}$ と $\{r_{ik}\}$ についての最適化

- 目的関数 J を、 $\{\mu_k\}$ と $\{r_{ik}\}$ について最小化する式は次のようになる

$$\begin{aligned}\mu_k &= \frac{1}{\sum_i r_{ik}} \sum_i r_{ik} \mathbf{x}_i \\ r_{ik} &= \begin{cases} 1 & (k = \arg \min_j \|\mathbf{x}_i - \mu_j\|^2 \text{ のとき}) \\ 0 & \text{それ以外} \end{cases}\end{aligned}$$

- μ_k を計算する式の中に r_{ik} が、 r_{ik} を計算する式の中に μ_k が入っている
- μ_k を求めるためには r_{ik} が、 r_{ik} を求めるためには μ_k が既知でなければならない
- 従って、 μ_k と r_{ik} の両方を同時に最適化することはできない
- どうすれば目的関数 J を、パラメータ $\{\mu_k\}$ と $\{r_{ik}\}$ の両方について最小化できるか?

K-Means 法によるクラスタリング

- $\{\mu_k\}$ と $\{r_{ik}\}$ についての最適化
 - μ_k と r_{ik} の両方を同時に最適化することはできない
 - μ_k と r_{ik} を交互に最適化すればよい
 - μ_k と r_{ik} のそれぞれを最適化する、2つのステップを交互に繰り返す
- μ_k と r_{ik} の最適化は、次のように行うことができる

1 $\mathcal{M} = \{\mu_k\}$ の初期値を設定

- N 個のデータ $\mathcal{D} = \{x_i\}$ を、ランダムなクラスタに割り振って、各クラスタの平均ベクトル $\{\mu_k\}$ を求める
- データ \mathcal{D} からランダムに選択した K 個のデータ点を、各クラスタの中心 μ_k ($k = 1, \dots, K$) とすることもできる

K-Means 法によるクラスタリング

- 2 第 1 フェーズでは、 μ_k を固定しつつ、 r_{ik} について J を最小化

$$r_{ik} = \begin{cases} 1 & (k = \arg \min_j ||\mathbf{x}_i - \mu_j||^2 \text{ のとき}) \\ 0 & \text{それ以外} \end{cases}$$

- 3 第 2 フェーズでは、 r_{ik} を固定しつつ、 μ_k について J を最小化

$$\mu_k = \frac{1}{\sum_i r_{ik}} \sum_i r_{ik} \mathbf{x}_i$$

- 4 (2) と (3) を、 μ_k と r_{ik} が収束するまで繰り返す

- 上記の 2 つのステップが、後述する EM アルゴリズムの E(Expectation) ステップと M(Maximization) ステップに対応する

K-Means 法によるクラスタリング

- データ点へのクラスタの再割り当てと、クラスタの平均ベクトルの再計算
- この2段階の処理を、クラスタの再割り当てが起こらなくなるまで(2段階の処理を行っても、データが属するクラスタが変化しなくなるまで)繰り返す
- 各フェーズで J の値が減少するので、アルゴリズムの収束が保証される
- 大域的最適解ではなく、局所最適解に収束する可能性はある

K-Means 法によるクラスタリング

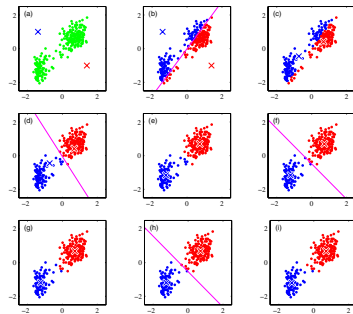


Figure 9.1 Illustration of the K -means algorithm using the re-scaled Old Faithful data set. (a) Green points denote the data set in a two-dimensional Euclidean space. The initial choices for centres μ_1 and μ_2 are shown by the red and blue crosses, respectively. (b) In the initial E step, each data point is assigned either to the red cluster or to the blue cluster, according to which cluster centre is nearer. This is equivalent to classifying the points according to which side of the perpendicular bisector of the two cluster centres, shown by the magenta line, they lie on. (c) In the subsequent M step, each cluster centre is re-computed to be the mean of the points assigned to the corresponding cluster. (d)–(i) show successive E and M steps through to final convergence of the algorithm.

図 1: K-Means アルゴリズムの動作

K-Means 法によるクラスタリング

Figure 9.2 Plot of the cost function J given by (9.1) after each E step (blue points) and M step (red points) of the K -means algorithm for the example shown in Figure 9.1. The algorithm has converged after the third M step, and the final EM cycle produces no changes in either the assignments or the prototype vectors.

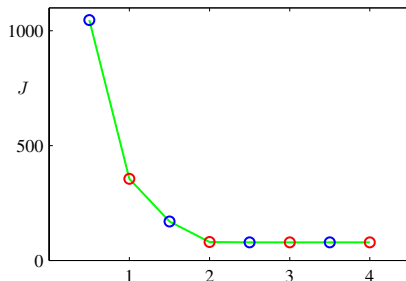


図 2: 目的関数 J の値の遷移

K-Means 法によるクラスタリング

- 素朴な実装では速度が遅いことがある
 - $\{r_{ik}\}$ の更新 (E ステップ) において、各データ点と、各平均ベクトルの全ての組み合わせ間の距離を計算する必要がある

$$r_{ik} = \begin{cases} 1 & (k = \arg \min_j \|\mathbf{x}_i - \boldsymbol{\mu}_j\|^2 \text{ のとき}) \\ 0 & \text{それ以外} \end{cases}$$

- K-Means の高速化について、これまでに様々な手法が提案されてきた
- 近接するデータ点同士が、同一の部分木に属するような木構造を採用する方法
- 距離の三角不等式を利用して、不必要な距離計算を避ける方法

逐次版の K-Means 法

- 逐次版の K-Means 法の導出

- これまでに紹介した K-Means 法は、利用する全てのデータ $\mathcal{D} = \{x_i\}$ が、最初から用意されていることが前提であった
- ここでは、Robbins-Monro 法を使って、オンラインのアルゴリズムを導出する

- Robbins-Monro 法

- 逐次学習アルゴリズムを導出するための手法
- 関数 $f(\theta^*) = 0$ を満たす根 θ^* を、逐次的に計算するための式を与える
- 同時分布 $p(z, \theta)$ に従う確率変数のペア θ, z について、関数 $f(\theta)$ は、 θ が与えられたときの z の条件付き期待値 $\mathbb{E}[z|\theta]$ として、定義される

$$f(\theta) \equiv \mathbb{E}[z|\theta] = \int zp(z|\theta)dz \quad (10)$$

- このとき、根 θ^* の逐次計算式は、次のように記述される

$$\theta^{(N)} = \theta^{(N-1)} - \eta_{N-1}z(\theta^{(N-1)}) \quad (11)$$

- η は学習率
- $z(\theta^{(N)})$ は、確率変数 θ が値 $\theta^{(N)}$ をとるときに観測される z の値

- Robbins-Monro 法を適用するための条件

- Robbins-Monro 法を使用するためには、満たすべき条件が幾つか存在する

1 z の条件付き分散 $\mathbb{E}[(z - f)^2|\theta]$ が、有限でなければならない

$$\mathbb{E}[(z - f)^2|\theta] = \int (z - f(\theta))^2 p(z|\theta) dz < \infty \quad (12)$$

2 $\theta > \theta^*$ では $f(\theta) > 0$ 、 $\theta < \theta^*$ では $f(\theta) < 0$ を仮定 (このように仮定しても一般性は失われない)

3 学習率の系列 $\{\eta_N\}$ は次の条件を満たす

$$\lim_{N \rightarrow \infty} \eta_N = 0 \quad (13)$$

$$\sum_{N=1}^{\infty} \eta_N = \infty \quad (14)$$

$$\sum_{N=1}^{\infty} \eta_N^2 < \infty \quad (15)$$

- 最初の条件は、推定系列 $\theta^{(N)}$ が、目標の根 θ^* に収束していくように、 θ の修正量を減らしていくことを保証
- 次の条件は、根 θ^* 以外の値に収束しないことを保証
- 最後の条件は、分散を有限に抑えることで、いつまで経っても収束しないことを防止

Figure 2.10 A schematic illustration of two correlated random variables z and θ , together with the regression function $f(\theta)$ given by the conditional expectation $\mathbb{E}[z|\theta]$. The Robbins-Monro algorithm provides a general sequential procedure for finding the root θ^* of such functions.

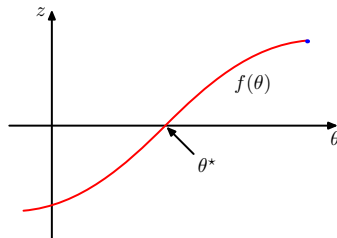


図 3: $f(\theta)$ のグラフ

逐次版の K-Means 法

● 逐次版の K-Means 法の導出

- 先程のパラメータ θ は、 $\{r_{ik}\}$ と $\{\mu_k\}$ である
- パラメータの最適解 $\{\mu_k^*\}$ は、以下の式を満たす

$$\left. \frac{\partial}{\partial \mu_k} J \right|_{\mu_k^*} = \left. \frac{\partial}{\partial \mu_k} \sum_{i=1}^N r_{ik} \|\mathbf{x}_i - \mu_k\|^2 \right|_{\mu_k^*} = 0 \quad (16)$$

これは以下の式と等価である ($N_k = \sum_i r_{ik}$)

$$\left. \frac{\partial}{\partial \mu_k} \frac{1}{N_k} \sum_{i=1}^N r_{ik} \|\mathbf{x}_i - \mu_k\|^2 \right|_{\mu_k^*} = 0 \quad (17)$$

これ以降、クラス k に属するデータのみを考えることにして、これを \mathbf{x}_j と書く ($j = 1, \dots, N_k$)

逐次版の K-Means 法

- このとき、上式から r_{ik} を省略でき、次のようになる

$$\left. \frac{\partial}{\partial \boldsymbol{\mu}_k} \frac{1}{N_k} \sum_{j=1}^{N_k} \|\mathbf{x}_j - \boldsymbol{\mu}_k\|^2 \right|_{\boldsymbol{\mu}_k^*} = 0 \quad (18)$$

- $N_k \rightarrow \infty$ の極限を取って次のように変形する

$$\begin{aligned} & \lim_{N_k \rightarrow \infty} \frac{\partial}{\partial \boldsymbol{\mu}_k} \frac{1}{N_k} \sum_j \|\mathbf{x}_j - \boldsymbol{\mu}_k\|^2 \\ &= \lim_{N_k \rightarrow \infty} \frac{1}{N_k} \frac{\partial}{\partial \boldsymbol{\mu}_k} \sum_j \|\mathbf{x}_j - \boldsymbol{\mu}_k\|^2 \\ &\simeq \mathbb{E} \left[\frac{\partial}{\partial \boldsymbol{\mu}_k} \|\mathbf{x} - \boldsymbol{\mu}_k\|^2 \middle| \boldsymbol{\mu}_k \right] \quad (\mathbf{x} \text{ はクラス } k \text{ に属する}) \\ &= \mathbb{E} [2(\mathbf{x} - \boldsymbol{\mu}_k) | \boldsymbol{\mu}_k] = f(\boldsymbol{\mu}_k) \end{aligned} \quad (19)$$

逐次版の K-Means 法

- これより、K-Means 法は、関数 $f(\mu_k^*) = 0$ をみたす根 μ_k^* を求める問題に帰着させられる
- 従って、Robbins-Monro 法を適用すると、 μ_k の逐次更新式は、以下のようになる

$$\mu_k^{\text{new}} = \mu_k^{\text{old}} - \eta_j(x_j - \mu_k^{\text{old}}) \quad (x_j \text{はクラス } k \text{ に属する}) \quad (20)$$

- η_j は学習率パラメータであり、一般に、 j の増加に伴って単調減少するように設定される
- クラス k に属するデータ x_j が 1 つずつ到着するときの、 μ_k のオンライン更新式が得られた

K-Means アルゴリズムの特徴

- K-Means アルゴリズムの特徴

- 各データを、**たった 1 つ**のクラスタにのみ割り当てる (**ハード割り当て**)
- あるクラスタの中心ベクトル μ_k に非常に近いデータ点があれば、複数のクラスタの中間領域にあるようなデータ点も存在する
- データ x_i が、クラスタ k に属するという結果を得たときに、そのクラスタに属することがほぼ確実なのか、それとも他のクラスタに割り振っても大差はないのかが、区別できない
- 後者のようなデータ点の場合、単一のクラスタへのハード割り当ては最適でない (不正確) かもしれない
- データ点を単一のクラスタに割り当てるのではなく、**各クラスタに属する確率**を、計算できれば良さそう
- 各クラスタへの割り当ての不明瞭さを反映できる
- このように曖昧さを含んだ割り当てを、**ソフト割り当て**という

K-Means アルゴリズムのまとめ

- K-Means アルゴリズムの目的

- 各データが属するクラスタを決定する (ハード割り当て)

- K-Means アルゴリズムで行っていること

- 各クラスタに対して、中心となるベクトル μ_k を考えた
- μ_k は、対応するクラスタに属する、全てのデータベクトルの平均として得られた
- 各データ点は、それと最も距離が近い μ_k に対応するクラスタ k に割り当てた
- 各データ点が属するクラスタを、 r_{ik} という変数 (1-of-K 符号化法) で表した

K-Means アルゴリズムのまとめ

- ここまでの話の流れ

- クラスタリングの問題を、目的関数 J の最小化として表現できた
- 目的関数 J を、 μ_k と r_{ik} の両方について一度に最適化することはできなかった
- そこで J を、 μ_k と r_{ik} について交互に最適化することを考えた
- 交互に行う最適化は、後ほど説明する EM アルゴリズムの、E ステップと M ステップに対応していた
- 逐次版の K-Means 法を、Robbins-Monro 法を使って導出した

目次

- 1 K-Means 法
- 2 混合ガウス分布
- 3 EM アルゴリズム
- 4 変分の導入
- 5 近似推論法
- 6 変分自己符号化器

混合ガウス分布

- 混合ガウス分布を導入する理由

- 曖昧さを含んだクラスタリング (ソフト割り当て) を実現するため
- 言い換えると、データに対して、各クラスタに属する確率が分かるようにするため

- 混合ガウス分布とは

- 各ガウス分布の線形の重ね合わせ

$$p(\boldsymbol{x}) = \sum_{k=1}^K \pi_k \mathcal{N}(\boldsymbol{x} | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k) \quad (21)$$

- 各ガウス分布 $\mathcal{N}(\boldsymbol{x} | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)$ は混合要素とよばれる
- 各ガウス分布は個別に、平均 $\boldsymbol{\mu}_k$ と共分散 $\boldsymbol{\Sigma}_k$ のパラメータをもつ

混合ガウス分布

- パラメータ π_k を **混合係数** といい、以下の条件を満たす

$$\sum_k \pi_k = 1 \quad (22)$$

これは、 $p(x)$ を x について積分すれば明らかである

$$\int p(x) dx = 1$$

$$\int \sum_k \pi_k \mathcal{N}(x | \mu_k, \Sigma_k) dx = 1$$

$$\sum_k \pi_k \int \mathcal{N}(x | \mu_k, \Sigma_k) dx = 1$$

$$\sum_k \pi_k = 1$$

混合ガウス分布

各ガウス分布 $\mathcal{N}(\boldsymbol{x}|\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)$ は、正規化されている

$$\forall k \quad \int \mathcal{N}(\boldsymbol{x}|\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k) d\boldsymbol{x} = 1$$

- $\mathcal{N}(\boldsymbol{x}|\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k) \geq 0$ であるので、 $p(\boldsymbol{x}) \geq 0$ となるための十分条件は、全ての k について、 $\pi_k \geq 0$ が成立することである

$$\forall k \in \{1, \dots, K\} \quad \pi_k \geq 0 \Rightarrow p(\boldsymbol{x}) \geq 0$$

- これと $\sum_k \pi_k = 1$ から、結局全ての π_k について以下が成り立つ

$$0 \leq \pi_k \leq 1 \tag{23}$$

混合ガウス分布

$$p(\mathbf{x}) = \sum_k \pi_k \mathcal{N}(\mathbf{x} | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)$$

$$\sum_k \pi_k = 1, \quad \forall k \quad 0 \leq \pi_k \leq 1$$

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

- 混合ガウス分布を決定づけるパラメータは、 $\boldsymbol{\pi} \equiv \{\pi_1, \pi_2, \dots, \pi_K\}$ 、 $\boldsymbol{\mu} \equiv \{\boldsymbol{\mu}_1, \boldsymbol{\mu}_2, \dots, \boldsymbol{\mu}_K\}$ 、 $\boldsymbol{\Sigma} \equiv \{\boldsymbol{\Sigma}_1, \boldsymbol{\Sigma}_2, \dots, \boldsymbol{\Sigma}_K\}$

混合ガウス分布

- 混合ガウス分布を導入する理由 (再確認)
 - データに対して、**各クラスタに属する確率**が分かるようにするため
- 問題設定
 - x の N 個の観測点で構成されるデータ集合 $\mathcal{D} = \{x_1, \dots, x_N\}$ ($x_i \in \mathbb{R}^D$)
 - データ集合 $\mathcal{D} = \{x_1, \dots, x_N\}$ を、 K 個のクラスタに分割
 - K は、**既知の定数**であるとする
 - k 番目のクラスタが、**平均 μ_k 、共分散行列 Σ_k の正規分布 $\mathcal{N}(x|\mu_k, \Sigma_k)$ で表現できる**とする
 - 各クラスタの分布 $\mathcal{N}(x|\mu_k, \Sigma_k)$ を、 π_k で重み付けして足し合わせた混合分布が、データ全体を表す分布である

- 最尤推定を試みる

- 最尤推定によって、混合ガウス分布のパラメータ π, μ, Σ が分かったとする
- このとき次のようにすれば、クラスタリングが可能
- 新たなデータ x が得られたとき、全ての $k(k = 1, \dots, K)$ について $\mathcal{N}(x|\mu_k, \Sigma_k)$ を計算する
- これを最大にするような k が、データ x が属するクラスタである

- 最尤推定を試みる

- 結論から先に言うと、いきなり最尤推定を試すと**失敗する**
- **最尤推定**によって、混合ガウス分布 $p(x)$ のパラメータ $\theta = \{\pi, \mu, \Sigma\}$ を求めている
- 尤度関数 $p(\mathcal{D}|\theta)$ を、パラメータ θ の関数とみなして、 θ について最大化することにより、 θ を求めるという考え方
- 尤度関数 $p(\mathcal{D}|\theta)$ は、パラメータ θ が与えられたときの、データの条件付き確率である
- パラメータを 1 つに決めたときに、データ \mathcal{D} が得られる確率

- 対数尤度関数 $\ln p(\mathcal{D}|\boldsymbol{\theta})$ は次のようになる

$$\begin{aligned} & \ln p(\mathcal{D}|\boldsymbol{\theta}) \\ = & \ln \prod_i p(\mathbf{x}_i|\boldsymbol{\theta}) \\ = & \ln \prod_i \left(\sum_k \pi_k \mathcal{N}(\mathbf{x}_i|\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k) \right) \\ = & \sum_i \ln \left(\sum_k \pi_k \mathcal{N}(\mathbf{x}_i|\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k) \right) \\ = & \sum_i \ln \left(\sum_k \pi_k \frac{1}{(2\pi)^{\frac{D}{2}}} \frac{1}{|\boldsymbol{\Sigma}_k|^{\frac{1}{2}}} \right. \\ & \left. \exp \left\{ -\frac{1}{2}(\mathbf{x}_i - \boldsymbol{\mu}_k)^T \boldsymbol{\Sigma}_k^{-1} (\mathbf{x}_i - \boldsymbol{\mu}_k) \right\} \right) \end{aligned} \quad (24)$$

混合ガウス分布

- 上式の最初の変形では、各データ $\mathcal{D} = \{\mathbf{x}_1, \dots, \mathbf{x}_N\}$ は、確率分布 $p(\mathbf{x})$ から、**独立に得られている**という仮定を用いた
- このようなデータ \mathcal{D} を、**i.i.d 標本**という (independently and identically distributed)
- 対数 \ln は単調増加関数であるため、対数を適用しても、関数の極値は変化しない
- 尤度関数 $p(\mathcal{D}|\theta)$ の最大化は、対数尤度関数 $\ln p(\mathcal{D}|\theta)$ の最大化と等価
- 重大な問題点
 - **対数関数 \ln の内部に、総和 (\sum) が入っている**
 - **log-sum の形状になっているため、これ以上式を簡単にできない!**
 - クラスタ数 $K = 1$ であれば、対数 \ln と、ガウス分布の指数 \exp が打ち消し合って、式が簡潔になる
- しかしここでは、このまま最尤推定を続けてみる

- パラメータ μ_k の最尤推定

- 対数尤度関数 $\ln p(\mathcal{D}|\theta)$ において、 θ はパラメータ (定数) で、 \mathcal{D} は変数であるが、実際は \mathcal{D} にはデータが入っているので、パラメータ θ を変数とみなす
- $\ln p(\mathcal{D}|\theta)$ をパラメータ μ_k で微分してみる
- これを 0 と等置することで、最適な μ_k が満たすべき式が得られる
- その前に、関数 f の対数の微分について、以下が成立することを確認しておく

$$(\ln f)' = \frac{f'}{f}, \quad f' = f \cdot (\ln f)' \quad (25)$$

- このとき次のようになる

$$\begin{aligned} & \frac{\partial}{\partial \boldsymbol{\mu}_k} \ln p(\mathcal{D} | \boldsymbol{\theta}) \\ &= \frac{\partial}{\partial \boldsymbol{\mu}_k} \sum_i \ln \left(\sum_k \pi_k \mathcal{N}(\mathbf{x}_i | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k) \right) \\ &= \sum_i \frac{\partial}{\partial \boldsymbol{\mu}_k} \ln \left(\sum_k \pi_k \mathcal{N}(\mathbf{x}_i | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k) \right) \\ &= \sum_i \frac{1}{\sum_k \pi_k \mathcal{N}(\mathbf{x}_i | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)} \frac{\partial}{\partial \boldsymbol{\mu}_k} \left(\sum_k \pi_k \mathcal{N}(\mathbf{x}_i | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k) \right) \\ &= \sum_i \frac{1}{\sum_k \pi_k \mathcal{N}(\mathbf{x}_i | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)} \frac{\partial}{\partial \boldsymbol{\mu}_k} \pi_k \mathcal{N}(\mathbf{x}_i | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k) \\ &= \sum_i \frac{1}{\sum_k \pi_k \mathcal{N}(\mathbf{x}_i | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)} \left(\frac{\partial}{\partial \boldsymbol{\mu}_k} \pi_k \right. \end{aligned}$$

$$\begin{aligned}
 & \frac{1}{(2\pi)^{\frac{D}{2}}} \frac{1}{|\Sigma_k|^{\frac{1}{2}}} \exp \left\{ -\frac{1}{2} (\mathbf{x}_i - \boldsymbol{\mu}_k)^T \Sigma_k^{-1} (\mathbf{x}_i - \boldsymbol{\mu}_k) \right\} \\
 = & \sum_i \frac{1}{\sum_k \pi_k \mathcal{N}(\mathbf{x}_i | \boldsymbol{\mu}_k, \Sigma_k)} \pi_k \frac{1}{(2\pi)^{\frac{D}{2}}} \frac{1}{|\Sigma_k|^{\frac{1}{2}}} \\
 & \frac{\partial}{\partial \boldsymbol{\mu}_k} \exp \left\{ -\frac{1}{2} (\mathbf{x}_i - \boldsymbol{\mu}_k)^T \Sigma_k^{-1} (\mathbf{x}_i - \boldsymbol{\mu}_k) \right\} \quad (26)
 \end{aligned}$$

ここで、以下のようにできる

$$\begin{aligned}
 & \frac{\partial}{\partial \boldsymbol{\mu}_k} \exp \left\{ -\frac{1}{2} (\mathbf{x}_i - \boldsymbol{\mu}_k)^T \Sigma_k^{-1} (\mathbf{x}_i - \boldsymbol{\mu}_k) \right\} \\
 = & \exp \left\{ -\frac{1}{2} (\mathbf{x}_i - \boldsymbol{\mu}_k)^T \Sigma_k^{-1} (\mathbf{x}_i - \boldsymbol{\mu}_k) \right\} \\
 & \frac{\partial}{\partial \boldsymbol{\mu}_k} \left(-\frac{1}{2} (\mathbf{x}_i - \boldsymbol{\mu}_k)^T \Sigma_k^{-1} (\mathbf{x}_i - \boldsymbol{\mu}_k) \right) \quad (27)
 \end{aligned}$$

更に、次が成立する

$$\begin{aligned}& \frac{\partial}{\partial \boldsymbol{\mu}_k} \left(-\frac{1}{2} (\mathbf{x}_i - \boldsymbol{\mu}_k)^T \boldsymbol{\Sigma}_k^{-1} (\mathbf{x}_i - \boldsymbol{\mu}_k) \right) \\&= -\frac{1}{2} \frac{\partial}{\partial \boldsymbol{\mu}_k} (\mathbf{x}_i^T \boldsymbol{\Sigma}_k^{-1} \mathbf{x}_i - \mathbf{x}_i^T \boldsymbol{\Sigma}_k^{-1} \boldsymbol{\mu}_k - \boldsymbol{\mu}_k^T \boldsymbol{\Sigma}_k^{-1} \mathbf{x}_i + \boldsymbol{\mu}_k^T \boldsymbol{\Sigma}_k^{-1} \boldsymbol{\mu}_k) \\&= -\frac{1}{2} \frac{\partial}{\partial \boldsymbol{\mu}_k} (-\mathbf{x}_i^T \boldsymbol{\Sigma}_k^{-1} \boldsymbol{\mu}_k - \mathbf{x}_i^T \boldsymbol{\Sigma}_k^{-1} \boldsymbol{\mu}_k + \boldsymbol{\mu}_k^T \boldsymbol{\Sigma}_k^{-1} \boldsymbol{\mu}_k) \\&= -\frac{1}{2} (-2\mathbf{x}_i^T \boldsymbol{\Sigma}_k^{-1} + 2\boldsymbol{\Sigma}_k^{-1} \boldsymbol{\mu}_k) \\&= \boldsymbol{\Sigma}_k^{-1} \mathbf{x}_i - \boldsymbol{\Sigma}_k^{-1} \boldsymbol{\mu}_k \\&= \boldsymbol{\Sigma}_k^{-1} (\mathbf{x}_i - \boldsymbol{\mu}_k)\end{aligned} \tag{28}$$

ここで、以下を用いた

$$\begin{aligned}& \boldsymbol{\mu}_k^T \boldsymbol{\Sigma}_k^{-1} \mathbf{x}_i \\&= (\boldsymbol{\mu}_k^T \boldsymbol{\Sigma}_k^{-1} \mathbf{x}_i)^T \quad (\because \text{スカラーであるため転置してもよい})\end{aligned}$$

$$\begin{aligned} &= \mathbf{x}_i^T (\boldsymbol{\Sigma}_k^{-1})^T (\boldsymbol{\mu}_k^T)^T \\ &= \mathbf{x}_i^T \boldsymbol{\Sigma}_k^{-1} \boldsymbol{\mu}_k \end{aligned} \quad (29)$$

共分散行列 $\boldsymbol{\Sigma}_k$ は対称行列 ($\boldsymbol{\Sigma}_k^T = \boldsymbol{\Sigma}_k$) であるため、以下が成立

$$(\boldsymbol{\Sigma}_k^{-1})^T = (\boldsymbol{\Sigma}_k^T)^{-1} = \boldsymbol{\Sigma}_k^{-1} \quad (30)$$

各項の微分は次のようになる

$$\begin{aligned} &\frac{\partial}{\partial \boldsymbol{\mu}_k} \mathbf{x}_i^T \boldsymbol{\Sigma}_k^{-1} \boldsymbol{\mu}_k \\ &= (\mathbf{x}_i^T \boldsymbol{\Sigma}_k^{-1})^T = (\boldsymbol{\Sigma}_k^{-1})^T (\mathbf{x}_i^T)^T = \boldsymbol{\Sigma}_k^{-1} \mathbf{x}_i \end{aligned} \quad (31)$$

$$\begin{aligned} &\frac{\partial}{\partial \boldsymbol{\mu}_k} \boldsymbol{\mu}_k^T \boldsymbol{\Sigma}_k^{-1} \boldsymbol{\mu}_k \\ &= \left(\boldsymbol{\Sigma}_k^{-1} + (\boldsymbol{\Sigma}_k^{-1})^T \right) \boldsymbol{\mu}_k = 2\boldsymbol{\Sigma}_k^{-1} \boldsymbol{\mu}_k \end{aligned} \quad (32)$$

結局、対数尤度関数 $\ln p(\mathcal{D}|\boldsymbol{\theta})$ の $\boldsymbol{\mu}_k$ による微分は

$$\begin{aligned} & \frac{\partial}{\partial \boldsymbol{\mu}_k} \ln p(\mathcal{D}|\boldsymbol{\theta}) \\ = & \sum_i \frac{1}{\sum_k \pi_k \mathcal{N}(\mathbf{x}_i|\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)} \pi_k \frac{1}{(2\pi)^{\frac{D}{2}}} \frac{1}{|\boldsymbol{\Sigma}_k|^{\frac{1}{2}}} \\ & \frac{\partial}{\partial \boldsymbol{\mu}_k} \exp \left\{ -\frac{1}{2} (\mathbf{x}_i - \boldsymbol{\mu}_k)^T \boldsymbol{\Sigma}_k^{-1} (\mathbf{x}_i - \boldsymbol{\mu}_k) \right\} \\ = & \sum_i \frac{1}{\sum_k \pi_k \mathcal{N}(\mathbf{x}_i|\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)} \pi_k \frac{1}{(2\pi)^{\frac{D}{2}}} \frac{1}{|\boldsymbol{\Sigma}_k|^{\frac{1}{2}}} \\ & \exp \left\{ -\frac{1}{2} (\mathbf{x}_i - \boldsymbol{\mu}_k)^T \boldsymbol{\Sigma}_k^{-1} (\mathbf{x}_i - \boldsymbol{\mu}_k) \right\} \boldsymbol{\Sigma}_k^{-1} (\mathbf{x}_i - \boldsymbol{\mu}_k) \\ = & \sum_i \frac{\pi_k \mathcal{N}(\mathbf{x}_i|\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)}{\sum_k \pi_k \mathcal{N}(\mathbf{x}_i|\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)} \boldsymbol{\Sigma}_k^{-1} (\mathbf{x}_i - \boldsymbol{\mu}_k) \\ = & 0 \end{aligned} \tag{33}$$

混合ガウス分布

- 未知のパラメータ π, μ, Σ が、分母と分子の双方に出現する複雑な式
- 直接この連立方程式を解いて、パラメータの最尤推定量を求めるのは難しそうである
- 勾配 $\nabla_{\mu_k} \ln p(\mathcal{D}|\theta)$ を利用した最適化も可能である
- この勾配の方向に、パラメータ μ_k を少しだけ更新する
- ここでは、**EM アルゴリズム**という別の手法を導出しようとしている
- x のほかに、**潜在変数**という仮想的な変数 z を導入することで、**簡単に解けるようになる**
- 上と似たような式が、後ほど登場する

混合ガウス分布

- ここまでの話の流れ

- 1 K-Means では、データを単一のクラスタに割り当てた (**ハード割り当て**)
- 2 データが属するクラスタだけではなく、より多くの情報 (各クラスタに属する確率) を手に入れたい
- 3 **ソフト割り当て**を実現するためには、クラスタリングを統計的機械学習 (確率分布) の観点から見直して、再定式化を行う必要があった
- 4 各クラスタをガウス分布として、データ全体を**混合ガウス分布**に当てはめることを考えた
- 5 混合ガウス分布のパラメータを、最尤推定により求めようとしたが、困難であることが分かった
- 6 そこで、**潜在変数**を導入して、最尤推定を簡単に解こうと考えている

- 潜在変数 z の導入

- 各データ x_i につき、1つのベクトル $z_i \in \mathbb{R}^K$ が対応しているとする
- z_i は、データ x_i が属するクラスタ を表現する

- 潜在変数 z の表現

- z_i は、 K 次元の二値確率変数 z の観測値である
- z の k 番目の要素を、 z_k と表すことにする
- 確率変数 z は、1-of- K 符号化法により表現されるとする
- 即ち、ある1つの $k \in \{1, \dots, K\}$ について $z_k = 1$ で、 $j \neq k$ に対し $z_j = 0$ となる
- $z_k (k = 1, \dots, K)$ は、 $z_k \in \{0, 1\}$ かつ $\sum_k z_k = 1$ をみたす
- ベクトル z は K 種類の状態を取る

- 潜在変数 z の例

- 例えば、データ点 x_i に対して z_i があるとする
- $z_{i1} = 1$ ($z_i = [1, 0, 0, \dots, 0]$) ならば、 x_i は 1 番目のクラスタ出身
- $z_{i2} = 1$ ($z_i = [0, 1, 0, \dots, 0]$) ならば、 x_i は 2 番目のクラスタ出身

- データ x_i が作られるまでの流れ

- z に関する確率分布 $p(z)$ から、 z_i がサンプルされる
- z が与えられた下での条件付き分布 $p(x|z)$ から、 x_i がサンプルされる
- 即ち、 x, z の同時分布は、周辺分布 $p(z)$ と、条件付き分布 $p(x|z)$ を用いて次のように書ける

$$p(x, z) = p(z)p(x|z) \quad (34)$$

- 潜在変数 z_i が最初に決められ、その z_i に応じて x_i が決まると考える
- z_i は実際には存在しない、仮想的なものである

- z_i は、実際に観測される x_i の裏側に潜んでいる

- $p(\mathbf{x})$ の表現 (予想)
 - $p(\mathbf{x})$ は混合ガウス分布になってほしい

$$p(\mathbf{x}) = \sum_k \pi_k \mathcal{N}(\mathbf{x} | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k) \quad (35)$$

- 周辺分布 $p(z)$ の定義
 - $p(z)$ は次のように定める

$$p(z) = \prod_k \pi_k^{z_k}, \quad p(z_k = 1) = \pi_k \quad (36)$$

- 但し π_k は混合係数であり、 $\sum_k \pi_k = 1, 0 \leq \pi_k \leq 1$ をみたす
- z の表現には 1-of-K 符号化法を使うため、左側のようにも書ける

混合ガウス分布

- 条件付き分布 $p(\mathbf{x}|\mathbf{z})$ の定義
 - $p(\mathbf{x}|\mathbf{z})$ は次のように定める

$$p(\mathbf{x}|\mathbf{z}) = \prod_k \mathcal{N}(\mathbf{x}|\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)^{z_k} \quad (37)$$

$$p(\mathbf{x}|z_k = 1) = \mathcal{N}(\mathbf{x}|\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k) \quad (38)$$

- $p(\mathbf{x})$ の導出

- \sum_z は、可能な全ての \mathbf{z} についての総和を取ること
- $\mathbf{z} = [1, 0, \dots, 0]^T, [0, 1, 0, \dots, 0]^T, \dots, [0, \dots, 0, 1]^T$ についての和
- これは、ベクトル \mathbf{z} の中で、1 である要素のインデックス k についての総和 \sum_k を取ることに相当

混合ガウス分布

- $p(\mathbf{x}, z) = p(z)p(\mathbf{x}|z)$ を、 z について周辺化すればよい

$$p(\mathbf{x}) = \sum_z p(\mathbf{x}, z) \quad (39)$$

$$= \sum_z p(z)p(\mathbf{x}|z) \quad (40)$$

$$= \sum_k p(z_k = 1)p(\mathbf{x}|z_k = 1) \quad (41)$$

$$= \sum_k \pi_k \mathcal{N}(\mathbf{x}|\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k) \quad (42)$$

- これは、混合ガウス分布と同じ形になっている
- 何が嬉しいのか?
 - 潜在変数を陽に含む表現 $p(\mathbf{x}, z) = p(z)p(\mathbf{x}|z)$ を得たことで、この同時分布を使った議論が可能になった

- $p(z_k = 1|\mathbf{x})$ の表現

- \mathbf{x} が与えられた下での、 z の条件付き確率
- 実は、 $p(z_k = 1|\mathbf{x})$ は、データ \mathbf{x} がクラス k に属する確率を表す
- 求めようとしているのは、この値である!

- $\gamma(z_k) = p(z_k = 1|\mathbf{x})$ とすると、ベイズの定理から次のように書ける

$$\gamma(z_k) \equiv p(z_k = 1|\mathbf{x}) \quad (43)$$

$$\begin{aligned} &= \frac{p(z_k = 1)p(\mathbf{x}|z_k = 1)}{\sum_j p(z_j = 1)p(\mathbf{x}|z_j = 1)} \\ &= \frac{\pi_k \mathcal{N}(\mathbf{x}|\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)}{\sum_j \pi_j \mathcal{N}(\mathbf{x}|\boldsymbol{\mu}_j, \boldsymbol{\Sigma}_j)} \end{aligned} \quad (44)$$

- $\sum_k \gamma(z_k) = 1$ であることに注意
- 潜在変数を導入したので、最尤推定について再度考えてみる

- 最尤推定を再挑戦

- パラメータ θ は、 $\pi \equiv \{\pi_1, \pi_2, \dots, \pi_K\}$ 、 $\mu \equiv \{\mu_1, \mu_2, \dots, \mu_K\}$ 、 $\Sigma \equiv \{\Sigma_1, \Sigma_2, \dots, \Sigma_K\}$ をまとめたもの
- 対数尤度関数 $\ln p(\mathcal{D}|\theta)$ は以下に示す通りであった

$$\begin{aligned} & \ln p(\mathcal{D}|\theta) \\ = & \sum_i \ln \left(\sum_k \pi_k \frac{1}{(2\pi)^{\frac{D}{2}}} \frac{1}{|\Sigma_k|^{\frac{1}{2}}} \right. \\ & \left. \exp \left\{ -\frac{1}{2} (x_i - \mu_k)^T \Sigma_k^{-1} (x_i - \mu_k) \right\} \right) \end{aligned} \quad (45)$$

- 尤度関数を、 π, μ, Σ のそれぞれについて最大化する
- ここでは、尤度関数を最大化する μ_k が、満たすべき条件を考える

混合ガウス分布

- $\ln p(\mathcal{D}|\boldsymbol{\theta})$ を $\boldsymbol{\mu}_k$ について偏微分して 0 と等置すると、以下を得る

$$\begin{aligned} & \frac{\partial}{\partial \boldsymbol{\mu}_k} \ln p(\mathcal{D}|\boldsymbol{\theta}) \\ &= \sum_i \frac{\pi_k \mathcal{N}(\mathbf{x}_i | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)}{\sum_k \pi_k \mathcal{N}(\mathbf{x}_i | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)} \boldsymbol{\Sigma}_k^{-1} (\mathbf{x}_i - \boldsymbol{\mu}_k) \end{aligned} \quad (46)$$

$$= \sum_i \gamma(z_{ik}) \boldsymbol{\Sigma}_k^{-1} (\mathbf{x}_i - \boldsymbol{\mu}_k) = 0 \quad (47)$$

- 途中までは、先程導出したものを利用
- 負担率 $\gamma(z_{ik})$ が現れていることに注意

混合ガウス分布

- 共分散行列 Σ_k が正則であると仮定して、両辺に左から掛けて整理する

$$\begin{aligned}\sum_i \gamma(z_{ik})(\mathbf{x}_i - \boldsymbol{\mu}_k) &= 0 \\ \Rightarrow \sum_i \gamma(z_{ik})\boldsymbol{\mu}_k &= \sum_i \gamma(z_{ik})\mathbf{x}_i \\ \Rightarrow \boldsymbol{\mu}_k \sum_i \gamma(z_{ik}) &= \sum_i \gamma(z_{ik})\mathbf{x}_i \\ \Rightarrow \boldsymbol{\mu}_k &= \frac{1}{\sum_i \gamma(z_{ik})} \sum_i \gamma(z_{ik})\mathbf{x}_i\end{aligned}\tag{48}$$

- これより、 $\boldsymbol{\mu}_k$ を導出する式が得られた

- K-Means 法との比較

- K-Means 法における、平均ベクトル μ_k の更新式と見比べてみる

$$\begin{aligned}\mu_k &= \frac{1}{\sum_i \gamma(z_{ik})} \sum_i \gamma(z_{ik}) \mathbf{x}_i \\ \mu_k &= \frac{1}{\sum_i r_{ik}} \sum_i r_{ik} \mathbf{x}_i\end{aligned}\tag{49}$$

- r_{ik} を、 $\gamma(z_{ik})$ に置き換えたものとなっている
- $\gamma(z_{ik})$ は、データ \mathbf{x}_i が、クラスタ k に属する確率である
- $\gamma(z_{ik})$ を、全てのデータ \mathbf{x}_i について足し合わせたもの $\sum_i \gamma(z_{ik})$ は、実質的に、 **k 番目のクラスタに割り当てられるデータの数**を表している (整数になるとは限らない)

混合ガウス分布

- そこで、K-Means 法のとおりと同じように、 N_k を次のように定める

$$N_k = \sum_i \gamma(z_{ik}) \quad (50)$$

このとき、 μ_k の式は

$$\mu_k = \frac{1}{N_k} \sum_i \gamma(z_{ik}) \mathbf{x}_i \quad (51)$$

- 例えば、 $\gamma(z_{ik})$ が 0, 1 のいずれかであれば、 $\sum_i \gamma(z_{ik})$ は、 k 番目のクラスタに属するデータの数と完全に一致
- クラスタ k に対応するガウス分布の平均 μ_k は、各データ \mathbf{x}_i の重み付き平均
- 重み因子は、事後確率 $p(z_k = 1 | \mathbf{x}_i) \equiv \gamma(z_{ik})$ である

混合ガウス分布

- $\gamma(z_{ik})$ は、 $\sum_k \gamma(z_{ik}) = 1$ となることから分かるように、 \mathbf{x}_i を生成するために、 k 番目のガウス分布が、どの程度貢献したかを表す
- 言い換えると、 k 番目のガウス分布が、 \mathbf{x}_i の出現を説明する度合いである
- この意味で、 $\gamma(z_{ik})$ のことを負担率 (Responsibility) という

- 尤度関数を最大化する Σ_k の導出

- μ_k の場合と同様に、対数尤度関数 $\ln p(\mathcal{D}|\theta)$ を、 Σ_k に関して微分して、0 と等置すればよい
- かなり導出が長くなるので注意

$$\begin{aligned} & \frac{\partial}{\partial \Sigma_k} \ln p(\mathcal{D}|\theta) \\ = & \frac{\partial}{\partial \Sigma_k} \sum_i \ln \left(\sum_k \pi_k \mathcal{N}(\mathbf{x}_i | \mu_k, \Sigma_k) \right) \\ = & \sum_i \frac{\partial}{\partial \Sigma_k} \ln \left(\sum_k \pi_k \mathcal{N}(\mathbf{x}_i | \mu_k, \Sigma_k) \right) \\ = & \sum_i \frac{1}{\sum_k \pi_k \mathcal{N}(\mathbf{x}_i | \mu_k, \Sigma_k)} \frac{\partial}{\partial \Sigma_k} \left(\sum_k \pi_k \mathcal{N}(\mathbf{x}_i | \mu_k, \Sigma_k) \right) \\ = & \sum_i \frac{1}{\sum_k \pi_k \mathcal{N}(\mathbf{x}_i | \mu_k, \Sigma_k)} \frac{\partial}{\partial \Sigma_k} \pi_k \mathcal{N}(\mathbf{x}_i | \mu_k, \Sigma_k) \end{aligned} \quad (52)$$

ここで、以下の部分を求める

$$\begin{aligned} & \frac{\partial}{\partial \Sigma_k} \pi_k \mathcal{N}(\mathbf{x}_i | \boldsymbol{\mu}_k, \Sigma_k) \\ = & \frac{\partial}{\partial \Sigma_k} \pi_k \frac{1}{(2\pi)^{\frac{D}{2}}} \frac{1}{|\Sigma_k|^{\frac{1}{2}}} \exp \left\{ -\frac{1}{2} (\mathbf{x}_i - \boldsymbol{\mu}_k)^T \Sigma_k^{-1} (\mathbf{x}_i - \boldsymbol{\mu}_k) \right\} \\ = & \pi_k \frac{1}{(2\pi)^{\frac{D}{2}}} \left(\frac{\partial}{\partial \Sigma_k} \frac{1}{|\Sigma_k|^{\frac{1}{2}}} \right. \\ & \left. \exp \left\{ -\frac{1}{2} (\mathbf{x}_i - \boldsymbol{\mu}_k)^T \Sigma_k^{-1} (\mathbf{x}_i - \boldsymbol{\mu}_k) \right\} \right) \end{aligned} \quad (53)$$

微分の中身は、 Σ_k についての合成関数となっている

$$\frac{\partial}{\partial \Sigma_k} \frac{1}{|\Sigma_k|^{\frac{1}{2}}} \exp \left\{ -\frac{1}{2} (\mathbf{x}_i - \boldsymbol{\mu}_k)^T \Sigma_k^{-1} (\mathbf{x}_i - \boldsymbol{\mu}_k) \right\}$$

混合ガウス分布

- 一般に、行列 \mathbf{X} についてのスカラー関数 $f(\mathbf{X}), g(\mathbf{X})$ があるとき、以下の連鎖律が成立

$$\frac{\partial}{\partial \mathbf{X}} f(\mathbf{X})g(\mathbf{X}) = f(\mathbf{X})\frac{\partial g(\mathbf{X})}{\partial \mathbf{X}} + g(\mathbf{X})\frac{\partial f(\mathbf{X})}{\partial \mathbf{X}} \quad (54)$$

これを利用して、先程の微分を求める

$$\begin{aligned} & \frac{\partial}{\partial \Sigma_k} \frac{1}{|\Sigma_k|^{\frac{1}{2}}} \exp \left\{ -\frac{1}{2}(\mathbf{x}_i - \boldsymbol{\mu}_k)^T \Sigma_k^{-1}(\mathbf{x}_i - \boldsymbol{\mu}_k) \right\} \\ = & \frac{1}{|\Sigma_k|^{\frac{1}{2}}} \frac{\partial}{\partial \Sigma_k} \exp \left\{ -\frac{1}{2}(\mathbf{x}_i - \boldsymbol{\mu}_k)^T \Sigma_k^{-1}(\mathbf{x}_i - \boldsymbol{\mu}_k) \right\} + \\ & \exp \left\{ -\frac{1}{2}(\mathbf{x}_i - \boldsymbol{\mu}_k)^T \Sigma_k^{-1}(\mathbf{x}_i - \boldsymbol{\mu}_k) \right\} \frac{\partial}{\partial \Sigma_k} \frac{1}{|\Sigma_k|^{\frac{1}{2}}} \end{aligned} \quad (55)$$

- 各項の微分を順番に求める

$$\begin{aligned}& \frac{\partial}{\partial \Sigma_k} \frac{1}{|\Sigma_k|^{\frac{1}{2}}} \\&= \frac{\partial}{\partial \Sigma_k} \exp \left(\ln \frac{1}{|\Sigma_k|^{\frac{1}{2}}} \right) \\&= \frac{\partial}{\partial \Sigma_k} \exp \left(-\frac{1}{2} \ln |\Sigma_k| \right) \\&= \exp \left(-\frac{1}{2} \ln |\Sigma_k| \right) \frac{\partial}{\partial \Sigma_k} \left(-\frac{1}{2} \ln |\Sigma_k| \right) \\&= \exp \left(-\frac{1}{2} \ln |\Sigma_k| \right) \left(-\frac{1}{2} \right) \frac{\partial}{\partial \Sigma_k} \ln |\Sigma_k| \\&= -\frac{1}{2} \exp \left(-\frac{1}{2} \ln |\Sigma_k| \right) (\Sigma_k^{-1})^T \\&= -\frac{1}{2} \frac{1}{|\Sigma_k|^{\frac{1}{2}}} \Sigma_k^{-1}\end{aligned} \tag{56}$$

また

$$\begin{aligned} & \frac{\partial}{\partial \Sigma_k} \exp \left\{ -\frac{1}{2} (\mathbf{x}_i - \boldsymbol{\mu}_k)^T \Sigma_k^{-1} (\mathbf{x}_i - \boldsymbol{\mu}_k) \right\} \\ = & \exp \left\{ -\frac{1}{2} (\mathbf{x}_i - \boldsymbol{\mu}_k)^T \Sigma_k^{-1} (\mathbf{x}_i - \boldsymbol{\mu}_k) \right\} \\ & \frac{\partial}{\partial \Sigma_k} \left(-\frac{1}{2} (\mathbf{x}_i - \boldsymbol{\mu}_k)^T \Sigma_k^{-1} (\mathbf{x}_i - \boldsymbol{\mu}_k) \right) \end{aligned} \quad (57)$$

であって

$$\begin{aligned} & \frac{\partial}{\partial \Sigma_k} \left(-\frac{1}{2} (\mathbf{x}_i - \boldsymbol{\mu}_k)^T \Sigma_k^{-1} (\mathbf{x}_i - \boldsymbol{\mu}_k) \right) \\ = & -\frac{1}{2} \frac{\partial}{\partial \Sigma_k} \left((\mathbf{x}_i - \boldsymbol{\mu}_k)^T \Sigma_k^{-1} (\mathbf{x}_i - \boldsymbol{\mu}_k) \right) \\ = & -\frac{1}{2} \frac{\partial}{\partial \Sigma_k} \text{Tr} \left((\mathbf{x}_i - \boldsymbol{\mu}_k)^T \Sigma_k^{-1} (\mathbf{x}_i - \boldsymbol{\mu}_k) \right) \end{aligned}$$

$$\begin{aligned}
 &= -\frac{1}{2} \frac{\partial}{\partial \Sigma_k} \text{Tr} \left(\Sigma_k^{-1} (\mathbf{x}_i - \boldsymbol{\mu}_k)(\mathbf{x}_i - \boldsymbol{\mu}_k)^T \right) \\
 &= -\frac{1}{2} \left(- \left(\Sigma_k^{-1} (\mathbf{x}_i - \boldsymbol{\mu}_k)(\mathbf{x}_i - \boldsymbol{\mu}_k)^T \Sigma_k^{-1} \right)^T \right) \\
 &= \frac{1}{2} \left(\Sigma_k^{-1} (\mathbf{x}_i - \boldsymbol{\mu}_k)(\mathbf{x}_i - \boldsymbol{\mu}_k)^T \Sigma_k^{-1} \right)^T \\
 &= \frac{1}{2} \left(\Sigma_k^{-1} \right)^T \left((\mathbf{x}_i - \boldsymbol{\mu}_k)(\mathbf{x}_i - \boldsymbol{\mu}_k)^T \right)^T \left(\Sigma_k^{-1} \right)^T \\
 &= \frac{1}{2} \Sigma_k^{-1} (\mathbf{x}_i - \boldsymbol{\mu}_k)(\mathbf{x}_i - \boldsymbol{\mu}_k)^T \Sigma_k^{-1} \tag{58}
 \end{aligned}$$

のように求まる

- 行列のトレースについて、一般に以下が成り立つことを利用している

$$\text{Tr}(\mathbf{X}\mathbf{Y}) = \text{Tr}(\mathbf{Y}\mathbf{X}) \tag{59}$$

$$\text{Tr}(\mathbf{X}\mathbf{Y}\mathbf{Z}) = \text{Tr}(\mathbf{Y}\mathbf{Z}\mathbf{X}) = \text{Tr}(\mathbf{Z}\mathbf{X}\mathbf{Y}) \tag{60}$$

混合ガウス分布

- また先程の微分では以下の公式を用いている

$$\frac{\partial}{\partial \mathbf{X}} \text{Tr}(\mathbf{X}^{-1} \mathbf{Y}) = -(\mathbf{X}^{-1} \mathbf{Y} \mathbf{X}^{-1})^T \quad (61)$$

- この公式を証明するためには、いくつかの段階を踏む必要がある
- まずは以下の微分公式の導出から始める

$$\frac{\partial}{\partial x} \mathbf{A} \mathbf{B}$$

上式の微分は、行列の積 $\mathbf{A} \mathbf{B}$ の i, k 成分について考えれば

$$\begin{aligned} & \frac{\partial}{\partial x} \sum_j A_{ij} B_{jk} \\ &= \sum_j \frac{\partial}{\partial x} A_{ij} B_{jk} \end{aligned}$$

$$\begin{aligned} &= \sum_j \left(\frac{\partial A_{ij}}{\partial x} B_{jk} + A_{ij} \frac{\partial B_{jk}}{\partial x} \right) \\ &= \sum_j \frac{\partial A_{ij}}{\partial x} B_{jk} + \sum_j A_{ij} \frac{\partial B_{jk}}{\partial x} \end{aligned} \quad (62)$$

であるから、結局

$$\frac{\partial}{\partial x} \mathbf{A} \mathbf{B} = \frac{\partial \mathbf{A}}{\partial x} \mathbf{B} + \mathbf{A} \frac{\partial \mathbf{B}}{\partial x} \quad (63)$$

となる

混合ガウス分布

- 上記の公式から、以下の公式を簡単に導ける

$$\begin{aligned}\frac{\partial}{\partial x} \mathbf{A}^{-1} \mathbf{A} &= \frac{\partial \mathbf{A}^{-1}}{\partial x} \mathbf{A} + \mathbf{A}^{-1} \frac{\partial \mathbf{A}}{\partial x} \\ \frac{\partial}{\partial x} \mathbf{I} &= \frac{\partial \mathbf{A}^{-1}}{\partial x} \mathbf{A} + \mathbf{A}^{-1} \frac{\partial \mathbf{A}}{\partial x} \\ 0 &= \frac{\partial \mathbf{A}^{-1}}{\partial x} \mathbf{A} + \mathbf{A}^{-1} \frac{\partial \mathbf{A}}{\partial x}\end{aligned}\tag{64}$$

これに右から \mathbf{A}^{-1} を掛ければ

$$\begin{aligned}0 &= \frac{\partial \mathbf{A}^{-1}}{\partial x} \mathbf{A} \mathbf{A}^{-1} + \mathbf{A}^{-1} \frac{\partial \mathbf{A}}{\partial x} \mathbf{A}^{-1} \\ 0 &= \frac{\partial \mathbf{A}^{-1}}{\partial x} \mathbf{I} + \mathbf{A}^{-1} \frac{\partial \mathbf{A}}{\partial x} \mathbf{A}^{-1} \\ 0 &= \frac{\partial \mathbf{A}^{-1}}{\partial x} + \mathbf{A}^{-1} \frac{\partial \mathbf{A}}{\partial x} \mathbf{A}^{-1}\end{aligned}$$

$$\frac{\partial}{\partial x} \mathbf{A}^{-1} = -\mathbf{A}^{-1} \frac{\partial \mathbf{A}}{\partial x} \mathbf{A}^{-1} \quad (65)$$

より

$$\frac{\partial}{\partial x} \mathbf{A}^{-1} = -\mathbf{A}^{-1} \frac{\partial \mathbf{A}}{\partial x} \mathbf{A}^{-1} \quad (66)$$

を得る (逆行列の微分公式)

- 逆行列 \mathbf{A}^{-1} の k, l 成分 $(\mathbf{A}^{-1})_{kl}$ については、以下のように書ける

$$\begin{aligned} & \frac{\partial}{\partial x} (\mathbf{A}^{-1})_{kl} \\ &= -\sum_{m,n} (\mathbf{A}^{-1})_{km} \left(\frac{\partial \mathbf{A}}{\partial x} \right)_{mn} (\mathbf{A}^{-1})_{ml} \\ &= -\sum_{m,n} (\mathbf{A}^{-1})_{km} \frac{\partial A_{mn}}{\partial x} (\mathbf{A}^{-1})_{ml} \end{aligned} \quad (67)$$

- この逆行列の微分公式を使えば、以下の微分公式を導出できる

$$\frac{\partial}{\partial \mathbf{X}} \text{Tr}(\mathbf{X}^{-1} \mathbf{Y})$$

行列 \mathbf{X} の i, j 要素による微分を考えれば

$$\begin{aligned} & \frac{\partial}{\partial X_{ij}} \text{Tr}(\mathbf{X}^{-1} \mathbf{Y}) \\ &= \frac{\partial}{\partial X_{ij}} \sum_{k,l} (\mathbf{X}^{-1})_{kl} Y_{lk} \\ &= \sum_{k,l} \frac{\partial}{\partial X_{ij}} ((\mathbf{X}^{-1})_{kl}) Y_{lk} \\ &= \sum_{k,l} \left(- \sum_{m,n} (\mathbf{X}^{-1})_{km} \frac{\partial X_{mn}}{\partial X_{ij}} (\mathbf{X}^{-1})_{ml} \right) Y_{lk} \\ &= \sum_{k,l} \left(- (\mathbf{X}^{-1})_{ki} (\mathbf{X}^{-1})_{jl} \right) Y_{lk} \end{aligned}$$

$$\begin{aligned} &= \sum_{k,l} \left(-(\mathbf{X}^{-1})_{jl} Y_{lk} (\mathbf{X}^{-1})_{ki} \right) \\ &= -(\mathbf{X}^{-1} \mathbf{Y} \mathbf{X}^{-1})_{ji} \\ &= -\left((\mathbf{X}^{-1} \mathbf{Y} \mathbf{X}^{-1})^T \right)_{ij} \end{aligned} \tag{68}$$

であるから

$$\frac{\partial}{\partial \mathbf{X}} \text{Tr}(\mathbf{X}^{-1} \mathbf{Y}) = -(\mathbf{X}^{-1} \mathbf{Y} \mathbf{X}^{-1})^T \tag{69}$$

を得られる

- 尤度関数を最大化する Σ_k の導出
 - これより、結局次のようになる

$$\begin{aligned} & \frac{\partial}{\partial \Sigma_k} \ln p(\mathcal{D} | \theta) \\ = & \sum_i \frac{1}{\sum_k \pi_k \mathcal{N}(\mathbf{x}_i | \boldsymbol{\mu}_k, \Sigma_k)} \frac{\partial}{\partial \Sigma_k} \pi_k \mathcal{N}(\mathbf{x}_i | \boldsymbol{\mu}_k, \Sigma_k) \\ = & \sum_i \frac{1}{\sum_k \pi_k \mathcal{N}(\mathbf{x}_i | \boldsymbol{\mu}_k, \Sigma_k)} \pi_k \frac{1}{(2\pi)^{\frac{D}{2}}} \\ & \left(\frac{\partial}{\partial \Sigma_k} \frac{1}{|\Sigma_k|^{\frac{1}{2}}} \exp \left\{ -\frac{1}{2} (\mathbf{x}_i - \boldsymbol{\mu}_k)^T \Sigma_k^{-1} (\mathbf{x}_i - \boldsymbol{\mu}_k) \right\} \right) \\ = & \sum_i \frac{1}{\sum_k \pi_k \mathcal{N}(\mathbf{x}_i | \boldsymbol{\mu}_k, \Sigma_k)} \pi_k \frac{1}{(2\pi)^{\frac{D}{2}}} \\ & \left(\frac{1}{|\Sigma_k|^{\frac{1}{2}}} \frac{\partial}{\partial \Sigma_k} \exp \left\{ -\frac{1}{2} (\mathbf{x}_i - \boldsymbol{\mu}_k)^T \Sigma_k^{-1} (\mathbf{x}_i - \boldsymbol{\mu}_k) \right\} + \right. \end{aligned}$$

$$\begin{aligned}
 & \exp \left\{ -\frac{1}{2} (\mathbf{x}_i - \boldsymbol{\mu}_k)^T \boldsymbol{\Sigma}_k^{-1} (\mathbf{x}_i - \boldsymbol{\mu}_k) \right\} \frac{\partial}{\partial \boldsymbol{\Sigma}_k} \frac{1}{|\boldsymbol{\Sigma}_k|^{\frac{1}{2}}} \\
 = & \sum_i \frac{1}{\sum_k \pi_k \mathcal{N}(\mathbf{x}_i | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)} \pi_k \frac{1}{(2\pi)^{\frac{D}{2}}} \\
 & \left(\frac{1}{|\boldsymbol{\Sigma}_k|^{\frac{1}{2}}} \exp \left\{ -\frac{1}{2} (\mathbf{x}_i - \boldsymbol{\mu}_k)^T \boldsymbol{\Sigma}_k^{-1} (\mathbf{x}_i - \boldsymbol{\mu}_k) \right\} \right. \\
 & \left(\frac{1}{2} \boldsymbol{\Sigma}_k^{-1} (\mathbf{x}_i - \boldsymbol{\mu}_k) (\mathbf{x}_i - \boldsymbol{\mu}_k)^T \boldsymbol{\Sigma}_k^{-1} \right) - \\
 & \left. \frac{1}{2} \frac{1}{|\boldsymbol{\Sigma}_k|^{\frac{1}{2}}} \boldsymbol{\Sigma}_k^{-1} \exp \left\{ -\frac{1}{2} (\mathbf{x}_i - \boldsymbol{\mu}_k)^T \boldsymbol{\Sigma}_k^{-1} (\mathbf{x}_i - \boldsymbol{\mu}_k) \right\} \right) \\
 = & \sum_i \frac{1}{\sum_k \pi_k \mathcal{N}(\mathbf{x}_i | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)} \pi_k \frac{1}{(2\pi)^{\frac{D}{2}}} \frac{1}{|\boldsymbol{\Sigma}_k|^{\frac{1}{2}}} \\
 & \exp \left\{ -\frac{1}{2} (\mathbf{x}_i - \boldsymbol{\mu}_k)^T \boldsymbol{\Sigma}_k^{-1} (\mathbf{x}_i - \boldsymbol{\mu}_k) \right\}
 \end{aligned}$$

$$\begin{aligned}
 & \left(\frac{1}{2} \Sigma_k^{-1} (\mathbf{x}_i - \boldsymbol{\mu}_k)(\mathbf{x}_i - \boldsymbol{\mu}_k)^T \Sigma_k^{-1} - \frac{1}{2} \Sigma_k^{-1} \right) \\
 = & \sum_i \frac{\pi_k \mathcal{N}(\mathbf{x}_i | \boldsymbol{\mu}_k, \Sigma_k)}{\sum_k \pi_k \mathcal{N}(\mathbf{x}_i | \boldsymbol{\mu}_k, \Sigma_k)} \\
 & \frac{1}{2} \Sigma_k^{-1} ((\mathbf{x}_i - \boldsymbol{\mu}_k)(\mathbf{x}_i - \boldsymbol{\mu}_k)^T \Sigma_k^{-1} - \mathbf{I}) \\
 = & \frac{1}{2} \sum_i \gamma(z_{ik}) \Sigma_k^{-1} ((\mathbf{x}_i - \boldsymbol{\mu}_k)(\mathbf{x}_i - \boldsymbol{\mu}_k)^T \Sigma_k^{-1} - \mathbf{I}) = 0 \quad (70)
 \end{aligned}$$

両辺に左右から Σ_k を掛けて、整理すれば

$$\begin{aligned}
 & \frac{1}{2} \sum_i \gamma(z_{ik}) ((\mathbf{x}_i - \boldsymbol{\mu}_k)(\mathbf{x}_i - \boldsymbol{\mu}_k)^T - \Sigma_k) = 0 \\
 \Rightarrow & \sum_i \gamma(z_{ik}) \Sigma_k = \sum_i \gamma(z_{ik}) (\mathbf{x}_i - \boldsymbol{\mu}_k)(\mathbf{x}_i - \boldsymbol{\mu}_k)^T \\
 \Rightarrow & \Sigma_k = \frac{1}{\sum_i \gamma(z_{ik})} \sum_i \gamma(z_{ik}) (\mathbf{x}_i - \boldsymbol{\mu}_k)(\mathbf{x}_i - \boldsymbol{\mu}_k)^T \quad (71)
 \end{aligned}$$

$$\Rightarrow \Sigma_k = \frac{1}{N_k} \sum_i \gamma(z_{ik})(\mathbf{x}_i - \boldsymbol{\mu}_k)(\mathbf{x}_i - \boldsymbol{\mu}_k)^T \quad (72)$$

- これより、 Σ_k を導出する式が得られた

- 尤度関数を最大化する π_k の導出
 - μ_k, Σ_k の場合と同様に、対数尤度関数 $\ln p(\mathcal{D}|\theta)$ を、 π_k に関して微分して、0 と等置すればよい
 - 但し、 $\sum_k \pi_k = 1$ という制約条件を考慮しなければならない
 - そのため、**ラグランジュの未定係数法**を用いる
- 以下を最大化する π_k を求める

$$\ln p(\mathcal{D}|\theta) + \lambda \left(\sum_k \pi_k - 1 \right) \quad (73)$$

- π_k で微分して 0 と等置すると、次のようになる

$$\begin{aligned} & \frac{\partial}{\partial \pi_k} \left(\ln p(\mathcal{D}|\boldsymbol{\theta}) + \lambda \left(\sum_k \pi_k - 1 \right) \right) \\ &= \frac{\partial}{\partial \pi_k} \left(\sum_i \ln \left(\sum_k \pi_k \mathcal{N}(\mathbf{x}_i | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k) \right) + \lambda \left(\sum_k \pi_k - 1 \right) \right) \\ &= \sum_i \frac{\mathcal{N}(\mathbf{x}_i | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)}{\sum_k \pi_k \mathcal{N}(\mathbf{x}_i | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)} + \lambda \end{aligned} \quad (74)$$

$$= \sum_i \frac{\gamma(z_{ik})}{\pi_k} + \lambda = 0 \quad (75)$$

両辺に π_k を掛けて

$$\sum_i \gamma(z_{ik}) + \lambda \pi_k = 0 \quad (76)$$

混合ガウス分布

k についての和を取ると

$$\sum_k \left(\sum_i \gamma(z_{ik}) + \lambda \pi_k \right) = 0 \quad (77)$$

$$\sum_i \sum_k \gamma(z_{ik}) + \lambda \sum_k \pi_k = 0 \quad (78)$$

$$\sum_i 1 + \lambda = 0 \quad (79)$$

$$N + \lambda = 0 \quad (80)$$

$$\therefore \lambda = -N \quad (81)$$

これより

$$\sum_i \gamma(z_{ik}) + (-N) \pi_k = 0 \quad (82)$$

$$\pi_k = \frac{1}{N} \sum_i \gamma(z_{ik}) = \frac{N_k}{N} \quad (83)$$

混合ガウス分布

ここで、以下が成立することに注意

$$N_k = \sum_i \gamma(z_{ik}), \quad 1 = \sum_k \gamma(z_{ik})$$

- これより、混合係数 π_k は、全ての要素における、クラスタ k の負担率 $\gamma(z_{ik})$ の平均である
- ここまでで、対数尤度関数 $\ln p(\mathcal{D}|\theta)$ を最大化するような、 μ_k, Σ_k, π_k の式が得られた

μ_k, Σ_k, π_k の更新式

$$\mu_k = \frac{1}{N_k} \sum_i \gamma(z_{ik}) \mathbf{x}_i \quad (84)$$

$$\Sigma_k = \frac{1}{N_k} \sum_i \gamma(z_{ik}) (\mathbf{x}_i - \mu_k)(\mathbf{x}_i - \mu_k)^T \quad (85)$$

$$\pi_k = \frac{N_k}{N} \quad (86)$$

$$N_k = \sum_i \gamma(z_{ik}) \quad (87)$$

$\gamma(z_{ik})$ の更新式

$$\gamma(z_{ik}) = \frac{\pi_k \mathcal{N}(\mathbf{x}_i | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)}{\sum_k \pi_k \mathcal{N}(\mathbf{x}_i | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)} \quad (88)$$

● 注意点

- $\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k, \pi_k$ の更新式は、これらのパラメータについての、**陽な解は与えていない**
- なぜなら、これらの更新式は全て、負担率 $\gamma(z_{ik})$ に依存しているため
- そしてその負担率 $\gamma(z_{ik})$ は、 $\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k, \pi_k$ の全てに依存する

- これらの更新式の意味

- 最尤推定の解を求めるための、**繰り返し手続きの存在**を示唆
- 即ち、 μ_k, Σ_k, π_k の初期化後に、(1) $\gamma(z_{ik})$ の更新と、(2) それを用いた μ_k, Σ_k, π_k の更新という、**2段階の処理を繰り返す**手続き
- これは、混合ガウス分布を確率モデルとして使ったときの、**EM アルゴリズム**となっている
- 混合ガウス分布に対する EM アルゴリズムは重要なので、次にまとめる

混合ガウス分布に対する EM アルゴリズム

- 目的は、混合ガウスモデルが与えられているとき、そのパラメータ (各ガウス分布の平均、分散、そして混合係数) について、尤度関数を最大化することである

- 1 平均 μ_k^{old} 、分散 Σ_k^{old} 、そして混合係数 π_k^{old} を初期化し、対数尤度 $\ln p(\mathcal{D}|\theta)$ の初期値を計算
- 2 **E ステップ**: 現在のパラメータを用いて、負担率 $\gamma(z_{ik})$ を計算

$$\gamma(z_{ik}) \leftarrow \frac{\pi_k^{\text{old}} \mathcal{N}(\mathbf{x}_i | \mu_k^{\text{old}}, \Sigma_k^{\text{old}})}{\sum_k \pi_k^{\text{old}} \mathcal{N}(\mathbf{x}_i | \mu_k^{\text{old}}, \Sigma_k^{\text{old}})} \quad (89)$$

3 M ステップ: 現在の負担率 $\gamma(z_{ik})$ を用いて、パラメータを更新

$$\boldsymbol{\mu}_k^{\text{new}} \leftarrow \frac{1}{N_k} \sum_i \gamma(z_{ik}) \boldsymbol{x}_i \quad (90)$$

$$\boldsymbol{\Sigma}_k^{\text{new}} \leftarrow \frac{1}{N_k} \sum_i \gamma(z_{ik}) (\boldsymbol{x}_i - \boldsymbol{\mu}_k^{\text{new}})(\boldsymbol{x}_i - \boldsymbol{\mu}_k^{\text{new}})^T \quad (91)$$

$$\pi_k^{\text{new}} \leftarrow \frac{N_k}{N} \quad (92)$$

但し

$$N_k = \sum_i \gamma(z_{ik}) \quad (93)$$

4 対数尤度 $\ln p(\mathcal{D}|\boldsymbol{\theta})$ を計算

$$\ln p(\mathcal{D}|\boldsymbol{\theta}) = \sum_i \ln \left(\sum_k \pi_k^{\text{new}} \mathcal{N}(\mathbf{x}_i | \boldsymbol{\mu}_k^{\text{new}}, \boldsymbol{\Sigma}_k^{\text{new}}) \right) \quad (94)$$

パラメータの変化量、あるいは対数尤度の変化量を見て、収束性を判定

5 収束基準を満たしていなければ、(2) に戻る

$$\boldsymbol{\mu}_k^{\text{old}} \leftarrow \boldsymbol{\mu}_k^{\text{new}}, \quad \boldsymbol{\Sigma}_k^{\text{old}} \leftarrow \boldsymbol{\Sigma}_k^{\text{new}}, \quad \pi_k^{\text{old}} \leftarrow \pi_k^{\text{new}} \quad (95)$$

- EM アルゴリズムの概要

- **E ステップ** (Expectation step) では、事後確率 $p(z_k = 1|x_i)$ 、即ち負担率 $\gamma(z_{ik})$ を計算
- **M ステップ** (Maximization step) では、事後確率を使って、各パラメータ μ_k, Σ_k, π_k を再計算

- EM アルゴリズムでの注意点

- 上記 (3) の M ステップにおける、各パラメータの計算順序に注意
- 最初に新しい平均値 μ_k^{new} を計算し、**その新しい平均値を使って**、新しい共分散行列 Σ_k^{new} を計算する
- E ステップと M ステップは、**対数尤度関数 $\ln p(\mathcal{D}|\theta)$ を増加させることが保証されている**
- EM アルゴリズムは、K-Means 法と比べて、収束までに必要な繰り返し回数と、各ステップでの計算量が非常に多くなる

混合ガウス分布

- 混合ガウス分布の良い初期値を見つけるために、最初に K-Means 法を実行し、その後に EM アルゴリズムを利用する、という方法がある
- K-Means 法により得られた平均ベクトル μ_k を、各ガウス分布の平均 μ_k の初期値とする
- 各クラスタに属するデータ点の**標本分散**を、共分散行列 Σ_k の初期値とする
- 各クラスタに属するデータ点の**割合**を、混合係数 π_k の初期値とする
- 一般に対数尤度には、多数の局所解が存在するため、**その中で最大のものの (大域的最適解) に収束するとは限らない**

混合ガウス分布のまとめ

- 混合ガウス分布を導入する理由

- 曖昧さを含んだ、各データ点のクラスタへの割り当て (ソフト割り当て) を実現するため
- 各データ点について、各クラスタに属する確率が分かるようにするため

- 混合ガウス分布による表現

- 各クラスタがガウス分布 $\mathcal{N}(x|\mu_k, \Sigma_k)$ に従うと仮定
- 各ガウス分布を、混合係数 π_k で重み付けして足し合わせることで、混合ガウス分布を作り、データ全体を表現する
- パラメータ θ は、各ガウス要素の平均 μ_k と共分散行列 Σ_k 、そして混合係数 π_k である

混合ガウス分布のまとめ

- ここまでの話の流れ

- 最尤推定 (対数尤度 $\ln p(\mathcal{D}|\theta)$ の最大化) によって、パラメータ θ を求める試みは失敗した
- 対数の中に総和が入っているせいで、対数尤度の式が複雑になっていた
- 潜在変数 z を導入し、 z に関する分布を考えたことで、混合ガウス分布に対する EM アルゴリズムを自然に導出した
- 事後確率 $p(z_k = 1|x_i)$ 即ち負担率 $\gamma(z_{ik})$ の計算と、パラメータ θ の更新という 2 段階の処理を、交互に繰り返していくアルゴリズムであった

目次

- 1 K-Means 法
- 2 混合ガウス分布
- 3 EM アルゴリズム**
- 4 変分の導入
- 5 近似推論法
- 6 変分自己符号化器

3 EM アルゴリズム

- EM アルゴリズムの解釈
- 混合ガウス分布の再解釈
- K-Means 法との関連
- 一般の EM アルゴリズム
- MAP 推定に対する EM アルゴリズム
- EM アルゴリズムの拡張
- EM アルゴリズムに関する補足
- 逐次型の EM アルゴリズム

EM アルゴリズムの解釈

● ここまでの話の流れ

- 1 ソフト割り当てを実現するために、確率モデル (混合ガウスモデル) を導入した
- 2 混合ガウス分布のパラメータを、最尤推定により直接求めるのは困難であった
- 3 潜在変数を導入して再度定式化を行い、混合ガウス分布に対する EM アルゴリズムを自然に導出した
- 4 EM アルゴリズムの中で、潜在変数は、負担率 (事後分布) の形で登場しただけであった ($\gamma(z_{ik}) = p(z_k = 1|x)$)

● これからの話の流れ

- 潜在変数が果たす重要な役割を明確にする
- そのうえで、混合ガウス分布の場合をもう一度見直す

- EM アルゴリズムの目的
 - 潜在変数をもつ確率モデルについて、パラメータの最尤解を求める
- 対数尤度関数の記述 (一般的な場合)
 - 全ての観測データをまとめた、データ行列を \mathbf{X} とする (第 i 行が \mathbf{x}_i^T)
 - 全ての潜在変数をまとめた行列を \mathbf{Z} とする (第 i 行が \mathbf{z}_i^T)
 - 確率モデルの全てのパラメータを、 $\boldsymbol{\theta}$ と表す
- 対数尤度関数は次のようになる

$$\ln p(\mathbf{X}|\boldsymbol{\theta}) = \ln \left(\sum_{\mathbf{Z}} p(\mathbf{X}, \mathbf{Z}|\boldsymbol{\theta}) \right) \quad (96)$$

- 潜在変数 z が連続変数の場合は

$$\ln p(\mathbf{X}|\boldsymbol{\theta}) = \ln \left(\int p(\mathbf{X}, \mathbf{Z}|\boldsymbol{\theta}) d\mathbf{Z} \right) \quad (97)$$

のように、単に総和を積分に置き換えればよい

- これ以降、離散潜在変数のみを扱うが、総和を積分に置き換えれば、ここでの議論は、連続潜在変数についても同様に成立
- 何が問題だったか
 - 対数の中に、潜在変数に関する総和が含まれる (log-sum の形)
 - 総和が存在するので、対数 \ln が、周辺分布 $p(\mathbf{X}, \mathbf{Z}|\boldsymbol{\theta})$ に直接作用することが妨げられる
 - その結果として、対数尤度関数が複雑な形となる

EM アルゴリズムの解釈

- 完全データと不完全データ
 - X だけでなく、 Z も観測できるとする
 - $\{X, Z\}$ の組を、**完全データ集合**という
 - 実際には X しか見えないので、実際の観測データ X は**不完全**である
- Z に関する知識は、潜在変数についての事後確率分布 $p(Z|X, \theta)$ のみからしか得られない

重要な仮定と考え方

- 1 完全データ対数尤度関数 $\ln p(\mathbf{X}, \mathbf{Z}|\boldsymbol{\theta})$ の最大化は、 $\ln p(\mathbf{X}|\boldsymbol{\theta})$ の最大化よりも、簡単であると仮定
- 2 $\ln p(\mathbf{X}|\boldsymbol{\theta})$ の代わりに、完全データ対数尤度関数 $\ln p(\mathbf{X}, \mathbf{Z}|\boldsymbol{\theta})$ を最大化したいが、 \mathbf{Z} に関する情報は $\ln p(\mathbf{Z}|\mathbf{X}, \boldsymbol{\theta})$ からしか得られない
- 3 そのため、完全データ対数尤度関数 $\ln p(\mathbf{X}, \mathbf{Z}|\boldsymbol{\theta})$ は使えない
- 4 そこで、事後確率分布 $p(\mathbf{Z}|\mathbf{X}, \boldsymbol{\theta})$ による、 $\ln p(\mathbf{X}, \mathbf{Z}|\boldsymbol{\theta})$ の期待値を最大化することを考える
- 5 これが、EM アルゴリズムの考え方である

- EM アルゴリズムへの落とし込み
 - パラメータ θ を適当に初期化する
 - **E ステップ**では、事後確率分布 $p(\mathbf{Z}|\mathbf{X}, \theta^{\text{old}})$ を、現在のパラメータ θ^{old} を使って求める
 - $p(\mathbf{Z}|\mathbf{X}, \theta^{\text{old}})$ を、M ステップでの期待値の計算に使う
 - **M ステップ**では、完全データ対数尤度関数 $\ln p(\mathbf{X}, \mathbf{Z}|\theta)$ の、事後確率分布 $p(\mathbf{Z}|\mathbf{X}, \theta^{\text{old}})$ に関する期待値 $\mathcal{Q}(\theta, \theta^{\text{old}})$ を計算

$$\mathcal{Q}(\theta, \theta^{\text{old}}) = \sum_{\mathbf{Z}} p(\mathbf{Z}|\mathbf{X}, \theta^{\text{old}}) \ln p(\mathbf{X}, \mathbf{Z}|\theta) \quad (98)$$

連続潜在変数の場合は次のようになる

$$\mathcal{Q}(\theta, \theta^{\text{old}}) = \int p(\mathbf{Z}|\mathbf{X}, \theta^{\text{old}}) \ln p(\mathbf{X}, \mathbf{Z}|\theta) d\mathbf{Z} \quad (99)$$

EM アルゴリズムの解釈

- 上式において、 $p(\mathbf{Z}|\mathbf{X}, \theta^{\text{old}})$ におけるパラメータ θ^{old} は、変数ではなく定数であることに注意
- 更に、 $Q(\theta, \theta^{\text{old}})$ を θ について最大化することで、新たなパラメータの推定値 θ^{new} を得る

$$\theta^{\text{new}} = \arg \max_{\theta} Q(\theta, \theta^{\text{old}}) \quad (100)$$

- 注意点
 - $Q(\theta, \theta^{\text{old}})$ において、対数 \ln は、同時分布 $p(\mathbf{X}, \mathbf{Z}|\theta)$ に直接作用していることに注意
 - これにより、期待値の計算が簡単になることが期待される
- なぜ事後確率分布 $p(\mathbf{Z}|\mathbf{X}, \theta)$ についての期待値なのか
 - 幾分恣意的にみえるが、後ほど、期待値を取ることの正当性が明らかになる

一般の EM アルゴリズム

- 観測変数 \mathbf{X} と、潜在変数 \mathbf{Z} の同時分布 $p(\mathbf{X}, \mathbf{Z}|\theta)$ が与えられているとする
- 目的は、尤度関数 $p(\mathbf{X}|\theta)$ を、パラメータ θ について最大化することである

- パラメータを θ^{old} に初期化する
- E ステップ**: 事後確率分布 $p(\mathbf{Z}|\mathbf{X}, \theta^{\text{old}})$ を計算する

- 3 **M ステップ**: 次式で与えられる θ^{new} を計算する

$$\theta^{\text{new}} = \arg \max_{\theta} Q(\theta, \theta^{\text{old}}) \quad (101)$$

但し

$$Q(\theta, \theta^{\text{old}}) = \sum_{\mathbf{Z}} p(\mathbf{Z}|\mathbf{X}, \theta^{\text{old}}) \ln p(\mathbf{X}, \mathbf{Z}|\theta) \quad (102)$$

- 4 対数尤度の変化量、あるいはパラメータの変化量をみて、収束性を判定
- 5 収束条件を満たしていなければ、(2) に戻る

$$\theta^{\text{old}} \leftarrow \theta^{\text{new}} \quad (103)$$

3 EM アルゴリズム

- EM アルゴリズムの解釈
- 混合ガウス分布の再解釈
- K-Means 法との関連
- 一般の EM アルゴリズム
- MAP 推定に対する EM アルゴリズム
- EM アルゴリズムの拡張
- EM アルゴリズムに関する補足
- 逐次型の EM アルゴリズム

混合ガウス分布の再解釈

- 先程の EM アルゴリズムの解釈で、混合ガウス分布を見直す
- これまでの話の流れ
 - 目的は、対数尤度関数 $\ln p(\mathbf{X}|\boldsymbol{\theta})$ の最大化であった
 - しかし、対数の中に総和が出現するため、最尤推定が困難であった
 - そこで、離散潜在変数 \mathbf{Z} を導入し、完全データ集合 $\{\mathbf{X}, \mathbf{Z}\}$ に関する尤度の最大化を考える

混合ガウス分布の再解釈

- 完全データ集合 $\{X, Z\}$ に関する尤度の最大化
 - 完全データ尤度関数 $p(X, Z|\theta)$ は次のようになる

$$\begin{aligned} & p(X, Z|\theta) \\ &= p(Z|\theta)p(X|Z, \theta) \\ &= \prod_i p(z_i|\theta)p(x_i|z_i, \theta) \\ &= \prod_i \left(\prod_k \pi_k^{z_{ik}} \right) \left(\prod_k \mathcal{N}(x_i|\mu_k, \Sigma_k)^{z_{ik}} \right) \\ &= \prod_i \prod_k \pi_k^{z_{ik}} \mathcal{N}(x_i|\mu_k, \Sigma_k)^{z_{ik}} \\ &= \prod_i \prod_k (\pi_k \mathcal{N}(x_i|\mu_k, \Sigma_k))^{z_{ik}} \end{aligned} \tag{104}$$

混合ガウス分布の再解釈

- ここで、データ点 x_i に対応する潜在変数を z_i 、また z_i の k 番目の要素を z_{ik} とする
- データ点 x_i, z_i は、 $p(X, Z|\theta)$ から独立にサンプルされているとする (このとき、要素ごとの積として書ける)
- 対数を取ると次のようになる

$$\begin{aligned} & \ln p(X, Z|\theta) \\ &= \ln \left(\prod_i \prod_k (\pi_k \mathcal{N}(x_i | \mu_k, \Sigma_k))^{z_{ik}} \right) \\ &= \sum_i \sum_k \ln ((\pi_k \mathcal{N}(x_i | \mu_k, \Sigma_k))^{z_{ik}}) \\ &= \sum_i \sum_k z_{ik} \ln (\pi_k \mathcal{N}(x_i | \mu_k, \Sigma_k)) \\ &= \sum_i \sum_k z_{ik} (\ln \pi_k + \ln \mathcal{N}(x_i | \mu_k, \Sigma_k)) \end{aligned} \tag{105}$$

混合ガウス分布の再解釈

- $\ln p(\mathbf{X}, \mathbf{Z}|\boldsymbol{\theta})$ を、元々最大化しようとしていた $\ln p(\mathbf{X}|\boldsymbol{\theta})$ と比較する

$$\ln p(\mathbf{X}|\boldsymbol{\theta}) = \sum_i \ln \left(\sum_k \pi_k \mathcal{N}(\mathbf{x}_i | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k) \right) \quad (106)$$

- $\ln p(\mathbf{X}, \mathbf{Z}|\boldsymbol{\theta})$ と $\ln p(\mathbf{X}|\boldsymbol{\theta})$ を比較すると、対数 \ln と、総和 \sum_k の、**順番が入れ替わっている**
- そして、対数 \ln が、ガウス分布 $\mathcal{N}(\mathbf{x}|\boldsymbol{\mu}, \boldsymbol{\Sigma})$ に**直接作用している**
- よって、 $\ln p(\mathbf{X}, \mathbf{Z}|\boldsymbol{\theta})$ の最大化は、 $\ln p(\mathbf{X}|\boldsymbol{\theta})$ の最大化よりも、**遥かに容易である** (そして、パラメータは**陽な形で解ける**)
- そこで、 $\ln p(\mathbf{X}, \mathbf{Z}|\boldsymbol{\theta})$ を最大化するようなパラメータを求めてみる

混合ガウス分布の再解釈

- $\ln p(\mathbf{X}, \mathbf{Z}|\boldsymbol{\theta})$ の $\boldsymbol{\mu}_k$ に関する最大化
 - 以下のように、 $\boldsymbol{\mu}_k$ で微分して 0 とおけば、簡単に解ける
 - ガウス分布の微分については、先程の EM アルゴリズムの導出時に求めたものを利用している

$$\begin{aligned} & \frac{\partial}{\partial \boldsymbol{\mu}_k} \ln p(\mathbf{X}, \mathbf{Z}|\boldsymbol{\theta}) \\ &= \frac{\partial}{\partial \boldsymbol{\mu}_k} \left(\sum_i \sum_k z_{ik} (\ln \pi_k + \ln \mathcal{N}(\mathbf{x}_i | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)) \right) \\ &= \sum_i \frac{\partial}{\partial \boldsymbol{\mu}_k} \left(\sum_k z_{ik} (\ln \pi_k + \ln \mathcal{N}(\mathbf{x}_i | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)) \right) \\ &= \sum_i \frac{\partial}{\partial \boldsymbol{\mu}_k} z_{ik} (\ln \pi_k + \ln \mathcal{N}(\mathbf{x}_i | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)) \\ &= \sum_i z_{ik} \frac{\partial}{\partial \boldsymbol{\mu}_k} \ln \mathcal{N}(\mathbf{x}_i | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k) \end{aligned}$$

混合ガウス分布の再解釈

$$\begin{aligned} &= \sum_i z_{ik} \frac{\partial}{\partial \boldsymbol{\mu}_k} \ln \left(\frac{1}{(2\pi)^{\frac{D}{2}}} \frac{1}{|\boldsymbol{\Sigma}_k|^{\frac{1}{2}}} \right. \\ &\quad \left. \exp \left\{ -\frac{1}{2} (\mathbf{x}_i - \boldsymbol{\mu}_k)^T \boldsymbol{\Sigma}_k^{-1} (\mathbf{x}_i - \boldsymbol{\mu}_k) \right\} \right) \\ &= \sum_i z_{ik} \frac{\partial}{\partial \boldsymbol{\mu}_k} \left(-\frac{D}{2} \ln 2\pi - \frac{1}{2} \ln |\boldsymbol{\Sigma}_k| \right. \\ &\quad \left. - \frac{1}{2} (\mathbf{x}_i - \boldsymbol{\mu}_k)^T \boldsymbol{\Sigma}_k^{-1} (\mathbf{x}_i - \boldsymbol{\mu}_k) \right) \\ &= \sum_i z_{ik} \frac{\partial}{\partial \boldsymbol{\mu}_k} \left(-\frac{1}{2} (\mathbf{x}_i - \boldsymbol{\mu}_k)^T \boldsymbol{\Sigma}_k^{-1} (\mathbf{x}_i - \boldsymbol{\mu}_k) \right) \\ &= \sum_i z_{ik} \boldsymbol{\Sigma}_k^{-1} (\mathbf{x}_i - \boldsymbol{\mu}_k) = 0 \end{aligned} \tag{107}$$

これより

$$\sum_i z_{ik} \boldsymbol{\Sigma}_k^{-1} \boldsymbol{\mu}_k = \sum_i z_{ik} \boldsymbol{\Sigma}_k^{-1} \mathbf{x}_i \tag{108}$$

混合ガウス分布の再解釈

であるから、両辺に左から Σ_k を掛けて

$$\begin{aligned}\sum_i z_{ik} \mu_k &= \sum_i z_{ik} \mathbf{x}_i \\ \mu_k \sum_i z_{ik} &= \sum_i z_{ik} \mathbf{x}_i \\ \mu_k &= \frac{1}{\sum_i z_{ik}} \sum_i z_{ik} \mathbf{x}_i\end{aligned}\tag{109}$$

のようになる

- 上式をみると、完全データ $\{X, Z\}$ について、 μ_k は陽な形で求まっていることが分かる
- 但し実際は Z が分からないので、 z_{ik} をどうにかして得る必要がある

混合ガウス分布の再解釈

- $\ln p(\mathbf{X}, \mathbf{Z}|\boldsymbol{\theta})$ の Σ_k に関する最大化
 - Σ_k について微分して 0 とおくと、次のようになる

$$\begin{aligned}& \frac{\partial}{\partial \Sigma_k} \ln p(\mathbf{X}, \mathbf{Z}|\boldsymbol{\theta}) \\&= \sum_i z_{ik} \frac{\partial}{\partial \Sigma_k} \ln \mathcal{N}(\mathbf{x}_i | \boldsymbol{\mu}_k, \Sigma_k) \\&= \sum_i z_{ik} \frac{\partial}{\partial \Sigma_k} \left(-\frac{1}{2} \ln |\Sigma_k| - \frac{1}{2} (\mathbf{x}_i - \boldsymbol{\mu}_k)^T \Sigma_k^{-1} (\mathbf{x}_i - \boldsymbol{\mu}_k) \right) \\&= \sum_i z_{ik} \left(-\frac{1}{2} (\Sigma_k^{-1})^T + \frac{1}{2} \Sigma_k^{-1} (\mathbf{x}_i - \boldsymbol{\mu}_k) (\mathbf{x}_i - \boldsymbol{\mu}_k)^T \Sigma_k^{-1} \right) \\&= \frac{1}{2} \sum_i z_{ik} \left(-\Sigma_k^{-1} + \Sigma_k^{-1} (\mathbf{x}_i - \boldsymbol{\mu}_k) (\mathbf{x}_i - \boldsymbol{\mu}_k)^T \Sigma_k^{-1} \right) \quad (110) \\&= 0\end{aligned}$$

となる

混合ガウス分布の再解釈

- ここで、以下の微分公式を用いた

$$\frac{\partial}{\partial \mathbf{X}} \ln |\mathbf{X}| = (\mathbf{X}^{-1})^T \quad (111)$$

- これより

$$\sum_i z_{ik} \Sigma_k^{-1} = \sum_i z_{ik} \Sigma_k^{-1} (\mathbf{x}_i - \boldsymbol{\mu}_k)(\mathbf{x}_i - \boldsymbol{\mu}_k)^T \Sigma_k^{-1} \quad (112)$$

であるから、両辺に左右から Σ_k を掛けて

$$\begin{aligned} \sum_i z_{ik} \Sigma_k &= \sum_i z_{ik} (\mathbf{x}_i - \boldsymbol{\mu}_k)(\mathbf{x}_i - \boldsymbol{\mu}_k)^T \\ \Sigma_k \sum_i z_{ik} &= \sum_i z_{ik} (\mathbf{x}_i - \boldsymbol{\mu}_k)(\mathbf{x}_i - \boldsymbol{\mu}_k)^T \\ \Sigma_k &= \frac{1}{\sum_i z_{ik}} \sum_i z_{ik} (\mathbf{x}_i - \boldsymbol{\mu}_k)(\mathbf{x}_i - \boldsymbol{\mu}_k)^T \end{aligned} \quad (113)$$

のようになる

混合ガウス分布の再解釈

- 上式をみても、やはり、完全データ $\{X, Z\}$ について、 Σ_k は陽な形で求まっていることが分かる

混合ガウス分布の再解釈

- $\ln p(\mathbf{X}, \mathbf{Z}|\boldsymbol{\theta})$ の π_k に関する最大化
 - $\sum_k \pi_k = 1$ という制約条件を考慮し、ラグランジュの未定乗数法で解く
 - 従って、以下の量を最大化する

$$\ln p(\mathbf{X}, \mathbf{Z}|\boldsymbol{\theta}) + \lambda \left(\sum_k \pi_k - 1 \right) \quad (114)$$

- π_k について微分して 0 とおくと、次のようになる

$$\begin{aligned} & \frac{\partial}{\partial \pi_k} \left(\ln p(\mathbf{X}, \mathbf{Z}|\boldsymbol{\theta}) + \lambda \left(\sum_k \pi_k - 1 \right) \right) \\ &= \frac{\partial}{\partial \pi_k} \left(\sum_i \sum_k z_{ik} (\ln \pi_k + \ln \mathcal{N}(\mathbf{x}_i | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)) + \lambda \left(\sum_k \pi_k - 1 \right) \right) \\ &= \sum_i z_{ik} \frac{\partial}{\partial \pi_k} \ln \pi_k + \lambda \end{aligned}$$

混合ガウス分布の再解釈

$$= \sum_i z_{ik} \frac{1}{\pi_k} + \lambda = 0 \quad (115)$$

- これより、両辺に π_k を掛けて

$$\sum_i z_{ik} + \lambda \pi_k = 0 \quad (116)$$

全ての k について総和を取ると

$$\begin{aligned} \sum_k \sum_i z_{ik} + \sum_k \lambda \pi_k &= 0 \\ \sum_i \left(\sum_k z_{ik} \right) + \lambda \sum_k \pi_k &= 0 \\ \sum_i 1 + \lambda &= 0 \\ N + \lambda &= 0 \\ \therefore \lambda &= -N \end{aligned} \quad (117)$$

混合ガウス分布の再解釈

- よって

$$\begin{aligned}\sum_i z_{ik} \frac{1}{\pi_k} - N &= 0 \\ \sum_i z_{ik} - N\pi_k &= 0 \\ N\pi_k &= \sum_i z_{ik} \\ \therefore \pi_k &= \frac{1}{N} \sum_i z_{ik}\end{aligned}\tag{118}$$

- π_k も、完全データ (特に潜在変数) が与えられていれば、陽な形で求まる
- EM アルゴリズムにおける μ_k, Σ_k, π_k の更新式は、ここで求めた式の z_{ik} を、負担率 $\gamma(z_{ik})$ にそのまま置き換えたものである

混合ガウス分布の再解釈

- 事後確率分布 $p(\mathbf{Z}|\mathbf{X}, \boldsymbol{\theta})$ に関する期待値の計算
 - 完全データ対数尤度関数 $\ln p(\mathbf{X}, \mathbf{Z}|\boldsymbol{\theta})$ の最大化は、陽な形で解けた
 - これらの全ての式には z_{ik} が登場したが、実際には潜在変数は分からないので、 z_{ik} を何かで代用しなければならない
 - 結局、完全データ対数尤度関数 $\ln p(\mathbf{X}, \mathbf{Z}|\boldsymbol{\theta})$ の、事後確率分布 $p(\mathbf{Z}|\mathbf{X}, \boldsymbol{\theta})$ に関する期待値を考えるしかない

混合ガウス分布の再解釈

- 事後確率分布は次のように書ける

$$\begin{aligned} & p(z_i | x_i, \theta) \\ = & \frac{p(x_i | z_i, \theta) p(z_i | \theta)}{p(x_i | \theta)} \end{aligned} \quad (119)$$

$$\begin{aligned} \propto & p(x_i | z_i, \theta) p(z_i | \theta) \\ & (\because p(x_i | \theta) \text{ は、} z_i \text{ には依存しない定数項}) \end{aligned} \quad (120)$$

$$\begin{aligned} = & \left(\prod_k \mathcal{N}(x_i | \mu_k, \Sigma_k)^{z_{ik}} \right) \left(\prod_k \pi_k^{z_{ik}} \right) \\ = & \prod_k (\pi_k \mathcal{N}(x_i | \mu_k, \Sigma_k))^{z_{ik}} \end{aligned} \quad (121)$$

以上より

$$p(z_i | x_i, \theta) \propto \prod_k (\pi_k \mathcal{N}(x_i | \mu_k, \Sigma_k))^{z_{ik}} \quad (122)$$

混合ガウス分布の再解釈

であるので、 $p(\mathbf{Z}|\mathbf{X}, \boldsymbol{\theta})$ は

$$p(\mathbf{Z}|\mathbf{X}, \boldsymbol{\theta}) \propto \prod_i \prod_k (\pi_k \mathcal{N}(\mathbf{x}_i | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k))^{z_{ik}} \quad (123)$$

- $p(z_i | \mathbf{x}_i, \boldsymbol{\theta})$ を等式で表すためには、 z_i で総和を取って 1 になる (確率としての条件を満たす) ように、**正規化すればよい**

$$p(z_i | \mathbf{x}_i, \boldsymbol{\theta}) = \frac{\prod_k (\pi_k \mathcal{N}(\mathbf{x}_i | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k))^{z_{ik}}}{\sum_{\mathbf{z}_i} \prod_k (\pi_k \mathcal{N}(\mathbf{x}_i | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k))^{z_{ik}}} \quad (124)$$

混合ガウス分布の再解釈

- まず、事後確率 $p(z_i | x_i, \theta)$ に関する、 z_{ik} の期待値を求めてみる

$$\begin{aligned} & \mathbb{E}_{z_i \sim p(z_i | x_i, \theta)} [z_{ik}] \\ &= \sum_{z_i} z_{ik} p(z_i | x_i, \theta) \\ &= \sum_{z_i} z_{ik} \frac{\prod_k (\pi_k \mathcal{N}(x_i | \mu_k, \Sigma_k))^{z_{ik}}}{\sum_{z_i} \prod_k (\pi_k \mathcal{N}(x_i | \mu_k, \Sigma_k))^{z_{ik}}} \end{aligned} \quad (125)$$

- ここで

$$\sum_{z_i} \prod_k (\pi_k \mathcal{N}(x_i | \mu_k, \Sigma_k))^{z_{ik}} = \sum_k (\pi_k \mathcal{N}(x_i | \mu_k, \Sigma_k)) \quad (126)$$

と書けることに注意する

- z_i は、**1-of-K 符号化法**で表現されている

混合ガウス分布の再解釈

- $\prod_k (\pi_k \mathcal{N}(\mathbf{x}_i | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k))^{z_{ik}}$ は、 $z_{ik} = 1$ の場合、 $j \neq k$ に対して $z_{ij} = 0$ であるから、 $\pi_k \mathcal{N}(\mathbf{x}_i | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)$ という単一の項として書ける
- 全ての z_i についての総和は、 z_i の中で、要素が 1 になるインデックス k についての総和を意味する
- また

$$\sum_{z_i} z_{ik} \prod_k (\pi_k \mathcal{N}(\mathbf{x}_i | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k))^{z_{ik}} = \pi_k \mathcal{N}(\mathbf{x}_i | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k) \quad (127)$$

であることにも注意する

- \sum_{z_i} の総和の中身は、 z_i が $z_{ik} = 1$ となるとき以外は、0 である (総和の中に z_{ik} があるため)
- 従って、 z_i が $z_{ik} = 1$ となるときの項 $\prod_k (\pi_k \mathcal{N}(\mathbf{x}_i | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k))^{z_{ik}}$ 、即ち $\pi_k \mathcal{N}(\mathbf{x}_i | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)$ だけが出現する

混合ガウス分布の再解釈

- これより

$$\begin{aligned} & \mathbb{E}_{\mathbf{z}_i \sim p(\mathbf{z}_i | \mathbf{x}_i, \boldsymbol{\theta})} [z_{ik}] \\ &= \frac{\sum_{\mathbf{z}_i} z_{ik} \prod_k (\pi_k \mathcal{N}(\mathbf{x}_i | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k))^{z_{ik}}}{\sum_{\mathbf{z}_i} \prod_k (\pi_k \mathcal{N}(\mathbf{x}_i | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k))^{z_{ik}}} \\ &= \frac{\pi_k \mathcal{N}(\mathbf{x}_i | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)}{\sum_k \pi_k \mathcal{N}(\mathbf{x}_i | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)} \equiv \gamma(z_{ik}) \end{aligned} \quad (128)$$

であるから、データ点 \mathbf{x}_i に対する、 k 番目のガウス要素の負担率に一致

混合ガウス分布の再解釈

- これより、事後確率 $p(\mathbf{Z}|\mathbf{X}, \boldsymbol{\theta})$ に関する、完全データ対数尤度関数 $\ln p(\mathbf{X}, \mathbf{Z}|\boldsymbol{\theta})$ の期待値は

$$\begin{aligned} & \mathbb{E}_{\mathbf{Z} \sim p(\mathbf{Z}|\mathbf{X}, \boldsymbol{\theta})} [\ln p(\mathbf{X}, \mathbf{Z}|\boldsymbol{\theta})] \\ &= \mathbb{E}_{\mathbf{Z} \sim p(\mathbf{Z}|\mathbf{X}, \boldsymbol{\theta})} \left[\sum_i \sum_k z_{ik} (\ln \pi_k + \ln \mathcal{N}(\mathbf{x}_i | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)) \right] \\ &= \sum_i \sum_k \mathbb{E}_{\mathbf{z}_i \sim p(\mathbf{z}_i|\mathbf{x}_i, \boldsymbol{\theta})} [z_{ik} (\ln \pi_k + \ln \mathcal{N}(\mathbf{x}_i | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k))] \quad (129) \\ &= \sum_i \sum_k \mathbb{E}_{\mathbf{z}_i \sim p(\mathbf{z}_i|\mathbf{x}_i, \boldsymbol{\theta})} [z_{ik}] (\ln \pi_k + \ln \mathcal{N}(\mathbf{x}_i | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)) \quad (130) \\ &= \sum_i \sum_k \gamma(z_{ik}) (\ln \pi_k + \ln \mathcal{N}(\mathbf{x}_i | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)) \quad (131) \end{aligned}$$

である

混合ガウス分布の再解釈

- これは $\ln p(\mathbf{X}, \mathbf{Z}|\boldsymbol{\theta})$ において、 z_{ik} を $\gamma(z_{ik})$ に置き換えたものと等しい

$$\ln p(\mathbf{X}, \mathbf{Z}|\boldsymbol{\theta}) = \sum_i \sum_k z_{ik} (\ln \pi_k + \ln \mathcal{N}(\mathbf{x}_i | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)) \quad (132)$$

- 先ほどは、 $\ln p(\mathbf{X}, \mathbf{Z}|\boldsymbol{\theta})$ を最大化するような、パラメータ $\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k, \pi_k$ の式を導出した
- これらの式について、 z_{ik} を $\gamma(z_{ik})$ に置き換えれば、そのまま期待値を最大化する式として使える
- $\gamma(z_{ik})$ に置き換えた式は、EM アルゴリズムにおける更新式と一致

混合ガウス分布の再解釈

- ここまでの話の流れ

- 1 対数尤度関数 $\ln p(\mathbf{X}|\boldsymbol{\theta})$ の最大化よりも、完全データ対数尤度関数 $\ln p(\mathbf{X}, \mathbf{Z}|\boldsymbol{\theta})$ の最大化の方が簡単であると仮定した
- 2 この仮定は、混合ガウス分布の場合について成り立っていた
- 3 $\ln p(\mathbf{X}|\boldsymbol{\theta})$ の代わりに、 $\ln p(\mathbf{X}, \mathbf{Z}|\boldsymbol{\theta})$ の最大化を考えた
- 4 しかし \mathbf{Z} に関する情報がないので、代わりに、事後確率分布 $p(\mathbf{Z}|\mathbf{X}, \boldsymbol{\theta})$ による、 $\ln p(\mathbf{X}, \mathbf{Z}|\boldsymbol{\theta})$ の期待値を最大化しようとするのが、EM アルゴリズムであった
- 5 混合ガウス分布の場合について実際に試すと、期待値の最大化によって、パラメータの更新式が再び導出できた

混合ガウス分布の再解釈

- これからの話の流れ
 - K-Means 法と、混合ガウス分布に対する EM アルゴリズムを比較する

3 EM アルゴリズム

- EM アルゴリズムの解釈
- 混合ガウス分布の再解釈
- K-Means 法との関連
- 一般の EM アルゴリズム
- MAP 推定に対する EM アルゴリズム
- EM アルゴリズムの拡張
- EM アルゴリズムに関する補足
- 逐次型の EM アルゴリズム

K-Means 法との関連

- K-Means 法と、混合ガウス分布に対する EM アルゴリズムの関係
 - K-Means 法では、各データ点は、ただ一つのクラスタに割り当てられる (ハード割り当て)
 - EM アルゴリズムでは、事後確率 $\gamma(z_{ik}) \equiv p(z_k = 1 | \mathbf{x}_i)$ に基づいて、各データをソフトに割り当てる (ソフト割り当て)
 - K-Means 法は、混合ガウス分布に対する EM アルゴリズムの、ある極限として得られる

● K-Means 法の導出

- 次のように、各ガウス分布の共分散行列が $\epsilon \mathbf{I}$ で与えられる、混合ガウスモデル $p(\mathbf{x}|\boldsymbol{\theta})$ を考える (ϵ は定数とする)

$$\begin{aligned} & p(\mathbf{x}|\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k) \\ &= p(\mathbf{x}|\boldsymbol{\mu}_k, \epsilon \mathbf{I}) \\ &= \mathcal{N}(\mathbf{x}|\boldsymbol{\mu}_k, \epsilon \mathbf{I}) \\ &= \frac{1}{(2\pi)^{\frac{D}{2}}} \frac{1}{|\epsilon \mathbf{I}|^{\frac{1}{2}}} \exp \left\{ -\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu}_k)^T (\epsilon \mathbf{I})^{-1} (\mathbf{x} - \boldsymbol{\mu}_k) \right\} \quad (133) \end{aligned}$$

$$= \frac{1}{(2\pi\epsilon)^{\frac{D}{2}}} \exp \left\{ -\frac{1}{2\epsilon}(\mathbf{x} - \boldsymbol{\mu}_k)^T (\mathbf{x} - \boldsymbol{\mu}_k) \right\} \quad (134)$$

$$\begin{aligned} & (\because |\epsilon \mathbf{I}|^{\frac{1}{2}} = (\epsilon^D |\mathbf{I}|)^{\frac{1}{2}} = (\epsilon^D)^{\frac{1}{2}} = \epsilon^{\frac{D}{2}}) \\ &= \frac{1}{(2\pi\epsilon)^{\frac{D}{2}}} \exp \left\{ -\frac{1}{2\epsilon} \|\mathbf{x} - \boldsymbol{\mu}_k\|^2 \right\} \quad (135) \end{aligned}$$

$$p(\mathbf{x}|\boldsymbol{\theta}) = \sum_k \pi_k p(\mathbf{x}|\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k) \quad (136)$$

$$= \sum_k \pi_k \frac{1}{(2\pi\epsilon)^{\frac{D}{2}}} \exp \left\{ -\frac{1}{2\epsilon} \|\mathbf{x} - \boldsymbol{\mu}_k\|^2 \right\} \quad (137)$$

- この混合ガウスモデルについて、EM アルゴリズムを実行する
- 最初に、データ点 \mathbf{x}_i に対する、 k 番目のガウス要素の負担率 $\gamma(z_{ik})$ を求めて、 $\epsilon \rightarrow 0$ についての極限を取ってみる

$$\begin{aligned} \gamma(z_{ik}) &= \frac{\pi_k \mathcal{N}(\mathbf{x}_i|\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)}{\sum_j \pi_j \mathcal{N}(\mathbf{x}_i|\boldsymbol{\mu}_j, \boldsymbol{\Sigma}_j)} \\ &= \frac{\pi_k \exp \left\{ -\frac{1}{2\epsilon} \|\mathbf{x} - \boldsymbol{\mu}_k\|^2 \right\}}{\sum_j \pi_j \exp \left\{ -\frac{1}{2\epsilon} \|\mathbf{x} - \boldsymbol{\mu}_j\|^2 \right\}} \end{aligned} \quad (138)$$

- 負担率は、以下のように変形できる

$$\begin{aligned} & \frac{\pi_k \exp \left\{ -\frac{1}{2\epsilon} \|\mathbf{x} - \boldsymbol{\mu}_k\|^2 \right\}}{\sum_j \pi_j \exp \left\{ -\frac{1}{2\epsilon} \|\mathbf{x} - \boldsymbol{\mu}_j\|^2 \right\}} \\ &= \left(\frac{\sum_j \pi_j \exp \left\{ -\frac{1}{2\epsilon} \|\mathbf{x} - \boldsymbol{\mu}_j\|^2 \right\}}{\pi_k \exp \left\{ -\frac{1}{2\epsilon} \|\mathbf{x} - \boldsymbol{\mu}_k\|^2 \right\}} \right)^{-1} \\ &= \left(\sum_j \frac{\pi_j}{\pi_k} \frac{(\exp \left\{ -\|\mathbf{x} - \boldsymbol{\mu}_j\|^2 \right\})^{\frac{1}{2\epsilon}}}{(\exp \left\{ -\|\mathbf{x} - \boldsymbol{\mu}_k\|^2 \right\})^{\frac{1}{2\epsilon}}} \right)^{-1} \\ &= \left(\sum_j \frac{\pi_j}{\pi_k} \left(\frac{\exp \left\{ -\|\mathbf{x} - \boldsymbol{\mu}_j\|^2 \right\}}{\exp \left\{ -\|\mathbf{x} - \boldsymbol{\mu}_k\|^2 \right\}} \right)^{\frac{1}{2\epsilon}} \right)^{-1} \\ &= \left(1 + \sum_{j \neq k} \frac{\pi_j}{\pi_k} \left(\frac{\exp \left\{ -\|\mathbf{x} - \boldsymbol{\mu}_j\|^2 \right\}}{\exp \left\{ -\|\mathbf{x} - \boldsymbol{\mu}_k\|^2 \right\}} \right)^{\frac{1}{2\epsilon}} \right)^{-1} \end{aligned} \tag{139}$$

- ここで、 k^* を次で定める

$$k^* = \arg \min_j \|\mathbf{x} - \boldsymbol{\mu}_j\|^2 = \arg \max_j (-\|\mathbf{x} - \boldsymbol{\mu}_j\|^2) \quad (140)$$

$k = k^*$ であるとき、以下の、 $\epsilon \rightarrow 0$ による極限

$$\lim_{\epsilon \rightarrow 0} \left(\sum_{j \neq k} \frac{\pi_j}{\pi_k} \left(\frac{\exp \{ -\|\mathbf{x} - \boldsymbol{\mu}_j\|^2 \}}{\exp \{ -\|\mathbf{x} - \boldsymbol{\mu}_k\|^2 \}} \right)^{\frac{1}{2\epsilon}} \right) \quad (141)$$

を考えると、全ての $j \neq k^*$ について

$$\frac{\exp \{ -\|\mathbf{x} - \boldsymbol{\mu}_j\|^2 \}}{\exp \{ -\|\mathbf{x} - \boldsymbol{\mu}_k\|^2 \}} < 1 \quad (142)$$

が成立するので

$$\lim_{\epsilon \rightarrow 0} \left(\sum_{j \neq k} \frac{\pi_j}{\pi_k} \left(\frac{\exp \{ -\|\mathbf{x} - \boldsymbol{\mu}_j\|^2 \}}{\exp \{ -\|\mathbf{x} - \boldsymbol{\mu}_k\|^2 \}} \right)^{\frac{1}{2\epsilon}} \right) = 0 \quad (143)$$

である

- 従って、 $k = k^*$ のとき

$$\begin{aligned} & \lim_{\epsilon \rightarrow 0} \gamma(z_{ik}) \\ &= \lim_{\epsilon \rightarrow 0} \frac{\pi_k \exp \left\{ -\frac{1}{2\epsilon} \|\mathbf{x} - \boldsymbol{\mu}_k\|^2 \right\}}{\sum_j \pi_j \exp \left\{ -\frac{1}{2\epsilon} \|\mathbf{x} - \boldsymbol{\mu}_j\|^2 \right\}} \\ &= \lim_{\epsilon \rightarrow 0} \left(1 + \sum_{j \neq k} \frac{\pi_j}{\pi_k} \left(\frac{\exp \{ -\|\mathbf{x} - \boldsymbol{\mu}_j\|^2 \}}{\exp \{ -\|\mathbf{x} - \boldsymbol{\mu}_k\|^2 \}} \right)^{\frac{1}{2\epsilon}} \right)^{-1} \end{aligned}$$

$$= (1 + 0)^{-1} = 1 \quad (144)$$

から、 $\gamma(z_{ik}) \rightarrow 1$ ($\epsilon \rightarrow 0$) がいえる

- $k \neq k^*$ のとき

$$1 = \sum_k \gamma(z_{ik}) = \gamma(z_{ik^*}) + \sum_{k \neq k^*} \gamma(z_{ik}) \quad (145)$$

であって、両辺の $\epsilon \rightarrow 0$ による極限を取れば

$$\begin{aligned} 1 &= \lim_{\epsilon \rightarrow 0} \left(\gamma(z_{ik^*}) + \sum_{k \neq k^*} \gamma(z_{ik}) \right) \\ \Rightarrow 1 &= \lim_{\epsilon \rightarrow 0} \gamma(z_{ik^*}) + \lim_{\epsilon \rightarrow 0} \sum_{k \neq k^*} \gamma(z_{ik}) \\ \Rightarrow 1 &= 1 + \lim_{\epsilon \rightarrow 0} \sum_{k \neq k^*} \gamma(z_{ik}) \end{aligned} \quad (146)$$

となるから

$$\lim_{\epsilon \rightarrow 0} \sum_{k \neq k^*} \gamma(z_{ik}) = 0 \quad (147)$$

が明らかに成立するほか、以下の不等式が

$$0 \leq \gamma(z_{ik}) \leq \sum_{k \neq k^*} \gamma(z_{ik}) \quad (148)$$

$\gamma(z_{ik}) \geq 0$ ゆえ成立するので ($\gamma(z_{ik})$ は確率値)、両辺の $\epsilon \rightarrow 0$ による極限を再び取れば

$$0 \leq \lim_{\epsilon \rightarrow 0} \gamma(z_{ik}) \leq \lim_{\epsilon \rightarrow 0} \sum_{k \neq k^*} \gamma(z_{ik}) = 0 \quad (149)$$

従って、 $k \neq k^*$ の場合は

$$\lim_{\epsilon \rightarrow 0} \gamma(z_{ik}) = 0 \quad (150)$$

である

K-Means 法との関連

- これより、データ点 x_i に関する負担率 $\gamma(z_{ik})$ は、1 に収束する k^* 番目の負担率 $\gamma(z_{ik^*})$ を除き、全て 0 に収束する

$$\gamma(z_{ik}) \equiv p(z_{ik} = 1 | x_i) = \begin{cases} 1 & (k = k^* \text{ の場合}) \\ 0 & (\text{それ以外の場合}) \end{cases} \quad (151)$$

- これは、 k^* 番目のクラスタに**確率 1 で属する**ということ、即ち、クラスタ k^* への**ハード割り当て**を意味する
- $k^* = \arg \min_j \|x - \mu_j\|^2$ であるから、結局、各データ点は、**平均ベクトル μ への二乗ユークリッド距離が最小となるクラスタ**に割り当てることになる

K-Means 法との関連

- $\gamma(z_{ik})$ を r_{ik} に置き換えれば、EM アルゴリズムにおける μ_k の更新式は、K-Means における平均ベクトルの更新式に帰着

$$\text{K-Means :} \quad \mu_k = \frac{1}{\sum_i r_{ik}} \sum_i r_{ik} \mathbf{x}_i \quad (152)$$

$$\text{EM アルゴリズム :} \quad \mu_k = \frac{1}{\sum_i \gamma(z_{ik})} \sum_i \gamma(z_{ik}) \mathbf{x}_i \quad (153)$$

- 従って、混合ガウスモデルの EM アルゴリズムにおいて、各ガウス分布の共分散行列を $\epsilon \mathbf{I}$ としたとき、 $\epsilon \rightarrow 0$ の極限を取ると、K-Means 法が得られる

- 期待完全データ対数尤度の計算

- $\mathbb{E}_{\mathbf{Z} \sim p(\mathbf{Z}|\mathbf{X}, \boldsymbol{\theta})} [\ln p(\mathbf{X}, \mathbf{Z}|\boldsymbol{\theta})]$ を計算する
- 完全データ対数尤度関数 $\ln p(\mathbf{X}, \mathbf{Z}|\boldsymbol{\theta})$ の、事後確率 $p(\mathbf{Z}|\mathbf{X}, \boldsymbol{\theta})$ による期待値
- 次のように計算する

$$\begin{aligned} & \mathbb{E}_{\mathbf{Z} \sim p(\mathbf{Z}|\mathbf{X}, \boldsymbol{\theta})} [\ln p(\mathbf{X}, \mathbf{Z}|\boldsymbol{\theta})] \\ &= \sum_i \sum_k \gamma(z_{ik}) (\ln \pi_k + \ln \mathcal{N}(\mathbf{x}_i | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)) \\ &= \sum_i \sum_k \gamma(z_{ik}) (\ln \pi_k + \ln \mathcal{N}(\mathbf{x}_i | \boldsymbol{\mu}_k, \epsilon \mathbf{I})) \\ &= \sum_i \sum_k \gamma(z_{ik}) \left(\ln \pi_k + \ln \left(\frac{1}{(2\pi\epsilon)^{\frac{D}{2}}} \exp \left\{ -\frac{1}{2\epsilon} \|\mathbf{x}_i - \boldsymbol{\mu}_k\|^2 \right\} \right) \right) \\ &= \sum_i \sum_k \gamma(z_{ik}) \left(\ln \pi_k - \frac{D}{2} \ln(2\pi\epsilon) - \frac{1}{2\epsilon} \|\mathbf{x}_i - \boldsymbol{\mu}_k\|^2 \right) \quad (154) \end{aligned}$$

- 両辺に ϵ を掛けると

$$\begin{aligned} & \epsilon \cdot \mathbb{E}_{\mathbf{Z} \sim p(\mathbf{Z}|\mathbf{X}, \boldsymbol{\theta})} [\ln p(\mathbf{X}, \mathbf{Z}|\boldsymbol{\theta})] \\ &= \sum_i \sum_k \gamma(z_{ik}) \left(\epsilon \ln \pi_k - \right. \\ & \quad \left. \frac{D}{2} \epsilon \ln(2\pi\epsilon) - \frac{1}{2} \|\mathbf{x}_i - \boldsymbol{\mu}_k\|^2 \right) \end{aligned} \quad (155)$$

- $\epsilon \rightarrow 0$ の極限を取ると

$$\gamma(z_{ik}) \rightarrow r_{ik}, \quad \epsilon \ln \pi_k \rightarrow 0, \quad \epsilon \ln(2\pi\epsilon) \rightarrow 0 \quad (156)$$

であるから

$$\begin{aligned} & \lim_{\epsilon \rightarrow 0} \epsilon \cdot \mathbb{E}_{\mathbf{Z} \sim p(\mathbf{Z}|\mathbf{X}, \boldsymbol{\theta})} [\ln p(\mathbf{X}, \mathbf{Z}|\boldsymbol{\theta})] \\ &= \sum_i \sum_k r_{ik} \left(-\frac{1}{2} \|\mathbf{x}_i - \boldsymbol{\mu}_k\|^2 \right) \end{aligned}$$

K-Means 法との関連

$$= -\frac{1}{2} \sum_i \sum_k r_{ik} \|\mathbf{x}_i - \boldsymbol{\mu}_k\|^2 \quad (157)$$

$$= -J \quad (158)$$

- よって、期待完全データ対数尤度 $\mathbb{E}_{\mathbf{Z}} [\ln p(\mathbf{X}, \mathbf{Z} | \boldsymbol{\theta})]$ の最大化は、
K-Means における目的関数 J の最小化と同等である

K-Means 法との関連

- その他のパラメータ

- K-Means 法では、各クラスタの分散は推定しない
- 実際に、混合ガウスモデルにおいて、各クラスタの共分散行列は ϵI で固定した

- 混合ガウスモデルの混合係数 π_k の更新式は、次のようであった

$$\pi_k = \frac{\sum_i \gamma(z_{ik})}{N} \quad (159)$$

$\epsilon \rightarrow 0$ の極限においては、 $\gamma(z_{ik}) \rightarrow r_{ik}$ であるから

$$\pi_k = \frac{\sum_i r_{ik}}{N} = \frac{N_k}{N} \quad (160)$$

- これは、 π_k の値を、 k 番目のクラスタに割り当てられる、データ数の割合に設定することを意味している
- π_k の値は K-Means 法においては、もはや何の意味も持たない

K-Means 法との関連

- ここまでの話の流れ

- K-Means 法は、混合ガウス分布に対する EM アルゴリズムの、ある極限として得られることが分かった

- これからの話の流れ

- $\ln p(\mathbf{X}|\boldsymbol{\theta})$ の代わりに、 $\ln p(\mathbf{X}, \mathbf{Z}|\boldsymbol{\theta})$ を最大化してもよい根拠を明らかにする
- $\ln p(\mathbf{X}, \mathbf{Z}|\boldsymbol{\theta})$ の、 $p(\mathbf{Z}|\mathbf{X}, \boldsymbol{\theta})$ による期待値を取る理由を明らかにする
- E ステップと M ステップが、確かに対数尤度関数 $\ln p(\mathbf{X}|\boldsymbol{\theta})$ を増加させることを証明する
- これらの解明のために、一般的な EM アルゴリズムの取り扱いについて調べる

3 EM アルゴリズム

- EM アルゴリズムの解釈
- 混合ガウス分布の再解釈
- K-Means 法との関連
- 一般の EM アルゴリズム
- MAP 推定に対する EM アルゴリズム
- EM アルゴリズムの拡張
- EM アルゴリズムに関する補足
- 逐次型の EM アルゴリズム

一般の EM アルゴリズム

- EM アルゴリズムの目的 (再掲)
 - 潜在変数をもつ確率モデルについて、パラメータの最尤解を求める
- 一般的な EM アルゴリズムの取り扱い
 - これまでは、混合ガウスモデルに対して、EM アルゴリズムを発見的に導いた
 - ここでは、EM アルゴリズムが、確かに尤度関数 $\ln p(\mathbf{X}|\boldsymbol{\theta})$ を極大化することを証明する
 - 後述する変分推論の基礎をなす部分
- 尤度関数 $p(\mathbf{X}|\boldsymbol{\theta})$ の記述
 - 全ての観測変数と、潜在変数をそれぞれ \mathbf{X}, \mathbf{Z} と表す
 - 確率モデルの全てのパラメータの組を、 $\boldsymbol{\theta}$ と表す

一般の EM アルゴリズム

- 同時確率分布を $p(\mathbf{X}, \mathbf{Z}|\boldsymbol{\theta})$ とすると、尤度関数は次のようになる

$$p(\mathbf{X}|\boldsymbol{\theta}) = \sum_{\mathbf{Z}} p(\mathbf{X}, \mathbf{Z}|\boldsymbol{\theta}) \quad (161)$$

- 連続潜在変数の場合は、次のように、総和を積分に置き換えればよい

$$p(\mathbf{X}|\boldsymbol{\theta}) = \int_{\mathbf{Z}} p(\mathbf{X}, \mathbf{Z}|\boldsymbol{\theta}) d\mathbf{Z} \quad (162)$$

- ここでは、連続潜在変数の場合を考える
- 重要な仮定**
 - $\ln p(\mathbf{X}|\boldsymbol{\theta})$ の最大化よりも、完全データ対数尤度関数 $\ln p(\mathbf{X}, \mathbf{Z}|\boldsymbol{\theta})$ の最大化の方が、容易である
 - 以前に見た尤度関数 $\ln p(\mathbf{X}|\boldsymbol{\theta})$ では、対数の中に総和が含まれており (**log-sum**)、複雑な形をしていた

一般の EM アルゴリズム

- Z についての情報を加えることで、尤度関数から log-sum の構造を消すことができた
- 対数 \ln がガウス分布に直接作用するようになったため、尤度関数の形が簡単になった

• EM アルゴリズムで行うこと

- $\ln p(\mathbf{X}|\boldsymbol{\theta})$ ではなく $\ln p(\mathbf{X}, \mathbf{Z}|\boldsymbol{\theta})$ を最適化しようとしたが、 Z に関する情報がないので、それはできない
- そこで、事後確率 $p(\mathbf{Z}|\mathbf{X}, \boldsymbol{\theta})$ による、 $\ln p(\mathbf{X}, \mathbf{Z}|\boldsymbol{\theta})$ の期待値 $\mathbb{E}_Z [\ln p(\mathbf{X}, \mathbf{Z}|\boldsymbol{\theta})]$ を最大化する
- これ以降の議論のために、**イェンセンの不等式**、**エントロピー**、**KL ダイバージェンス**について確認しておく

一般の EM アルゴリズム

- イェンセンの不等式

- 凸関数 $f(x)$ は、任意の点集合 $\{x_i\}$ について以下を満たす

$$f\left(\sum_i \lambda_i x_i\right) \leq \sum_i \lambda_i f(x_i) \quad (163)$$

- ここで、 $\lambda_i \geq 0$, $\sum_i \lambda_i = 1$ であるとする
- λ_i を、値 $\{x_i\}$ を取る離散確率変数 x 上の確率分布 $p(x)$ と解釈すると

$$\begin{aligned} f\left(\sum_i p(x_i) x_i\right) &\leq \sum_i p(x_i) f(x_i) \\ f(\mathbb{E}[x]) &\leq \mathbb{E}[f(x)] \end{aligned} \quad (164)$$

一般の EM アルゴリズム

- x が連続変数であれば、イェンセンの不等式は次のように書ける

$$f\left(\int \boldsymbol{x} p(\boldsymbol{x}) d\boldsymbol{x}\right) \leq \int f(\boldsymbol{x}) p(\boldsymbol{x}) d\boldsymbol{x} \quad (165)$$

例えば、 $f(x) = -\ln x$ は凸関数であるから

$$-\ln\left(\int \boldsymbol{x} p(\boldsymbol{x}) d\boldsymbol{x}\right) \leq \int (-\ln \boldsymbol{x}) p(\boldsymbol{x}) d\boldsymbol{x} \quad (166)$$

よって

$$\ln\left(\int \boldsymbol{x} p(\boldsymbol{x}) d\boldsymbol{x}\right) \geq \int (\ln \boldsymbol{x}) p(\boldsymbol{x}) d\boldsymbol{x} \quad (167)$$

一般の EM アルゴリズム

- エントロピー

- 確率分布 $p(\boldsymbol{x})$ について、エントロピーは以下で定義される

$$H[p] = - \int p(\boldsymbol{x}) \ln p(\boldsymbol{x}) d\boldsymbol{x} \quad (168)$$

- エントロピーは、確率分布 $p(\boldsymbol{x})$ を入力として、上記の量を返す、**汎関数** (Functional) である
- 汎関数とは、入力として関数を取り、出力として汎関数の値を返すものである

一般の EM アルゴリズム

● KL ダイバージェンス

- 確率分布 $p(\boldsymbol{x})$ と $q(\boldsymbol{x})$ の間の、カルバック-ライブラーダイバージェンスを、 $\text{KL}(p||q)$ と表す
- 確率分布 $p(\boldsymbol{x})$ と $q(\boldsymbol{x})$ の間の、(擬似的な) 距離を表す指標である

$$\text{KL}(p||q) = - \int p(\boldsymbol{x}) \ln \left\{ \frac{q(\boldsymbol{x})}{p(\boldsymbol{x})} \right\} d\boldsymbol{x} \quad (169)$$

- $\text{KL}(p||q) \geq 0$ であり、等号成立は $p(\boldsymbol{x}) = q(\boldsymbol{x})$ のときに限る
- 2つの分布が完全に同一であれば、KL ダイバージェンスは 0 で最小値を取る
- また厳密には距離ではないため、対称性は成立しない
- 従って、一般に $\text{KL}(p||q) \neq \text{KL}(q||p)$ となる

一般の EM アルゴリズム

- $\ln p(\mathbf{X}|\boldsymbol{\theta})$ の分解

- EM アルゴリズムについて考察するために、まずは $\ln p(\mathbf{X}|\boldsymbol{\theta})$ を分解してみよう
- 潜在変数についての分布を $q(\mathbf{Z})$ とおく
- $q(\mathbf{Z})$ の設定の仕方によらず、 $\ln p(\mathbf{X}|\boldsymbol{\theta})$ を次のように分解できる

$$\ln p(\mathbf{X}|\boldsymbol{\theta}) = \mathcal{L}(q, \boldsymbol{\theta}) + \text{KL}(q||p) \quad (170)$$

- $\mathcal{L}(q, \boldsymbol{\theta})$ は、分布 $q(\mathbf{Z})$ の汎関数であり、かつパラメータ $\boldsymbol{\theta}$ の関数である
- $\text{KL}(q||p)$ は、確率分布 $q(\mathbf{Z})$ と $p(\mathbf{X}|\boldsymbol{\theta})$ の間の、**KL ダイバージェンス**である

一般の EM アルゴリズム

- 分解は次のように行える

$$\begin{aligned}\ln p(\mathbf{X}|\boldsymbol{\theta}) &= \underbrace{\left(\sum_{\mathbf{Z}} q(\mathbf{Z}) \right)}_{=1} \ln p(\mathbf{X}|\boldsymbol{\theta}) \\&= \sum_{\mathbf{Z}} q(\mathbf{Z}) \ln p(\mathbf{X}|\boldsymbol{\theta}) \\&= \sum_{\mathbf{Z}} q(\mathbf{Z}) \ln \frac{p(\mathbf{X}, \mathbf{Z}|\boldsymbol{\theta})}{p(\mathbf{Z}|\mathbf{X}, \boldsymbol{\theta})} \\&= \sum_{\mathbf{Z}} q(\mathbf{Z}) \ln \left(\frac{p(\mathbf{X}, \mathbf{Z}|\boldsymbol{\theta})}{q(\mathbf{Z})} \frac{q(\mathbf{Z})}{p(\mathbf{Z}|\mathbf{X}, \boldsymbol{\theta})} \right) \\&= \sum_{\mathbf{Z}} q(\mathbf{Z}) \ln \frac{p(\mathbf{X}, \mathbf{Z}|\boldsymbol{\theta})}{q(\mathbf{Z})} - \sum_{\mathbf{Z}} q(\mathbf{Z}) \ln \frac{p(\mathbf{Z}|\mathbf{X}, \boldsymbol{\theta})}{q(\mathbf{Z})} \\&= \mathcal{L}(q, \boldsymbol{\theta}) + \text{KL}(q||p)\end{aligned}\tag{171}$$

一般の EM アルゴリズム

- ここで $\mathcal{L}(q, \theta)$ と $\text{KL}(q||p)$ は以下のように定義した

$$\mathcal{L}(q, \theta) = \sum_{\mathbf{Z}} q(\mathbf{Z}) \ln \frac{p(\mathbf{X}, \mathbf{Z} | \theta)}{q(\mathbf{Z})} \quad (172)$$

$$\text{KL}(q||p) = - \sum_{\mathbf{Z}} q(\mathbf{Z}) \ln \frac{p(\mathbf{Z} | \mathbf{X}, \theta)}{q(\mathbf{Z})} \quad (173)$$

- $\text{KL}(q||p) \geq 0$ ゆえ、以下の不等式を得る

$$\mathcal{L}(q, \theta) \leq \ln p(\mathbf{X} | \theta) \quad (174)$$

- $\mathcal{L}(q, \theta)$ は、 $q(\mathbf{Z}), \theta$ によらず、常に $\ln p(\mathbf{X} | \theta)$ の下界をなす
- 下界について確認した後に、EM アルゴリズムの各ステップについて見ていく

一般の EM アルゴリズム

Figure 9.11 Illustration of the decomposition given by (9.70), which holds for any choice of distribution $q(\mathbf{Z})$. Because the Kullback-Leibler divergence satisfies $\text{KL}(q||p) \geq 0$, we see that the quantity $\mathcal{L}(q, \theta)$ is a lower bound on the log likelihood function $\ln p(\mathbf{X}|\theta)$.

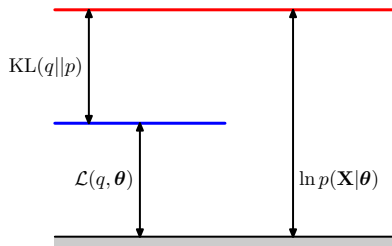


図 4: $\ln p(\mathbf{X}|\theta)$ の分解

下界と上界の定義

- 下界と上界の定義

- **下界**と**上界**は、次のように定義される (実数の集合を \mathbb{R} とする)

実数 $a \in \mathbb{R}$ が、実数の部分集合 $A \subset \mathbb{R}$ の上界である
 \Leftrightarrow 集合 A に属する任意の要素 x は、 a 以下である
 $\Leftrightarrow \forall x \in A, x \leq a$ (175)

実数 $a \in \mathbb{R}$ が、実数の部分集合 $A \subset \mathbb{R}$ の下界である
 \Leftrightarrow 集合 A に属する任意の要素 x は、 a 以上である
 $\Leftrightarrow \forall x \in A, a \leq x$ (176)

- 上記の定義は、実数の集合 \mathbb{R} だけでなく、一般の半順序集合について成り立つ
- 例えば、集合 A を $A = \{x | 0 \leq x\}$ と定義すると、0 以下である実数、例えば -1 、 -3 などは、いずれも集合 A の下界となる

一般の EM アルゴリズム

- EM アルゴリズムの概要

- EM アルゴリズムでは、 $\ln p(\mathbf{X}|\boldsymbol{\theta})$ の最尤解を求めるために、**E ステップ**と **M ステップ**の二段階の処理を、交互に繰り返す
- パラメータの現在値を $\boldsymbol{\theta}^{\text{old}}$ とする

- **E ステップ**

- E ステップでは、下界 $\mathcal{L}(q, \boldsymbol{\theta}^{\text{old}})$ を、 $\boldsymbol{\theta}^{\text{old}}$ を固定しながら、 $q(\mathbf{Z})$ について最大化する
- この問題は、 $\ln p(\mathbf{X}|\boldsymbol{\theta})$ の分解をみれば簡単に解ける

$$\begin{aligned} & \ln p(\mathbf{X}|\boldsymbol{\theta}^{\text{old}}) \\ = & \mathcal{L}(q, \boldsymbol{\theta}^{\text{old}}) + \text{KL}(q||p) \end{aligned} \tag{177}$$

$$\begin{aligned} = & \mathcal{L}(q, \boldsymbol{\theta}^{\text{old}}) + \text{KL}(q(\mathbf{Z})||p(\mathbf{Z}|\mathbf{X}, \boldsymbol{\theta}^{\text{old}})) \\ & (\text{E ステップ前}) \end{aligned} \tag{178}$$

一般の EM アルゴリズム

- 上式において、左辺の $\ln p(\mathbf{X}|\boldsymbol{\theta}^{\text{old}})$ は、 q には依存しない定数である
- 従って、 q について $\mathcal{L}(q, \boldsymbol{\theta}^{\text{old}})$ を最大化するためには、 $\text{KL}(q||p)$ を最小化するしかない
- $\text{KL}(q||p)$ を最小化するためには、 $q(\mathbf{Z}) = p(\mathbf{Z}|\mathbf{X}, \boldsymbol{\theta}^{\text{old}})$ とおいて、 $\text{KL}(q||p) = 0$ とすればよい ($\text{KL}(q||p) \geq 0$ であるから、最小値は 0)
- このとき、下界 $\mathcal{L}(q, \boldsymbol{\theta}^{\text{old}})$ は、対数尤度 $\ln p(\mathbf{X}|\boldsymbol{\theta}^{\text{old}})$ に一致する

$$\ln p(\mathbf{X}|\boldsymbol{\theta}^{\text{old}}) \tag{179}$$

$$\begin{aligned} &= \mathcal{L}(q, \boldsymbol{\theta}^{\text{old}}) + \text{KL}(q(\mathbf{Z})||p(\mathbf{Z}|\mathbf{X}, \boldsymbol{\theta}^{\text{old}})) \\ &= \mathcal{L}(q, \boldsymbol{\theta}^{\text{old}}) + \text{KL}(p(\mathbf{Z}|\mathbf{X}, \boldsymbol{\theta}^{\text{old}})||p(\mathbf{Z}|\mathbf{X}, \boldsymbol{\theta}^{\text{old}})) \end{aligned} \tag{180}$$

$$\begin{aligned} &= \mathcal{L}(q, \boldsymbol{\theta}^{\text{old}}) \\ &\quad (\text{E ステップ後}) \end{aligned} \tag{181}$$

一般の EM アルゴリズム

- 次の図 5 には E ステップの概要が示されている
- $KL(q||p) = 0$ となるように q を調節している
- 青線（図 5 の青線）で示されている下界 $\mathcal{L}(q, \theta^{\text{old}})$ が、赤線（図 5 の赤線）で示されている対数尤度 $\ln p(\mathbf{X}|\theta^{\text{old}})$ のところまで、持ち上げられている

一般の EM アルゴリズム

Figure 9.12 Illustration of the E step of the EM algorithm. The q distribution is set equal to the posterior distribution for the current parameter values θ^{old} , causing the lower bound to move up to the same value as the log likelihood function, with the KL divergence vanishing.

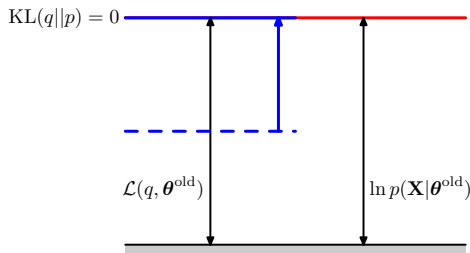


図 5: EM アルゴリズムの E ステップ

一般の EM アルゴリズム

● M ステップ

- M ステップでは、下界 $\mathcal{L}(q, \theta)$ を、分布 $q(\mathbf{Z})$ を固定しながら、 θ について最大化し、新たなパラメータ θ^{new} を得る
- M ステップは下界 \mathcal{L} を増加させるが、 $\text{KL}(q||p) \geq 0$ であるから、対数尤度 $\ln p(\mathbf{X}|\theta)$ も必然的に増加する

$$\begin{aligned} & \ln p(\mathbf{X}|\theta) \\ = & \mathcal{L}(q, \theta) + \text{KL}(q||p) \end{aligned} \tag{182}$$

$$= \mathcal{L}(q, \theta) + \text{KL}(q(\mathbf{Z})||p(\mathbf{Z}|\mathbf{X}, \theta)) \tag{183}$$

$$\begin{aligned} = & \mathcal{L}(q, \theta) + \text{KL}(p(\mathbf{Z}|\mathbf{X}, \theta^{\text{old}})||p(\mathbf{Z}|\mathbf{X}, \theta)) \\ & \text{(M ステップ前)} \end{aligned} \tag{184}$$

- 分布 $q(\mathbf{Z}) = p(\mathbf{Z}|\mathbf{X}, \theta^{\text{old}})$ は、古いパラメータ θ^{old} によって決められており、**M ステップの間は固定**されている

一般の EM アルゴリズム

- $KL(q||p)$ は、 $q(\mathbf{Z}) = p(\mathbf{Z}|\mathbf{X}, \boldsymbol{\theta}^{\text{old}})$ と $p(\mathbf{Z}|\mathbf{X}, \boldsymbol{\theta})$ との KL ダイバージェンスである
- $q(\mathbf{Z}) = p(\mathbf{Z}|\mathbf{X}, \boldsymbol{\theta}^{\text{old}})$ と、M ステップ後の新しい事後分布 $p(\mathbf{Z}|\mathbf{X}, \boldsymbol{\theta}^{\text{new}})$ とは一致しないため、 $KL(q||p) > 0$ となる

$$\begin{aligned} & \ln p(\mathbf{X}|\boldsymbol{\theta}^{\text{new}}) \\ = & \mathcal{L}(q, \boldsymbol{\theta}^{\text{new}}) + KL(p(\mathbf{Z}|\mathbf{X}, \boldsymbol{\theta}^{\text{old}})||p(\mathbf{Z}|\mathbf{X}, \boldsymbol{\theta}^{\text{new}})) \quad (185) \\ & \text{(M ステップ後)} \end{aligned}$$

- 対数尤度の増加量は、下界 \mathcal{L} の増加量よりも大きくなる (図 6)

$$\begin{aligned} & \ln p(\mathbf{X}|\boldsymbol{\theta}^{\text{new}}) - \ln p(\mathbf{X}|\boldsymbol{\theta}^{\text{old}}) \\ = & \mathcal{L}(q, \boldsymbol{\theta}^{\text{new}}) + KL(p(\mathbf{Z}|\mathbf{X}, \boldsymbol{\theta}^{\text{old}})||p(\mathbf{Z}|\mathbf{X}, \boldsymbol{\theta}^{\text{new}})) - \\ & \mathcal{L}(q, \boldsymbol{\theta}^{\text{old}}) \quad (186) \end{aligned}$$

$$\begin{aligned} = & (\mathcal{L}(q, \boldsymbol{\theta}^{\text{new}}) - \mathcal{L}(q, \boldsymbol{\theta}^{\text{old}})) + \\ & KL(p(\mathbf{Z}|\mathbf{X}, \boldsymbol{\theta}^{\text{old}})||p(\mathbf{Z}|\mathbf{X}, \boldsymbol{\theta}^{\text{new}})) \quad (187) \end{aligned}$$

$$\geq \mathcal{L}(q, \boldsymbol{\theta}^{\text{new}}) - \mathcal{L}(q, \boldsymbol{\theta}^{\text{old}}) \quad (188)$$

一般の EM アルゴリズム

- 次の図 6 には M ステップの概要が示されている
- 下界 $\mathcal{L}(q, \theta)$ を、 $q(\mathbf{Z})$ を固定しつつ、 θ について最大化している
- 青の点線で示されている下界 $\mathcal{L}(q, \theta^{\text{old}})$ が、青の実線で示されている下界 $\mathcal{L}(q, \theta^{\text{new}})$ へと、持ち上げられている
- 赤の点線で示される対数尤度 $\ln p(\mathbf{X}|\theta^{\text{old}})$ は、赤の実線で示される対数尤度 $\ln p(\mathbf{X}|\theta^{\text{new}})$ へと、持ち上げられている
- 新たに生じた $\text{KL}(q||p)$ によって、対数尤度の増加量は、下界 \mathcal{L} の増加量よりも大きくなっている

一般の EM アルゴリズム

Figure 9.13 Illustration of the M step of the EM algorithm. The distribution $q(\mathbf{Z})$ is held fixed and the lower bound $\mathcal{L}(q, \theta)$ is maximized with respect to the parameter vector θ to give a revised value θ^{new} . Because the KL divergence is nonnegative, this causes the log likelihood $\ln p(\mathbf{X}|\theta)$ to increase by at least as much as the lower bound does.

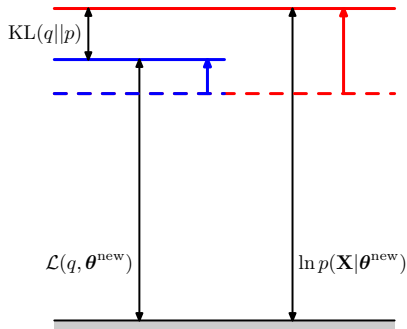


図 6: EM アルゴリズムの M ステップ

一般の EM アルゴリズム

- M ステップで最大化される量

- M ステップでは下界 $\mathcal{L}(q, \theta)$ を、 q を固定しつつ θ について最大化する
- M ステップで最大化するのは、**E ステップ後の下界** $\mathcal{L}(q, \theta)$ であり、これは次のように表せる ($q(\mathbf{Z}) = p(\mathbf{X}, \mathbf{Z}|\theta^{\text{old}})$ である)

$$\begin{aligned} & \mathcal{L}(q, \theta) \\ = & \sum_{\mathbf{Z}} q(\mathbf{Z}) \ln \frac{p(\mathbf{X}, \mathbf{Z}|\theta)}{q(\mathbf{Z})} \\ = & \sum_{\mathbf{Z}} p(\mathbf{X}, \mathbf{Z}|\theta^{\text{old}}) \ln \frac{p(\mathbf{X}, \mathbf{Z}|\theta)}{p(\mathbf{X}, \mathbf{Z}|\theta^{\text{old}})} \end{aligned} \quad (189)$$

$$\begin{aligned} = & \sum_{\mathbf{Z}} p(\mathbf{X}, \mathbf{Z}|\theta^{\text{old}}) \ln p(\mathbf{X}, \mathbf{Z}|\theta) - \\ & \sum_{\mathbf{Z}} p(\mathbf{X}, \mathbf{Z}|\theta^{\text{old}}) \ln p(\mathbf{X}, \mathbf{Z}|\theta^{\text{old}}) \end{aligned} \quad (190)$$

$$= \mathcal{Q}(\theta, \theta^{\text{old}}) + \text{Const.} \quad (191)$$

一般の EM アルゴリズム

- 定数項は、単に分布 $q(\mathbf{Z}) = p(\mathbf{X}, \mathbf{Z}|\boldsymbol{\theta}^{\text{old}})$ のエントロピーであって、 $\boldsymbol{\theta}$ には依存しないため無視できる
- M ステップで最大化される量は、結局、完全データ対数尤度関数 $\ln p(\mathbf{X}, \mathbf{Z}|\boldsymbol{\theta})$ の、事後確率分布 $p(\mathbf{Z}|\mathbf{X}, \boldsymbol{\theta}^{\text{old}})$ による期待値 $Q(\boldsymbol{\theta}, \boldsymbol{\theta}^{\text{old}})$ である
- 最適化しようとしているパラメータ $\boldsymbol{\theta}$ は、**対数の中にしか現れない**
- 同時分布 $p(\mathbf{X}, \mathbf{Z}|\boldsymbol{\theta})$ に対して**対数が直接作用**するので、同時分布が例えばガウス分布であれば、対数と指数が打ち消されて、簡単な形になる
- その結果として、不完全データ対数尤度関数 $\ln p(\mathbf{X}|\boldsymbol{\theta})$ の最適化よりも、**非常に単純な手続きとなる**

一般の EM アルゴリズム

- E ステップのまとめ

- 下界 $\mathcal{L}(q, \theta^{\text{old}})$ を、 θ^{old} を固定しつつ、 q について最大化する
- これは、単に $q(\mathbf{Z}) = p(\mathbf{Z}|\mathbf{X}, \theta^{\text{old}})$ とすればよい
- 即ち、 $p(\mathbf{Z}|\mathbf{X}, \theta^{\text{old}})$ を計算するだけである

- M ステップのまとめ

- 下界 $\mathcal{L}(q, \theta)$ を、 q を固定しつつ、 θ について最大化する
- これは、期待値 $\mathcal{Q}(\theta, \theta^{\text{old}})$ を最大化するような、パラメータ θ を求めることに相当

● 疑問に対する答え

- $\ln p(\mathbf{X}|\boldsymbol{\theta})$ の代わりに、 $\ln p(\mathbf{X}, \mathbf{Z}|\boldsymbol{\theta})$ を最大化してよい根拠
- そして、 $\ln p(\mathbf{X}, \mathbf{Z}|\boldsymbol{\theta})$ の、 $p(\mathbf{Z}|\mathbf{X}, \boldsymbol{\theta})$ による期待値を取る理由
- 期待値を取る操作は、式の導出の中で、極めて自然に現れた
- 期待値 $Q(\boldsymbol{\theta}, \boldsymbol{\theta}^{\text{old}})$ の最大化は、 $\mathcal{L}(q, \boldsymbol{\theta})$ の最大化と等価である
- $\mathcal{L}(q, \boldsymbol{\theta})$ は、 q や $\boldsymbol{\theta}$ によらず、常に $\ln p(\mathbf{X}|\boldsymbol{\theta})$ の下界である
- 下界を最大化することは、 $\ln p(\mathbf{X}|\boldsymbol{\theta})$ を徐々に大きくしていくことにつながる (図 5 と図 6 を参照)
- これらより、 $\ln p(\mathbf{X}, \mathbf{Z}|\boldsymbol{\theta})$ の期待値を最適化させることは、 $\ln p(\mathbf{X}|\boldsymbol{\theta})$ を最適化させることと等価

一般の EM アルゴリズム

- パラメータの更新によって $\ln p(\mathbf{X}|\boldsymbol{\theta})$ が常に大きくなることの補足
 - 以下のように式変形を行う

$$\begin{aligned} & (\text{M ステップ後の } \ln p(\mathbf{X}|\boldsymbol{\theta})) - (\text{E ステップ後の } \ln p(\mathbf{X}|\boldsymbol{\theta})) \\ &= \ln p(\mathbf{X}|\boldsymbol{\theta}^{\text{new}}) - \ln p(\mathbf{X}|\boldsymbol{\theta}^{\text{old}}) \\ &= \ln \frac{p(\mathbf{X}|\boldsymbol{\theta}^{\text{new}})}{p(\mathbf{X}|\boldsymbol{\theta}^{\text{old}})} \\ &= \underbrace{\sum_{\mathbf{Z}} p(\mathbf{Z}|\mathbf{X}, \boldsymbol{\theta}^{\text{old}})}_{=1} \ln \frac{p(\mathbf{X}|\boldsymbol{\theta}^{\text{new}})}{p(\mathbf{X}|\boldsymbol{\theta}^{\text{old}})} \\ &= \sum_{\mathbf{Z}} p(\mathbf{Z}|\mathbf{X}, \boldsymbol{\theta}^{\text{old}}) \ln \frac{p(\mathbf{X}, \mathbf{Z}|\boldsymbol{\theta}^{\text{new}})}{p(\mathbf{Z}|\mathbf{X}, \boldsymbol{\theta}^{\text{new}})} \frac{p(\mathbf{Z}|\mathbf{X}, \boldsymbol{\theta}^{\text{old}})}{p(\mathbf{X}, \mathbf{Z}|\boldsymbol{\theta}^{\text{old}})} \\ &= \sum_{\mathbf{Z}} p(\mathbf{Z}|\mathbf{X}, \boldsymbol{\theta}^{\text{old}}) \ln \frac{p(\mathbf{X}, \mathbf{Z}|\boldsymbol{\theta}^{\text{new}})}{p(\mathbf{X}, \mathbf{Z}|\boldsymbol{\theta}^{\text{old}})} \frac{p(\mathbf{Z}|\mathbf{X}, \boldsymbol{\theta}^{\text{old}})}{p(\mathbf{Z}|\mathbf{X}, \boldsymbol{\theta}^{\text{new}})} \end{aligned}$$

一般の EM アルゴリズム

$$\begin{aligned} &= \sum_{\mathbf{Z}} p(\mathbf{Z}|\mathbf{X}, \boldsymbol{\theta}^{\text{old}}) \ln \frac{p(\mathbf{X}, \mathbf{Z}|\boldsymbol{\theta}^{\text{new}})}{p(\mathbf{X}, \mathbf{Z}|\boldsymbol{\theta}^{\text{old}})} + \\ &\quad \sum_{\mathbf{Z}} p(\mathbf{Z}|\mathbf{X}, \boldsymbol{\theta}^{\text{old}}) \ln \frac{p(\mathbf{Z}|\mathbf{X}, \boldsymbol{\theta}^{\text{old}})}{p(\mathbf{Z}|\mathbf{X}, \boldsymbol{\theta}^{\text{new}})} \\ &= \sum_{\mathbf{Z}} p(\mathbf{Z}|\mathbf{X}, \boldsymbol{\theta}^{\text{old}}) \ln p(\mathbf{X}, \mathbf{Z}|\boldsymbol{\theta}^{\text{new}}) - \\ &\quad \sum_{\mathbf{Z}} p(\mathbf{Z}|\mathbf{X}, \boldsymbol{\theta}^{\text{old}}) \ln p(\mathbf{X}, \mathbf{Z}|\boldsymbol{\theta}^{\text{old}}) - \\ &\quad \sum_{\mathbf{Z}} p(\mathbf{Z}|\mathbf{X}, \boldsymbol{\theta}^{\text{old}}) \ln \frac{p(\mathbf{Z}|\mathbf{X}, \boldsymbol{\theta}^{\text{old}})}{p(\mathbf{Z}|\mathbf{X}, \boldsymbol{\theta}^{\text{new}})} \\ &= \mathcal{Q}(\boldsymbol{\theta}^{\text{old}}, \boldsymbol{\theta}^{\text{new}}) - \mathcal{Q}(\boldsymbol{\theta}^{\text{old}}, \boldsymbol{\theta}^{\text{old}}) + \\ &\quad \text{KL} (p(\mathbf{Z}|\mathbf{X}, \boldsymbol{\theta}^{\text{old}}) || p(\mathbf{Z}|\mathbf{X}, \boldsymbol{\theta}^{\text{new}})) \end{aligned} \tag{192}$$

$$\geq \mathcal{Q}(\boldsymbol{\theta}^{\text{old}}, \boldsymbol{\theta}^{\text{new}}) - \mathcal{Q}(\boldsymbol{\theta}^{\text{old}}, \boldsymbol{\theta}^{\text{old}}) \tag{193}$$

$$\geq 0$$

一般の EM アルゴリズム

- 最後の変形は、M ステップでは $Q(\theta, \theta^{\text{old}})$ を、 θ について最大化しているから、 $Q(\theta^{\text{old}}, \theta^{\text{new}}) \geq Q(\theta^{\text{old}}, \theta^{\text{old}})$ であることを利用
- 更新によって $\ln p(\mathbf{X}|\theta)$ は、収束していない限り常に大きくなる

一般の EM アルゴリズム

- $\ln p(\mathbf{X}|\boldsymbol{\theta})$ の分解の導出の補足
 - イェンセンの不等式を用いて導出してみよう

$$\begin{aligned}\ln p(\mathbf{X}|\boldsymbol{\theta}) &= \ln \sum_{\mathbf{Z}} p(\mathbf{X}, \mathbf{Z}|\boldsymbol{\theta}) \\ &= \ln \sum_{\mathbf{Z}} q(\mathbf{Z}) \frac{p(\mathbf{X}, \mathbf{Z}|\boldsymbol{\theta})}{q(\mathbf{Z})} \\ &\geq \sum_{\mathbf{Z}} q(\mathbf{Z}) \ln \frac{p(\mathbf{X}, \mathbf{Z}|\boldsymbol{\theta})}{q(\mathbf{Z})} \\ &= \mathcal{L}(q, \boldsymbol{\theta})\end{aligned}\tag{194}$$

- 不等式の部分でイェンセンの不等式 $\log(\mathbb{E}[x]) \leq \mathbb{E}[\log x]$ を用いた

一般の EM アルゴリズム

- これより、 $\ln p(\mathbf{X}|\boldsymbol{\theta})$ と $\mathcal{L}(q, \boldsymbol{\theta})$ の差を調べると

$$\begin{aligned} & \ln p(\mathbf{X}|\boldsymbol{\theta}) - \mathcal{L}(q, \boldsymbol{\theta}) \\ = & \ln p(\mathbf{X}|\boldsymbol{\theta}) - \sum_{\mathbf{Z}} q(\mathbf{Z}) \ln \frac{p(\mathbf{X}, \mathbf{Z}|\boldsymbol{\theta})}{q(\mathbf{Z})} \\ = & \underbrace{\sum_{\mathbf{Z}} q(\mathbf{Z}) \ln p(\mathbf{X}|\boldsymbol{\theta})}_{=1} - \sum_{\mathbf{Z}} q(\mathbf{Z}) \ln \frac{p(\mathbf{X}, \mathbf{Z}|\boldsymbol{\theta})}{q(\mathbf{Z})} \\ = & \sum_{\mathbf{Z}} q(\mathbf{Z}) \ln p(\mathbf{X}|\boldsymbol{\theta}) - \sum_{\mathbf{Z}} q(\mathbf{Z}) \ln \frac{p(\mathbf{Z}|\mathbf{X}, \boldsymbol{\theta})p(\mathbf{X}|\boldsymbol{\theta})}{q(\mathbf{Z})} \\ = & \sum_{\mathbf{Z}} q(\mathbf{Z}) \ln p(\mathbf{X}|\boldsymbol{\theta}) - \\ & \sum_{\mathbf{Z}} q(\mathbf{Z}) (\ln p(\mathbf{Z}|\mathbf{X}, \boldsymbol{\theta}) + \ln p(\mathbf{X}|\boldsymbol{\theta}) - \ln q(\mathbf{Z})) \\ = & - \sum_{\mathbf{Z}} q(\mathbf{Z}) (\ln p(\mathbf{Z}|\mathbf{X}, \boldsymbol{\theta}) - \ln q(\mathbf{Z})) \end{aligned}$$

一般の EM アルゴリズム

$$\begin{aligned} &= - \sum_{\mathbf{Z}} q(\mathbf{Z}) \ln \frac{p(\mathbf{Z}|\mathbf{X}, \boldsymbol{\theta})}{q(\mathbf{Z})} \\ &= \text{KL}(q||p) \end{aligned} \tag{195}$$

ゆえ、 $\text{KL}(q||p)$ となることが分かったので

$$\ln p(\mathbf{X}|\boldsymbol{\theta}) = \mathcal{L}(q, \boldsymbol{\theta}) + \text{KL}(q||p) \tag{196}$$

のように分解できることが分かる

一般の EM アルゴリズム

- パラメータ空間での図示

- EM アルゴリズムは、パラメータ空間でも視覚化できる (図 7)
- **赤の実線**は、最大化したい対象である、不完全データ対数尤度関数 $\ln p(\mathbf{X}|\boldsymbol{\theta})$ を表す

- E ステップ

- パラメータの初期値 $\boldsymbol{\theta}^{\text{old}}$ から始めて、最初の E ステップでは、潜在変数の事後確率分布 $p(\mathbf{Z}|\mathbf{X}, \boldsymbol{\theta})$ を計算
- このとき、**青の実線**で示す下界 $\mathcal{L}(q, \boldsymbol{\theta}^{\text{old}})$ が q について更新され、下界 \mathcal{L} は、 $\ln p(\mathbf{X}|\boldsymbol{\theta})$ と $\boldsymbol{\theta}^{\text{old}}$ において一致する
- 下界 \mathcal{L} の曲線は、 $\boldsymbol{\theta}^{\text{old}}$ において $\ln p(\mathbf{X}|\boldsymbol{\theta})$ と**接する**ことに注意する
- 下界 \mathcal{L} と対数尤度 $\ln p(\mathbf{X}|\boldsymbol{\theta})$ は、 $\boldsymbol{\theta}^{\text{old}}$ において**同じ勾配を持つ**

一般の EM アルゴリズム

- M ステップ

- 下界 \mathcal{L} が凹関数で、唯一の最大値をもつとする (例えば混合ガウスモデル)
- M ステップでは、下界 $\mathcal{L}(q, \theta)$ が θ について最大化されて、パラメータ θ^{new} が得られる

- 続く E ステップ

- 続く E ステップでは、緑の実線で示した下界 $\mathcal{L}(q, \theta^{\text{new}})$ が計算される
- 下界 $\mathcal{L}(q, \theta^{\text{new}})$ は、 $\ln p(\mathbf{X}|\theta)$ と θ^{new} で接する

一般の EM アルゴリズム

- 勾配が等しくなることについての証明
 - 以下の式の、 θ による微分を考えれば明らか

$$\begin{aligned} & \left. \frac{\partial}{\partial \theta} \ln p(\mathbf{X}|\theta) \right|_{\theta^{\text{old}}} \\ &= \left. \frac{\partial}{\partial \theta} \mathcal{L}(q, \theta) \right|_{\theta^{\text{old}}} + \left. \frac{\partial}{\partial \theta} \text{KL}(q||p) \right|_{\theta^{\text{old}}} \\ &= \left. \frac{\partial}{\partial \theta} \mathcal{L}(q, \theta) \right|_{\theta^{\text{old}}} \end{aligned} \tag{197}$$

- E ステップによって $\text{KL}(q||p)$ が最小化されるので、 θ による勾配も当然 0 になるはずである
- このとき、 $\ln p(\mathbf{X}|\theta)$ と $\mathcal{L}(q, \theta)$ の、 θ^{old} における微分値が等しくなる
- 従って、 θ^{old} において両者は接することが分かる
- 直感的には、次のように考えればよい

一般の EM アルゴリズム

- 両者が接していなければ、交差しているはずである
- このとき、対数尤度 $\ln p(\mathbf{X}|\boldsymbol{\theta})$ が、下界 \mathcal{L} を上回る ($\mathcal{L}(q, \boldsymbol{\theta}) > \ln p(\mathbf{X}|\boldsymbol{\theta})$) ような $\boldsymbol{\theta}$ が存在する
- これは、 $\text{KL}(q||p) < 0$ となる可能性があることを示し、従って有り得ないので、両者は接しているはず

一般の EM アルゴリズム

Figure 9.14 The EM algorithm involves alternately computing a lower bound on the log likelihood for the current parameter values and then maximizing this bound to obtain the new parameter values. See the text for a full discussion.

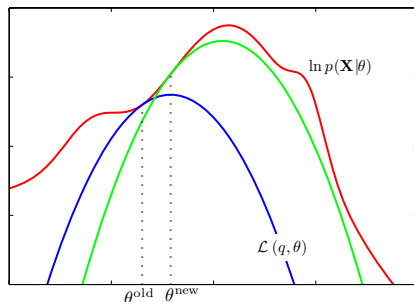


図 7: EM アルゴリズムの手続き

一般の EM アルゴリズム

- **i.i.d 標本**である場合

- データ点 x_i と、対応する潜在変数 z_i が、同一の確率分布 $p(x, z)$ から独立に得られている場合
- 以下のように同時分布 $p(\mathbf{X}, \mathbf{Z})$ を分解できる

$$p(\mathbf{X}, \mathbf{Z}) = \prod_i p(x_i, z_i) \quad (198)$$

- 従って、E ステップで計算される事後確率 $p(\mathbf{Z}|\mathbf{X}, \theta)$ は次のようになる

$$\begin{aligned} p(\mathbf{Z}|\mathbf{X}, \theta) &= \frac{p(\mathbf{X}, \mathbf{Z}|\theta)}{p(\mathbf{X}|\theta)} \\ &= \frac{p(\mathbf{X}, \mathbf{Z}|\theta)}{\sum_{\mathbf{Z}} p(\mathbf{X}, \mathbf{Z}|\theta)} \\ &= \frac{\prod_i p(x_i, z_i|\theta)}{\sum_{\mathbf{Z}} \prod_i p(x_i, z_i|\theta)} \end{aligned}$$

$$\begin{aligned} &= \frac{\prod_i p(\mathbf{x}_i, \mathbf{z}_i | \boldsymbol{\theta})}{\prod_i \sum_{\mathbf{z}} p(\mathbf{x}_i, \mathbf{z}_i | \boldsymbol{\theta})} \\ &= \frac{\prod_i p(\mathbf{x}_i, \mathbf{z}_i | \boldsymbol{\theta})}{\prod_i p(\mathbf{x}_i | \boldsymbol{\theta})} \\ &= \prod_i \frac{p(\mathbf{x}_i, \mathbf{z}_i | \boldsymbol{\theta})}{p(\mathbf{x}_i | \boldsymbol{\theta})} \\ &= \prod_i p(\mathbf{z}_i | \mathbf{x}_i, \boldsymbol{\theta}) \end{aligned} \tag{199}$$

各データ点に対する事後確率 $p(\mathbf{z}_i | \mathbf{x}_i, \boldsymbol{\theta})$ の積として、 $p(\mathbf{Z} | \mathbf{X}, \boldsymbol{\theta})$ を表現できた

- 例えば混合ガウスモデルであれば、データ点 \mathbf{x}_i に対する各ガウス分布の負担率は、データ \mathbf{x}_i とガウス分布のパラメータ $\boldsymbol{\theta}$ にのみ依存し、他のデータ点には依存しないことを示している

一般の EM アルゴリズム

- ここまでの話の流れ

- 一般的な EM アルゴリズムの取り扱いを調べた
- EM アルゴリズムに対する次の疑問を解決した
 - $\ln p(\mathbf{X}|\boldsymbol{\theta})$ の代わりに、 $\ln p(\mathbf{X}, \mathbf{Z}|\boldsymbol{\theta})$ を最大化してもよい根拠
 - $\ln p(\mathbf{X}, \mathbf{Z}|\boldsymbol{\theta})$ の、 $p(\mathbf{Z}|\mathbf{X}, \boldsymbol{\theta})$ による期待値を取る理由
 - E ステップと M ステップが、対数尤度関数 $\ln p(\mathbf{X}|\boldsymbol{\theta})$ を増加させる理由

- これからの話の流れ

- MAP 推定に対する EM アルゴリズムの適用を考える
- EM アルゴリズムの拡張 (一般化 EM アルゴリズム) について簡単に触れる
- 混合ガウスモデルについて、逐次型の EM アルゴリズムを導出する

3 EM アルゴリズム

- EM アルゴリズムの解釈
- 混合ガウス分布の再解釈
- K-Means 法との関連
- 一般の EM アルゴリズム
- MAP 推定に対する EM アルゴリズム
- EM アルゴリズムの拡張
- EM アルゴリズムに関する補足
- 逐次型の EM アルゴリズム

MAP 推定に対する EM アルゴリズム

- 事後分布の対数 $\ln p(\boldsymbol{\theta}|\mathbf{X})$ の最大化
 - 今までは、尤度関数 $\ln p(\mathbf{X}|\boldsymbol{\theta})$ の最適化を考えてきた
 - 即ち、**最尤推定に対する EM アルゴリズム**を考えてきた
- パラメータの事前分布 $p(\boldsymbol{\theta})$ を導入したモデルであれば、最尤推定だけでなく **MAP 推定**に対しても、EM アルゴリズムを使える
- MAP 推定とは、次式のように、事後分布 $p(\boldsymbol{\theta}|\mathbf{X})$ を最大化するパラメータ $\boldsymbol{\theta}^*$ を求める問題である

$$\boldsymbol{\theta}^* = \arg \max_{\boldsymbol{\theta}} p(\boldsymbol{\theta}|\mathbf{X}) \quad (200)$$

$$= \arg \max_{\boldsymbol{\theta}} \frac{p(\mathbf{X}, \boldsymbol{\theta})}{p(\mathbf{X})} \quad (201)$$

$$= \arg \max_{\boldsymbol{\theta}} \frac{p(\mathbf{X}|\boldsymbol{\theta})p(\boldsymbol{\theta})}{p(\mathbf{X})} \quad (202)$$

$$= \arg \max_{\boldsymbol{\theta}} p(\mathbf{X}|\boldsymbol{\theta})p(\boldsymbol{\theta}) \quad (203)$$

MAP 推定に対する EM アルゴリズム

- 事後分布の対数 $\ln p(\boldsymbol{\theta}|\mathbf{X})$ は

$$\ln p(\boldsymbol{\theta}|\mathbf{X}) = \ln \frac{p(\mathbf{X}, \boldsymbol{\theta})}{p(\mathbf{X})} \quad (204)$$

$$\begin{aligned} &= \ln \frac{p(\mathbf{X}|\boldsymbol{\theta})p(\boldsymbol{\theta})}{p(\mathbf{X})} \\ &= \ln p(\mathbf{X}|\boldsymbol{\theta}) + \ln p(\boldsymbol{\theta}) - p(\mathbf{X}) \end{aligned} \quad (205)$$

$$= \mathcal{L}(q, \boldsymbol{\theta}) + \text{KL}(q||p) + \ln p(\boldsymbol{\theta}) - p(\mathbf{X}) \quad (206)$$

$$\geq \mathcal{L}(q, \boldsymbol{\theta}) + \ln p(\boldsymbol{\theta}) - p(\mathbf{X}) \quad (207)$$

$$= \mathcal{L}(q, \boldsymbol{\theta}) + \ln p(\boldsymbol{\theta}) + \text{Const.} \quad (208)$$

- $\ln p(\mathbf{X})$ は定数とみなせるから、 $\ln p(\boldsymbol{\theta}|\mathbf{X})$ の最大化は、結局 $\mathcal{L}(q, \boldsymbol{\theta}) + \ln p(\boldsymbol{\theta})$ の最大化に相当する

MAP 推定に対する EM アルゴリズム

- MAP 推定に対する EM アルゴリズム

- **E ステップ**では、パラメータ θ を固定しつつ、 q について $\mathcal{L}(q, \theta)$ を最大化する
- q は下界 $\mathcal{L}(q, \theta)$ にしか現れないので、**通常の EM アルゴリズムと全く同様**である
- **M ステップ**では、分布 q を固定しつつ、パラメータ θ について $\mathcal{L}(q, \theta) + \ln p(\theta)$ を最大化する
- 事前分布の項 $\ln p(\theta)$ が現れているが、大抵は、通常の最尤推定に関する EM アルゴリズムと、少ししか変わらない

3 EM アルゴリズム

- EM アルゴリズムの解釈
- 混合ガウス分布の再解釈
- K-Means 法との関連
- 一般の EM アルゴリズム
- MAP 推定に対する EM アルゴリズム
- **EM アルゴリズムの拡張**
- EM アルゴリズムに関する補足
- 逐次型の EM アルゴリズム

EM アルゴリズムの拡張

- EM アルゴリズムに対する懸念

- EM アルゴリズムは、潜在的に困難である尤度関数 $\ln p(\mathbf{X}|\boldsymbol{\theta})$ の最大化を、**E ステップ**と **M ステップ**の2つに分解してくれる
- この2つのステップは多くの場合、実装が単純になる
- 但し、複雑なモデルに対しては、2つのどちらかのステップが、**依然として手に負えないかもしれない**

- 一般化 EM アルゴリズム

- 手に負えない M ステップに対処するためのアルゴリズム
- M ステップで、下界 $\mathcal{L}(q, \boldsymbol{\theta})$ を $\boldsymbol{\theta}$ について**最大化するのは諦める**代わりに、下界 $\mathcal{L}(q, \boldsymbol{\theta})$ を**少しでも増加させるように**、 $\boldsymbol{\theta}$ を更新する
- $\mathcal{L}(q, \boldsymbol{\theta})$ は、**常に**尤度関数 $\ln p(\mathbf{X}|\boldsymbol{\theta})$ の下界であるから、 \mathcal{L} を押し上げることは、尤度関数の増加につながる

EM アルゴリズムの拡張

- M ステップで制限付きの最適化を行うことができる
- パラメータ θ を幾つかのグループに分割
- 各グループに属するパラメータを、他のグループに属するパラメータを固定しながら、順番に最適化していく

3 EM アルゴリズム

- EM アルゴリズムの解釈
- 混合ガウス分布の再解釈
- K-Means 法との関連
- 一般の EM アルゴリズム
- MAP 推定に対する EM アルゴリズム
- EM アルゴリズムの拡張
- EM アルゴリズムに関する補足
- 逐次型の EM アルゴリズム

パラメータについての補足

- パラメータ θ についての補足

- 任意の θ について、下界 $\mathcal{L}(q, \theta)$ は q について**唯一の最大点**をもつ
- それは事後分布 $q(\mathbf{Z}) = p(\mathbf{Z}|\mathbf{X}, \theta)$ である
- またこのとき、下界 \mathcal{L} は対数尤度関数 $\ln p(\mathbf{X}|\theta)$ に一致する
- θ が下界 $\mathcal{L}(q, \theta)$ の大域的最適解に収束するなら、そのような θ は、対数尤度関数 $\ln p(\mathbf{X}|\theta)$ の大域的最適解でもある
- 任意の下界 $\mathcal{L}(q, \theta)$ の任意の極大点は、 $\ln p(\mathbf{X}|\theta)$ の極大点でもある

3 EM アルゴリズム

- EM アルゴリズムの解釈
- 混合ガウス分布の再解釈
- K-Means 法との関連
- 一般の EM アルゴリズム
- MAP 推定に対する EM アルゴリズム
- EM アルゴリズムの拡張
- EM アルゴリズムに関する補足
- 逐次型の EM アルゴリズム

逐次型の EM アルゴリズムの例

- 混合ガウスモデルに対する逐次型の EM アルゴリズム
 - E ステップでは、事後確率分布 $p(Z|X, \theta)$ を計算する
 - データが **i.i.d 集合** であれば、次のように、各データ点ごとの事後確率 $p(z_i|x_i, \theta)$ の積として分解できる

$$p(Z|X, \theta) = \prod_i p(z_i|x_i, \theta) \quad (209)$$

- このとき、**全てのデータ点** X に対して事後確率を求める必要がある
- これを、**1つのデータ点** についてだけ事後確率を求めるように変更する
- M ステップでも、1つのデータ点に対して求めた事後確率だけを使って、パラメータを逐次的に更新するように変更を加える
- 混合ガウスモデルであれば、逐次的な更新式を導出することが可能

逐次型の EM アルゴリズムの例

- 従って、全てのデータ点に対する事後確率を使って、パラメータを再計算する必要がない
- これらの変更によって、**逐次版の EM アルゴリズム**を導出できる
- 混合ガウスモデルに対する逐次型の EM アルゴリズムの導出
 - データ点 x_m について、事後確率 (負担率) $\gamma(z_{mk})$ を更新したとする
 - 新しい負担率を $\gamma^{\text{new}}(z_{mk})$ 、以前の負担率を $\gamma^{\text{old}}(z_{mk})$ とする
 - d を次のようにおく (前後の負担率の差)

$$d = \gamma^{\text{new}}(z_{mk}) - \gamma^{\text{old}}(z_{mk}) \quad (210)$$

- N_k^{new} を次のようにおく (クラス k に属するデータの、実質的な個数)

$$N_k^{\text{new}} = N_k^{\text{old}} + \gamma^{\text{new}}(z_{mk}) - \gamma^{\text{old}}(z_{mk}) = N_k^{\text{old}} + d \quad (211)$$

逐次型の EM アルゴリズムの例

- 以前の平均 $\boldsymbol{\mu}_k^{\text{old}}$ 、共分散行列 $\boldsymbol{\Sigma}_k^{\text{old}}$ 、混合係数 π_k^{old} を、以下のように書くことにする

$$\boldsymbol{\mu}_k^{\text{old}} = \frac{1}{N_k^{\text{old}}} \sum_i \gamma^{\text{old}}(z_{ik}) \mathbf{x}_i \quad (212)$$

$$\boldsymbol{\Sigma}_k^{\text{old}} = \frac{1}{N_k^{\text{old}}} \sum_i \gamma^{\text{old}}(z_{ik}) (\mathbf{x}_i - \boldsymbol{\mu}_k^{\text{old}})(\mathbf{x}_i - \boldsymbol{\mu}_k^{\text{old}})^T \quad (213)$$

$$\pi_k^{\text{old}} = \frac{N_k^{\text{old}}}{N} \quad (214)$$

但し

$$N_k^{\text{old}} = \sum_i \gamma^{\text{old}}(z_{ik}) \quad (215)$$

逐次型の EM アルゴリズムの例

- 平均の更新式は

$$\boldsymbol{\mu}_k^{\text{new}} = \frac{1}{N_k^{\text{new}}} \left(\sum_{i \neq m} \gamma^{\text{old}}(z_{ik}) \mathbf{x}_i + \gamma^{\text{new}}(z_{mk}) \mathbf{x}_m \right) \quad (216)$$

$$= \frac{1}{N_k^{\text{new}}} (N_k^{\text{old}} \boldsymbol{\mu}_k^{\text{old}} - \gamma^{\text{old}}(z_{mk}) \mathbf{x}_m + \gamma^{\text{new}}(z_{mk}) \mathbf{x}_m)$$

$$= \frac{1}{N_k^{\text{new}}} ((N_k^{\text{new}} - \gamma^{\text{new}}(z_{mk}) + \gamma^{\text{old}}(z_{mk})) \boldsymbol{\mu}_k^{\text{old}} - \gamma^{\text{old}}(z_{mk}) \mathbf{x}_m + \gamma^{\text{new}}(z_{mk}) \mathbf{x}_m)$$

$$= \boldsymbol{\mu}_k^{\text{old}} + \frac{1}{N_k^{\text{new}}} (-(\gamma^{\text{new}}(z_{mk}) - \gamma^{\text{old}}(z_{mk})) \boldsymbol{\mu}_k^{\text{old}} + (\gamma^{\text{new}}(z_{mk}) - \gamma^{\text{old}}(z_{mk})) \mathbf{x}_m)$$

$$= \boldsymbol{\mu}_k^{\text{old}} + \frac{\gamma^{\text{new}}(z_{mk}) - \gamma^{\text{old}}(z_{mk})}{N_k^{\text{new}}} (\mathbf{x}_m - \boldsymbol{\mu}_k^{\text{old}}) \quad (217)$$

逐次型の EM アルゴリズムの例

$$= \boldsymbol{\mu}_k^{\text{old}} + \frac{d}{N_k^{\text{new}}} (\mathbf{x}_m - \boldsymbol{\mu}_k^{\text{old}}) \quad (218)$$

- 分散の更新式は

$$\boldsymbol{\Sigma}_k^{\text{new}} = \frac{1}{N_k^{\text{new}}} \sum_i \gamma^{\text{new}}(z_{ik}) \mathbf{x}_i \mathbf{x}_i^T - \boldsymbol{\mu}_k^{\text{new}} (\boldsymbol{\mu}_k^{\text{new}})^T \quad (219)$$

$$= \frac{1}{N_k^{\text{new}}} \left(\sum_{i \neq m} \gamma^{\text{old}}(z_{ik}) \mathbf{x}_i \mathbf{x}_i^T + \gamma^{\text{new}}(z_{mk}) \mathbf{x}_m \mathbf{x}_m^T \right) - \boldsymbol{\mu}_k^{\text{new}} (\boldsymbol{\mu}_k^{\text{new}})^T \quad (220)$$

$$= \frac{1}{N_k^{\text{new}}} \left(\left(\sum_i \gamma^{\text{old}}(z_{ik}) \mathbf{x}_i \mathbf{x}_i^T - \gamma^{\text{old}}(z_{mk}) \mathbf{x}_m \mathbf{x}_m^T \right) + \gamma^{\text{new}}(z_{mk}) \mathbf{x}_m \mathbf{x}_m^T \right) - \boldsymbol{\mu}_k^{\text{new}} (\boldsymbol{\mu}_k^{\text{new}})^T \quad (221)$$

逐次型の EM アルゴリズムの例

$$\begin{aligned} &= \frac{1}{N_k^{\text{new}}} \left(N_k^{\text{old}} \left(\boldsymbol{\Sigma}_k^{\text{old}} + \boldsymbol{\mu}_k^{\text{old}} (\boldsymbol{\mu}_k^{\text{old}})^T \right) - \right. \\ &\quad \left. \gamma^{\text{old}}(z_{mk}) \mathbf{x}_m \mathbf{x}_m^T + \gamma^{\text{new}}(z_{mk}) \mathbf{x}_m \mathbf{x}_m^T \right) - \\ &\quad \boldsymbol{\mu}_k^{\text{new}} (\boldsymbol{\mu}_k^{\text{new}})^T \end{aligned} \quad (222)$$

$$\begin{aligned} &= \frac{1}{N_k^{\text{new}}} \left(N_k^{\text{old}} \boldsymbol{\Sigma}_k^{\text{old}} + N_k^{\text{old}} \boldsymbol{\mu}_k^{\text{old}} (\boldsymbol{\mu}_k^{\text{old}})^T + \right. \\ &\quad \left(\gamma^{\text{new}}(z_{mk}) - \gamma^{\text{old}}(z_{mk}) \right) \mathbf{x}_m \mathbf{x}_m^T - \\ &\quad \left. N_k^{\text{new}} \boldsymbol{\mu}_k^{\text{new}} (\boldsymbol{\mu}_k^{\text{new}})^T \right) \\ &= \frac{1}{N_k^{\text{new}}} \left(N_k^{\text{old}} \boldsymbol{\Sigma}_k^{\text{old}} + N_k^{\text{old}} \boldsymbol{\mu}_k^{\text{old}} (\boldsymbol{\mu}_k^{\text{old}})^T + \right. \\ &\quad \left. d\mathbf{x}_m \mathbf{x}_m^T - N_k^{\text{new}} \boldsymbol{\mu}_k^{\text{new}} (\boldsymbol{\mu}_k^{\text{new}})^T \right) \end{aligned} \quad (223)$$

逐次型の EM アルゴリズムの例

ここで

$$\begin{aligned} & N_k^{\text{new}} \boldsymbol{\mu}_k^{\text{new}} (\boldsymbol{\mu}_k^{\text{new}})^T \\ = & (N_k^{\text{new}} \boldsymbol{\mu}_k^{\text{old}} + d(\mathbf{x}_m - \boldsymbol{\mu}_k^{\text{old}})) (\boldsymbol{\mu}_k^{\text{new}})^T \\ = & \frac{1}{N_k^{\text{new}}} (N_k^{\text{new}} \boldsymbol{\mu}_k^{\text{old}} + d(\mathbf{x}_m - \boldsymbol{\mu}_k^{\text{old}})) \\ & (N_k^{\text{new}} \boldsymbol{\mu}_k^{\text{old}} + d(\mathbf{x}_m - \boldsymbol{\mu}_k^{\text{old}}))^T \\ = & N_k^{\text{new}} \boldsymbol{\mu}_k^{\text{old}} (\boldsymbol{\mu}_k^{\text{old}})^T + 2\boldsymbol{\mu}_k^{\text{old}} d(\mathbf{x}_m - \boldsymbol{\mu}_k^{\text{old}})^T + \\ & \frac{1}{N_k^{\text{new}}} d^2 (\mathbf{x}_m - \boldsymbol{\mu}_k^{\text{old}}) (\mathbf{x}_m - \boldsymbol{\mu}_k^{\text{old}})^T \end{aligned} \quad (224)$$

$$\begin{aligned} = & (N_k^{\text{old}} + d) \boldsymbol{\mu}_k^{\text{old}} (\boldsymbol{\mu}_k^{\text{old}})^T + \\ & 2d \boldsymbol{\mu}_k^{\text{old}} \mathbf{x}_m^T - 2d \boldsymbol{\mu}_k^{\text{old}} (\boldsymbol{\mu}_k^{\text{old}})^T + \\ & \frac{1}{N_k^{\text{new}}} d^2 (\mathbf{x}_m - \boldsymbol{\mu}_k^{\text{old}}) (\mathbf{x}_m - \boldsymbol{\mu}_k^{\text{old}})^T \end{aligned} \quad (225)$$

逐次型の EM アルゴリズムの例

であるから

$$\begin{aligned} & N_k^{\text{old}} \boldsymbol{\mu}_k^{\text{old}} (\boldsymbol{\mu}_k^{\text{old}})^T + d \mathbf{x}_m \mathbf{x}_m^T - N_k^{\text{new}} \boldsymbol{\mu}_k^{\text{new}} (\boldsymbol{\mu}_k^{\text{new}})^T \\ = & N_k^{\text{old}} \boldsymbol{\mu}_k^{\text{old}} (\boldsymbol{\mu}_k^{\text{old}})^T + d \mathbf{x}_m \mathbf{x}_m^T - \\ & (N_k^{\text{old}} + d) \boldsymbol{\mu}_k^{\text{old}} (\boldsymbol{\mu}_k^{\text{old}})^T - \\ & 2d \boldsymbol{\mu}_k^{\text{old}} \mathbf{x}_m^T + 2d \boldsymbol{\mu}_k^{\text{old}} (\boldsymbol{\mu}_k^{\text{old}})^T - \\ & \frac{1}{N_k^{\text{new}}} d^2 (\mathbf{x}_m - \boldsymbol{\mu}_k^{\text{old}}) (\mathbf{x}_m - \boldsymbol{\mu}_k^{\text{old}})^T \\ = & d \mathbf{x}_m \mathbf{x}_m^T + d \boldsymbol{\mu}_k^{\text{old}} (\boldsymbol{\mu}_k^{\text{old}})^T - 2d \boldsymbol{\mu}_k^{\text{old}} \mathbf{x}_m^T - \\ & \frac{1}{N_k^{\text{new}}} d^2 (\mathbf{x}_m - \boldsymbol{\mu}_k^{\text{old}}) (\mathbf{x}_m - \boldsymbol{\mu}_k^{\text{old}})^T \\ = & d (\mathbf{x}_m \mathbf{x}_m^T + \boldsymbol{\mu}_k^{\text{old}} (\boldsymbol{\mu}_k^{\text{old}})^T - 2 \boldsymbol{\mu}_k^{\text{old}} \mathbf{x}_m^T) - \\ & \frac{1}{N_k^{\text{new}}} d^2 (\mathbf{x}_m - \boldsymbol{\mu}_k^{\text{old}}) (\mathbf{x}_m - \boldsymbol{\mu}_k^{\text{old}})^T \\ = & d (\mathbf{x}_m - \boldsymbol{\mu}_k^{\text{old}}) (\mathbf{x}_m - \boldsymbol{\mu}_k^{\text{old}})^T - \end{aligned}$$

逐次型の EM アルゴリズムの例

$$\begin{aligned} & \frac{1}{N_k^{\text{new}}} d^2 (\mathbf{x}_m - \boldsymbol{\mu}_k^{\text{old}}) (\mathbf{x}_m - \boldsymbol{\mu}_k^{\text{old}})^T \\ = & \frac{d}{N_k^{\text{new}}} (N_k^{\text{new}} - d) (\mathbf{x}_m - \boldsymbol{\mu}_k^{\text{old}}) (\mathbf{x}_m - \boldsymbol{\mu}_k^{\text{old}})^T \\ = & \frac{d}{N_k^{\text{new}}} N_k^{\text{old}} (\mathbf{x}_m - \boldsymbol{\mu}_k^{\text{old}}) (\mathbf{x}_m - \boldsymbol{\mu}_k^{\text{old}})^T \end{aligned} \quad (226)$$

以上より

$$\begin{aligned} \boldsymbol{\Sigma}_k^{\text{new}} &= \frac{1}{N_k^{\text{new}}} \left(N_k^{\text{old}} \boldsymbol{\Sigma}_k^{\text{old}} + N_k^{\text{old}} \boldsymbol{\mu}_k^{\text{old}} (\boldsymbol{\mu}_k^{\text{old}})^T + \right. \\ & \quad \left. d \mathbf{x}_m \mathbf{x}_m^T - N_k^{\text{new}} \boldsymbol{\mu}_k^{\text{new}} (\boldsymbol{\mu}_k^{\text{new}})^T \right) \\ &= \frac{1}{N_k^{\text{new}}} \left(N_k^{\text{old}} \boldsymbol{\Sigma}_k^{\text{old}} + \right. \\ & \quad \left. \frac{d}{N_k^{\text{new}}} N_k^{\text{old}} (\mathbf{x}_m - \boldsymbol{\mu}_k^{\text{old}}) (\mathbf{x}_m - \boldsymbol{\mu}_k^{\text{old}})^T \right) \end{aligned}$$

逐次型の EM アルゴリズムの例

$$\begin{aligned} &= \frac{N_k^{\text{old}}}{N_k^{\text{new}}} \left(\boldsymbol{\Sigma}_k^{\text{old}} + \frac{d}{N_k^{\text{new}}} \right. \\ &\quad \left. (\mathbf{x}_m - \boldsymbol{\mu}_k^{\text{old}}) (\mathbf{x}_m - \boldsymbol{\mu}_k^{\text{old}})^T \right) \\ &= \frac{N_k^{\text{old}}}{N_k^{\text{new}}} \left(\boldsymbol{\Sigma}_k^{\text{old}} + \frac{\gamma^{\text{new}}(z_{ik}) - \gamma^{\text{old}}(z_{ik})}{N_k^{\text{new}}} \right. \\ &\quad \left. (\mathbf{x}_m - \boldsymbol{\mu}_k^{\text{old}}) (\mathbf{x}_m - \boldsymbol{\mu}_k^{\text{old}})^T \right) \end{aligned} \quad (227)$$

- 混合係数の更新式は

$$\pi_k^{\text{new}} = \frac{N_k^{\text{new}}}{N} = \frac{N_k^{\text{old}} + \gamma^{\text{new}}(z_{mk}) - \gamma^{\text{old}}(z_{mk})}{N} \quad (228)$$

逐次型の EM アルゴリズムの例

- 混合ガウスモデルにおける逐次型の EM アルゴリズム
 - 上記より、パラメータ μ_k, Σ_k, π_k の逐次更新式が得られた

$$\mu_k^{\text{new}} = \mu_k^{\text{old}} + \frac{\gamma^{\text{new}}(z_{mk}) - \gamma^{\text{old}}(z_{mk})}{N_k^{\text{new}}} (\mathbf{x}_m - \mu_k^{\text{old}}) \quad (229)$$

$$\Sigma_k^{\text{new}} = \frac{N_k^{\text{old}}}{N_k^{\text{new}}} \left(\Sigma_k^{\text{old}} + \frac{\gamma^{\text{new}}(z_{ik}) - \gamma^{\text{old}}(z_{ik})}{N_k^{\text{new}}} (\mathbf{x}_m - \mu_k^{\text{old}}) (\mathbf{x}_m - \mu_k^{\text{old}})^T \right) \quad (230)$$

$$\pi_k^{\text{new}} = \frac{N_k^{\text{new}}}{N} = \frac{N_k^{\text{old}} + \gamma^{\text{new}}(z_{mk}) - \gamma^{\text{old}}(z_{mk})}{N} \quad (231)$$

但し

$$N_k^{\text{new}} = N_k^{\text{old}} + \gamma^{\text{new}}(z_{mk}) - \gamma^{\text{old}}(z_{mk}) \quad (232)$$

逐次型の EM アルゴリズムの例

- 到着したデータ x_m について、E ステップで負担率 $\gamma(z_{mk})$ を求めた後に、M ステップで (上記の更新式を用いて) パラメータを更新することを、交互に繰り返せばよい
- 逐次型の EM アルゴリズムの特徴
 - x_m が新しく到着したデータであれば、 $\gamma^{\text{old}}(z_{mk}) = 0$ とする
 - E ステップと M ステップの計算に必要な時間は、データ点の総数とは無関係に決まる
 - パラメータの更新は、全データについての処理を待たずに、各データ点についての処理の後に行われる
 - そのため、逐次型の EM アルゴリズムは、従来のバッチ型に比べて、**速く収束する**

逐次型の EM アルゴリズムの例

- ここまでの話の流れ
 - MAP 推定に対する EM アルゴリズムについて考えた
 - EM アルゴリズムの拡張 (一般化 EM アルゴリズム) について簡単に触れた
 - 混合ガウスモデルについて、逐次型の EM アルゴリズムを導出した

EM アルゴリズムのまとめ

- EM アルゴリズムの目的

- 潜在変数をもつ確率モデルについて、パラメータの最尤解を求める

- EM アルゴリズムで行っていること

- 対数尤度 $\ln p(\mathbf{X}|\boldsymbol{\theta})$ の直接の最適化が困難であっても、E ステップと M ステップという 2 段階の簡単な手続きに分割し、交互に繰り返すことで最適化できるようにする
- 完全データ対数尤度 $\ln p(\mathbf{X}, \mathbf{Z}|\boldsymbol{\theta})$ の事後確率 $\ln p(\mathbf{Z}|\mathbf{X}, \boldsymbol{\theta})$ による期待値 $Q(\boldsymbol{\theta}, \boldsymbol{\theta}^{\text{old}})$ の最大化を行う
- 期待値の最大化は、 $\mathcal{L}(q, \boldsymbol{\theta})$ の最大化と等価である
- $\mathcal{L}(q, \boldsymbol{\theta})$ は $\ln p(\mathbf{X}|\boldsymbol{\theta})$ の下界であるから、 \mathcal{L} の最大化は、対数尤度 $\ln p(\mathbf{X}|\boldsymbol{\theta})$ の最大化に相当

EM アルゴリズムのまとめ

- ここまでの話の流れ

- 発見的に導出した、混合ガウスモデルに対する EM アルゴリズムも、期待値の最大化という考え方で解釈可能であった
- K-Means 法は、混合ガウスモデルに対する EM アルゴリズムの一種の極限として得られた
- 一般的な EM アルゴリズムの取り扱いについて調べた
- 最尤推定だけでなく、MAP 推定に対しても EM アルゴリズムを適用できた
- 混合ガウスモデルに対する、逐次版の EM アルゴリズムを導出した

目次

- 1 K-Means 法
- 2 混合ガウス分布
- 3 EM アルゴリズム
- 4 変分の導入**
- 5 近似推論法
- 6 変分自己符号化器

- 4 変分の導入
 - EM アルゴリズムが困難な場合
 - 変分法
 - 変分法で解ける問題の例
 - 変分法のまとめ

EM アルゴリズムが困難な場合

- EM アルゴリズムで行う計算

- **E ステップ**では、潜在変数の事後確率分布 $p(\mathbf{Z}|\mathbf{X}, \boldsymbol{\theta}^{\text{old}})$ を計算
- **M ステップ**では、完全データ対数尤度 $\ln p(\mathbf{X}, \mathbf{Z}|\boldsymbol{\theta})$ の期待値を計算

$$Q(\boldsymbol{\theta}, \boldsymbol{\theta}^{\text{old}}) = \sum_{\mathbf{Z}} p(\mathbf{Z}|\mathbf{X}, \boldsymbol{\theta}^{\text{old}}) \ln p(\mathbf{X}, \mathbf{Z}|\boldsymbol{\theta}) \quad (233)$$

$$Q(\boldsymbol{\theta}, \boldsymbol{\theta}^{\text{old}}) = \int_{\mathbf{Z}} p(\mathbf{Z}|\mathbf{X}, \boldsymbol{\theta}^{\text{old}}) \ln p(\mathbf{X}, \mathbf{Z}|\boldsymbol{\theta}) d\mathbf{Z} \quad (234)$$

そして、 $Q(\boldsymbol{\theta}, \boldsymbol{\theta}^{\text{old}})$ を最大化するパラメータ $\boldsymbol{\theta}^{\text{new}}$ を求める

$$\boldsymbol{\theta}^{\text{new}} = \arg \max_{\boldsymbol{\theta}} Q(\boldsymbol{\theta}, \boldsymbol{\theta}^{\text{old}}) \quad (235)$$

EM アルゴリズムが困難な場合

● EM アルゴリズムの困難さ

- 実際に扱うモデルでは、事後分布 $p(\mathbf{Z}|\mathbf{X}, \boldsymbol{\theta}^{\text{old}})$ の計算や、事後分布に従った期待値 $Q(\boldsymbol{\theta}, \boldsymbol{\theta}^{\text{old}})$ の計算が、**不可能であることが多い**
- 隠れ変数の**次元が高すぎる**
- 事後分布が**複雑な形をしていて**、期待値を解析的に計算できない
- **連続変数**であれば、積分が閉形式の解を持たないかもしれない
- 空間の次元の問題や、被積分項の複雑さから、数値積分すら困難かもしれない
- **離散変数**であれば、期待値を計算するためには、**潜在変数の可能な全ての組み合わせについての和を取る**必要がある
- 隠れ変数の次元が高くなると、組み合わせ数が指数的に増大する
- 計算量が大きすぎて、期待値の厳密な計算がもはや不可能

EM アルゴリズムが困難な場合

- 近似法

- EM アルゴリズムが困難であるとき、何らかの方法で近似しなければならない
- 近似法は、確率的な近似と、決定的な近似の2つに分けられる

- 確率的な近似

- マルコフ連鎖モンテカルロ法などの手法がある
- 無限の計算資源があれば、厳密な結果が得られる
- 実際には計算量が有限であるため、得られる解は近似解となる

- 決定的な近似

- 事後分布 $p(\mathbf{Z}|\mathbf{X}, \boldsymbol{\theta})$ を解析的に近似する
- 事後分布に対して、何らかの仮定をおく
- 例えば、単純な項の積として分解できる、あるいは、(ガウス分布などの特別な) パラメトリックな分布であるといった仮定

- 4 変分の導入
 - EM アルゴリズムが困難な場合
 - 変分法
 - 変分法で解ける問題の例
 - 変分法のまとめ

- ここで扱う近似法
 - **変分推論法** (Variational inference) あるいは**変分ベイズ法** (Variational Bayes) について扱う
- **変分推論** (Variational inference)
 - 18 世紀のオイラー、ラグランジュらによる**変分法** (Calculus of variations) に起源をもつ
 - まずは、変分法について説明をしていく

- 関数と汎関数の違い

- 通常関数は、入力として値をとり、出力として関数の値を返す
- 通常関数は、値から値への写像である
- 関数の導関数は、入力値を微小に変えたときに、出力の関数値がどの程度変わるかを表す
- 汎関数 (Functional) とは、入力として関数を取り、出力として汎関数の値を返す
- 汎関数は、関数から値への写像である
- 汎関数微分 (Functional derivative) とは、入力関数が微小に変わったときに、出力の汎関数値がどの程度変わるかを表す
- 汎関数の微分を、変分という

- 汎関数の例

- エントロピー $H[p]$ は、確率分布 $p(x)$ を入力として、以下の量を返す汎関数

$$H[p] = - \int p(x) \ln p(x) dx \quad (236)$$

- 汎関数の最適化

- 多くの問題は、**汎関数の値を最適化する問題**として定式化できる
- 汎関数の最適化とは、**可能な全ての入力関数の中から**、汎関数の値を最大化、あるいは最小化するような**関数を選び出す**ことである
- 通常の最適化では、可能な全てのパラメータ (入力値) の中から、関数を最大化、あるいは最小化するような 1 つのパラメータを選び出す
- 次は、いよいよ**変分**の計算について説明する

- 変分法

- 通常の微分を使えば、ある関数 $y(x)$ を最大化 (最小化) するような x の値が求められる
- **変分法**を使えば、汎関数 $F[y]$ を最大化 (最小化) するような、関数 $y(x)$ が求められる
- 従って、可能な全ての関数 $y(x)$ の中から、 $F[y]$ を最大 (最小) にするような関数が得られる

- 変分法によって解ける問題の例

- 2 点を結ぶ最短経路は? (答えは直線)
- 最速降下曲線は? (答えはサイクロイド)
- **エントロピーが最大**になるような確率分布は? (答えは**ガウス分布**)

- 通常の微分の表現

- 関数 $y(x + \epsilon)$ のテイラー展開は次のように記述できた

$$y(x + \epsilon) = \sum_{n=0}^{\infty} \frac{y^{(n)}(x)}{n!} \epsilon^n \quad (237)$$

$$= y(x) + \frac{dy}{dx} \epsilon + \frac{1}{2!} \frac{d^2 y}{dx^2} \epsilon^2 + \frac{1}{3!} \frac{d^3 y}{dx^3} \epsilon^3 + \dots \quad (238)$$

$$= y(x) + \frac{dy}{dx} \epsilon + O(\epsilon^2) \quad (239)$$

- これより微分 dy/dx は、次のように求められる
- 変数 x に微小な変化 ϵ を加え、このときの関数値 $y(x + \epsilon)$ を ϵ の累乗形として表現する
- 最後に $\epsilon \rightarrow 0$ の極限をとればよい

$$\frac{dy}{dx} = \lim_{\epsilon \rightarrow 0} \frac{y(x + \epsilon) - y(x)}{\epsilon} \quad (240)$$

- 多変数関数 $y(x_1, \dots, x_D)$ の偏微分の表現
 - 多変数関数 $y(x_1, \dots, x_D)$ のテイラー展開は次のように記述できた

$$\mathbf{D}^n = \left(\epsilon_1 \frac{\partial y}{\partial x_1} + \dots + \epsilon_D \frac{\partial y}{\partial x_D} \right)^n \quad (241)$$

上記のような演算子 \mathbf{D} を考えれば

$$\begin{aligned} & y(x_1 + \epsilon_1, \dots, x_D + \epsilon_D) \\ &= \sum_{n=0}^{\infty} \frac{1}{n!} (\mathbf{D}^n y)(x_1, \dots, x_D) \end{aligned} \quad (242)$$

$$\begin{aligned} &= y(x_1, \dots, x_D) + \sum_{i=1}^D \frac{\partial y}{\partial x_i} \epsilon_i + \frac{1}{2!} \sum_{i=1}^D \sum_{j=1}^D \frac{\partial^2 y}{\partial x_i \partial x_j} \epsilon_i \epsilon_j + \\ & \quad \frac{1}{3!} \sum_{i=1}^D \sum_{j=1}^D \sum_{k=1}^D \frac{\partial^3 y}{\partial x_i \partial x_j \partial x_k} \epsilon_i \epsilon_j \epsilon_k + \dots \end{aligned} \quad (243)$$

であるから

$$\begin{aligned} & y(x_1 + \epsilon_1, \dots, x_D + \epsilon_D) \\ = & y(x_1, \dots, x_D) + \sum_{i=1}^D \frac{\partial y}{\partial x_i} \epsilon_i + O(\epsilon^2) \end{aligned} \quad (244)$$

- これより偏微分 $\partial y / \partial x_i$ は、次のように求められる

$$\begin{aligned} \frac{\partial y}{\partial x_i} = \lim_{\epsilon_i \rightarrow 0} \frac{1}{\epsilon_i} & (y(x_1, \dots, x_{i-1}, x_i + \epsilon_i, x_{i+1}, \dots, x_D) - \\ & y(x_1, \dots, x_{i-1}, x_i, x_{i+1}, \dots, x_D)) \end{aligned} \quad (245)$$

- 変分の表現

- 多少不正確だが、変分をどのように定義すればよいか考えてみる
- ここで、各 x_i に対する関数の値 $z_i = y(x_i)$ を個別の変数とみなして、次の関数 $F(z_1, \dots, z_D)$ について考えてみよう

$$\begin{aligned} & F(z_1 + \epsilon\eta(x_1), \dots, z_D + \epsilon\eta(x_D)) \\ &= F(z_1, \dots, z_D) + \sum_{i=1}^D \frac{\partial F}{\partial z_i} \epsilon\eta(x_i) + O(\epsilon^2) \end{aligned} \quad (246)$$

$z_i = y(x_i)$ を代入してみると

$$\begin{aligned} & F(y(x_1) + \epsilon\eta(x_1), \dots, y(x_D) + \epsilon\eta(x_D)) \\ &= F(y(x_1), \dots, y(x_D)) + \sum_{i=1}^D \frac{\partial F}{\partial y(x_i)} \epsilon\eta(x_i) + O(\epsilon^2) \end{aligned} \quad (247)$$

変分法

- ここで $D \rightarrow \infty$ の極限を取り、 x_1, \dots, x_D が、ある連続した区間 $[a, b]$ に含まれる、全ての実数を表すことにする
- このとき $y(x_1), \dots, y(x_D)$ は、実数の区間 $[a, b]$ で定義される連続関数 $y(x)$ として書けることが分かる
- 同様に $y(x_1) + \epsilon\eta(x_1), \dots, y(x_D) + \epsilon\eta(x_D)$ は、実数の区間 $[a, b]$ で定義される連続関数 $y(x) + \epsilon\eta(x)$ として、まとめることができる
- 関数 $\eta(x)$ も、実数の区間 $[a, b]$ で定義される連続関数
- $\epsilon\eta(x)$ は、 $y(x)$ に加わる摂動として、考えることができる

- 関数 F は、関数 $y(x)$ や $y(x) + \epsilon\eta(x)$ を入力として受け取る、汎関数 $F[y]$ として解釈できるから、次のように書ける

$$F(y(x_1) + \epsilon\eta(x_1), \dots, y(x_D) + \epsilon\eta(x_D)) = F[y(x) + \epsilon\eta(x)] \quad (248)$$

$$F(y(x_1), \dots, y(x_D)) = F[y(x)] \quad (249)$$

- 以下の項は、入力を $y(x)$ に摂動を加えて $y(x) + \epsilon\eta(x)$ へと微小に変化させたときの、汎関数の ($F[y(x)]$ から $F[y(x) + \epsilon\eta(x)]$ への) 変化量を表している

$$\sum_{i=1}^D \frac{\partial F}{\partial y(x_i)} \epsilon\eta(x_i) \quad (250)$$

- 点 x_i における汎関数 F の変化量を、 x_1, \dots, x_D の範囲について、即ち、実数の区間 $[a, b]$ について足し合わせていると解釈する

- $D \rightarrow \infty$ のとき、 x_1, \dots, x_D は区間 $[a, b]$ における全ての実数を表すから、総和を積分に置き換えられそうである
- 汎関数の微分 $\frac{\delta F}{\delta y(x)}$ を使えば、次のように書ける

$$\begin{aligned} & \sum_{i=1}^D \frac{\partial F}{\partial y(x_i)} \epsilon \eta(x_i) \\ \Rightarrow & \int_a^b \frac{\delta F}{\delta y(x)} \epsilon \eta(x) dx = \epsilon \int_a^b \frac{\delta F}{\delta y(x)} \eta(x) dx \end{aligned} \quad (251)$$

- 結局、変分 $\frac{\delta F}{\delta y(x)}$ は次のように定義できる

$$F[y(x) + \epsilon \eta(x)] = F[y(x)] + \epsilon \int_a^b \frac{\delta F}{\delta y(x)} \eta(x) dx + O(\epsilon^2) \quad (252)$$

変分法

- $F[y]$ は、区間 $[a, b]$ で定義される関数 y を受け取るとする
- 変分 $\delta F / \delta y$ は、入力関数 $y(x)$ に、任意の微小な変化 $\epsilon \eta(x)$ を加えたときの、汎関数 $F[y]$ の変化量として定義できる
- $\eta(x)$ は x についての任意の関数

Figure D.1 A functional derivative can be defined by considering how the value of a functional $F[y]$ changes when the function $y(x)$ is changed to $y(x) + \epsilon\eta(x)$ where $\eta(x)$ is an arbitrary function of x .

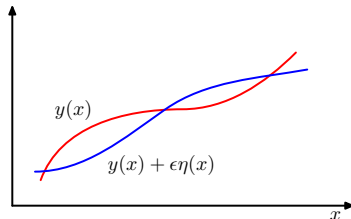


図 8: $y(x)$ と $y(x) + \epsilon\eta(x)$ の表現

- 変分法の例

- 次の図 9 を使って、実際に変分を求めてみよう

- 汎関数 $F[y]$ は、次のように定義されたとする

$$F[y] = \int_a^b y(x) dx \quad (253)$$

- 汎関数の値 $F[y(x)], F[y(x) + \epsilon\eta(x)]$ は次のようになる

$$F[y(x)] = \int_a^b y(x) dx \quad (254)$$

$$F[y(x) + \epsilon\eta(x)] = \int_a^b (y(x) + \epsilon\eta(x)) dx \quad (255)$$

- $F[y(x) + \epsilon\eta(x)]$ は次のように分解できる

$$F[y(x) + \epsilon\eta(x)] = \int_a^b y(x)dx + \epsilon \int_a^b \eta(x)dx \quad (256)$$

$$= F[y(x)] + \epsilon \int_a^b \eta(x)dx \quad (257)$$

- ここで、変分の定義式は

$$F[y(x) + \epsilon\eta(x)] = F[y(x)] + \epsilon \int_a^b \frac{\delta F}{\delta y(x)} \eta(x)dx + O(\epsilon^2) \quad (258)$$

であったので、上の2つの式を見比べれば、変分 $\delta F/\delta y$ は結局

$$\frac{\delta F}{\delta y(x)} = 1 \quad (259)$$

となることが分かる

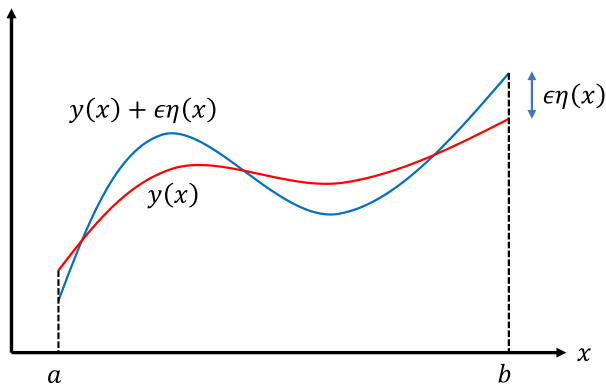


図 9: 区間 $[a, b]$ で定義された関数 $y(x)$ の表現

● 汎関数の最適化

- 汎関数 $F[y]$ が最大 (最小) となるとき、関数 $y(x)$ の微小な変化に対して、汎関数は変化しないはず
- 即ち、汎関数が最大 (最小) となるとき、 $F[y(x) + \epsilon\eta(x)] = F[y(x)]$ が成り立つ
- 従って、変分の定義式から、以下が成り立つ

$$\int_a^b \frac{\delta F}{\delta y(x)} \eta(x) dx = 0 \quad (260)$$

- 上式は任意の $\eta(x)$ について成立しなければならない
- 従って、変分 $\delta F / \delta y$ は、任意の x について 0 とならなければならない
- 汎関数 $F[y]$ が最大 (最小) となるとき、 $\delta F / \delta y = 0$ が成立することが分かった (通常の微分と同じ)

- 変分法の例

- 様々な汎関数について、変分を導出してみよう
- また、その汎関数が最大 (最小) となるときに成り立つ条件を、導出してみよう

- 汎関数の例 (1)

- $y(x)$ とその微分 $y'(x) = dy/dx$ 、そして x によって決まる関数 $G(y(x), y'(x), x)$ があるとする
- 汎関数 $F[y]$ を、 $G(y(x), y'(x), x)$ を区間 $[a, b]$ にわたって積分した結果を出力する関数として、次のように定める

$$F[y] = \int_a^b G(y(x), y'(x), x) dx \quad (261)$$

- 積分区間は無限であってもよいとする ($a = -\infty, b = \infty$ でもよい)

変分法

- $y(x)$ に摂動 $\epsilon\eta(x)$ を加えたときの、汎関数の値 $F[y(x) + \epsilon\eta(x)]$ を使って、変分 $\delta F/\delta y$ を調べてみる

$$F[y(x) + \epsilon\eta(x)] = \int_a^b G(y(x) + \epsilon\eta(x), y'(x) + \epsilon\eta'(x), x) dx \quad (262)$$

ここで、被積分項のテーラー展開を考えれば

$$\begin{aligned} & G(y(x) + \epsilon\eta(x), y'(x) + \epsilon\eta'(x), x) \\ = & G(y(x), y'(x), x) + \frac{\partial G}{\partial y} \epsilon\eta(x) + \\ & \frac{\partial G}{\partial y'} \epsilon\eta'(x) + \frac{\partial G}{\partial x} \cdot 0 + O(\epsilon^2) \end{aligned} \quad (263)$$

$$= G(y(x), y'(x), x) + \epsilon \left(\frac{\partial G}{\partial y} \eta(x) + \frac{\partial G}{\partial y'} \eta'(x) \right) + O(\epsilon^2) \quad (264)$$

であるから

$$\begin{aligned}
 & F[y(x) + \epsilon \eta(x)] \\
 = & \int_a^b G(y(x) + \epsilon \eta(x), y'(x) + \epsilon \eta'(x), x) dx \\
 = & \int_a^b \left(G(y(x), y'(x), x) + \right. \\
 & \left. \epsilon \left(\frac{\partial G}{\partial y} \eta(x) + \frac{\partial G}{\partial y'} \eta'(x) \right) + O(\epsilon^2) \right) dx \quad (265)
 \end{aligned}$$

$$\begin{aligned}
 = & \int_a^b G(y(x), y'(x), x) dx + \\
 & \epsilon \int_a^b \left(\frac{\partial G}{\partial y} \eta(x) + \frac{\partial G}{\partial y'} \eta'(x) \right) dx + O(\epsilon^2) \quad (266)
 \end{aligned}$$

$$= F[y(x)] + \epsilon \int_a^b \left(\frac{\partial G}{\partial y} \eta(x) + \frac{\partial G}{\partial y'} \eta'(x) \right) dx + O(\epsilon^2) \quad (267)$$

ここで

$$\begin{aligned} & \int_a^b \left(\frac{\partial G}{\partial y} \eta(x) + \frac{\partial G}{\partial y'} \eta'(x) \right) dx \\ = & \int_a^b \left(\frac{\partial G}{\partial y} \eta(x) \right) dx + \int_a^b \left(\frac{\partial G}{\partial y'} \eta'(x) \right) dx \end{aligned} \quad (268)$$

$$\begin{aligned} = & \int_a^b \left(\frac{\partial G}{\partial y} \eta(x) \right) dx + \\ & \left[\frac{\partial G}{\partial y'} \eta(x) \right]_a^b - \int_a^b \frac{d}{dx} \left(\frac{\partial G}{\partial y'} \right) \eta(x) dx \end{aligned} \quad (269)$$

$$= \left[\frac{\partial G}{\partial y'} \eta(x) \right]_a^b + \int_a^b \left(\frac{\partial G}{\partial y} - \frac{d}{dx} \left(\frac{\partial G}{\partial y'} \right) \right) \eta(x) dx \quad (270)$$

である

- 途中の式変形では、部分積分を使っていることに注意
- いま、積分区間の両端において、 $y(x)$ の値は固定されているとする
- これを**固定端条件**という (図 10)
- このとき、 $\eta(a) = \eta(b) = 0$ であるから、上式の最初の項が消えて

$$\begin{aligned} & \int_a^b \left(\frac{\partial G}{\partial y} \eta(x) + \frac{\partial G}{\partial y'} \eta'(x) \right) dx \\ &= \left[\frac{\partial G}{\partial y'} \eta(x) \right]_a^b + \int_a^b \left(\frac{\partial G}{\partial y} - \frac{d}{dx} \left(\frac{\partial G}{\partial y'} \right) \right) \eta(x) dx \quad (271) \end{aligned}$$

$$= \int_a^b \left(\frac{\partial G}{\partial y} - \frac{d}{dx} \left(\frac{\partial G}{\partial y'} \right) \right) \eta(x) dx \quad (272)$$

のようになる

- 従って、摂動を加えたときの汎関数の値 $F[y(x) + \epsilon\eta(x)]$ は

$$\begin{aligned} & F[y(x) + \epsilon\eta(x)] \\ = & F[y(x)] + \epsilon \int_a^b \left(\frac{\partial G}{\partial y} - \frac{d}{dx} \left(\frac{\partial G}{\partial y'} \right) \right) \eta(x) dx + O(\epsilon^2) \quad (273) \end{aligned}$$

となる

- 上式を、変分の定義式と比べれば

$$F[y(x) + \epsilon\eta(x)] = F[y(x)] + \epsilon \int_a^b \frac{\delta F}{\delta y(x)} \eta(x) dx + O(\epsilon^2) \quad (274)$$

変分 $\delta F/\delta y$ は次のように書ける

$$\frac{\delta F}{\delta y(x)} = \frac{\partial G}{\partial y} - \frac{d}{dx} \left(\frac{\partial G}{\partial y'} \right) \quad (275)$$

- 汎関数 $F[y]$ が最大 (最小) になるとき、変分 $\delta F/\delta y$ が 0 になる
- 従って、汎関数が最大 (最小) になるとき、以下の方程式が成り立つ

$$\frac{\partial G}{\partial y} - \frac{d}{dx} \left(\frac{\partial G}{\partial y'} \right) = 0 \quad (276)$$

- これをオイラー-ラグランジュ方程式という
- オイラー-ラグランジュ方程式は、次のような考え方で導出することもできる
- $F[y]$ が最大 (最小) であれば、摂動 $\epsilon\eta(x)$ によって $y(x)$ が少し変化しても、 $F[y]$ の値は変化しないはず
- 従って、 $F[y]$ が最大 (最小) であるとき、 $F[y]$ の ϵ による微分は 0 になるはず

- これを数式で表現すると、次のようになる

$$\left. \frac{\partial F[y]}{\partial \epsilon} \right|_{\epsilon=0} = 0 \quad (277)$$

左辺は通常の偏微分であり、これを計算すると

$$\begin{aligned} & \frac{\partial F[y]}{\partial \epsilon} \\ = & \frac{\partial}{\partial \epsilon} \int_a^b G(y, y', x) dx \end{aligned} \quad (278)$$

$$= \int_a^b \frac{\partial}{\partial \epsilon} G(y, y', x) dx \quad (279)$$

$$= \int_a^b \left(\frac{\partial G}{\partial y} \frac{\partial y}{\partial \epsilon} + \frac{\partial G}{\partial y'} \frac{\partial y'}{\partial \epsilon} + \frac{\partial G}{\partial x} \frac{\partial x}{\partial \epsilon} \right) dx \quad (280)$$

$$= \int_a^b \left(\frac{\partial G}{\partial y} \eta(x) + \frac{\partial G}{\partial y'} \eta'(x) \right) dx \quad (281)$$

$$= \left[\frac{\partial G}{\partial y'} \eta(x) \right]_a^b + \int_a^b \left(\frac{\partial G}{\partial y} - \frac{d}{dx} \left(\frac{\partial G}{\partial y'} \right) \right) \eta(x) dx \quad (282)$$

($\because \eta(a) = \eta(b) = 0$)

$$= \int_a^b \left(\frac{\partial G}{\partial y} - \frac{d}{dx} \left(\frac{\partial G}{\partial y'} \right) \right) \eta(x) dx \quad (283)$$

$$= 0$$

- 上の式変形では、 $y = y(x) + \epsilon \eta(x)$ であるから

$$\frac{\partial y}{\partial \epsilon} = \eta(x) \quad (284)$$

$$\frac{\partial y'}{\partial \epsilon} = \frac{\partial}{\partial \epsilon} \left(\frac{\partial y}{\partial x} \right) = \frac{\partial}{\partial \epsilon} (y'(x) + \epsilon \eta'(x)) = \eta'(x) \quad (285)$$

$$\frac{\partial x}{\partial \epsilon} = 0 \quad (286)$$

が成立することを利用している

- 任意の $\eta(x)$ について、上式が恒等的に成り立つためには

$$\frac{\partial G}{\partial y} - \frac{d}{dx} \left(\frac{\partial G}{\partial y'} \right) = 0 \quad (287)$$

でなければならないことが分かり、先程と同様に、オイラー-ラグランジュ方程式を得る

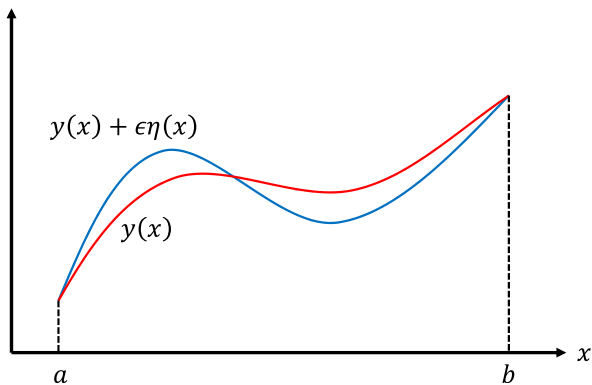


図 10: 制約条件を含んでいる場合の表現

- 汎関数の例 (2)

- 上では $G(y(x), y'(x), x)$ について考えて、変分を導出した
- $y(x)$ と x のみによって決まり、 $y'(x)$ には依存しない関数 $G(y(x), x)$ を考えよう
- 汎関数 $F[y]$ は、先程と同様に以下で表されとする

$$F[y] = \int_a^b G(y(x), x) dx \quad (288)$$

- このとき変分 $\delta F / \delta y$ を求めるのは、非常に簡単である
- 先程の式に、 $\partial G / \partial y' = 0$ を代入すれば直ちに得られる

$$\frac{\delta F}{\delta y(x)} = \frac{\partial G}{\partial y} \quad (289)$$

- あるいは以下のように書ける

$$\frac{\delta}{\delta y(x)} \int_a^b G(y(x), x) dx = \frac{\partial}{\partial y} G(y(x), x) \quad (290)$$

- $F[y]$ が最大 (最小) であるとき、以下のオイラー-ラグランジュ方程式が成り立つ

$$\frac{\partial G}{\partial y} = 0 \quad (291)$$

- 汎関数の例 (3)

- 今度は、 $y'(x)$ と x のみによって決まり、 $y(x)$ には依存しない関数 $G(y'(x), x)$ を考えよう
- この場合も変分 $\delta F / \delta y$ を求めるのは簡単である
- $G(y(x), y'(x), x)$ の変分の式に、 $\partial G / \partial y = 0$ を代入すればよい

$$\frac{\delta F}{\delta y(x)} = -\frac{d}{dx} \left(\frac{\partial G}{\partial y'} \right) \quad (292)$$

- オイラー-ラグランジュ方程式は次のようになる

$$-\frac{d}{dx} \left(\frac{\partial G}{\partial y'} \right) = 0 \quad (293)$$

● 汎関数の例 (4)

- $y(x)$ と $y'(x)$ によって決まる関数 $G(y(x), y'(x))$ を考えよう
- このときのオイラー-ラグランジュ方程式を導出してみよう
- $G(y(x), y'(x))$ を x で微分すれば

$$\begin{aligned} & \frac{d}{dx} G(y, y') \\ = & \frac{\partial}{\partial y} G(y, y') \frac{dy}{dx} + \frac{\partial}{\partial y'} G(y, y') \frac{dy'}{dx} \end{aligned} \quad (294)$$

$$= y' \frac{\partial}{\partial y} G(y, y') + \frac{\partial}{\partial y'} G(y, y') \frac{dy'}{dx} \quad (295)$$

となるから

$$y' \frac{\partial}{\partial y} G(y, y') = \frac{d}{dx} G(y, y') - \frac{\partial}{\partial y'} G(y, y') \frac{dy'}{dx} \quad (296)$$

また、オイラー-ラグランジュ方程式の両辺に y' を掛けたものは

$$y' \frac{\partial}{\partial y} G(y, y') - y' \frac{d}{dx} \left(\frac{\partial}{\partial y'} G(y, y') \right) = 0 \quad (297)$$

これらを連立させて

$$y' \frac{d}{dx} \left(\frac{\partial}{\partial y'} G(y, y') \right) = \frac{d}{dx} G(y, y') - \frac{\partial}{\partial y'} G(y, y') \frac{dy'}{dx} \quad (298)$$

$$y' \frac{d}{dx} \left(\frac{\partial}{\partial y'} G(y, y') \right) + \frac{\partial}{\partial y'} G(y, y') \frac{dy'}{dx} = \frac{d}{dx} G(y, y') \quad (299)$$

$$\frac{d}{dx} \left(y' \cdot \frac{\partial}{\partial y'} G(y, y') \right) = \frac{d}{dx} G(y, y') \quad (300)$$

$$\int \left(\frac{d}{dx} \left(y' \cdot \frac{\partial}{\partial y'} G(y, y') \right) \right) dx = \int \left(\frac{d}{dx} G(y, y') \right) dx + C \quad (301)$$

$$G(y, y') = y' \cdot \frac{\partial}{\partial y'} G(y, y') + C \quad (302)$$

となるので、結局オイラー-ラグランジュ方程式は

$$G - y' \frac{\partial G}{\partial y'} = \text{Const.} \quad (303)$$

と書ける

- 4 変分の導入
 - EM アルゴリズムが困難な場合
 - 変分法
 - 変分法で解ける問題の例
 - 変分法のまとめ

変分法で解ける問題の例

- 変分法で解ける問題の例 (1)

- 2 点 $P(0, 0)$ 、 $Q(a, b)$ を結ぶ最短経路は?
- 2 点を結ぶ経路 $y = f(x) (0 \leq x \leq a)$ の長さ l は、次のようになる

$$l = \int_0^a \sqrt{1 + \left(\frac{dy}{dx}\right)^2} dx = \int_0^a \sqrt{1 + y'^2} dx \quad (304)$$

- 被積分項が y' のみの関数となっていることが分かる

変分法で解ける問題の例

- $G(y'(x), x)$ の場合の公式を使えば、 l の $y = f(x)$ による変分が求まる

$$\frac{\delta l}{\delta f(x)} = -\frac{d}{dx} \left(\frac{\partial}{\partial y'} \sqrt{1 + y'^2} \right) \quad (305)$$

$$= -\frac{d}{dx} \left(\frac{1}{2} \frac{1}{\sqrt{1 + y'^2}} \frac{\partial}{\partial y'} (1 + y'^2) \right) \quad (306)$$

$$= -\frac{d}{dx} \frac{y'}{\sqrt{1 + y'^2}} \quad (307)$$

- l を最小化するような $y = f(x)$ は、上式の変分を 0 と等置すれば

$$-\frac{d}{dx} \frac{y'}{\sqrt{1 + y'^2}} = 0 \quad (308)$$

$$\therefore \frac{y'}{\sqrt{1 + y'^2}} = \text{Const.} \quad (309)$$

変分法で解ける問題の例

- これは、 y' が定数であることを意味している
- 従って、 $y = C_0x + C_1$ と書ける
- 以上より、2 点間を結ぶ最短経路は直線である
- $y = f(x)$ の形について、具体的な仮定は特に置いていないことに注意
- 変分法では、関数そのものを最適化する
- 従って、関数の具体的な形については、特に仮定する必要がない

4 変分の導入

- EM アルゴリズムが困難な場合
- 変分法
- 変分法で解ける問題の例
- 変分法のまとめ

変分法のまとめ

- 変分のまとめ

- これまでの計算で、次の変分が明らかとなった

$$\begin{aligned} \frac{\delta}{\delta y(x)} \int G(y(x), y'(x), x) dx = \\ \frac{\partial}{\partial y} G(y(x), y'(x), x) - \frac{d}{dx} \left(\frac{\partial}{\partial y'} G(y(x), y'(x), x) \right) \end{aligned} \quad (310)$$

$$\frac{\delta}{\delta y(x)} \int G(y(x), x) dx = \frac{\partial}{\partial y} G(y(x), x) \quad (311)$$

$$\frac{\delta}{\delta y(x)} \int G(y'(x), x) dx = -\frac{d}{dx} \left(\frac{\partial}{\partial y'} G(y'(x), x) \right) \quad (312)$$

変分法のまとめ

- ここまでの話の流れ

- 1 変分の定義や、変分の計算法について調査した
- 2 **変分** (汎関数の微分) とは、入力関数が微小に変化したときの、出力値の変化量として定義される
- 3 汎関数が特定の形で表せるとき、変分がどのようなになるか計算した
- 4 汎関数が最大 (最小) になるとき、**オイラー-ラグランジュ方程式**が成立した
- 5 変分法を用いて、2 点間を結ぶ最短経路が**直線**になることを確認した

- これからの話の流れ

- 変分最適化を、どのように推論問題に適用するのかについて調べていく

目次

- 1 K-Means 法
- 2 混合ガウス分布
- 3 EM アルゴリズム
- 4 変分の導入
- 5 近似推論法**
- 6 変分自己符号化器

5 近似推論法

- 変分推論
- 座標降下法
- MAP 推定および最尤推定と変分推論
- 変分近似による弊害
- 離散潜在変数の変分推論
- 連続潜在変数の変分推論
- 正則化との関わり
- 雑多な話題

- 変分推論が必要だった理由

- 潜在変数に関する事後分布 $p(\mathbf{Z}|\mathbf{X}, \theta)$ の計算は、困難であることが多い
- どのような場合に困難になるのか、次の図 11 に示す
- $p(\mathbf{Z}|\mathbf{X}, \theta)$ の厳密な計算は諦める代わりに、別の確率分布で近似したい
- 別の確率分布で近似するとき、単純な項の積として表現できるといった、何らかの仮定を置く

Figure 19.1: Intractable inference problems in deep learning are usually the result of interactions between latent variables in a structured graphical model. These can be due to edges directly connecting one latent variable to another, or due to longer paths that are activated when the child of a V-structure is observed. *(Left)* A **semi-restricted Boltzmann machine** (Osindero and Hinton, 2008) with connections between hidden units. These direct connections between latent variables make the posterior distribution intractable due to large cliques of latent variables. *(Center)* A deep Boltzmann machine, organized into layers of variables without intra-layer connections, still has an intractable posterior distribution due to the connections between layers. *(Right)* This directed model has interactions between latent variables when the visible variables are observed, because every two latent variables are co-parents. Some probabilistic models are able to provide tractable inference over the latent variables despite having one of the graph structures depicted above. This is possible if the conditional probability distributions are chosen to introduce additional independences beyond those described by the graph. For example, probabilistic PCA has the graph structure shown in the right, yet still has simple inference

図 12: 事後分布 $p(\mathbf{Z}|\mathbf{X}, \boldsymbol{\theta})$ の計算が困難な場合

- 事後分布 $p(\mathbf{Z}|\mathbf{X}, \theta)$ の計算が困難な場合 1
 - グラフィカルモデルにおいて、**潜在変数間の相互作用がある**場合、計算が困難になる
 - 左側の**半制限付きボルツマンマシン**では、全ての潜在変数の組み合わせ間で、接続がある
 - 従って、潜在変数間に**依存関係**が存在し、事後分布の計算が手に負えない
- 事後分布 $p(\mathbf{Z}|\mathbf{X}, \theta)$ の計算が困難な場合 2
 - 中央は、層間の結合がない潜在変数の層で構成される、**深層ボルツマンマシン**を表す
 - 潜在変数の層間の結合があるため、事後分布の計算が手に負えない

5 近似推論法

- 変分推論
- 座標降下法
- MAP 推定および最尤推定と変分推論
- 変分近似による弊害
- 離散潜在変数の変分推論
- 連続潜在変数の変分推論
- 正則化との関わり
- 雑多な話題

- 変分推論の目的

- 同時分布 $p(X, Z)$ が分かっているときに、事後分布 $p(Z|X)$ の近似を求めたい

- 注意点

- X は、観測データ $\{x_1, \dots, x_N\}$ をまとめた行列 (i 行ベクトルが x_i^T)
- X の定義は、これまでと全く同様である
- Z は、データ X とは異なり、観測することができない潜在変数を全てまとめた行列
- Z の定義は、これまでとは異なる
- 今まで、各データ x_i に対して存在する変数 z_i のみを潜在変数として扱ってきた
- これ以降、確率モデルのパラメータも潜在変数に含めることにする

- 即ち、潜在変数には、各データ x_i の裏側に潜んでいて表には現れない、データの実体 z_i (今までに考えてきた潜在変数) が含まれるかもしれない
- また、各データ x_i について存在する本質的な情報 z_i だけでなく、モデル全体の構造に組み込まれた、何らかの変数 (データに結び付いているのではなく、モデルに対して存在するパラメータ) も含まれるかもしれない
- 要するに、 Z は、**データ X としては観測できない変数を一般に指す**
- EM アルゴリズムにおける潜在変数とは、**意味が異なっている**ので注意
- 事後分布 $p(Z|X)$ の近似
 - 事後分布 $p(Z|X)$ を解析的に導出するのは困難である
 - $p(Z|X)$ を直接計算するのは諦めて、別の分布 $q(Z)$ で近似する
 - $q(Z)$ は、 $p(Z|X)$ よりは簡単に求められることが前提

- 分布 $q(\mathbf{Z})$ は、 $q(\mathbf{Z}|\mathbf{X})$ のように、データ \mathbf{X} を明示的に書くこともある
(両方を適宜使い分けているので注意)
- $q(\mathbf{Z})$ を、 $p(\mathbf{Z}|\mathbf{X})$ にできるだけ近づけたい
- そのためには、KL ダイバージェンス $\text{KL}(q(\mathbf{Z})||p(\mathbf{Z}|\mathbf{X}))$ を最小化すればよい
- KL ダイバージェンス $\text{KL}(q(\mathbf{Z})||p(\mathbf{Z}|\mathbf{X}))$ の最小化
 - 当然であるが、KL ダイバージェンス $\text{KL}(q(\mathbf{Z})||p(\mathbf{Z}|\mathbf{X}))$ を最小化するためには、KL ダイバージェンス $\text{KL}(q(\mathbf{Z})||p(\mathbf{Z}|\mathbf{X}))$ を計算できなければならない
 - KL ダイバージェンス $\text{KL}(q(\mathbf{Z})||p(\mathbf{Z}|\mathbf{X}))$ を計算するためには、 $p(\mathbf{Z}|\mathbf{X})$ が分かっている必要がある
 - しかし $p(\mathbf{Z}|\mathbf{X})$ を計算するのは困難であったため、KL ダイバージェンスは求められない

- これは本末転倒である
- そこで、KL ダイバージェンスの最小化を、別の問題で置き換えることにする
- 後述するように、KL ダイバージェンスの最小化は、変分下界の最大化と等価である
- これを示すために、まずは周辺分布の対数 $\ln p(\mathbf{X})$ を 2 つの項に分解してみよう
- $p(\mathbf{X})$ は、確率モデルからデータ \mathbf{X} が生起する確率である
- データからみたモデルの好みと解釈できるから、 $p(\mathbf{X})$ をモデルエビデンスという

- 周辺分布の対数 $\ln p(\mathbf{X})$ の分解
 - EM アルゴリズムのときと同様に計算でき、次のようになる

$$\ln p(\mathbf{X}) = \mathcal{L}(q) + \text{KL}(q||p) \quad (313)$$

- 但し、 $\mathcal{L}(q)$ と $\text{KL}(q||p)$ は次のように定義した

$$\mathcal{L}(q) = \int q(\mathbf{Z}) \ln \frac{p(\mathbf{X}, \mathbf{Z})}{q(\mathbf{Z})} d\mathbf{Z} \quad (314)$$

$$\text{KL}(q||p) = - \int q(\mathbf{Z}) \ln \frac{p(\mathbf{Z}|\mathbf{X})}{q(\mathbf{Z})} d\mathbf{Z} \quad (315)$$

- ここでは連続潜在変数について考えるが、離散潜在変数であれば、積分を \mathbf{Z} に関する総和に置き換えればよい

- 下界 $\mathcal{L}(q)$ を最適化する動機

- $\mathcal{L}(q)$ はエビデンスの対数 $\ln p(\mathbf{X})$ の下界であるから、エビデンス下界 (Evidence lower bound, ELBO) ともいう
- または、変分下界、変分下限、負の変分自由エネルギー (Variational free energy) などという
- $\ln p(\mathbf{X})$ は q には依存しないため、定数項とみなせる
- 従って、 $\mathcal{L}(q)$ を q について最大化することは、 $\text{KL}(q||p)$ の最小化に相当
- このとき、分布 $q(\mathbf{Z})$ を真の事後分布 $p(\mathbf{Z}|\mathbf{X})$ に近づけられる
- $q(\mathbf{Z}) = q(\mathbf{Z}|\mathbf{X}) = p(\mathbf{Z}|\mathbf{X})$ が分かれば、データ \mathbf{X} から、潜在変数やパラメータ \mathbf{Z} が得られる
- これにより、KL ダイバージェンス $\text{KL}(q||p)$ の最小化という問題は、変分下界 $\mathcal{L}(q)$ の最大化に読み替えてよいことが分かった

- 下界 $\mathcal{L}(q)$ の最適化

- EM アルゴリズムのときと同じように、下界 $\mathcal{L}(q)$ を、分布 $q(\mathbf{Z})$ について最大化する
- これは、KL ダイバージェンス $\text{KL}(q||p)$ を最小化することと等価である
- 従って、もし $q(\mathbf{Z})$ を任意の分布にしていれば、 $q(\mathbf{Z}) = p(\mathbf{Z}|\mathbf{X})$ において、KL ダイバージェンスを 0 にすればよい
- しかしここでは、真の事後分布 $p(\mathbf{Z}|\mathbf{X})$ を求めることは不可能という状況を想定していた

- 分布 $q(\mathbf{Z})$ の近似

- 計算コストを削減するために、 $q(\mathbf{Z})$ の形をある程度制限する
- 制限したクラスの $q(\mathbf{Z})$ の中で、KL ダイバージェンス $\text{KL}(q||p)$ を最小化するものを探す

● 変分推論の目的

- 分布のクラスを制限することで、 $q(\mathbf{Z})$ を計算可能にすること
- 表現力が豊かな分布のクラスを使えば、真の事後分布 $p(\mathbf{Z}|\mathbf{X})$ を良く近似できる
- 計算可能な分布のクラスの中で、可能な限り豊かな表現力を持つものを選びたい
- 表現力が豊かな分布を使うことは、真の事後分布を、精度良く近似することにつながるのであって、過学習は起こらない

- 分布 $q(\mathbf{Z})$ のクラスを制限する方法
 - 例えば分布 $q(\mathbf{Z})$ を、**パラメトリックな分布に限定**することができる
 - 即ち、パラメータベクトル ω によって $q(\mathbf{Z}|\omega)$ と記述されるような、分布に制限する
 - 分布 $q(\mathbf{Z})$ を、ガウス分布などの、何らかの特別なパラメトリックな分布と仮定することに相当

- クラスを制限する別の方法 (平均場近似)
 - 分布 $q(\mathbf{Z})$ のクラスを制限する別の方法として、平均場近似がある
 - 潜在変数 \mathbf{Z} を、 M 個の互いに排反なグループ $\{\mathbf{Z}_1, \dots, \mathbf{Z}_M\}$ に分割
 - 分布 $q(\mathbf{Z})$ が、これらのグループによって分解されると仮定

$$q(\mathbf{Z}) = \prod_{i=1}^M q_i(\mathbf{Z}_i) \quad (316)$$

- 分布 q について、これ以上の仮定はしない
- 従って、各因子 $q_i(\mathbf{Z}_i)$ の関数形については、何の制限も課さない
- 平均場近似とは、元々は物理学における用語である

- 下界 $\mathcal{L}(q)$ の最大化

- $q(\mathbf{Z}) = \prod_i q_i(\mathbf{Z}_i)$ と分解できるような分布 $q(\mathbf{Z})$ の中で、**下界 $\mathcal{L}(q)$ を最大にするもの**を探す
- $\mathcal{L}(q)$ を $q(\mathbf{Z})$ について最大化するために、 $\mathcal{L}(q)$ を各因子 $q_i(\mathbf{Z}_i)$ について**順番に最大化**していくことが考えられる
- $\mathcal{L}(q)$ に $q(\mathbf{Z}) = \prod_i q_i(\mathbf{Z}_i)$ を代入して、因子の一つ $q_j(\mathbf{Z}_j)$ に関する**依存項**を抜き出してみよう

$$\mathcal{L}(q) = \int q(\mathbf{Z}) \ln \frac{p(\mathbf{X}, \mathbf{Z})}{q(\mathbf{Z})} d\mathbf{Z} \quad (317)$$

$$= \int \left(\prod_i q_i(\mathbf{Z}_i) \right) \ln \frac{p(\mathbf{X}, \mathbf{Z})}{q(\mathbf{Z})} d\mathbf{Z} \quad (318)$$

$$= \int \prod_i q_i(\mathbf{Z}_i) \left(\ln p(\mathbf{X}, \mathbf{Z}) - \sum_i \ln q_i(\mathbf{Z}_i) \right) d\mathbf{Z} \quad (319)$$

$$= \int \prod_i q_i(\mathbf{Z}_i) (\ln p(\mathbf{X}, \mathbf{Z})) d\mathbf{Z} - \int \prod_i q_i(\mathbf{Z}_i) \left(\sum_i \ln q_i(\mathbf{Z}_i) \right) d\mathbf{Z} \quad (320)$$

ここで第 1 項は

$$\begin{aligned} & \int \prod_i q_i(\mathbf{Z}_i) (\ln p(\mathbf{X}, \mathbf{Z})) d\mathbf{Z} \\ &= \int q_j(\mathbf{Z}_j) \left(\prod_{i \neq j} q_i(\mathbf{Z}_i) \right) (\ln p(\mathbf{X}, \mathbf{Z})) d\mathbf{Z} \end{aligned} \quad (321)$$

$d\mathbf{Z} = d\mathbf{Z}_1 d\mathbf{Z}_2 \cdots d\mathbf{Z}_M$ であるから

$$= \int q_j(\mathbf{Z}_j) \left(\prod_{i \neq j} q_i(\mathbf{Z}_i) \right) (\ln p(\mathbf{X}, \mathbf{Z})) d\mathbf{Z}_1 d\mathbf{Z}_2 \cdots d\mathbf{Z}_M \quad (322)$$

$$= \int q_j(\mathbf{Z}_j) (\ln p(\mathbf{X}, \mathbf{Z})) \left(\prod_{i \neq j} q_i(\mathbf{Z}_i) \right) \left(\prod_{i \neq j} d\mathbf{Z}_i \right) d\mathbf{Z}_j \quad (323)$$

$$= \int q_j(\mathbf{Z}_j) \left(\int (\ln p(\mathbf{X}, \mathbf{Z})) \prod_{i \neq j} q_i(\mathbf{Z}_i) d\mathbf{Z}_i \right) d\mathbf{Z}_j \quad (324)$$

$$= \int q_j(\mathbf{Z}_j) \ln \tilde{p}(\mathbf{X}, \mathbf{Z}_j) d\mathbf{Z}_j \quad (325)$$

- 但し、新しい分布 $\tilde{p}(\mathbf{X}, \mathbf{Z}_j)$ は以下の式で定義した (積分の結果であるため、定数項が出現する)

$$\ln \tilde{p}(\mathbf{X}, \mathbf{Z}_j) = \mathbb{E}_{i \neq j} [\ln p(\mathbf{X}, \mathbf{Z})] + \text{Const.} \quad (326)$$

- 記法 $\mathbb{E}_{i \neq j}$ は、 $i \neq j$ をみたす全ての分布 $q_i(\mathbf{Z}_i)$ による、期待値を取ることを表す

$$\mathbb{E}_{i \neq j} [\ln p(\mathbf{X}, \mathbf{Z})] = \int (\ln p(\mathbf{X}, \mathbf{Z})) \prod_{i \neq j} q_i(\mathbf{Z}_i) d\mathbf{Z} \quad (327)$$

- 第2項は

$$\begin{aligned} & \int \prod_i q_i(\mathbf{Z}_i) \left(\sum_i \ln q_i(\mathbf{Z}_i) \right) d\mathbf{Z} \\ &= \sum_i \int \prod_i q_i(\mathbf{Z}_i) (\ln q_i(\mathbf{Z}_i)) d\mathbf{Z} \end{aligned} \quad (328)$$

$$\begin{aligned} &= \sum_i \int \prod_k q_k(\mathbf{Z}_k) (\ln q_i(\mathbf{Z}_i)) d\mathbf{Z}_1 d\mathbf{Z}_2 \cdots d\mathbf{Z}_M \\ &= \int \prod_k q_k(\mathbf{Z}_k) (\ln q_j(\mathbf{Z}_j)) d\mathbf{Z}_1 d\mathbf{Z}_2 \cdots d\mathbf{Z}_M + \end{aligned} \quad (329)$$

$$\sum_{i \neq j} \int \prod_k q_k(\mathbf{Z}_k) (\ln q_i(\mathbf{Z}_i)) d\mathbf{Z}_1 d\mathbf{Z}_2 \cdots d\mathbf{Z}_M \quad (330)$$

但し

$$\int \prod_k q_k(\mathbf{Z}_k) (\ln q_j(\mathbf{Z}_j)) d\mathbf{Z}_1 d\mathbf{Z}_2 \cdots d\mathbf{Z}_M \quad (331)$$

$$= \int q_j(\mathbf{Z}_j) \ln q_j(\mathbf{Z}_j) \left(\prod_{k \neq j} q_k(\mathbf{Z}_k) \right) \left(\prod_{k \neq j} d\mathbf{Z}_k \right) d\mathbf{Z}_j \quad (332)$$

$$= \int q_j(\mathbf{Z}_j) \ln q_j(\mathbf{Z}_j) \left(\int \prod_{k \neq j} q_k(\mathbf{Z}_k) d\mathbf{Z}_k \right) d\mathbf{Z}_j \quad (333)$$

$$= \int q_j(\mathbf{Z}_j) \ln q_j(\mathbf{Z}_j) \prod_{k \neq j} \underbrace{\left(\int q_k(\mathbf{Z}_k) d\mathbf{Z}_k \right)}_{=1} d\mathbf{Z}_j \quad (334)$$

$$= \int q_j(\mathbf{Z}_j) \ln q_j(\mathbf{Z}_j) d\mathbf{Z}_j \quad (335)$$

であるほか

$$\begin{aligned} & \sum_{i \neq j} \int \prod_k q_k(\mathbf{Z}_k) (\ln q_i(\mathbf{Z}_i)) d\mathbf{Z}_1 d\mathbf{Z}_2 \cdots d\mathbf{Z}_M \\ &= \sum_{i \neq j} \int q_i(\mathbf{Z}_i) q_j(\mathbf{Z}_j) \left(\prod_{k \neq i, j} \ln q_k(\mathbf{Z}_k) \right) \\ & \quad (\ln q_i(\mathbf{Z}_i)) \left(\prod_{k \neq i, j} d\mathbf{Z}_k \right) d\mathbf{Z}_i d\mathbf{Z}_j \quad (336) \\ &= \sum_{i \neq j} \underbrace{\left(\int q_j(\mathbf{Z}_j) d\mathbf{Z}_j \right)}_{=1} \left(\int \prod_{k \neq i, j} \ln q_k(\mathbf{Z}_k) d\mathbf{Z}_k \right) \end{aligned}$$

$$\int \ln q_i(\mathbf{Z}_i) q_i(\mathbf{Z}_i) d\mathbf{Z}_i \quad (337)$$

$$= \sum_{i \neq j} \prod_{k \neq i, j} \underbrace{\left(\int \ln q_k(\mathbf{Z}_k) d\mathbf{Z}_k \right)}_{=\text{Const.}} \underbrace{\int \ln q_i(\mathbf{Z}_i) q_i(\mathbf{Z}_i) d\mathbf{Z}_i}_{=\text{Const.}} \quad (338)$$

$$= \sum_{i \neq j} \text{Const.} = \text{Const.} \quad (339)$$

となるから結局

$$\begin{aligned} & \int \prod_i q_i(\mathbf{Z}_i) \left(\sum_i \ln q_i(\mathbf{Z}_i) \right) d\mathbf{Z} \\ &= \int q_j(\mathbf{Z}_j) \ln q_j(\mathbf{Z}_j) d\mathbf{Z}_j + \text{Const.} \end{aligned} \quad (340)$$

- 従って、下界 $\mathcal{L}(q)$ から $q_j(\mathbf{Z}_j)$ に依存する項を取り出すと

$$\mathcal{L}(q) = \int q_j(\mathbf{Z}_j) \ln \tilde{p}(\mathbf{X}, \mathbf{Z}_j) d\mathbf{Z}_j - \int q_j(\mathbf{Z}_j) \ln q_j(\mathbf{Z}_j) d\mathbf{Z}_j + \text{Const.} \quad (341)$$

但し

$$\ln \tilde{p}(\mathbf{X}, \mathbf{Z}_j) = \mathbb{E}_{i \neq j} [\ln p(\mathbf{X}, \mathbf{Z})] = \int \ln p(\mathbf{X}, \mathbf{Z}) \prod_{i \neq j} q_i(\mathbf{Z}_i) d\mathbf{Z}_i \quad (342)$$

- $\mathcal{L}(q)$ を、 $i \neq j$ である全ての $q_i(\mathbf{Z}_i)$ について固定した上で、 $q_j(\mathbf{Z}_j)$ について最大化することになる
- $q_j(\mathbf{Z}_j)$ について可能な全ての分布の中で、 $\mathcal{L}(q)$ を最大にするようなものを選ぶ

- $\mathcal{L}(q)$ は次のように変形できる

$$\mathcal{L}(q) = \int q_j(\mathbf{Z}_j) \ln \frac{\tilde{p}(\mathbf{X}, \mathbf{Z}_j)}{q_j(\mathbf{Z}_j)} d\mathbf{Z}_j + \text{Const.} \quad (343)$$

$$= -\text{KL}(q_j(\mathbf{Z}_j) \parallel \tilde{p}(\mathbf{X}, \mathbf{Z}_j)) + \text{Const.} \quad (344)$$

- これより、 $\mathcal{L}(q)$ は $q_j(\mathbf{Z}_j)$ と $\tilde{p}(\mathbf{X}, \mathbf{Z}_j)$ の間の、負の KL ダイバージェンスとなっている
- そして、 $\mathcal{L}(q)$ の $q_j(\mathbf{Z}_j)$ に関する最大化は、**KL ダイバージェンスの最小化**と等価
- KL ダイバージェンスを最小にするためには、 $q_j(\mathbf{Z}_j) = \tilde{p}(\mathbf{X}, \mathbf{Z}_j)$ とすればよい
- 従って、 $q_j(\mathbf{Z}_j)$ の最適解は、次のように書ける

$$\ln q_j^*(\mathbf{Z}_j) = \mathbb{E}_{i \neq j}[\ln p(\mathbf{X}, \mathbf{Z})] + \text{Const.} \quad (345)$$

- 下界 $\mathcal{L}(q)$ を最大化する $\ln q_j(\mathbf{Z}_j)$ の解
 - $q_j(\mathbf{Z}_j)$ の最適解は次のように書けた

$$\ln q_j^*(\mathbf{Z}_j) = \mathbb{E}_{i \neq j} [\ln p(\mathbf{X}, \mathbf{Z})] + \text{Const.} \quad (346)$$

- 上式は次のことを意味している
- 因子 $q_j(\mathbf{Z}_j)$ の最適解の対数 $\ln q_j^*(\mathbf{Z}_j)$ は、観測データ \mathbf{X} と潜在変数 \mathbf{Z} の同時分布の対数 $\ln p(\mathbf{X}, \mathbf{Z})$ を考え、 $i \neq j$ である他の因子 $q_i(\mathbf{Z}_i)$ について期待値を取ったものである
- 定数項は、正規化することで得られるので、結局次のようになる

$$q_j^*(\mathbf{Z}_j) = \frac{\exp(\mathbb{E}_{i \neq j} [\ln p(\mathbf{X}, \mathbf{Z})])}{\int \exp(\mathbb{E}_{i \neq j} [\ln p(\mathbf{X}, \mathbf{Z})]) d\mathbf{Z}_j} \quad (347)$$

- 正規化定数は必要に応じて計算すればよいので、取り敢えず無視できる

- 最適解の式は、 $q(\mathbf{Z})$ の分解の数だけ得られるので、 $\{q_i(\mathbf{Z}_i)\}$ に関する M 本の連立方程式となる
- この方程式は、分布 $q(\mathbf{Z})$ が M 個の因子に分解されるという仮定の下で、下界 $\mathcal{L}(q)$ の最大値が満たすべき条件である
- $\ln q_j^*(\mathbf{Z}_j)$ の右辺は、 $i \neq j$ である $q_i(\mathbf{Z}_i)$ の期待値に依存するため、 $q_j^*(\mathbf{Z}_j)$ を陽に求めることができない
- そこで、下界 $\mathcal{L}(q)$ は次のように最適化される (**重要**)
- $i \neq j$ である全ての $q_i(\mathbf{Z}_i)$ を**固定**した状態で、 $q_j(\mathbf{Z}_j)$ を最適化することを、全ての $j = 1, \dots, M$ について繰り返す手続きを、**座標降下法**という

- 下界 $\mathcal{L}(q)$ の最適化 (座標降下法)

- 1 全ての因子 $q_j(\mathbf{Z}_j)$ を適当に初期化する

- 2 各因子を、以下の式を使って更新する

$$\ln q_j(\mathbf{Z}_j) = \mathbb{E}_{i \neq j} [\ln p(\mathbf{X}, \mathbf{Z})] + \text{Const.} \quad (348)$$

$$\mathbb{E}_{i \neq j} [\ln p(\mathbf{X}, \mathbf{Z})] = \int \ln p(\mathbf{X}, \mathbf{Z}) \prod_{i \neq j} q_i(\mathbf{Z}_i) d\mathbf{Z}_i \quad (349)$$

即ち、因子 $q_j(\mathbf{Z}_j)$ を、他の全ての因子の現在の値 $q_i(\mathbf{Z}_i)$ を使って改良する

- 3 (2) を、下界 $\mathcal{L}(q)$ が収束するまで繰り返す

- ここまでの話の流れ

- 1 $\ln p(\mathbf{X}) = \mathcal{L}(q) + \text{KL}(q(\mathbf{Z})||p(\mathbf{Z}|\mathbf{X}))$ であるから、エビデンス下界 $\mathcal{L}(q)$ を q について最大化すれば、 $\text{KL}(q||p) = 0$ とでき、従って $q(\mathbf{Z}) = p(\mathbf{Z}|\mathbf{X})$ を得る
- 2 $q(\mathbf{Z}) = p(\mathbf{Z}|\mathbf{X})$ が分かれば、データ \mathbf{X} から、潜在変数やパラメータ \mathbf{Z} が得られる (パラメータは潜在変数 \mathbf{Z} に含まれている)
- 3 しかし、事後分布 $q(\mathbf{Z}) = p(\mathbf{Z}|\mathbf{X})$ は計算不可能なので、何らかの方法で近似するしかない
- 4 近似するといっても、計算可能でなければならないので、 $q(\mathbf{Z})$ の形には、通常何らかの制限を課す
- 5 $q(\mathbf{Z})$ を、パラメトリックな分布 $q(\mathbf{Z}|\omega)$ と仮定することがある

- 6 または、 $q(\mathbf{Z})$ を、 $\prod_i q_i(\mathbf{Z}_i)$ のように分解できるとする (平均場近似)
 - 7 平均場近似を行うとき、各因子 $q_j(\mathbf{Z}_j)$ の最適解は、 $\ln q_j^*(\mathbf{Z}_j) = \mathbb{E}_{i \neq j} [\ln p(\mathbf{X}, \mathbf{Z})] + \text{Const.}$ であった
 - 8 全ての因子 $\{q_j(\mathbf{Z}_j)\}$ を同時に最適化することはできない
 - 9 下界 $\mathcal{L}(q)$ を、各因子 $q_j(\mathbf{Z}_j)$ について順番に最適化することはできる
- これからの話の流れ
 - MAP 推定と最尤推定は、変分推論の特殊な場合であることを確認する

5 近似推論法

- 変分推論
- 座標降下法
- MAP 推定および最尤推定と変分推論
- 変分近似による弊害
- 離散潜在変数の変分推論
- 連続潜在変数の変分推論
- 正則化との関わり
- 雑多な話題

- MAP 推定および最尤推定の変分推論からの導出
 - 変分推論で得たいのは、潜在変数に関する事後分布 $p(\mathbf{Z}|\mathbf{X})$ である
 - この変分推論の特殊ケースが、MAP 推定や最尤推定であることを導く
 - 変分推論では、エビデンス下界 $\mathcal{L}(q)$ を q について最大化し、従って $\text{KL}(q||p)$ を最小化する
 - いま、分布 $q(\mathbf{Z})$ が、次のデルタ関数であるとする
 - 特定の値 $\mathbf{Z} = \hat{\mathbf{Z}}$ についてのみ確率が非零になる、無限に鋭い分布

$$q(\mathbf{Z}) = \delta(\mathbf{Z} - \hat{\mathbf{Z}}) \quad (350)$$

- このとき、KL ダイバージェンス $\text{KL}(q||p)$ は次のようになる

$$\begin{aligned} & \text{KL}(q||p) \\ \equiv & \text{KL}(q(\mathbf{Z})||p(\mathbf{Z}|\mathbf{X})) \\ = & - \int q(\mathbf{Z}) \ln \frac{p(\mathbf{Z}|\mathbf{X})}{q(\mathbf{Z})} d\mathbf{Z} \\ = & - \int q(\mathbf{Z}) \ln p(\mathbf{Z}|\mathbf{X}) d\mathbf{Z} + \int q(\mathbf{Z}) \ln q(\mathbf{Z}) d\mathbf{Z} \quad (351) \end{aligned}$$

$$\begin{aligned} = & - \int \delta(\mathbf{Z} - \hat{\mathbf{Z}}) \ln p(\mathbf{Z}|\mathbf{X}) d\mathbf{Z} + \\ & \int \delta(\mathbf{Z} - \hat{\mathbf{Z}}) \ln \delta(\mathbf{Z} - \hat{\mathbf{Z}}) d\mathbf{Z} \quad (352) \end{aligned}$$

$$= - \ln p(\hat{\mathbf{Z}}|\mathbf{X}) + \ln \delta(\hat{\mathbf{Z}} - \hat{\mathbf{Z}}) \quad (353)$$

$$= - \ln p(\hat{\mathbf{Z}}|\mathbf{X}) + \text{Const.} \quad (354)$$

- これを $q(\mathbf{Z})$ について最小化することは、 $\hat{\mathbf{Z}}$ について最小化することに相当する
- よって $\hat{\mathbf{Z}}$ の最適解 $\hat{\mathbf{Z}}^*$ は

$$\hat{\mathbf{Z}}^* = \arg \min_{\hat{\mathbf{Z}}} \left(-\ln p(\hat{\mathbf{Z}}|\mathbf{X}) \right) \quad (355)$$

$$= \arg \max_{\hat{\mathbf{Z}}} \ln p(\hat{\mathbf{Z}}|\mathbf{X}) \quad (356)$$

$$= \arg \max_{\hat{\mathbf{Z}}} p(\hat{\mathbf{Z}}|\mathbf{X}) \quad (357)$$

$$= \arg \max_{\hat{\mathbf{Z}}} p(\mathbf{X}|\hat{\mathbf{Z}})p(\hat{\mathbf{Z}}) \quad (358)$$

となるから、**MAP 推定** (**最大事後確率推定**) の式と一致する

- これより、MAP 推定は、KL ダイバージェンス $\text{KL}(q||p)$ の最小化、従って、**エビデンス下界 $\mathcal{L}(q)$ の最大化と等価**である

- 更に、 $p(\hat{\mathbf{Z}}) = \text{Const.}$ 、即ち $\hat{\mathbf{Z}}$ に関する事前の情報がないとすると

$$\begin{aligned}\hat{\mathbf{Z}}^* &= \arg \max_{\hat{\mathbf{Z}}} p(\mathbf{X}|\hat{\mathbf{Z}})p(\hat{\mathbf{Z}}) \\ &= \arg \max_{\hat{\mathbf{Z}}} p(\mathbf{X}|\hat{\mathbf{Z}})\end{aligned}\tag{359}$$

となるから、これは**最尤推定**の式に一致する

- 下界 $\mathcal{L}(q)$ の式から導くことも可能である

$$\begin{aligned}\mathcal{L}(q) &= \int q(\mathbf{Z}) \ln \frac{p(\mathbf{X}, \mathbf{Z})}{q(\mathbf{Z})} d\mathbf{Z} \\ &= \int q(\mathbf{Z}) \ln p(\mathbf{X}, \mathbf{Z}) d\mathbf{Z} - \int q(\mathbf{Z}) \ln q(\mathbf{Z}) d\mathbf{Z} \quad (360)\end{aligned}$$

$$\begin{aligned}&= \int \delta(\mathbf{Z} - \hat{\mathbf{Z}}) \ln p(\mathbf{X}, \mathbf{Z}) d\mathbf{Z} - \\ &\quad \int \delta(\mathbf{Z} - \hat{\mathbf{Z}}) \ln \delta(\mathbf{Z} - \hat{\mathbf{Z}}) d\mathbf{Z} \quad (361)\end{aligned}$$

$$= \ln p(\mathbf{X}, \hat{\mathbf{Z}}) - \ln \delta(\hat{\mathbf{Z}} - \hat{\mathbf{Z}}) \quad (362)$$

$$= \ln p(\mathbf{X}, \hat{\mathbf{Z}}) + \text{Const.} \quad (363)$$

$$= \ln p(\mathbf{X} | \hat{\mathbf{Z}}) p(\hat{\mathbf{Z}}) + \text{Const.} \quad (364)$$

- これより、下界 $\mathcal{L}(q)$ の最大化は、 $\hat{\mathbf{Z}}$ に関する $p(\hat{\mathbf{Z}} | \mathbf{X}) \propto p(\mathbf{X} | \hat{\mathbf{Z}}) p(\hat{\mathbf{Z}})$ の最大化、**即ち MAP 推定と等価**である

MAP 推定の例

- スパース符号化

- ここでは MAP 推定の例として、スパース符号化を扱う (変分推論の例ではない)
- データ x に対する潜在変数 z に、スパース性を持たせる
- そのために、スパース性を導く事前分布 (ラプラス事前分布) を、潜在変数に用いる
- データ x を D 次元、また潜在変数 z を K 次元とする
- データ x に対応する潜在変数を z_i と表し、 z_i の第 k 成分を z_{ik} とする

$$p(z_{ik}|\lambda) = \frac{\lambda}{2} \exp(-\lambda|z_{ik}|) \quad (365)$$

$$p(z_i|\lambda) = \prod_k p(z_{ik}) = \prod_k \frac{\lambda}{2} \exp(-\lambda|z_{ik}|) \quad (366)$$

MAP 推定の例

- データに関する事後分布 $p(\mathbf{x}_i | \mathbf{z}_i, \mathbf{W}, \mathbf{b}, \beta)$ を、次で定義する

$$\begin{aligned} & p(\mathbf{x}_i | \mathbf{z}_i, \mathbf{W}, \mathbf{b}, \beta) \\ &= \mathcal{N}(\mathbf{x}_i | \mathbf{W}\mathbf{z}_i + \mathbf{b}, \beta^{-1}\mathbf{I}) \end{aligned} \quad (367)$$

$$\begin{aligned} &= \frac{1}{(2\pi)^{\frac{D}{2}}} \frac{1}{|\beta^{-1}\mathbf{I}|^{\frac{1}{2}}} \exp \left(-\frac{1}{2} (\mathbf{x}_i - (\mathbf{W}\mathbf{z}_i + \mathbf{b}))^T \right. \\ &\quad \left. (\beta^{-1}\mathbf{I})^{-1} (\mathbf{x}_i - (\mathbf{W}\mathbf{z}_i + \mathbf{b})) \right) \end{aligned} \quad (368)$$

$$\begin{aligned} &= \frac{1}{(2\pi\beta^{-1})^{\frac{D}{2}}} \exp \left(-\frac{\beta}{2} (\mathbf{x}_i - (\mathbf{W}\mathbf{z}_i + \mathbf{b}))^T \right. \\ &\quad \left. (\mathbf{x}_i - (\mathbf{W}\mathbf{z}_i + \mathbf{b})) \right) \end{aligned} \quad (369)$$

- データ \mathbf{x}_i は、対応する潜在変数 \mathbf{z}_i に、線形変換 $\mathbf{W}\mathbf{z}_i + \mathbf{b}$ を施し、更に分散 $\beta^{-1}\mathbf{I}$ のガウスノイズを足すことで、生成されると考える

MAP 推定の例

- パラメータについての補足

- 今回は、 $b = 0$ において**バイアス**を**無視**したものを考える
- λ と、精度 β は**ハイパーパラメータ**であり、予め決められているとする
- そこでこれ以降、次のように確率分布を表現する

$$p(z_i|\lambda) = p(z_i) \quad (370)$$

$$p(x_i|z_i, \mathbf{W}, \mathbf{b}, \beta) = p(x_i|z_i, \mathbf{W}) \quad (371)$$

- \mathbf{X} を、全てのデータ $\{x_i\}$ を集めた行列とする (第 i 行ベクトルが x_i^T)
- \mathbf{Z} についても同様に、全ての潜在変数 $\{z_i\}$ を集めた行列として定める (第 i 行ベクトルが z_i^T)

- スパース符号化の学習

- 事後分布 $p(z_i|x_i)$ は、**表現することすら困難**であるため、最尤推定による手法 (EM アルゴリズムなど) は利用できない

MAP 推定の例

- そこで、最尤推定の代わりに **MAP 推定** を利用することで、最適なパラメータが得られる
- 最大化するのは、以下の事後分布 $p(\mathbf{Z}|\mathbf{X}, \mathbf{W})$ である

$$p(\mathbf{Z}|\mathbf{X}, \mathbf{W}) = \frac{p(\mathbf{X}, \mathbf{Z}|\mathbf{W})}{p(\mathbf{X})} \quad (372)$$

$$\propto p(\mathbf{X}, \mathbf{Z}|\mathbf{W}) \quad (373)$$

$$= p(\mathbf{X}|\mathbf{Z}, \mathbf{W})p(\mathbf{Z}) \quad (374)$$

$$= \left(\prod_i p(\mathbf{x}_i|\mathbf{z}_i, \mathbf{W}) \right) \left(\prod_i p(\mathbf{z}_i) \right) \quad (375)$$

$$= \prod_i p(\mathbf{x}_i|\mathbf{z}_i, \mathbf{W})p(\mathbf{z}_i) \quad (376)$$

MAP 推定の例

- 対数を取って最大化してもよいので、最大化する量は結局

$$\ln p(\mathbf{Z}|\mathbf{X}, \mathbf{W}) = \ln \left(\prod_i p(\mathbf{x}_i|\mathbf{z}_i, \mathbf{W})p(\mathbf{z}_i) \right) \quad (377)$$

$$= \sum_i (\ln p(\mathbf{x}_i|\mathbf{z}_i, \mathbf{W}) + \ln p(\mathbf{z}_i)) \quad (378)$$

$$= \sum_i \ln p(\mathbf{x}_i|\mathbf{z}_i, \mathbf{W}) + \sum_i \ln p(\mathbf{z}_i) \quad (379)$$

各項を求めると

$$\begin{aligned} & \sum_i \ln p(\mathbf{x}_i|\mathbf{z}_i, \mathbf{W}) \\ = & \sum_i \ln \mathcal{N}(\mathbf{x}_i|\mathbf{W}\mathbf{z}_i, \beta^{-1}\mathbf{I}) \\ = & \sum_i \ln \left(\frac{1}{(2\pi\beta^{-1})^{\frac{D}{2}}} \exp \left(-\frac{\beta}{2} (\mathbf{x}_i - \mathbf{W}\mathbf{z}_i)^T (\mathbf{x}_i - \mathbf{W}\mathbf{z}_i) \right) \right) \end{aligned}$$

MAP 推定の例

$$\begin{aligned} &= \sum_i \left(-\frac{D}{2} \ln 2\pi - \frac{D}{2} \ln \beta^{-1} - \right. \\ &\quad \left. \frac{\beta}{2} (\mathbf{x}_i - \mathbf{W} \mathbf{z}_i)^T (\mathbf{x}_i - \mathbf{W} \mathbf{z}_i) \right) \\ &= \sum_i \left(-\frac{\beta}{2} (\mathbf{x}_i - \mathbf{W} \mathbf{z}_i)^T (\mathbf{x}_i - \mathbf{W} \mathbf{z}_i) \right) + \text{Const.} \\ &= -\frac{\beta}{2} \sum_i (\mathbf{x}_i - \mathbf{W} \mathbf{z}_i)^T (\mathbf{x}_i - \mathbf{W} \mathbf{z}_i) + \text{Const.} \\ &= -\frac{\beta}{2} \sum_i (\mathbf{X} - \mathbf{Z} \mathbf{W}^T)_{i:} ((\mathbf{X} - \mathbf{Z} \mathbf{W}^T)_{i:})^T + \text{Const.} \quad (380) \end{aligned}$$

$$= -\frac{\beta}{2} \sum_{i,j} ((\mathbf{X} - \mathbf{Z} \mathbf{W}^T) \odot (\mathbf{X} - \mathbf{Z} \mathbf{W}^T))_{i,j} + \text{Const.} \quad (381)$$

$$= -\frac{\beta}{2} \sum_{i,j} (\mathbf{X} - \mathbf{Z} \mathbf{W}^T)_{i,j}^2 + \text{Const.} \quad (382)$$

MAP 推定の例

$$= -\frac{\beta}{2} \|\mathbf{X} - \mathbf{Z}\mathbf{W}^T\|_2^2 + \text{Const.} \quad (383)$$

また

$$\begin{aligned} & \sum_i \ln p(\mathbf{z}_i) \\ &= \sum_i \ln \prod_j \left(\frac{\lambda}{2} \exp(-\lambda |z_{ij}|) \right) \\ &= \sum_i \sum_j \left(\ln \frac{\lambda}{2} - \lambda |z_{ij}| \right) \\ &= -\lambda \sum_{i,j} |z_{ij}| + \text{Const.} \end{aligned} \quad (384)$$

$$= -\lambda \sum_{i,j} |\mathbf{Z}_{i,j}| + \text{Const.} \quad (385)$$

$$= -\lambda \|\mathbf{Z}\|_1 + \text{Const.} \quad (386)$$

MAP 推定の例

のようになる

- \odot はアダマール積、 $A_{i:}$ は行列 A の第 i 行目のベクトルを表す
- また、行列のノルム $\|A\|_p$ は次で定義される

$$\|A\|_p = \left(\sum_{i,j} |A_{i,j}|^p \right)^{\frac{1}{p}} \quad (387)$$

$p = 1, p = 2$ のときは次のようになる

$$\|A\|_1 = \sum_{i,j} |A_{i,j}|, \quad \|A\|_2 = \sqrt{\sum_{i,j} A_{i,j}^2} \quad (388)$$

- これより、 $\|A\|_1$ は、**行列 A の各要素の絶対値の和**を表す
- また $\|A\|_2$ は、行列 A の各要素の二乗和の平方根を表し、**フロベニウスノルム $\|A\|_F$** ともよばれる

MAP 推定の例

- 従って、 $\|\mathbf{A}\|_2^2$ は、**行列 \mathbf{A} の各要素の二乗和**である
- 上記から

$$\begin{aligned} & \ln p(\mathbf{Z}|\mathbf{X}, \mathbf{W}) \\ &= \sum_i \ln p(\mathbf{x}_i|z_i, \mathbf{W}) + \sum_i \ln p(z_i) \\ &= -\frac{\beta}{2}\|\mathbf{X} - \mathbf{Z}\mathbf{W}^T\|_2^2 - \lambda\|\mathbf{Z}\|_1 + \text{Const.} \end{aligned} \quad (389)$$

となるので、 $\ln p(\mathbf{Z}|\mathbf{X}, \mathbf{W})$ の最大化は、以下の**関数 $J(\mathbf{Z}, \mathbf{W})$ の最小化**である

$$J(\mathbf{Z}, \mathbf{W}) = \|\mathbf{X} - \mathbf{Z}\mathbf{W}^T\|_2^2 + \|\mathbf{Z}\|_1 \quad (390)$$

- スパース符号化の学習方法のまとめ
 - 関数 $J(\mathbf{Z}, \mathbf{W})$ を、 \mathbf{Z} と \mathbf{W} について**交互に最小化**すればよい

MAP 推定の例

- $p(\mathbf{Z}|\mathbf{X}, \mathbf{W})$ の最大化は、 $p(\mathbf{X}|\mathbf{Z}, \mathbf{W})p(\mathbf{Z})$ の最大化、即ち $J(\mathbf{Z}, \mathbf{W})$ の最小化と等価であった
- 従って、 $J(\mathbf{Z}, \mathbf{W})$ の各パラメータ \mathbf{Z}, \mathbf{W} についての最小化は、事後確率を大きくする方向に働く
- 関数 $J(\mathbf{Z}, \mathbf{W})$ をもう一度みてみよう

$$J(\mathbf{Z}, \mathbf{W}) = \|\mathbf{X} - \mathbf{Z}\mathbf{W}^T\|_2^2 + \|\mathbf{Z}\|_1 \quad (391)$$

- 第1項は、明らかに再構成誤差を表している
- 第2項は、データ \mathbf{X} の表現 \mathbf{Z} がスパースになるように付加した、正則化項である
- $\mathbf{Z}\mathbf{W}^T$ は、データ \mathbf{X} の内部表現 \mathbf{Z} に、重み \mathbf{W} を掛けることで、データ \mathbf{X} を復元したもの

5 近似推論法

- 変分推論
- 座標降下法
- MAP 推定および最尤推定と変分推論
- 変分近似による弊害
- 離散潜在変数の変分推論
- 連続潜在変数の変分推論
- 正則化との関わり
- 雑多な話題

- これまでの話の流れ

- 1 変分推論から、MAP 推定と最尤推定が導出できることを確認した
- 2 MAP 推定は、 $q(\mathbf{Z})$ を無限に鋭い確率分布 (デルタ関数) として、 \mathbf{Z} が特定の値しか取らないと仮定した場合であった
- 3 最尤推定は、MAP 推定において、 \mathbf{Z} の事前確率分布を設けない場合であった
- 4 MAP 推定の例として、スパース符号化を扱った

- これからの話の流れ

- 話題を変えて、分解 $q(\mathbf{Z}) = \prod_i q_i(\mathbf{Z}_i)$ によって $p(\mathbf{Z}|\mathbf{X})$ を近似するときの、弊害を調べる

- $q(\mathbf{Z})$ の分解による近似の性質
 - 変分推論では、真の事後分布 $q(\mathbf{Z}) = p(\mathbf{Z}|\mathbf{X})$ を、分解により近似する
 - 分解で近似することによって、**どのような不正確さが生じるのか?**
- ガウス分布の分解による近似
 - ガウス分布を、**分解されたガウス分布**で近似することを考えてみよう
 - 分解による近似で、どのような問題が起こるのか考えてみよう
 - 2つの変数 $\mathbf{z} = (z_1, z_2)$ 間には、**相関がある**とする
 - \mathbf{z} はガウス分布 $p(\mathbf{z}) = \mathcal{N}(\mathbf{z}|\boldsymbol{\mu}, \boldsymbol{\Lambda}^{-1})$ に従っているとする

$$\boldsymbol{\mu} = \begin{bmatrix} \mu_1 \\ \mu_2 \end{bmatrix}, \quad \boldsymbol{\Lambda} = \begin{bmatrix} \Lambda_{11} & \Lambda_{12} \\ \Lambda_{21} & \Lambda_{22} \end{bmatrix} \quad (392)$$

- 精度行列 $\boldsymbol{\Lambda}$ は対称行列であるから、 $\Lambda_{12} = \Lambda_{21}$

- この分布 $p(\mathbf{z})$ を、分解されたガウス分布 $q(\mathbf{z}) = q_1(z_1)q_2(z_2)$ で近似する
- 各因子 $q_1(z_1), q_2(z_2)$ の関数形については何の仮定も置いていないことに注意

- 最適な因子 $q_1(z_1), q_2(z_2)$ の計算

- 最適な因子 $q_1^*(z_1)$ を、先程の結果を使って求める

$$\ln q_j(\mathbf{Z}_j) = \mathbb{E}_{i \neq j} [\ln p(\mathbf{X}, \mathbf{Z})] + \text{Const.} \quad (393)$$

- 従って、 $q_1^*(z_1)$ を計算する式は次のようになる

$$\ln q_1^*(z_1) = \mathbb{E}_{z_2} [\ln p(\mathbf{z})] + \text{Const.} \quad (394)$$

$$\mathbb{E}_{z_2} [\ln p(\mathbf{z})] = \int \ln p(\mathbf{z}) q_2(z_2) dz_2 \quad (395)$$

- 上式の右辺では、 z_1 に依存する項だけを考えればよい

- z_1 の関数を求めようとしているため
- z_1 に依存しない項は、全て定数項 (正規化定数) に含まれてしまうため
- 従って $q_1^*(z_1)$ は

$$\ln q_1^*(z_1) = \mathbb{E}_{z_2} [\ln p(\mathbf{z})] + \text{Const.} \quad (396)$$

$$= \mathbb{E}_{z_2} [\ln \mathcal{N}(\mathbf{z} | \boldsymbol{\mu}, \boldsymbol{\Lambda}^{-1})] + \text{Const.} \quad (397)$$

但し

$$\begin{aligned} & \ln \mathcal{N}(\mathbf{z} | \boldsymbol{\mu}, \boldsymbol{\Lambda}^{-1}) \\ = & \ln \left(\frac{1}{(2\pi)^{\frac{2}{2}}} \frac{1}{|\boldsymbol{\Lambda}^{-1}|^{\frac{1}{2}}} \exp \left(-\frac{1}{2} (\mathbf{z} - \boldsymbol{\mu})^T (\boldsymbol{\Lambda}^{-1})^{-1} (\mathbf{z} - \boldsymbol{\mu}) \right) \right) \\ = & \ln \left(\frac{1}{2\pi} \frac{1}{|\boldsymbol{\Lambda}|^{-\frac{1}{2}}} \exp \left(-\frac{1}{2} (\mathbf{z} - \boldsymbol{\mu})^T \boldsymbol{\Lambda} (\mathbf{z} - \boldsymbol{\mu}) \right) \right) \\ = & -\frac{1}{2} (\mathbf{z} - \boldsymbol{\mu})^T \boldsymbol{\Lambda} (\mathbf{z} - \boldsymbol{\mu}) + \text{Const.} \end{aligned} \quad (398)$$

z_1 に依存する項だけを取り出せば

$$\begin{aligned}
 & -\frac{1}{2}(\mathbf{z} - \boldsymbol{\mu})^T \boldsymbol{\Lambda}(\mathbf{z} - \boldsymbol{\mu}) \\
 = & -\frac{1}{2} [z_1 - \mu_1, z_2 - \mu_2] \begin{bmatrix} \Lambda_{11} & \Lambda_{12} \\ \Lambda_{21} & \Lambda_{22} \end{bmatrix} \begin{bmatrix} z_1 - \mu_1 \\ z_2 - \mu_2 \end{bmatrix} \\
 = & -\frac{1}{2} \left[\left\{ (\Lambda_{11}(z_1 - \mu_1) + \Lambda_{21}(z_2 - \mu_2)) \right\} (z_1 - \mu_1) + \right. \\
 & \left. \left\{ \Lambda_{12}(z_1 - \mu_1) + \Lambda_{22}(z_2 - \mu_2) \right\} (z_2 - \mu_2) \right] \\
 = & -\frac{1}{2} (\Lambda_{11}(z_1 - \mu_1)^2 + 2\Lambda_{12}(z_1 - \mu_1)(z_2 - \mu_2) + \\
 & \Lambda_{22}(z_2 - \mu_2)^2) \quad (\because \Lambda_{21} = \Lambda_{12}) \\
 = & -\frac{1}{2} \Lambda_{11}(z_1 - \mu_1)^2 - \Lambda_{12}(z_1 - \mu_1)(z_2 - \mu_2) + \text{Const.} \quad (399)
 \end{aligned}$$

これを代入して

$$\begin{aligned} & \ln \mathcal{N}(\mathbf{z} | \boldsymbol{\mu}, \boldsymbol{\Lambda}^{-1}) \\ = & -\frac{1}{2} \Lambda_{11} (z_1 - \mu_1)^2 - \Lambda_{12} (z_1 - \mu_1) (z_2 - \mu_2) + \text{Const.} \quad (400) \end{aligned}$$

従って

$$\begin{aligned} & \ln q_1^*(z_1) \\ = & \mathbb{E}_{z_2} [\ln \mathcal{N}(\mathbf{z} | \boldsymbol{\mu}, \boldsymbol{\Lambda}^{-1})] + \text{Const.} \\ = & \mathbb{E}_{z_2} \left[-\frac{1}{2} \Lambda_{11} (z_1 - \mu_1)^2 - \Lambda_{12} (z_1 - \mu_1) (z_2 - \mu_2) \right] + \text{Const.} \\ = & -\frac{1}{2} \Lambda_{11} z_1^2 + \Lambda_{11} \mu_1 z_1 - \Lambda_{12} z_1 (\mathbb{E}[z_2] - \mu_2) + \text{Const.} \quad (401) \end{aligned}$$

- これより、 $q_1^*(z_1)$ は次のように書ける

$$\begin{aligned} & q_1^*(z_1) \\ \propto & \exp \left(-\frac{1}{2} \Lambda_{11} z_1^2 + \Lambda_{11} \mu_1 z_1 - \Lambda_{12} z_1 (\mathbb{E}[z_2] - \mu_2) \right) \\ = & \exp \left(-\frac{1}{2} \Lambda_{11} \left(z_1 - (\mu_1 - \Lambda_{11}^{-1} \Lambda_{12} (\mathbb{E}[z_2] - \mu_2)) \right)^2 + \right. \\ & \left. \frac{1}{2} \Lambda_{11} (\mu_1 - \Lambda_{11}^{-1} \Lambda_{12} (\mathbb{E}[z_2] - \mu_2))^2 \right) \\ \propto & \exp \left(-\frac{1}{2 \Lambda_{11}^{-1}} \left(z_1 - (\mu_1 - \Lambda_{11}^{-1} \Lambda_{12} (\mathbb{E}[z_2] - \mu_2)) \right)^2 \right) \\ = & \mathcal{N}(z_1 | m_1, \Lambda_{11}^{-1}) \end{aligned} \tag{402}$$

但し m_1 は次のようにおいた

$$m_1 = \mu_1 - \Lambda_{11}^{-1} \Lambda_{12} (\mathbb{E}[z_2] - \mu_2) \tag{403}$$

- 対称性から、 $q_2^*(z_2)$ も次のように求められる

$$\ln q_2^*(z_2) = \mathbb{E}_{z_2}[\ln p(z)] + \text{Const.} \quad (404)$$

$$= \mathbb{E}_{z_2}[\ln \mathcal{N}(z|\boldsymbol{\mu}, \boldsymbol{\Lambda}^{-1})] + \text{Const.} \quad (405)$$

$$q_2^*(z_2) = \mathcal{N}(z_2|m_2, \Lambda_{22}^{-1}) \quad (406)$$

但し m_2 は次のようにおいた

$$m_2 = \mu_2 - \Lambda_{22}^{-1} \Lambda_{21}(\mathbb{E}[z_1] - \mu_1) \quad (407)$$

- これより $q_1^*(z_1), q_2^*(z_2)$ は

$$q_1^*(z_1) = \mathcal{N}(z_1|m_1, \Lambda_{11}^{-1}) \quad (408)$$

$$m_1 = \mu_1 - \Lambda_{11}^{-1} \Lambda_{12}(\mathbb{E}[z_2] - \mu_2) \quad (409)$$

$$q_2^*(z_2) = \mathcal{N}(z_2|m_2, \Lambda_{22}^{-1}) \quad (410)$$

$$m_2 = \mu_2 - \Lambda_{22}^{-1} \Lambda_{21}(\mathbb{E}[z_1] - \mu_1) \quad (411)$$

- $\mathbb{E}[z_2]$ は、 z_2 の $q_2(z_2) = \mathcal{N}(z_2|m_2, \Lambda_{22}^{-1})$ による平均であるから、
 $\mathbb{E}[z_2] = m_2$ である (同様に、 $\mathbb{E}[z_1] = m_1$)
- これより m_1, m_2 を連立させれば

$$m_1 = \mu_1 - \Lambda_{11}^{-1} \Lambda_{12} (m_2 - \mu_2) \quad (412)$$

$$m_2 = \mu_2 - \Lambda_{22}^{-1} \Lambda_{21} (m_1 - \mu_1) \quad (413)$$

であるから、 m_1 について解けば

$$m_1 = \mu_1 - \Lambda_{11}^{-1} \Lambda_{12} (m_2 - \mu_2) \quad (414)$$

$$= \mu_1 - \Lambda_{11}^{-1} \Lambda_{12} (\mu_2 - \Lambda_{22}^{-1} \Lambda_{21} (m_1 - \mu_1) - \mu_2) \quad (415)$$

$$= \mu_1 + \Lambda_{11}^{-1} \Lambda_{12} \Lambda_{22}^{-1} \Lambda_{21} (m_1 - \mu_1) \quad (416)$$

$$= \mu_1 + \Lambda_{11}^{-1} \Lambda_{22}^{-1} \Lambda_{12}^2 (m_1 - \mu_1) \quad (417)$$

$$= (1 - \Lambda_{11}^{-1} \Lambda_{22}^{-1} \Lambda_{12}^2) \mu_1 + \Lambda_{11}^{-1} \Lambda_{22}^{-1} \Lambda_{12}^2 m_1 \quad (418)$$

従って

$$(1 - \Lambda_{11}^{-1} \Lambda_{22}^{-1} \Lambda_{12}^2) \mu_1 = (1 - \Lambda_{11}^{-1} \Lambda_{22}^{-1} \Lambda_{12}^2) m_1 \quad (419)$$

精度行列 Λ が正則であれば

$$\Lambda_{11} \Lambda_{12} - \Lambda_{12} \Lambda_{21} \neq 0 \quad (420)$$

$$\Rightarrow \Lambda_{11} \Lambda_{12} - \Lambda_{12}^2 \neq 0 \quad (421)$$

$$\Rightarrow (\Lambda_{11} \Lambda_{12}) (1 - \Lambda_{11}^{-1} \Lambda_{12}^{-1} \Lambda_{12}^2) \neq 0 \quad (422)$$

$$\Rightarrow (1 - \Lambda_{11}^{-1} \Lambda_{12}^{-1} \Lambda_{12}^2) \neq 0 \quad (423)$$

が成立するから、分布 $p(z)$ が非特異 (精度行列が正則) ならば

$$m_1 = \mu_1 \quad (424)$$

が唯一の解であるほか、対称性から、 m_2 についても以下を得る

$$m_2 = \mu_2 \quad (425)$$

- ゆえに、 $q_1^*(z_1), q_2^*(z_2)$ は

$$q_1^*(z_1) = \mathcal{N}(z_1 | \mu_1, \Lambda_{11}^{-1}) \quad (426)$$

$$q_2^*(z_2) = \mathcal{N}(z_2 | \mu_2, \Lambda_{22}^{-1}) \quad (427)$$

- 注意点 1

- $q_1^*(z_1)$ は、 $q_2^*(z_2)$ を使って計算される $p(z)$ の期待値 $\mathbb{E}[z_2]$ に依存 (逆も成り立つ)
- $q_1^*(z_1)$ と、 $q_2^*(z_2)$ は相互に依存しているため、2つを同時に求めることはできない
- その代わりに、次のように最適化すればよい
- $q_1(z_1), q_2(z_2)$ を適当に初期化したあと、 $q_1^*(z_1), q_2^*(z_2)$ の式を使って、交互に $q_1(z_1)$ と $q_2(z_2)$ を更新していく (収束するまでこれを繰り返す)

- 注意点 2

変分推論

- $q_1(z_1), q_2(z_2)$ の具体的な関数形については、何の仮定も置かなかった
- $q_i^*(z_i)$ がガウス分布だという仮定は置いていないが、 $\text{KL}(q||p)$ を最適化する変分推論によって、結果的にガウス分布が得られた

- $KL(q||p)$ の最適化と $KL(p||q)$ の最適化の比較

- 上記の結果は、 $KL(q||p)$ の最適化 (エビデンス下界 $\mathcal{L}(q)$ の最適化) によって得た
- $KL(q||p)$ ではなく、 $KL(p||q)$ を最適化したらどうなるか?
- 変分推論ではない、もう一つの近似推論の方法である、**EP 法**で使われる考え方
- $q(\mathbf{Z})$ を $p(\mathbf{Z}|\mathbf{X})$ に近づけたいのであれば、 $KL(q||p)$ と $KL(p||q)$ のどちらを最小化してもよいはず
- なぜなら、KL ダイバージェンスは、確率分布間の (擬似的な) 距離を表すため

- $KL(p||q)$ の最適化

- $q(\mathbf{Z})$ が平均場近似によって分解できるとき、 $KL(p||q)$ を最適化したい

- KL ダイバージェンス $\text{KL}(p||q)$ は、次のように書ける

$$\text{KL}(p||q) \equiv \text{KL}(p(\mathbf{Z}|\mathbf{X})||q(\mathbf{Z})) \quad (428)$$

$$= - \int p(\mathbf{Z}|\mathbf{X}) \ln \frac{q(\mathbf{Z})}{p(\mathbf{Z}|\mathbf{X})} d\mathbf{Z} \quad (429)$$

$$= - \int p(\mathbf{Z}|\mathbf{X}) (\ln q(\mathbf{Z}) - \ln p(\mathbf{Z}|\mathbf{X})) d\mathbf{Z} \quad (430)$$

$$= - \int p(\mathbf{Z}|\mathbf{X}) \ln q(\mathbf{Z}) d\mathbf{Z} - \underbrace{\int p(\mathbf{Z}|\mathbf{X}) \ln p(\mathbf{Z}|\mathbf{X}) d\mathbf{Z}}_{q \text{ には依存しない定数項}} \quad (431)$$

$$= - \int p(\mathbf{Z}|\mathbf{X}) \ln \prod_i q_i(\mathbf{Z}_i) d\mathbf{Z} + \text{Const.} \quad (432)$$

$$= - \int p(\mathbf{Z}|\mathbf{X}) \sum_i \ln q_i(\mathbf{Z}_i) d\mathbf{Z} + \text{Const.} \quad (433)$$

$$= - \sum_i \int p(\mathbf{Z}|\mathbf{X}) \ln q_i(\mathbf{Z}_i) d\mathbf{Z} + \text{Const.} \quad (434)$$

定数項は、 $p(\mathbf{Z}|\mathbf{X})$ のエントロピーであり、 q には依存しない

- 各因子 $q_j(\mathbf{Z}_j)$ について $\text{KL}(p||q)$ を最適化したい
- このとき、 $i \neq j$ となる、全ての $q_i(\mathbf{Z}_i)$ は**固定する**
- $q_j(\mathbf{Z}_j)$ に依存する項を取り出せば、次のようになる

$$\begin{aligned} & \sum_i \int p(\mathbf{Z}|\mathbf{X}) \ln q_i(\mathbf{Z}_i) d\mathbf{Z} \\ &= \int p(\mathbf{Z}|\mathbf{X}) \ln q_j(\mathbf{Z}_j) d\mathbf{Z} \end{aligned} \quad (435)$$

$$= \int p(\mathbf{Z}|\mathbf{X}) \ln q_j(\mathbf{Z}_j) d\mathbf{Z}_1 d\mathbf{Z}_2 \cdots d\mathbf{Z}_M \quad (436)$$

$$= \int p(\mathbf{Z}|\mathbf{X}) \ln q_j(\mathbf{Z}_j) d\mathbf{Z}_j \left(\prod_{i \neq j} d\mathbf{Z}_i \right) \quad (437)$$

$$= \int \left(\int p(\mathbf{Z}|\mathbf{X}) \prod_{i \neq j} d\mathbf{Z}_i \right) \ln q_j(\mathbf{Z}_j) d\mathbf{Z}_j \quad (438)$$

- $q_j(\mathbf{Z}_j)$ は確率分布であるから、以下の条件を満たさなければならない

$$\int q_j(\mathbf{Z}_j) d\mathbf{Z}_j = 1 \quad (\text{規格化条件}) \quad (439)$$

$$q_j(\mathbf{Z}_j) \geq 0 \quad (440)$$

- 従って $\text{KL}(p||q)$ を $q_j(\mathbf{Z}_j)$ について最適化するとき、**ラグランジュの未定乗数法**を使って、規格化条件を組み込む必要がある

- $q_j(\mathbf{Z}_j) \geq 0$ という条件は、 $\ln q_i(\mathbf{Z}_i)$ という項が既にあるから、何もしなくても常に満たされる (ラグランジュ関数に、制約条件を改めて取り入れる必要がない)
- 結局、ラグランジュ汎関数 $\mathcal{L}[q_j]$ は、次のようになる

$$\begin{aligned} \mathcal{L}[q_j] = & - \int \left(\int p(\mathbf{Z}|\mathbf{X}) \prod_{i \neq j} d\mathbf{Z}_i \right) \ln q_j(\mathbf{Z}_j) d\mathbf{Z}_j + \\ & \lambda \left(\int q_j(\mathbf{Z}_j) d\mathbf{Z}_j - 1 \right) \end{aligned} \quad (441)$$

- 上記は $q_j(\mathbf{Z}_j)$ についての汎関数となっていることに注意
- 次の公式を使って、 $\mathcal{L}[q_j]$ を変分最適化する

$$\frac{\delta}{\delta y(x)} \int G(y(x), x) dx = \frac{\partial}{\partial y} G(y(x), x) \quad (442)$$

- 従って

$$\begin{aligned}
 & \frac{\delta}{\delta q_j(\mathbf{Z}_j)} \mathcal{L}[q_j] \\
 = & -\frac{\delta}{\delta q_j(\mathbf{Z}_j)} \int \left(\int p(\mathbf{Z}|\mathbf{X}) \prod_{i \neq j} d\mathbf{Z}_i \right) \ln q_j(\mathbf{Z}_j) d\mathbf{Z}_j + \\
 & \frac{\delta}{\delta q_j(\mathbf{Z}_j)} \lambda \left(\int q_j(\mathbf{Z}_j) d\mathbf{Z}_j - 1 \right) \quad (443)
 \end{aligned}$$

$$\begin{aligned}
 = & -\frac{\partial}{\partial q_j} \left(\int p(\mathbf{Z}|\mathbf{X}) \prod_{i \neq j} d\mathbf{Z}_i \right) \ln q_j(\mathbf{Z}_j) + \\
 & \lambda \frac{\partial}{\partial q_j} q_j(\mathbf{Z}_j) \quad (444)
 \end{aligned}$$

$$= - \left(\int p(\mathbf{Z}|\mathbf{X}) \prod_{i \neq j} d\mathbf{Z}_i \right) \frac{1}{q_j(\mathbf{Z}_j)} + \lambda = 0 \quad (445)$$

- これより、未定乗数 λ は

$$\lambda q_j(\mathbf{Z}_j) = \left(\int p(\mathbf{Z}|\mathbf{X}) \prod_{i \neq j} d\mathbf{Z}_i \right) \frac{1}{q_j(\mathbf{Z}_j)} \quad (446)$$

$$\Rightarrow \int \lambda q_j(\mathbf{Z}_j) d\mathbf{Z}_j = \int \underbrace{\left(\int p(\mathbf{Z}|\mathbf{X}) \prod_{i \neq j} d\mathbf{Z}_i \right)}_{=p(\mathbf{Z}_j|\mathbf{X})} d\mathbf{Z}_j \quad (447)$$

$$\Rightarrow \lambda \underbrace{\int q_j(\mathbf{Z}_j) d\mathbf{Z}_j}_{=1} = \underbrace{\int p(\mathbf{Z}_j|\mathbf{X}) d\mathbf{Z}_j}_{=1} \quad (448)$$

$$\Rightarrow \lambda = 1 \quad (449)$$

- 結局、最適解 $q_j^*(\mathbf{Z}_j)$ は次のようになる

$$- \left(\int p(\mathbf{Z}|\mathbf{X}) \prod_{i \neq j} d\mathbf{Z}_i \right) \frac{1}{q_j(\mathbf{Z}_j)} + \lambda = 0 \quad (450)$$

$$\Rightarrow q_j^*(\mathbf{Z}_j) = \underbrace{\int p(\mathbf{Z}|\mathbf{X}) \prod_{i \neq j} d\mathbf{Z}_i}_{=p(\mathbf{Z}_j|\mathbf{X})} \quad (451)$$

$$\Rightarrow q_j^*(\mathbf{Z}_j) = p(\mathbf{Z}_j|\mathbf{X}) \quad (452)$$

- $q_j^*(\mathbf{Z}_j)$ の最適解は、 $p(\mathbf{Z}|\mathbf{X})$ を、 $i \neq j$ である全ての \mathbf{Z}_i について周辺化した分布
- これは閉じた解であり、繰り返しを必要としない

- 最適な因子 $q_1(z_1), q_2(z_2)$ の計算

- 今回は $p(\mathbf{z}|\mathbf{x}) = p(\mathbf{z})$ の場合を考えており、かつ $p(\mathbf{z}) = \mathcal{N}(\mathbf{z}|\boldsymbol{\mu}, \boldsymbol{\Lambda}^{-1})$ であった
- 従って、 $q_1^*(z_1)$ は、 $p(\mathbf{z})$ を z_2 について周辺化すればよいから

$$\begin{aligned} & q_1^*(z_1) \\ &= \int p(\mathbf{z}) d\mathbf{z}_2 \\ &= \int \mathcal{N}(\mathbf{z}|\boldsymbol{\mu}, \boldsymbol{\Lambda}^{-1}) d\mathbf{z}_2 \end{aligned} \tag{453}$$

$$= \frac{1}{2\pi} \frac{1}{|\boldsymbol{\Lambda}|^{-\frac{1}{2}}} \int \exp\left(-\frac{1}{2}(\mathbf{z} - \boldsymbol{\mu})^T \boldsymbol{\Lambda}(\mathbf{z} - \boldsymbol{\mu})\right) d\mathbf{z}_2 \tag{454}$$

- ここで、指数の内側を、積分変数 z_2 に依存する項と、そうでない項に分ける

$$\begin{aligned} & -\frac{1}{2}(\mathbf{z} - \boldsymbol{\mu})^T \boldsymbol{\Lambda}(\mathbf{z} - \boldsymbol{\mu}) \\ = & -\frac{1}{2}(\Lambda_{11}(z_1 - \mu_1)^2 + \\ & 2\Lambda_{12}(z_1 - \mu_1)(z_2 - \mu_2) + \Lambda_{22}(z_2 - \mu_2)^2) \quad (455) \end{aligned}$$

$$\begin{aligned} = & -\frac{1}{2}\Lambda_{11}(z_1 - \mu_1)^2 + \Lambda_{12}(z_1 - \mu_1)\mu_2 + \\ & -\frac{1}{2}(2\Lambda_{12}(z_1 - \mu_1)z_2 + \Lambda_{22}(z_2 - \mu_2)^2) \quad (456) \end{aligned}$$

そして

$$\begin{aligned} & -\frac{1}{2}(2\Lambda_{12}(z_1 - \mu_1)z_2 + \Lambda_{22}(z_2 - \mu_2)^2) \\ = & -\frac{1}{2}(\Lambda_{22}z_2^2 - 2\Lambda_{22}\mu_2z_2 + 2\Lambda_{12}(z_1 - \mu_1)z_2 + \Lambda_{22}\mu_2^2) \quad (457) \end{aligned}$$

$$= -\frac{1}{2} (\Lambda_{22} z_2^2 - 2 (\Lambda_{22} \mu_2 - \Lambda_{12} (z_1 - \mu_1)) z_2 + \Lambda_{22} \mu_2^2) \quad (458)$$

$$= -\frac{1}{2} \left(\Lambda_{22} (z_2 - \Lambda_{22}^{-1} (\Lambda_{22} \mu_2 - \Lambda_{12} (z_1 - \mu_1)))^2 - \Lambda_{22} (\Lambda_{22}^{-1} (\Lambda_{22} \mu_2 - \Lambda_{12} (z_1 - \mu_1)))^2 + \Lambda_{22} \mu_2^2 \right) \quad (459)$$

$$= -\frac{1}{2} \left(\Lambda_{22} (z_2 - \Lambda_{22}^{-1} m)^2 - \Lambda_{22}^{-1} m^2 + \Lambda_{22} \mu_2^2 \right) \quad (460)$$

ゆえ ($m = \Lambda_{22} \mu_2 - \Lambda_{12} (z_1 - \mu_1)$ とおいた)

$$\begin{aligned} & -\frac{1}{2} (\mathbf{z} - \boldsymbol{\mu})^T \boldsymbol{\Lambda} (\mathbf{z} - \boldsymbol{\mu}) \\ = & -\frac{1}{2} \Lambda_{11} (z_1 - \mu_1)^2 + \Lambda_{12} (z_1 - \mu_1) \mu_2 - \frac{1}{2} \Lambda_{22} \mu_2^2 - \\ & \frac{1}{2} \Lambda_{22} (z_2 - \Lambda_{22}^{-1} m)^2 + \frac{1}{2} \Lambda_{22}^{-1} m^2 \end{aligned} \quad (461)$$

- これより、積分変数 z_2 の依存項だけを取り出せたので

$$\begin{aligned}
 & \int \exp \left(-\frac{1}{2} (z - \mu)^T \Lambda (z - \mu) \right) dz_2 \\
 = & \exp \left(-\frac{1}{2} \Lambda_{11} (z_1 - \mu_1)^2 + \Lambda_{12} (z_1 - \mu_1) \mu_2 - \right. \\
 & \quad \left. \frac{1}{2} \Lambda_{22} \mu_2^2 + \frac{1}{2} \Lambda_{22}^{-1} m^2 \right) \\
 & \int \exp \left(-\frac{1}{2} \Lambda_{22} (z_2 - \Lambda_{22}^{-1} m)^2 \right) dz_2 \quad (462)
 \end{aligned}$$

であって、右側の積分は、中身が (正規化されていない) ガウス分布であるから

$$\int \exp \left(-\frac{1}{2} \Lambda_{22} (z_2 - \Lambda_{22}^{-1} m)^2 \right) dz_2$$

$$\begin{aligned}
 &= (2\pi\Lambda_{22}^{-1})^{\frac{1}{2}} \cdot \int \frac{1}{(2\pi\Lambda_{22}^{-1})^{\frac{1}{2}}} \exp\left(-\frac{1}{2}\Lambda_{22}(z_2 - \Lambda_{22}^{-1}m)^2\right) dz_2 \\
 &= (2\pi\Lambda_{22}^{-1})^{\frac{1}{2}} \tag{463}
 \end{aligned}$$

となって、 z_2 を積分により消去できる

- また指数の残りの部分から、 z_1 に依存する項だけを取り出して

$$\begin{aligned}
 &-\frac{1}{2}\Lambda_{11}(z_1 - \mu_1)^2 + \Lambda_{12}(z_1 - \mu_1)\mu_2 - \frac{1}{2}\Lambda_{22}\mu_2^2 + \frac{1}{2}\Lambda_{22}^{-1}m^2 \\
 = &-\frac{1}{2}\Lambda_{11}(z_1 - \mu_1)^2 + \Lambda_{12}(z_1 - \mu_1)\mu_2 - \frac{1}{2}\Lambda_{22}\mu_2^2 + \\
 &\frac{1}{2}\Lambda_{22}^{-1}(\Lambda_{22}\mu_2 - \Lambda_{12}(z_1 - \mu_1))^2 \\
 = &-\frac{1}{2}\Lambda_{11}(z_1 - \mu_1)^2 + \Lambda_{12}(z_1 - \mu_1)\mu_2 - \frac{1}{2}\Lambda_{22}\mu_2^2 + \\
 &\frac{1}{2}\Lambda_{22}\mu_2^2 - \mu_2\Lambda_{12}(z_1 - \mu_1) +
 \end{aligned}$$

$$\frac{1}{2} \Lambda_{22}^{-1} \Lambda_{12}^2 (z_1 - \mu_1)^2 \quad (464)$$

$$= -\frac{1}{2} (\Lambda_{11} - \Lambda_{22}^{-1} \Lambda_{12}^2) (z_1 - \mu_1)^2 \quad (465)$$

$$= -\frac{1}{2} (\Lambda_{11} - \Lambda_{22}^{-1} \Lambda_{12}^2) z_1^2 + \frac{1}{2} (\Lambda_{11} - \Lambda_{22}^{-1} \Lambda_{12}^2) \mu_1 z_1 + \text{Const.} \quad (466)$$

- これより結局、 z_2 による積分は次のようになる

$$\begin{aligned} & \int \exp \left(-\frac{1}{2} (z - \mu)^T \Lambda (z - \mu) \right) dz_2 \\ &= (2\pi \Lambda_{22}^{-1})^{\frac{1}{2}} \exp \left(-\frac{1}{2} (\Lambda_{11} - \Lambda_{22}^{-1} \Lambda_{12}^2) (z_1 - \mu_1)^2 \right) \end{aligned} \quad (467)$$

- 従って、 $q_1^*(z_1)$ は次のようになる

$$= \frac{q_1^*(z_1)}{2\pi |\mathbf{\Lambda}|^{-\frac{1}{2}}} \int \exp\left(-\frac{1}{2}(\mathbf{z} - \boldsymbol{\mu})^T \mathbf{\Lambda}(\mathbf{z} - \boldsymbol{\mu})\right) dz_2 \quad (468)$$

$$= \frac{1}{2\pi} \frac{1}{(\Lambda_{11}\Lambda_{22} - \Lambda_{12}^2)^{-\frac{1}{2}}} (2\pi\Lambda_{22}^{-1})^{\frac{1}{2}} \exp\left(-\frac{1}{2}(\Lambda_{11} - \Lambda_{22}^{-1}\Lambda_{12}^2)(z_1 - \mu_1)^2\right) \quad (469)$$

$$= \frac{1}{(2\pi)^{\frac{1}{2}}} \frac{1}{(\Lambda_{11} - \Lambda_{22}^{-1}\Lambda_{12}^2)^{-\frac{1}{2}}} \exp\left(-\frac{1}{2}(\Lambda_{11} - \Lambda_{22}^{-1}\Lambda_{12}^2)(z_1 - \mu_1)^2\right) \quad (470)$$

$$= \mathcal{N}(z_1 | \mu_1, (\Lambda_{11} - \Lambda_{22}^{-1}\Lambda_{12}^2)^{-1}) \quad (471)$$

- 共分散行列 Σ を、精度行列 Λ を使って次のように定めれば

$$\begin{pmatrix} \Sigma_{11} & \Sigma_{12} \\ \Sigma_{21} & \Sigma_{22} \end{pmatrix}^{-1} = \begin{pmatrix} \Lambda_{11} & \Lambda_{12} \\ \Lambda_{21} & \Lambda_{22} \end{pmatrix} \quad (\Sigma_{12} = \Sigma_{21}) \quad (472)$$

次が成り立つから

$$\begin{pmatrix} \Sigma_{11} & \Sigma_{12} \\ \Sigma_{21} & \Sigma_{22} \end{pmatrix} \begin{pmatrix} \Lambda_{11} & \Lambda_{12} \\ \Lambda_{21} & \Lambda_{22} \end{pmatrix} = \mathbf{I} \quad (473)$$

各成分に注目すれば

$$\begin{cases} \Sigma_{11}\Lambda_{11} + \Sigma_{12}\Lambda_{21} = 1 \\ \Sigma_{11}\Lambda_{12} + \Sigma_{12}\Lambda_{22} = 0 \end{cases} \quad (474)$$

これを Σ_{11} について解けば

$$\Sigma_{11}\Lambda_{11} + (-\Lambda_{22}^{-1}\Lambda_{12}\Sigma_{11})\Lambda_{21} = 1 \quad (475)$$

$$\Rightarrow \Sigma_{11}(\Lambda_{11} - \Lambda_{22}^{-1}\Lambda_{12}\Lambda_{21}) = 1 \quad (476)$$

$$\Rightarrow \Sigma_{11} (\Lambda_{11} - \Lambda_{22}^{-1} \Lambda_{12}^2) = 1 \quad (477)$$

$$\Rightarrow \Sigma_{11} = (\Lambda_{11} - \Lambda_{22}^{-1} \Lambda_{12}^2)^{-1} \quad (478)$$

- これから、 $q_1^*(z_1)$ は次のようにも書ける

$$\begin{aligned} & q_1^*(z_1) \\ = & \frac{1}{(2\pi)^{\frac{1}{2}}} \frac{1}{(\Lambda_{11} - \Lambda_{22}^{-1} \Lambda_{12}^2)^{-\frac{1}{2}}} \\ & \exp\left(-\frac{1}{2} (\Lambda_{11} - \Lambda_{22}^{-1} \Lambda_{12}^2) (z_1 - \mu_1)^2\right) \end{aligned} \quad (479)$$

$$= \frac{1}{(2\pi)^{\frac{1}{2}}} \frac{1}{\Sigma_{11}^{\frac{1}{2}}} \exp\left(-\frac{1}{2} \Sigma_{11}^{-1} (z_1 - \mu_1)^2\right) \quad (480)$$

$$= \mathcal{N}(z_1 | \mu_1, \Sigma_{11}) \quad (481)$$

- $q_2^*(z_2)$ は、対称性から次のようになる

$$q_2^*(z_2) = \mathcal{N}(z_2 | \mu_2, (\Lambda_{22} - \Lambda_{11}^{-1} \Lambda_{12}^2)^{-1}) = \mathcal{N}(z_2 | \mu_2, \Sigma_{22}) \quad (482)$$

- 2つの解の比較

- $\text{KL}(q||p)$ の最小化によって次の解を得た

$$q_1^*(z_1) = \mathcal{N}(z_1 | \mu_1, \Lambda_{11}^{-1}) \quad (483)$$

$$q_2^*(z_2) = \mathcal{N}(z_2 | \mu_2, \Lambda_{22}^{-1}) \quad (484)$$

- $\text{KL}(p||q)$ の最小化では、次の解を得た

$$q_1^*(z_1) = \mathcal{N}(z_1|\mu_1, \Sigma_{11}) \quad (485)$$

$$\Sigma_{11} = (\Lambda_{11} - \Lambda_{22}^{-1}\Lambda_{12}^2)^{-1} \quad (486)$$

$$q_2^*(z_2) = \mathcal{N}(z_2|\mu_2, \Sigma_{22}) \quad (487)$$

$$\Sigma_{22} = (\Lambda_{22} - \Lambda_{11}^{-1}\Lambda_{12}^2)^{-1} \quad (488)$$

- $p(z) = \mathcal{N}(z|\mu, \Sigma)$ の平均は $\mu = [\mu_1, \mu_2]^T$ であったので、いずれの場合も、平均は正しく捉えている
- しかし、両者の間では、分散が異なっている
- また、変数 z_1 と z_2 の間の相関は消えてなくなっている
- 両者の違いを、次の図 13 に示す
- 緑色の線が真の分布 $p(z)$ を表す
- 左側の赤線は、 $\text{KL}(q||p)$ の最小化によって得られた分布 $q(z)$
- 右側の赤線は、 $\text{KL}(p||q)$ の最小化によって得られた分布 $q(z)$

Figure 10.2 Comparison of the two alternative forms for the Kullback-Leibler divergence. The green contours corresponding to 1, 2, and 3 standard deviations for a correlated Gaussian distribution $p(\mathbf{z})$ over two variables z_1 and z_2 , and the red contours represent the corresponding levels for an approximating distribution $q(\mathbf{z})$ over the same variables given by the product of two independent univariate Gaussian distributions whose parameters are obtained by minimization of (a) the Kullback-Leibler divergence $\text{KL}(q\|p)$, and (b) the reverse Kullback-Leibler divergence $\text{KL}(p\|q)$.

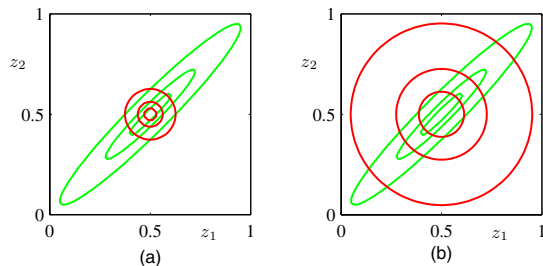


図 13: KL ダイバージェンスの 2 つの形の比較

● 2つの解の比較

- $KL(q||p)$ の最小化で得られる $q(z)$ は、分散が小さくなる方向に制御されていて、それと直交する方向の分散は、大きく過小評価されている
- 分解による近似では、一般に事後分布 $p(Z|X)$ をコンパクトに近似しすぎる
- $KL(p||q)$ の最小化で得られる $q(z)$ は、分散が大きくなる方向に制御されていて、それと直交する方向の分散は、過大評価されている
- 非常に低い確率しか持たないはずの領域にも、多くの確率質量が割り当てられている
- 変分推論では、計算コストの観点から $KL(q||p)$ の方を用いる
- $KL(q||p)$ の計算には、 q に関する期待値の評価が含まれるので、 q を単純な形に制限することで、必要な期待値を単純化できる
- $KL(p||q)$ の計算には、真の事後分布 p に関する期待値の計算が必要である

- 違いが生じる理由

- KL ダイバージェンス $KL(q||p)$ は次のようであった

$$KL(q||p) = - \int q(\mathbf{Z}) \ln \frac{p(\mathbf{Z}|\mathbf{X})}{q(\mathbf{Z})} d\mathbf{Z} \quad (489)$$

- $KL(q||p)$ が大きくなる主要因は、 $p(\mathbf{Z})$ がほとんど 0 で、 $q(\mathbf{Z})$ はそうではない領域
- 従って、 $KL(q||p)$ を最小化すると、 $q(\mathbf{Z})$ は、 $p(\mathbf{Z})$ が小さい領域を避けるようになる
- また、 $KL(p||q)$ は次のようであった

$$KL(p||q) = - \int p(\mathbf{Z}|\mathbf{X}) \ln \frac{q(\mathbf{Z})}{p(\mathbf{Z}|\mathbf{X})} d\mathbf{Z} \quad (490)$$

- $\text{KL}(p||q)$ が大きくなる主要因は、 $q(\mathbf{Z})$ がほとんど 0 で、 $p(\mathbf{Z})$ はそうではない領域
- 従って、 $\text{KL}(p||q)$ を最小化すると、 $q(\mathbf{Z})$ は、 $p(\mathbf{Z})$ が 0 でない領域にも、必ず確率を持たせるようになる
- 別の分布について、この両者の振舞いの違いを観察してみよう

- 多峰性のある分布を、単峰の分布で近似する場合
 - $KL(q||p)$ を最小化する変分近似では、多数ある峰のうちの 1 つを再現
 - $KL(p||q)$ を最小化する変分近似では、全ての峰を平均したような分布が得られる
- 多峰性のある分布を平均してしまうと、予測性能の悪化をもたらす
- これらの比較を次の図 14 に示す

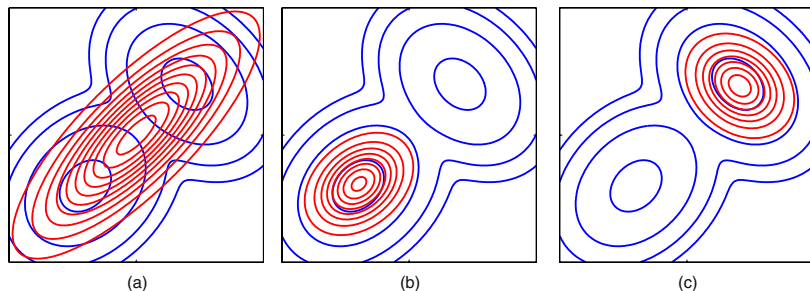


Figure 10.3 Another comparison of the two alternative forms for the Kullback-Leibler divergence. (a) The blue contours show a bimodal distribution $p(\mathbf{Z})$ given by a mixture of two Gaussians, and the red contours correspond to the single Gaussian distribution $q(\mathbf{Z})$ that best approximates $p(\mathbf{Z})$ in the sense of minimizing the Kullback-Leibler divergence $KL(p||q)$. (b) As in (a) but now the red contours correspond to a Gaussian distribution $q(\mathbf{Z})$ found by numerical minimization of the Kullback-Leibler divergence $KL(q||p)$. (c) As in (b) but showing a different local minimum of the Kullback-Leibler divergence.

図 14: KL ダイバージェンスの 2 つの形の別の比較

- ここまでの話の流れ

- 1 分解 $q(\mathbf{Z}) = \prod_i q_i(\mathbf{Z}_i)$ を使ったエビデンス下界 $\mathcal{L}(q)$ の最適化は、 $\text{KL}(q(\mathbf{Z})||p(\mathbf{Z}|\mathbf{X}))$ の最小化と等価である
- 2 $\text{KL}(q||p)$ と、 $\text{KL}(p||q)$ を最小化する変分近似を、2 変数のガウス分布を例として試した
- 3 $\text{KL}(q||p)$ の最小化を使って求めた $q(\mathbf{Z})$ は、事後分布 $p(\mathbf{Z}|\mathbf{X})$ を **コンパクトに近似**する傾向にあった
- 4 $\text{KL}(p||q)$ の最小化によって求めた $q(\mathbf{Z})$ は、事後分布 $p(\mathbf{Z}|\mathbf{X})$ を **大きく捉えて近似**する傾向にあった

- これからの話の流れ

- 変分推論の具体的な例について更に見ていく

変分推論

- 離散潜在変数のモデル (二値スパース符号化モデル) に変分推論を適用する
- 連続潜在変数の場合は、簡単な確率モデルを使って、変分推論を試してみる

5 近似推論法

- 変分推論
- 座標降下法
- MAP 推定および最尤推定と変分推論
- 変分近似による弊害
- 離散潜在変数の変分推論
- 連続潜在変数の変分推論
- 正則化との関わり
- 雑多な話題

離散潜在変数の変分推論

- 離散潜在変数をもつ変分推論の概要

- ここでは単純な場合を扱う
- データ \mathbf{x}_i に対する潜在変数 z_i は、各要素が二値であるとする
- データ \mathbf{x} は D 次元、潜在変数 z は K 次元とする
- 分布 $q(z_i|\mathbf{x}_i)$ は、平均場近似によって、次のように分解できるとする

$$q(z_i|\mathbf{x}_i) = \prod_k q(z_{ik}|\mathbf{x}_i) \quad (491)$$

- 潜在変数は二値であるから、 $q(z_{ik}=1|\mathbf{x}_i) = \widehat{z_{ik}}$ と書くことにする
- このとき、 $q(z_{ik}|\mathbf{x}_i)$ は次のようになる

$$q(z_{ik}|\mathbf{x}_i) = \widehat{z_{ik}}^{z_{ik}} (1 - \widehat{z_{ik}})^{(1-z_{ik})} \quad (492)$$

- 各パラメータ $\widehat{z_{ik}}$ について、下界 $\mathcal{L}(q)$ を順番に最適化することを、 $\mathcal{L}(q)$ が収束するまで繰り返し行う

離散潜在変数の変分推論

- 即ち、以下の不動点方程式を、各 $\widehat{z_{ik}}$ について繰り返し解く

$$\frac{\partial}{\partial \widehat{z_{ik}}} \mathcal{L}(q) = 0 \quad (493)$$

- 離散潜在変数の場合は、単なる標準的な最適化問題を解くことになる
- 二値スパース符号化モデル
 - 二値スパース符号化モデルでのデータ生成過程は、次のようになる
 - 各データ x_i には、潜在変数 $z_i \in \{0, 1\}^K$ が対応する
 - z_i に対し、重み W を用いて線形変換を施し、更にガウスノイズを足し合わせることで、データ x_i が生成される
 - 潜在変数 z_i は、 K 次元ベクトルであり、その各要素は 0 または 1 である

離散潜在変数の変分推論

- 従って、確率分布は次のようになる

$$p(z_{ik} = 1) = \sigma(b_{ik}) \quad (494)$$

$$p(\mathbf{x}_i | \mathbf{z}_i) = \mathcal{N}(\mathbf{x}_i | \mathbf{W} \mathbf{z}_i, \boldsymbol{\beta}^{-1}) \quad (495)$$

- $\sigma(\cdot)$ はシグモイド関数、 $\mathbf{b}_i = [b_{i1}, \dots, b_{iK}]^T$ は学習可能な**バイアス**、 \mathbf{W} は学習可能な**重み行列**、 $\boldsymbol{\beta}$ は学習可能な**対角精度行列**である
- $\boldsymbol{\beta} = \text{diag}(\beta_1, \beta_2, \dots, \beta_D)$ と書ける
- $p(\mathbf{z}_i)$ は次のように計算できる

$$p(\mathbf{z}_i) = \prod_k p(z_{ik}) \quad (496)$$

$$= \prod_k (\sigma(b_{ik}))^{z_{ik}} (1 - \sigma(b_{ik}))^{1-z_{ik}} \quad (497)$$

$$= \prod_k (\sigma(b_{ik}))^{z_{ik}} (\sigma(-b_{ik}))^{1-z_{ik}} \quad (498)$$

離散潜在変数の変分推論

- 事後分布 $p(z_i|x_i)$ は複雑であり、表現することも、計算することもできない
- 従って、最尤推定の手法 (EM アルゴリズム) により学習することはできない
- 例えばバイアス b_{ik} に関する微分を考えると

$$\begin{aligned} & \frac{\partial}{\partial b_{ik}} \ln p(z_i|x_i) \\ = & \frac{\partial}{\partial b_{ik}} \ln \frac{p(x_i, z_i)}{p(x_i)} \end{aligned} \quad (499)$$

$$= \frac{\partial}{\partial b_{ik}} (\ln p(x_i, z_i) - \ln p(x_i)) \quad (500)$$

であって、第 1 項の $p(x_i, z_i)$ は容易に計算できるが、第 2 項について考えると

$$\frac{\partial}{\partial b_{ik}} \ln p(x_i)$$

$$= \frac{1}{p(\mathbf{x}_i)} \frac{\partial}{\partial b_{ik}} p(\mathbf{x}_i) \quad (501)$$

$$= \frac{1}{p(\mathbf{x}_i)} \frac{\partial}{\partial b_{ik}} \sum_{\mathbf{z}_i} p(\mathbf{x}_i, \mathbf{z}_i) \quad (502)$$

$$= \frac{1}{p(\mathbf{x}_i)} \frac{\partial}{\partial b_{ik}} \sum_{\mathbf{z}_i} p(\mathbf{z}_i) p(\mathbf{x}_i | \mathbf{z}_i) \quad (503)$$

$$= \frac{1}{p(\mathbf{x}_i)} \sum_{\mathbf{z}_i} p(\mathbf{x}_i | \mathbf{z}_i) \frac{\partial}{\partial b_{ik}} p(\mathbf{z}_i) \quad (504)$$

$$= \sum_{\mathbf{z}_i} \frac{1}{p(\mathbf{x}_i)} \frac{p(\mathbf{x}_i, \mathbf{z}_i)}{p(\mathbf{z}_i)} \frac{\partial}{\partial b_{ik}} p(\mathbf{z}_i) \quad (505)$$

$$= \sum_{\mathbf{z}_i} p(\mathbf{z}_i | \mathbf{x}_i) \frac{1}{p(\mathbf{z}_i)} \frac{\partial}{\partial b_{ik}} p(\mathbf{z}_i) \quad (506)$$

$$= \sum_{\mathbf{z}_i} p(\mathbf{z}_i | \mathbf{x}_i) \frac{\partial}{\partial b_{ik}} \ln p(\mathbf{z}_i) \quad (507)$$

$$= \mathbb{E}_{\mathbf{z}_i \sim p(\mathbf{z}_i | \mathbf{x}_i)} \left[\frac{\partial}{\partial b_{ik}} \ln p(\mathbf{z}_i) \right] \quad (508)$$

何とか $p(\mathbf{z}_i)$ と $p(\mathbf{x}_i | \mathbf{z}_i)$ を使って計算できないかと試行錯誤したが、結局、 $p(\mathbf{z}_i | \mathbf{x}_i)$ に関する期待値の計算が必要になった

- $p(\mathbf{x}_i, \mathbf{z}_i)$ と、 $p(\mathbf{z}_i | \mathbf{x}_i)$ のグラフ構造を次の図 15 に示す
- 図では、 $\mathbf{h} = \mathbf{z}$ 、 $\mathbf{v} = \mathbf{x}$ のように読み替える必要がある

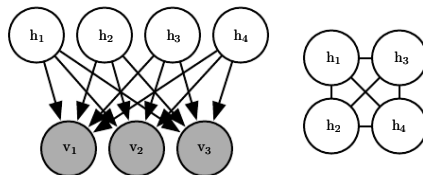


Figure 19.2: The graph structure of a binary sparse coding model with four hidden units. *(Left)* The graph structure of $p(\mathbf{h}, \mathbf{v})$. Note that the edges are directed, and that every two hidden units are co-parents of every visible unit. *(Right)* The graph structure of $p(\mathbf{h} \mid \mathbf{v})$. In order to account for the active paths between co-parents, the posterior distribution needs an edge between all of the hidden units.

図 15: 二値スパース符号化モデルのグラフ構造 (4 つの潜在変数をもつ場合)

- 変分推論による二値スパース符号化モデルの学習

- 最尤推定の手法を諦める代わりに、変分推論を使うことで、困難を解決できる
- 分布 $p(\mathbf{z}_i|\mathbf{x}_i)$ を $q(\mathbf{z}_i)$ で表現し、更に平均場近似を行う
- 但し $\mathbf{z}_i = [z_{i1}, z_{i2}, \dots, z_{iK}]^T$ とする

$$q(\mathbf{z}_i|\mathbf{x}_i) = \prod_k q(z_{ik}|\mathbf{x}_i) \quad (509)$$

- 潜在変数の各要素は二値であるから、各 $q(z_{ik}|\mathbf{x}_i)$ はベルヌーイ分布とすればよい
- 即ち、 $q(z_{ik} = 1|\mathbf{x}_i) = \widehat{z_{ik}}$ とする
- $\widehat{z_{ik}} \neq 0, 1$ という制約を課すことで、 $\ln \widehat{z_{ik}}$ を計算できる
- このようにすれば、潜在変数 $z_{ik}, z_{il} (k \neq l)$ 間の相関を断ち切ることができる
- 先程の図 15 の右側から、潜在変数間の辺を、全て消し去ることになる

離散潜在変数の変分推論

- これで、平均場近似により、因数分解可能な q を表現することができる

$$q(\mathbf{z}_i | \mathbf{x}_i) = \prod_k q(z_{ik} | \mathbf{x}_i) \quad (510)$$

$$= \prod_k \widehat{z}_{ik}^{z_{ik}} (1 - \widehat{z}_{ik})^{1-z_{ik}} \quad (511)$$

$$q(\mathbf{Z} | \mathbf{X}) = \prod_i q(\mathbf{z}_i | \mathbf{x}_i) \quad (512)$$

$$= \prod_i \prod_k \widehat{z}_{ik}^{z_{ik}} (1 - \widehat{z}_{ik})^{1-z_{ik}} \quad (513)$$

- ソフトウェア上では、丸め誤差などによって \widehat{z}_{ik} が 0 や 1 になり、計算を続行できなくなるかもしれない
- これを回避するためには、パラメータ \tilde{z}_i を使って二値スパース符号化モデルを学習させる

離散潜在変数の変分推論

- そして、 $\hat{z}_i = \sigma(\tilde{z}_i)$ の関係によって、 $\hat{z}_i = [\hat{z}_{i1}, \dots, \hat{z}_{iK}]^T$ を得るようにする
- $\ln \hat{z}_{ik} = \ln \sigma(\tilde{z}_{ik}) = -\zeta(-\tilde{z}_{ik})$ によって、コンピュータ上で安全に $\ln \hat{z}_{ik}$ を計算できる ($\zeta(\cdot)$ はソフトプラス関数)
- 変分推論のために、まずはエビデンス下界 $\mathcal{L}(q)$ を計算する

$$\begin{aligned} & \mathcal{L}(q) \\ = & \sum_{\mathbf{Z}} q(\mathbf{Z}|\mathbf{X}) \ln \frac{p(\mathbf{X}, \mathbf{Z})}{q(\mathbf{Z}|\mathbf{X})} \end{aligned} \quad (514)$$

$$= \sum_{\mathbf{Z}} q(\mathbf{Z}|\mathbf{X}) \ln \frac{p(\mathbf{Z})p(\mathbf{X}|\mathbf{Z})}{q(\mathbf{Z}|\mathbf{X})} \quad (515)$$

$$\begin{aligned} = & \sum_{\mathbf{Z}} q(\mathbf{Z}|\mathbf{X}) (\ln p(\mathbf{Z}) + \ln p(\mathbf{X}|\mathbf{Z}) - \ln q(\mathbf{Z}|\mathbf{X})) \quad (516) \\ = & \sum_{\mathbf{Z}} q(\mathbf{Z}|\mathbf{X}) (\ln p(\mathbf{Z}) - \ln q(\mathbf{Z}|\mathbf{X})) + \end{aligned}$$

$$\sum_{\mathbf{Z}} q(\mathbf{Z}|\mathbf{X}) \ln p(\mathbf{X}|\mathbf{Z}) \quad (517)$$

- $q(\mathbf{Z}|\mathbf{X})$ についての期待値を取っていることに注意する
- 全ての \mathbf{Z} についての和を取ればよいので、第 1 項は次のようになる

$$\begin{aligned} & \sum_{\mathbf{Z}} q(\mathbf{Z}|\mathbf{X}) (\ln p(\mathbf{Z}) - \ln q(\mathbf{Z}|\mathbf{X})) \\ = & \sum_i \sum_k \sum_{z_{ik}} q(z_{ik}|\mathbf{x}_i) (\ln p(z_{ik}) - \ln q(z_{ik}|\mathbf{x}_i)) \end{aligned} \quad (518)$$

$$\begin{aligned} = & \sum_i \sum_k q(z_{ik} = 1|\mathbf{x}_i) (\ln p(z_{ik} = 1) - \ln q(z_{ik} = 1|\mathbf{x}_i)) + \\ & q(z_{ik} = 0|\mathbf{x}_i) (\ln p(z_{ik} = 0) - \ln q(z_{ik} = 0|\mathbf{x}_i)) \end{aligned} \quad (519)$$

$$\begin{aligned} = & \sum_i \sum_k \{ \widehat{z_{ik}} (\ln \sigma(b_{ik}) - \ln \widehat{z_{ik}}) + \\ & (1 - \widehat{z_{ik}}) (\ln (1 - \sigma(b_{ik})) - \ln (1 - \widehat{z_{ik}})) \} \end{aligned} \quad (520)$$

$$= \sum_i \sum_k \{ \widehat{z}_{ik} (\ln \sigma(b_{ik}) - \ln \widehat{z}_{ik}) + (1 - \widehat{z}_{ik}) (\ln \sigma(-b_{ik}) - \ln (1 - \widehat{z}_{ik})) \} \quad (521)$$

- また第 2 項は、次のようになる

$$\begin{aligned} & \sum_{\mathbf{Z}} q(\mathbf{Z}|\mathbf{X}) \ln p(\mathbf{X}|\mathbf{Z}) \\ = & \sum_i \sum_{\mathbf{z}_i} q(\mathbf{z}_i|\mathbf{x}_i) \ln p(\mathbf{x}_i|\mathbf{z}_i) \end{aligned} \quad (522)$$

$$= \sum_i \sum_{\mathbf{z}_i} q(\mathbf{z}_i|\mathbf{x}_i) \ln \mathcal{N}(\mathbf{x}_i|\mathbf{W}\mathbf{z}_i, \boldsymbol{\beta}^{-1}) \quad (523)$$

離散潜在変数の変分推論

- 但し、 $\ln \mathcal{N}(\mathbf{x}_i | \mathbf{W} \mathbf{z}_i, \boldsymbol{\beta}^{-1})$ は次のように分解できる

$$\begin{aligned} & \ln \mathcal{N}(\mathbf{x}_i | \mathbf{W} \mathbf{z}_i, \boldsymbol{\beta}^{-1}) \\ = & \ln \left(\frac{1}{(2\pi)^{\frac{D}{2}}} \frac{1}{|\boldsymbol{\beta}^{-1}|^{\frac{1}{2}}} \right. \\ & \left. \exp \left(-\frac{1}{2} (\mathbf{x}_i - \mathbf{W} \mathbf{z}_i)^T \boldsymbol{\beta} (\mathbf{x}_i - \mathbf{W} \mathbf{z}_i) \right) \right) \end{aligned} \quad (524)$$

$$\begin{aligned} = & -\frac{D}{2} \ln 2\pi - \frac{1}{2} \ln |\boldsymbol{\beta}^{-1}| - \\ & \frac{1}{2} (\mathbf{x}_i - \mathbf{W} \mathbf{z}_i)^T \boldsymbol{\beta} (\mathbf{x}_i - \mathbf{W} \mathbf{z}_i) \end{aligned} \quad (525)$$

離散潜在変数の変分推論

- ここで、 $\beta = \text{diag}(\beta_1, \beta_2, \dots, \beta_D)$ であるので

$$\begin{aligned}\ln |\beta^{-1}| &= \ln |\beta|^{-1} = -\ln |\beta| \\ &= -\ln \prod_{j=1}^D \beta_j = -\sum_{j=1}^D \ln \beta_j\end{aligned}\quad (526)$$

となるほか

$$(\mathbf{x}_i - \mathbf{W} \mathbf{z}_i)^T \beta (\mathbf{x}_i - \mathbf{W} \mathbf{z}_i) = \sum_{j=1}^D \beta_j (x_{ij} - \mathbf{W}_{j:} \mathbf{z}_i)^2 \quad (527)$$

であるから

$$\begin{aligned}&\ln \mathcal{N}(\mathbf{x}_i | \mathbf{W} \mathbf{z}_i, \beta^{-1}) \\ &= -\frac{1}{2} \sum_j \ln 2\pi + \frac{1}{2} \sum_j \ln \beta_j - \frac{1}{2} \sum_j \beta_j (x_{ij} - \mathbf{W}_{j:} \mathbf{z}_i)^2\end{aligned}\quad (528)$$

$$= \frac{1}{2} \sum_j \left(-\ln 2\pi + \ln \beta_j - \beta_j (x_{ij} - \mathbf{W}_{j:} \mathbf{z}_i)^2 \right) \quad (529)$$

$$= \frac{1}{2} \sum_j \left(\ln \frac{\beta_j}{2\pi} - \beta_j (x_{ij} - \mathbf{W}_{j:} \mathbf{z}_i)^2 \right) \quad (530)$$

$$= \frac{1}{2} \sum_j \left(\ln \frac{\beta_j}{2\pi} - \beta_j \left(x_{ij} - \sum_k W_{jk} z_{ik} \right)^2 \right) \quad (531)$$

- 上式の対数を外すと、確率分布の積になっていることが分かる

$$\begin{aligned} & \mathcal{N}(\mathbf{x}_i | \mathbf{W} \mathbf{z}_i, \boldsymbol{\beta}^{-1}) \\ &= \exp \left(\sum_j \left(\frac{1}{2} \ln \frac{\beta_j}{2\pi} - \frac{1}{2} \beta_j (x_{ij} - \mathbf{W}_{j:} \mathbf{z}_i)^2 \right) \right) \end{aligned} \quad (532)$$

$$= \prod_j \exp \left(\frac{1}{2} \ln \frac{\beta_j}{2\pi} - \frac{1}{2} \beta_j (x_{ij} - \mathbf{W}_{j:} \mathbf{z}_i)^2 \right) \quad (533)$$

$$= \prod_j \exp\left(\frac{1}{2} \ln \frac{\beta_j}{2\pi}\right) \exp\left(-\frac{1}{2} \beta_j (x_{ij} - \mathbf{W}_{j:} \mathbf{z}_i)^2\right) \quad (534)$$

$$= \prod_j \left(\frac{\beta_j}{2\pi}\right)^{\frac{1}{2}} \exp\left(-\frac{1}{2} \beta_j (x_{ij} - \mathbf{W}_{j:} \mathbf{z}_i)^2\right) \quad (535)$$

$$= \prod_j \left(\frac{1}{(2\pi\beta_j^{-1})^{\frac{1}{2}}} \exp\left(-\frac{1}{2} \beta_j (x_{ij} - \mathbf{W}_{j:} \mathbf{z}_i)^2\right)\right) \quad (536)$$

$$= \prod_j \mathcal{N}(x_{ij} | \mathbf{W}_{j:} \mathbf{z}_i, \beta_j^{-1}) \quad (537)$$

- 精度行列 β は対角行列であるから、 x_i の各成分は無相関である
- 従って、 $\mathcal{N}(x_i | \mathbf{W} \mathbf{z}_i, \beta^{-1})$ は、各成分 x_{ij} についての確率 $\mathcal{N}(x_{ij} | \mathbf{W}_{j:} \mathbf{z}_i, \beta_j^{-1})$ の積として、記述できる

離散潜在変数の変分推論

- さて、第2項は

$$\begin{aligned} & \sum_{\mathbf{Z}} q(\mathbf{Z}|\mathbf{X}) \ln p(\mathbf{X}|\mathbf{Z}) \\ &= \sum_i \sum_{\mathbf{z}_i} q(\mathbf{z}_i|\mathbf{x}_i) \ln \mathcal{N}(\mathbf{x}_i|\mathbf{W}\mathbf{z}_i, \beta^{-1}) \end{aligned} \quad (538)$$

$$= \sum_i \mathbb{E}_{\mathbf{z}_i \sim q(\mathbf{z}_i|\mathbf{x}_i)} [\ln \mathcal{N}(\mathbf{x}_i|\mathbf{W}\mathbf{z}_i, \beta^{-1})] \quad (539)$$

$$= \sum_i \mathbb{E}_{\mathbf{z}_i} \left[\frac{1}{2} \sum_j \left(\ln \frac{\beta_j}{2\pi} - \beta_j \left(x_{ij} - \sum_k W_{jk} z_{ik} \right)^2 \right) \right] \quad (540)$$

$$= \frac{1}{2} \sum_i \sum_j \mathbb{E}_{\mathbf{z}_i} \left[\ln \frac{\beta_j}{2\pi} - \beta_j \left(x_{ij} - \sum_k W_{jk} z_{ik} \right)^2 \right] \quad (541)$$

$$= \frac{1}{2} \sum_i \sum_j \left(\ln \frac{\beta_j}{2\pi} - \beta_j \mathbb{E}_{\mathbf{z}_i} \left[\left(x_{ij} - \sum_k W_{jk} z_{ik} \right)^2 \right] \right) \quad (542)$$

ここで

$$\begin{aligned} & \mathbb{E}_{\mathbf{z}_i} \left[\left(x_{ij} - \sum_k W_{jk} z_{ik} \right)^2 \right] \\ = & \mathbb{E}_{\mathbf{z}_i} \left[x_{ij}^2 - 2x_{ij} \sum_k W_{jk} z_{ik} + \left(\sum_k W_{jk} z_{ik} \right)^2 \right] \end{aligned} \quad (543)$$

$$= x_{ij}^2 - 2x_{ij} \mathbb{E}_{\mathbf{z}_i} \left[\sum_k W_{jk} z_{ik} \right] + \mathbb{E}_{\mathbf{z}_i} \left[\left(\sum_k W_{jk} z_{ik} \right)^2 \right] \quad (544)$$

各項を順番に計算すると

$$\mathbb{E}_{\mathbf{z}_i} \left[\sum_k W_{jk} z_{ik} \right]$$

$$= \sum_k W_{jk} \mathbb{E}_{\mathbf{z}_i} [z_{ik}] \quad (545)$$

$$= \sum_k W_{jk} \mathbb{E}_{z_{ik}} [z_{ik}] \quad (546)$$

$$= \sum_k W_{jk} (q(z_{ik} = 1 | \mathbf{x}_i) \cdot 1 + q(z_{ik} = 0 | \mathbf{x}_i) \cdot 0) \quad (547)$$

$$= \sum_k W_{jk} \widehat{z_{ik}} \quad (548)$$

また

$$\begin{aligned} & \mathbb{E}_{\mathbf{z}_i} \left[\left(\sum_k W_{jk} z_{ik} \right)^2 \right] \\ &= \mathbb{E}_{\mathbf{z}_i} \left[\sum_k \sum_l W_{jk} W_{jl} z_{ik} z_{il} \right] \end{aligned} \quad (549)$$

$$= \mathbb{E}_{\mathbf{z}_i} \left[\sum_k \left(W_{jk}^2 z_{ik}^2 + \sum_{l \neq k} (W_{jk} W_{jl} z_{ik} z_{il}) \right) \right] \quad (550)$$

$$= \sum_k \left(W_{jk}^2 \mathbb{E}_{\mathbf{z}_i} [z_{ik}^2] + \sum_{l \neq k} W_{jk} W_{jl} \mathbb{E}_{\mathbf{z}_i} [z_{ik} z_{il}] \right) \quad (551)$$

$$\begin{aligned} & \mathbb{E}_{\mathbf{z}_i} [z_{ik}^2] \\ &= \mathbb{E}_{z_{ik}} [z_{ik}^2] \end{aligned} \quad (552)$$

$$= q(z_{ik} = 1 | \mathbf{x}_i) \cdot 1^2 + q(z_{ik} = 0 | \mathbf{x}_i) \cdot 0^2 \quad (553)$$

$$= \widehat{z_{ik}} \quad (554)$$

$$\begin{aligned} & \mathbb{E}_{\mathbf{z}_i} [z_{ik} z_{il}] \\ &= \mathbb{E}_{\mathbf{z}_i} [z_{ik}] \mathbb{E}_{\mathbf{z}_i} [z_{il}] \quad (\because z_{ik} \text{ と } z_{il} \text{ は独立であるため}) \end{aligned} \quad (555)$$

$$= \mathbb{E}_{z_{ik}} [z_{ik}] \mathbb{E}_{z_{il}} [z_{il}] \quad (556)$$

$$= \widehat{z_{ik}} \widehat{z_{il}} \quad (557)$$

従って

$$\begin{aligned} & \mathbb{E}_{\mathbf{z}_i} \left[\left(x_{ij} - \sum_k W_{jk} z_{ik} \right)^2 \right] \\ &= x_{ij}^2 - 2x_{ij} \mathbb{E}_{\mathbf{z}_i} \left[\sum_k W_{jk} z_{ik} \right] + \mathbb{E}_{\mathbf{z}_i} \left[\left(\sum_k W_{jk} z_{ik} \right)^2 \right] \end{aligned} \quad (558)$$

$$\begin{aligned} &= x_{ij}^2 - 2x_{ij} \sum_k W_{jk} \widehat{z_{ik}} + \\ & \quad \sum_k \left(W_{jk}^2 \widehat{z_{ik}} + \sum_{l \neq k} W_{jk} W_{jl} \widehat{z_{ik}} \widehat{z_{il}} \right) \end{aligned} \quad (559)$$

これより

$$\begin{aligned} & \sum_{\mathbf{Z}} q(\mathbf{Z}|\mathbf{X}) \ln p(\mathbf{X}|\mathbf{Z}) \\ &= \frac{1}{2} \sum_i \sum_j \left(\ln \frac{\beta_j}{2\pi} - \beta_j \mathbb{E}_{\mathbf{z}_i} \left[\left(x_{ij} - \sum_k W_{jk} z_{ik} \right)^2 \right] \right) \end{aligned} \quad (560)$$

$$\begin{aligned} &= \frac{1}{2} \sum_i \sum_j \left(\ln \frac{\beta_j}{2\pi} - \beta_j \left(x_{ij}^2 - 2x_{ij} \sum_k W_{jk} \widehat{z}_{ik} + \right. \right. \\ &\quad \left. \left. \sum_k \left(W_{jk}^2 \widehat{z}_{ik} + \sum_{l \neq k} W_{jk} W_{jl} \widehat{z}_{ik} \widehat{z}_{il} \right) \right) \right) \end{aligned} \quad (561)$$

- よって、エビデンス下界 $\mathcal{L}(q)$ は

$$\begin{aligned} & \mathcal{L}(q) \\ = & \sum_{\mathbf{Z}} q(\mathbf{Z}|\mathbf{X}) (\ln p(\mathbf{Z}) - \ln q(\mathbf{Z}|\mathbf{X})) + \\ & \sum_{\mathbf{Z}} q(\mathbf{Z}|\mathbf{X}) \ln p(\mathbf{X}|\mathbf{Z}) \end{aligned} \quad (562)$$

$$\begin{aligned} = & \sum_i \sum_k \{ \widehat{z}_{ik} (\ln \sigma(b_{ik}) - \ln \widehat{z}_{ik}) + \\ & (1 - \widehat{z}_{ik}) (\ln \sigma(-b_{ik}) - \ln (1 - \widehat{z}_{ik})) \} + \\ & \frac{1}{2} \sum_i \sum_j \left(\ln \frac{\beta_j}{2\pi} - \beta_j \left(x_{ij}^2 - 2x_{ij} \sum_k W_{jk} \widehat{z}_{ik} + \right. \right. \\ & \left. \left. \sum_k \left(W_{jk}^2 \widehat{z}_{ik} + \sum_{l \neq k} W_{jk} W_{jl} \widehat{z}_{ik} \widehat{z}_{il} \right) \right) \right) \end{aligned} \quad (563)$$

離散潜在変数の変分推論

- この式にはあまり美的魅力がない (Somewhat unappealing aesthetically) が、事後分布よりは計算しやすい
- 事後分布 $\ln p(\mathbf{Z}|\mathbf{X})$ を最大化する代わりに、上記の下界 $\mathcal{L}(q)$ を q について最大化することができる
- データ x_i に対するパラメータ $\{\widehat{z}_{i1}, \dots, \widehat{z}_{iK}\}$ をまとめて、ベクトル $\widehat{\mathbf{z}}_i$ として記述する
- パラメータ $\widehat{\mathbf{z}}_i = [\widehat{z}_{i1}, \dots, \widehat{z}_{iK}]^T$ は、データ x_i に対応する潜在変数 z の各要素が、1 となる確率を集めたベクトルである
- $\widehat{z}_{ik} = q(z_{ik} = 1|x_i)$ と定義されていることに注意
- よってベクトル $\widehat{\mathbf{z}}_i$ は、データ x_i に対応する二値のスパース符号である
- 勾配上昇法を利用しない理由
 - データ \mathbf{X} と、潜在変数 \mathbf{Z} についての勾配上昇法を用いれば、学習することが可能

離散潜在変数の変分推論

- 但し、その方法では、各 x_i について平均場パラメータ \hat{z}_i を保管する必要がある
 - 各事例について、動的に更新されるベクトルが必要であるため、そのアルゴリズムを、数十億もの事例に対して適用することは困難である
 - また、収束するまで繰り返し計算を行うため、データ x から、パラメータ \hat{z}_i を素早く抽出することができない
 - 現実にデプロイされるときは、 \hat{z}_i をリアルタイムで計算できなければならない
-
- 不動点方程式による平均場パラメータ \hat{z}_i の推定
 - 勾配上昇法の代わりに、不動点方程式を使ってパラメータ \hat{z}_{ik} を素早く推定できる
 - $\nabla_{\hat{z}_i} \mathcal{L}(q) = 0$ をみtas、 \hat{z}_i の極大値を見つけ出す

離散潜在変数の変分推論

- \widehat{z}_i の全ての成分について同時に解くことはできないので、各成分 \widehat{z}_{ik} について繰り返し解く
- 即ち、各パラメータ \widehat{z}_{ik} について、下界 $\mathcal{L}(q)$ を順番に最適化する手続きを、 $\mathcal{L}(q)$ の収束基準を満たすまで繰り返す

$$\frac{\partial}{\partial \widehat{z}_{ik}} \mathcal{L}(q) = 0 \quad (564)$$

- 平均場不動点方程式を導くためには、 $\mathcal{L}(q)$ を \widehat{z}_{ik} で微分する必要がある

$$\begin{aligned} & \frac{\partial}{\partial \widehat{z}_{ik}} \mathcal{L}(q) \\ = & \frac{\partial}{\partial \widehat{z}_{ik}} \sum_i \sum_k \{ \widehat{z}_{ik} (\ln \sigma(b_{ik}) - \ln \widehat{z}_{ik}) + \\ & (1 - \widehat{z}_{ik}) (\ln \sigma(-b_{ik}) - \ln (1 - \widehat{z}_{ik})) \} + \end{aligned}$$

$$\frac{\partial}{\partial \widehat{z_{ik}}} \frac{1}{2} \sum_i \sum_j \left(\ln \frac{\beta_j}{2\pi} - \beta_j \left(x_{ij}^2 - 2x_{ij} \sum_k W_{jk} \widehat{z_{ik}} + \sum_k \left(W_{jk}^2 \widehat{z_{ik}} + \sum_{l \neq k} W_{jk} W_{jl} \widehat{z_{ik}} \widehat{z_{il}} \right) \right) \right) \quad (565)$$

- 前半部分は

$$\frac{\partial}{\partial \widehat{z_{ik}}} \sum_i \sum_k \{ \widehat{z_{ik}} (\ln \sigma(b_{ik}) - \ln \widehat{z_{ik}}) + (1 - \widehat{z_{ik}}) (\ln \sigma(-b_{ik}) - \ln (1 - \widehat{z_{ik}})) \} \quad (566)$$

$$= \frac{\partial}{\partial \widehat{z_{ik}}} \{ \widehat{z_{ik}} (\ln \sigma(b_{ik}) - \ln \widehat{z_{ik}}) + (1 - \widehat{z_{ik}}) (\ln \sigma(-b_{ik}) - \ln (1 - \widehat{z_{ik}})) \} \quad (567)$$

$$= \frac{\partial}{\partial \widehat{z_{ik}}} \widehat{z_{ik}} (\ln \sigma(b_{ik}) - \ln \widehat{z_{ik}}) +$$

$$\frac{\partial}{\partial \widehat{z_{ik}}} (1 - \widehat{z_{ik}}) (\ln \sigma(-b_{ik}) - \ln (1 - \widehat{z_{ik}})) \quad (568)$$

$$\begin{aligned} = & (\ln \sigma(b_{ik}) - \ln \widehat{z_{ik}}) + \widehat{z_{ik}} \left(-\frac{1}{\widehat{z_{ik}}} \right) + \\ & (- (\ln \sigma(-b_{ik}) - \ln (1 - \widehat{z_{ik}}))) + \\ & (1 - \widehat{z_{ik}}) \frac{1}{1 - \widehat{z_{ik}}} \end{aligned} \quad (569)$$

$$= \ln \sigma(b_{ik}) - \ln \widehat{z_{ik}} - 1 - \ln \sigma(-b_{ik}) + \ln(1 - \widehat{z_{ik}}) + 1 \quad (570)$$

$$= \ln \sigma(b_{ik}) - \ln \widehat{z_{ik}} - \ln \sigma(-b_{ik}) + \ln(1 - \widehat{z_{ik}}) \quad (571)$$

$$= -\zeta(-b_{ik}) - \ln \widehat{z_{ik}} + \zeta(b_{ik}) + \ln(1 - \widehat{z_{ik}}) \quad (572)$$

$$= (\zeta(b_{ik}) - \zeta(-b_{ik})) - \ln \widehat{z_{ik}} + \ln(1 - \widehat{z_{ik}}) \quad (573)$$

$$= b_{ik} - \ln \widehat{z_{ik}} + \ln(1 - \widehat{z_{ik}}) \quad (574)$$

離散潜在変数の変分推論

- ここで、シグモイド関数 $\sigma(\cdot)$ と、ソフトプラス関数 $\zeta(\cdot)$ に関する、以下の公式を用いた

$$\ln \sigma(x) = -\zeta(-x) \quad (575)$$

$$\zeta(x) - \zeta(-x) = x \quad (576)$$

- 後半部分は

$$\begin{aligned} & \frac{\partial}{\partial \widehat{z_{ik}}} \frac{1}{2} \sum_i \sum_j \left(\ln \frac{\beta_j}{2\pi} - \beta_j \left(x_{ij}^2 - 2x_{ij} \sum_k W_{jk} \widehat{z_{ik}} + \right. \right. \\ & \quad \left. \left. \sum_k \left(W_{jk}^2 \widehat{z_{ik}} + \sum_{l \neq k} W_{jk} W_{jl} \widehat{z_{ik}} \widehat{z_{il}} \right) \right) \right) \quad (577) \\ &= \frac{1}{2} \sum_j \beta_j \frac{\partial}{\partial \widehat{z_{ik}}} \left(2x_{ij} \sum_k W_{jk} \widehat{z_{ik}} - \right. \end{aligned}$$

$$\sum_k \left(W_{jk}^2 \widehat{z_{ik}} + \sum_{l \neq k} W_{jk} W_{jl} \widehat{z_{ik}} \widehat{z_{il}} \right) \quad (578)$$

$$\begin{aligned} &= \frac{1}{2} \sum_j \beta_j \frac{\partial}{\partial \widehat{z_{ik}}} \left(2x_{ij} \sum_k W_{jk} \widehat{z_{ik}} - \right. \\ &\quad \left. \left(W_{jk}^2 \widehat{z_{ik}} + \sum_{l \neq k} W_{jk} W_{jl} \widehat{z_{ik}} \widehat{z_{il}} \right) - \right. \\ &\quad \left. \sum_{m \neq k} \left(W_{jm}^2 \widehat{z_{im}} + \sum_{l \neq m} W_{jm} W_{jl} \widehat{z_{im}} \widehat{z_{il}} \right) \right) \quad (579) \end{aligned}$$

$$\begin{aligned} &= \frac{1}{2} \sum_j \beta_j \left(2x_{ij} W_{jk} - W_{jk}^2 - \right. \\ &\quad \left. \sum_{l \neq k} W_{jk} W_{jl} \widehat{z_{il}} - \sum_{m \neq k} W_{jm} W_{jk} \widehat{z_{im}} \right) \quad (580) \end{aligned}$$

$$= \frac{1}{2} \sum_j \beta_j \left(2x_{ij}W_{jk} - W_{jk}^2 - 2 \sum_{l \neq k} W_{jk}W_{jl}\widehat{z_{il}} \right) \quad (581)$$

$$= \sum_j \beta_j \left(x_{ij}W_{jk} - \frac{1}{2}W_{jk}^2 - \sum_{l \neq k} W_{jk}W_{jl}\widehat{z_{il}} \right) \quad (582)$$

$$= \sum_j x_{ij}\beta_j W_{jk} - \frac{1}{2} \sum_j W_{jk}\beta_j W_{jk} - \sum_j \beta_j \sum_{l \neq k} W_{jk}W_{jl}\widehat{z_{il}} \quad (583)$$

$$= \sum_j x_{ij}\beta_j W_{jk} - \frac{1}{2} \sum_j W_{jk}\beta_j W_{jk} - \sum_{l \neq k} \left(\sum_j W_{jl}\beta_j W_{jk} \right) \widehat{z_{il}} \quad (584)$$

$$= \mathbf{x}_i^T \beta \mathbf{W}_{:k} - \frac{1}{2} \mathbf{W}_{:k}^T \beta \mathbf{W}_{:k} - \sum_{l \neq k} \mathbf{W}_{:l}^T \beta \mathbf{W}_{:k} \widehat{z}_{il} \quad (585)$$

- これより、 $\mathcal{L}(q)$ の \widehat{z}_{ik} による微分は次のようになる

$$\begin{aligned} & \frac{\partial}{\partial \widehat{z}_{ik}} \mathcal{L}(q) \\ = & b_{ik} - \ln \widehat{z}_{ik} + \ln(1 - \widehat{z}_{ik}) + \\ & \mathbf{x}_i^T \beta \mathbf{W}_{:k} - \frac{1}{2} \mathbf{W}_{:k}^T \beta \mathbf{W}_{:k} - \sum_{l \neq k} \mathbf{W}_{:l}^T \beta \mathbf{W}_{:k} \widehat{z}_{il} \quad (586) \end{aligned}$$

- これを 0 と等置して、 \widehat{z}_{ik} について解くと次のようになる

$$\ln \widehat{z}_{ik} - \ln(1 - \widehat{z}_{ik}) = b_{ik} + \mathbf{x}_i^T \boldsymbol{\beta} \mathbf{W}_{:k} - \frac{1}{2} \mathbf{W}_{:k}^T \boldsymbol{\beta} \mathbf{W}_{:k} - \sum_{l \neq k} \mathbf{W}_{:l}^T \boldsymbol{\beta} \mathbf{W}_{:k} \widehat{z}_{il} \quad (587)$$

右辺を A とおけば

$$\begin{aligned} \ln \widehat{z}_{ik} - \ln(1 - \widehat{z}_{ik}) &= A & (588) \\ \Rightarrow \ln \frac{\widehat{z}_{ik}}{1 - \widehat{z}_{ik}} &= A \\ \Rightarrow \frac{\widehat{z}_{ik}}{1 - \widehat{z}_{ik}} &= \exp A \\ \Rightarrow \frac{1}{1 - \widehat{z}_{ik}} - 1 &= \exp A \\ \Rightarrow 1 - \widehat{z}_{ik} &= \frac{1}{1 + \exp A} \end{aligned}$$

$$\begin{aligned}\Rightarrow \widehat{z_{ik}} &= 1 - \frac{1}{1 + \exp A} \\ \Rightarrow \widehat{z_{ik}} &= \frac{\exp A}{1 + \exp A} \\ \Rightarrow \widehat{z_{ik}} &= \frac{1}{1 + \exp(-A)}\end{aligned}\tag{589}$$

$$\Rightarrow \widehat{z_{ik}} = \sigma(A)\tag{590}$$

従って、不動点方程式は

$$\widehat{z_{ik}} = \sigma \left(b_{ik} + \mathbf{x}_i^T \boldsymbol{\beta} \mathbf{W}_{:k} - \frac{1}{2} \mathbf{W}_{:k}^T \boldsymbol{\beta} \mathbf{W}_{:k} - \sum_{l \neq k} \mathbf{W}_{:l}^T \boldsymbol{\beta} \mathbf{W}_{:k} \widehat{z_{il}} \right)\tag{591}$$

- 不動点方程式の観察

離散潜在変数の変分推論

- 不動点方程式は次で表された

$$\widehat{z}_{ik} = \sigma \left(b_{ik} + \mathbf{x}_i^T \beta \mathbf{W}_{:k} - \frac{1}{2} \mathbf{W}_{:k}^T \beta \mathbf{W}_{:k} - \sum_{l \neq k} \mathbf{W}_{:l}^T \beta \mathbf{W}_{:k} \widehat{z}_{il} \right) \quad (592)$$

- 第2項 $\mathbf{x}_i^T \beta \mathbf{W}_{:k}$ は、潜在変数のユニット k に対する入力
- 第3項 $-\frac{1}{2} \mathbf{W}_{:k}^T \beta \mathbf{W}_{:k}$ は、隠れユニット k から自身への入力
- 第4項 $-\sum_{l \neq k} \mathbf{W}_{:l}^T \beta \mathbf{W}_{:k} \widehat{z}_{il}$ は、他の隠れユニット $l \neq k$ から隠れユニット k への入力
- これより、平均場不動点方程式は、回帰結合型ニューラルネットワーク (RNN) との関係があることが分かる
- 隠れユニット k と l は、それらの重みベクトル $\mathbf{W}_{:l}$ と $\mathbf{W}_{:k}$ が互いに同調するとき (似たような重みを持っているとき) に、互いに抑制し合う

離散潜在変数の変分推論

- 即ち、2つの隠れユニット k, l が、共に入力を説明するとき (入力から同じような表現を抽出するとき)、**入力を最もよく説明するユニットのみがアクティブ**になる (強く活性化される)
- これはユニット間の競合の一形態である
- 従って、実際には多峰性の事後分布かもしれないが、そのうちの1つのみが選択される (図 14 の (b) と (c) 参照)
- 不動点方程式を更に以下のように変形する

$$\widehat{z}_{ik} = \sigma \left(b_{ik} + \left(\mathbf{x}_i - \sum_{l \neq k} \mathbf{W}_{:l} \widehat{z}_{il} \right)^T \boldsymbol{\beta} \mathbf{W}_{:k} - \frac{1}{2} \mathbf{W}_{:k}^T \boldsymbol{\beta} \mathbf{W}_{:k} \right) \quad (593)$$

- これより、ユニット k への入力は、 \mathbf{x}_i ではなく $\mathbf{x}_i - \sum_{l \neq k} \mathbf{W}_{:l} \widehat{z}_{il}$ であるとみなせる

離散潜在変数の変分推論

- ユニット k への入力は、他の全てのユニットによる x_i の再構成と、実際の入力 x_i との誤差である
 - ユニット k は、この残差誤差を符号化していると分かるので、スパース符号化は、**反復自己符号化器**とみなせる
 - スパース符号化では、入力 x_i の符号化 ($\widehat{z_{ik}}$ の計算) と、復号 ($\sum_{l \neq k} \mathbf{W}_{:l} \widehat{z_{il}}$ の計算) を繰り返す
 - この反復のたびに、再構成の誤差を修正していく
-
- **ダンピング**
 - 1つのユニットの更新則を (不動点方程式として) 導出した
 - 複数のユニットを同時に更新することは、二値スパース符号化モデルでは通常できない
 - 但し、**ダンピング**という発見的手法を使えば可能になる
 - 各要素 $\widehat{z_{ik}}$ についての最適値を計算し、その値の変化の方向に、他の要素 $\widehat{z_{il}}$ を小さいステップで動かす

離散潜在変数の変分推論

- 下界 $\mathcal{L}(q)$ が増加することはもはや保証されないが、多くの場合はうまくいく

離散潜在変数の変分推論のまとめ

● ここまでの話の流れ

- 1 離散潜在変数における変分推論の具体例として、二値スパース符号化モデルをみた
- 2 事後分布 $p(Z|X)$ が複雑になるので、最尤推定 (EM アルゴリズム) が使えない
- 3 代わりに、別の分布 $q(Z)$ を使って事後分布を近似することにした
- 4 エビデンス下界 $\mathcal{L}(q)$ を苦勞して求めた
- 5 更に、下界 $\mathcal{L}(q)$ を (各パラメータについて) 最大化するための、不動点方程式を導出した
- 6 不動点方程式を観察し、回帰結合型ニューラルネットワークや、自己符号化器との関係を考えて

● これからの話の流れ

離散潜在変数の変分推論のまとめ

- 連続潜在変数に対する変分推論を、簡単な確率モデルを使って、試してみよう

5 近似推論法

- 変分推論
- 座標降下法
- MAP 推定および最尤推定と変分推論
- 変分近似による弊害
- 離散潜在変数の変分推論
- 連続潜在変数の変分推論
- 正則化との関わり
- 雑多な話題

連続潜在変数の変分推論

- 連続潜在変数をもつ変分推論の概要

- 平均場近似を行う場合は、以下の式によって最適な因子 $q_j^*(Z_j|X)$ が得られる

$$\ln q_j^*(Z_j|X) = \mathbb{E}_{i \neq j} [\ln p(X, Z)] + \text{Const.} \quad (594)$$

$$\mathbb{E}_{i \neq j} [\ln p(X, Z)] = \int \ln p(X, Z) \prod_{i \neq j} q_i(Z_i|X) dZ_i \quad (595)$$

$$q_j^*(Z_j|X) = \frac{\exp(\mathbb{E}_{i \neq j} [\ln p(X, Z)])}{\int \exp(\mathbb{E}_{i \neq j} [\ln p(X, Z)]) dZ_j} \quad (596)$$

- 上式は、下界 $\mathcal{L}(q)$ を最大化する q であり、従って連続潜在変数の場合の不動点方程式とみなせる
- 各因子 $q_j(Z_j|X)$ を、上式を用いて順番に更新していく ($i \neq j$ である全ての q_i を固定した状態で、各 q_j について下界を最適化する)

連続潜在変数の変分推論

- このステップを、下界 $\mathcal{L}(q)$ が収束するまで繰り返し行う (座標降下法)
- 不動点方程式は、下界が最適な値に収束するかどうかには関係なく、 q_j の最適解が取る関数形を提供してくれる

- 扱う確率モデルの表現

- ここでは次のような確率モデルを対象として、変分推論を扱う

$$p(\mathbf{z}_i) = \mathcal{N}(\mathbf{z}_i | 0, \mathbf{I}) \quad (597)$$

$$p(x_i | \mathbf{z}_i) = \mathcal{N}(x_i | \mathbf{w}^T \mathbf{z}_i, 1) \quad (598)$$

- 1次元のデータ $x \in \mathbb{R}$ に対して、2次元の潜在変数 $\mathbf{z} \in \mathbb{R}^2$ が存在する
- 同時分布 $p(x_i, \mathbf{z}_i) = p(x_i | \mathbf{z}_i)p(\mathbf{z}_i)$ を \mathbf{z}_i で積分消去すれば、 x_i についての単なるガウス分布となる

- 真の事後分布 $p(\mathbf{z}_i | x_i)$ の計算

連続潜在変数の変分推論

- 正規化定数を見捨てて次のように計算できる

$$\begin{aligned} & p(\mathbf{z}_i | x_i) \\ \propto & p(\mathbf{z}_i | x_i) p(x_i) \\ = & p(x_i, \mathbf{z}_i) \\ = & p(\mathbf{z}_i) p(x_i | \mathbf{z}_i) \\ = & \mathcal{N}(\mathbf{z}_i | 0, \mathbf{I}) \mathcal{N}(x_i | \mathbf{w}^T \mathbf{z}_i, 1) \\ \propto & \exp\left(-\frac{1}{2} \mathbf{z}_i^T \mathbf{z}_i\right) \exp\left(-\frac{1}{2} (x_i - \mathbf{w}^T \mathbf{z}_i)^T (x_i - \mathbf{w}^T \mathbf{z}_i)\right) \\ = & \exp\left(-\frac{1}{2} (z_{i1}^2 + z_{i2}^2)\right) \\ & \exp\left(-\frac{1}{2} (x_i - w_1 z_{i1} - w_2 z_{i2})^T (x_i - w_1 z_{i1} - w_2 z_{i2})\right) \\ = & \exp\left(-\frac{1}{2} (z_{i1}^2 + z_{i2}^2 + x_i^2 + w_1^2 z_{i1}^2 + w_2^2 z_{i2}^2 - \right. \end{aligned}$$

$$2x_i w_1 z_{i1} - 2x_i w_2 z_{i2} + 2w_1 w_2 z_{i1} z_{i2} \Big) \quad (599)$$

- 正規化項を C とすれば次のように書ける

$$\begin{aligned} p(z_i | x_i) \\ = C \exp \left(-\frac{1}{2} (z_{i1}^2 + z_{i2}^2 + x_i^2 + w_1^2 z_{i1}^2 + w_2^2 z_{i2}^2 - \right. \\ \left. 2x_i w_1 z_{i1} - 2x_i w_2 z_{i2} + 2w_1 w_2 z_{i1} z_{i2}) \right) \end{aligned} \quad (600)$$

- z_{i1} と z_{i2} を乗算する項が存在する
- 従って、真の事後分布は、 z_{i1} と z_{i2} のみの因子には分解できないことが分かる

- 平均場近似の計算

- 平均場近似を次のように表現する

$$q(\mathbf{z}_i | x_i) = q_1(z_{i1} | x_i) q_2(z_{i2} | x_i) \quad (601)$$

- q_2 を固定した状態で、最適な $q_1^*(z_{i1} | x_i)$ を求める

$$\begin{aligned} & \ln q_1^*(z_{i1} | x_i) \\ = & \mathbb{E}_{z_{i2} \sim q_2(z_{i2} | x_i)} [\ln p(\mathbf{z}_i, x_i)] + \text{Const.} \end{aligned} \quad (602)$$

$$\begin{aligned} = & \mathbb{E}_{z_{i2}} \left[\ln \left(C \exp \left(-\frac{1}{2} (z_{i1}^2 + z_{i2}^2 + x_i^2 + w_1^2 z_{i1}^2 + w_2^2 z_{i2}^2 - \right. \right. \right. \\ & \left. \left. \left. 2x_i w_1 z_{i1} - 2x_i w_2 z_{i2} + 2w_1 w_2 z_{i1} z_{i2} \right) \right) \right] + \text{Const.} \\ = & \mathbb{E}_{z_{i2}} \left[-\frac{1}{2} (z_{i1}^2 + z_{i2}^2 + x_i^2 + w_1^2 z_{i1}^2 + w_2^2 z_{i2}^2 - \right. \\ & \left. 2x_i w_1 z_{i1} - 2x_i w_2 z_{i2} + 2w_1 w_2 z_{i1} z_{i2}) \right] + \text{Const.} \end{aligned} \quad (603)$$

連続潜在変数の変分推論

- $z_{i2} \sim q_2(z_{i2}|x_i)$ による期待値を取っている
- $q_2(z_{i2}|x_i)$ から得る必要があるのは、結局 $\mathbb{E}_{z_{i2}}[z_{i2}]$ と、 $\mathbb{E}_{z_{i2}}[z_{i2}^2]$ の2つだけである
- $\langle z_{i2} \rangle = \mathbb{E}_{z_{i2}}[z_{i2}]$ 、 $\langle z_{i2}^2 \rangle = \mathbb{E}_{z_{i2}}[z_{i2}^2]$ と書くことにする
- このとき次式が得られる

$$\begin{aligned} & \ln q_1^*(z_{i1}|x_i) \\ = & -\frac{1}{2} \left(z_{i1}^2 + \langle z_{i2}^2 \rangle + x_i^2 + w_1^2 z_{i1}^2 + w_2^2 \langle z_{i2}^2 \rangle - \right. \\ & \quad \left. 2x_i w_1 z_{i1} - 2x_i w_2 \langle z_{i2} \rangle + 2w_1 w_2 z_{i1} \langle z_{i2} \rangle \right) + \\ & \quad \text{Const.} \end{aligned} \tag{604}$$

連続潜在変数の変分推論

- これより、最適な $q_1^*(z_{i1}|x_i)$ は**ガウス分布**の形であると分かる

$$\begin{aligned} & q_1^*(z_{i1}|x_i) \\ = & C \exp \left(-\frac{1}{2} (z_{i1}^2 + \langle z_{i2}^2 \rangle + x_i^2 + w_1^2 z_{i1}^2 + w_2^2 \langle z_{i2}^2 \rangle - \right. \\ & \quad \left. 2x_i w_1 z_{i1} - 2x_i w_2 \langle z_{i2} \rangle + 2w_1 w_2 z_{i1} \langle z_{i2} \rangle) \right) \\ \propto & \exp \left(-\frac{1}{2} (z_{i1}^2 + \langle z_{i2}^2 \rangle + x_i^2 + w_1^2 z_{i1}^2 + w_2^2 \langle z_{i2}^2 \rangle - \right. \\ & \quad \left. 2x_i w_1 z_{i1} - 2x_i w_2 \langle z_{i2} \rangle + 2w_1 w_2 z_{i1} \langle z_{i2} \rangle) \right) \quad (605) \end{aligned}$$

- 対称性から、最適な $q_2^*(z_{i2}|x_i)$ も**ガウス分布**であることが分かる
- ガウス分布同士の積もガウス分布になるので、結局 $q(z_i|x_i) = q_1(z_{i1}|x_i)q_2(z_{i2}|x_i)$ はガウス分布である

連続潜在変数の変分推論

- $q(z_i|x_i)$ が 2 つの因子に分解できるとは仮定したが、**各因子の関数形については全く仮定していないことに注意**
- ガウス分布は、下界 $\mathcal{L}(q)$ を q について変分最適化する過程で、**自然に出現した**

連続潜在変数の変分推論のまとめ

- ここまでの話の流れ

- 1 連続潜在変数における変分推論の例として、簡単な確率モデルを扱った
- 2 分布 $q(\mathbf{Z}|\mathbf{X})$ を、 $\prod_i q_i(\mathbf{Z}_i|\mathbf{X})$ のように因数分解できるという仮定 (平均場近似) のみを置いた
- 3 各因子 $q_i(\mathbf{Z}_i|\mathbf{X})$ の関数形については全く仮定を置かなかった
- 4 変分推論によって、最適解となる q_i の関数形が自然に導出できた

- これからの話の流れ

- 変分推論と、正則化との関係を見てみよう

5 近似推論法

- 変分推論
- 座標降下法
- MAP 推定および最尤推定と変分推論
- 変分近似による弊害
- 離散潜在変数の変分推論
- 連続潜在変数の変分推論
- 正則化との関わり
- 雑多な話題

- 正則化との関係

- 変分下界 $\mathcal{L}(q)$ を次のように分解してみよう

$$\begin{aligned}\mathcal{L}(q) &= \int q(\mathbf{Z}) \ln \frac{p(\mathbf{X}, \mathbf{Z})}{q(\mathbf{Z})} d\mathbf{Z} \\&= \int q(\mathbf{Z}) \ln \frac{p(\mathbf{X}|\mathbf{Z})p(\mathbf{Z})}{q(\mathbf{Z})} d\mathbf{Z} \\&= \int q(\mathbf{Z}) \left(\ln p(\mathbf{X}|\mathbf{Z}) + \ln \frac{p(\mathbf{Z})}{q(\mathbf{Z})} \right) d\mathbf{Z} \\&= \int q(\mathbf{Z}) \ln p(\mathbf{X}|\mathbf{Z}) d\mathbf{Z} + \int q(\mathbf{Z}) \ln \frac{p(\mathbf{Z})}{q(\mathbf{Z})} d\mathbf{Z} \\&= \int q(\mathbf{Z}) \ln p(\mathbf{X}|\mathbf{Z}) d\mathbf{Z} - \text{KL}(q||p)\end{aligned}\tag{606}$$

正則化との関わり

- 下界 $\mathcal{L}(q)$ を最大化するとき、第 1 項を**最大化**し、第 2 項 $\text{KL}(q||p)$ を**最小化**する必要がある
- 第 1 項は、観測データの対数尤度 $\ln p(\mathbf{X}|\mathbf{Z})$ の、近似事後分布 $q(\mathbf{Z})$ についての期待値である
- この項の最大化は、**最尤推定**の考え方と似ている
- $q(\mathbf{Z})$ を、 $\ln(\mathbf{X}|\mathbf{Z})$ が最大となる \mathbf{Z}^* でのみ非零となる、デルタ関数 $\delta(\mathbf{Z} - \mathbf{Z}^*)$ とすれば、最大化できる
- 第 2 項は、近似事後分布 $q(\mathbf{Z})$ と、潜在変数に関する事前分布 $p(\mathbf{Z})$ との KL ダイバージェンス
- これを小さくしようとすることは、近似分布 $q(\mathbf{Z})$ が事前分布 $p(\mathbf{Z})$ から、大きく離れることを防ぐことにあたる
- 第 1 項で、近似分布が際限なく尖ろうとすることに対して、第 2 項では、**予め定められた事前分布 $p(\mathbf{Z})$ を通じて、ペナルティを与える**

正則化との関わり

- 第 2 項は、MAP 推定において、パラメータの事前分布 $p(\mathbf{Z})$ を導入することに相当する
- これより、第 2 項は正則化項として作用していることが分かる

正則化との関わりのおまとめ

- ここまでの話の流れ

- 1 変分下界 $\mathcal{L}(q)$ の式を変形し、2つの項に分解した
- 2 最尤推定に相当する項と、パラメータに関する事前分布の導入による **正則化** に相当する項とに、分かれていることを確認した

- これからの話の流れ

- 教科書で触れられている雑多な話題を扱う
- 近似推論が、推論アルゴリズムの精度に影響することについて考える
- 学習による近似推論の手法を考える

5 近似推論法

- 変分推論
- 座標降下法
- MAP 推定および最尤推定と変分推論
- 変分近似による弊害
- 離散潜在変数の変分推論
- 連続潜在変数の変分推論
- 正則化との関わり
- 雑多な話題

- 近似推論が精度に与える影響

- 近似推論では、以下の下界 $\mathcal{L}(q)$ を q について最適化する

$$\ln p(\mathbf{X}) = \mathcal{L}(q) + \text{KL}(q||p) \quad (607)$$

$$\mathcal{L}(q) = \int q(\mathbf{Z}) \ln \frac{p(\mathbf{X}, \mathbf{Z})}{q(\mathbf{Z})} d\mathbf{Z} \quad (608)$$

- 下界 $\mathcal{L}(q)$ を次のように分解する

$$\mathcal{L}(q) = \int q(\mathbf{Z}) \ln p(\mathbf{X}, \mathbf{Z}) d\mathbf{Z} - \int q(\mathbf{Z}) \ln q(\mathbf{Z}) d\mathbf{Z} \quad (609)$$

$$= \int q(\mathbf{Z}) \ln p(\mathbf{X}, \mathbf{Z}) d\mathbf{Z} + H[q] \quad (610)$$

$$= \mathbb{E}_{\mathbf{Z} \sim q(\mathbf{Z})} [\ln p(\mathbf{X}, \mathbf{Z})] + H[q] \quad (611)$$

$$= \mathbb{E}_{\mathbf{Z} \sim q(\mathbf{Z})} [\ln p(\mathbf{X}|\mathbf{Z})p(\mathbf{Z})] + H[q] \quad (612)$$

- 分布に含まれるパラメータを明示的に記述する

学習と推論の相互作用

- このとき、モデリングされる関数は $p(\mathbf{X}|\mathbf{Z}, \theta)$ 、 $p(\mathbf{Z}|\theta)$ と表せるので、 $p(\mathbf{X}|\mathbf{Z}) = p(\mathbf{X}, \mathbf{Z}|\theta) = p(\mathbf{X}|\mathbf{Z}, \theta)p(\mathbf{Z}|\theta)$ と書ける
- 従って、下界 $\mathcal{L}(q)$ は次のようになる

$$\mathcal{L}(q, \theta) = \mathbb{E}_{\mathbf{Z}} [\ln p(\mathbf{X}, \mathbf{Z}|\theta)] + H[q] \quad (613)$$

- 下界 $\mathcal{L}(q, \theta)$ をパラメータ θ について増加させるとする
- これは、 $\mathbb{E}_{\mathbf{Z}} [\ln p(\mathbf{X}, \mathbf{Z}|\theta)] = \mathbb{E}_{\mathbf{Z}} [\ln p(\mathbf{Z}|\mathbf{X}, \theta)] + \text{Const.}$ を、 q を固定しつつ、 θ について増加させることに繋がる
- このとき、 $q(\mathbf{Z})$ が高い確率となる \mathbf{Z} について $p(\mathbf{Z}|\mathbf{X}, \theta)$ が増大し、 $q(\mathbf{Z})$ が低い確率となる \mathbf{Z} について $p(\mathbf{Z}|\mathbf{X}, \theta)$ が減少する

$$\begin{aligned} \mathbb{E}_{\mathbf{Z}} [\ln p(\mathbf{X}, \mathbf{Z}|\theta)] &= \mathbb{E}_{\mathbf{Z}} [\ln p(\mathbf{Z}|\mathbf{X}, \theta)p(\mathbf{X})] \\ &= \mathbb{E}_{\mathbf{Z}} [\ln p(\mathbf{Z}|\mathbf{X}, \theta)] + \mathbb{E}_{\mathbf{Z}} [\ln p(\mathbf{X})] \\ &= \mathbb{E}_{\mathbf{Z}} [\ln p(\mathbf{Z}|\mathbf{X}, \theta)] + \ln p(\mathbf{X}) \\ &= \mathbb{E}_{\mathbf{Z}} [\ln p(\mathbf{Z}|\mathbf{X}, \theta)] + \text{Const.} \end{aligned} \quad (614)$$

学習と推論の相互作用

- これより、近似仮定が**自己充足的予言**となることが分かる
- 自己充足的予言とは、ある事象や状況に関する判断や思い込みが原因となり、その結果として、その判断や思い込みが現実化することである
- 即ち、予め近似事後分布 $q(\mathbf{Z})$ に何らかの仮定を置いて訓練することで、**結果としてその仮定に沿うような分布**が得られてしまう
- 単峰性の近似事後分布 $q(\mathbf{Z})$ を使って訓練することを考える
- 得られる事後分布 $q(\mathbf{Z}) = p(\mathbf{Z}|\mathbf{X})$ は、厳密な推論によって得られた $p(\mathbf{Z}|\mathbf{X})$ と比べて、単峰性がはるかに強くなっている
- 事後分布が多峰性であるという仮定が、**推定結果にも現れてしまう**
- 変分推論によってモデルに課される損害
 - モデルを訓練した後に、 $\ln p(\mathbf{X}|\boldsymbol{\theta})$ を推定し、下界 $\mathcal{L}(q, \boldsymbol{\theta})$ との**差が十分に小さい** ($\mathcal{L}(q, \boldsymbol{\theta}) \simeq \ln p(\mathbf{X}|\boldsymbol{\theta})$) ことを確認する

学習と推論の相互作用

- このとき、ある特定のパラメータ θ について、変分近似が正確であるといえる
- 但し、一般的に変分近似が正確であるとはいえない
- また、変分推論が、学習過程に殆ど害を及ぼさないとも結論付けてはいけない
- 変分近似による真の損害は、 $\ln p(\mathbf{X}|\theta)$ を最大にする θ^* が分からなければ、測定することができない
- ある θ について $\mathcal{L}(q, \theta) \simeq \ln p(\mathbf{X}|\theta)$ が成立しても、 $\mathcal{L}(q, \theta) \simeq \ln p(\mathbf{X}|\theta) \ll \ln p(\mathbf{X}|\theta^*)$ であるかもしれない
- あるパラメータ θ については変分近似は正確だが、実際には近似が全く上手く行っていないことになる

学習と推論の相互作用

- また $\max_q \mathcal{L}(q, \theta^*) \ll \ln p(\mathbf{X}|\theta^*)$ であれば、 θ^* のときの事後分布は、制限された q にとって複雑すぎるため、訓練では θ^* を見つけることができない
- この問題は、変分推論とは別の方法で θ^* が求まる場合にしか、検知することができない

- 近似推論の実行の学習

- 不動点方程式や勾配に基づく最適化を、繰り返し行う近似推論は、計算コストが非常に高い
- 近似推論による最適化処理は、入力 X から、近似事後分布 $q^*(Z) = \arg \max_q \mathcal{L}(X, q)$ へと写像する関数 f とみなせる
- これにより、反復的な最適化処理を単なる関数とみなして、関数 $\hat{f}(X, \theta)$ を、ニューラルネットワークで近似することができる

- Wake-Sleep アルゴリズム

- データ X から潜在変数 Z を推論するときの困難は、正しい潜在変数 Z が分からず、従って教師あり訓練集合が得られないことである
- データ X から潜在変数 Z への写像は、モデル $(p(X|Z, \theta), p(Z|\theta))$ の選択) に依存し、更に θ の変化に伴って変わり続ける

学習による近似推論

- Wake-Sleep アルゴリズムでは、 X と Z の両方のサンプルを、モデル分布から抽出する
- モデルの分布が高い確率を示すような X に対してしか、推論ネットワークを学習できない
- モデルの分布が、データの分布に似ていないとき、推論ネットワークはデータに似たサンプルを学習できない
- 推論学習の他の形式
 - 変分自己符号化器は、生成モデリングで主要なアプローチとなった
 - 推論ネットワークは、単に下界 $\mathcal{L}(q)$ を定義するために使用される
 - ネットワークのパラメータは、下界 $\mathcal{L}(q)$ が増加するように適用される

目次

- 1 K-Means 法
- 2 混合ガウス分布
- 3 EM アルゴリズム
- 4 変分の導入
- 5 近似推論法
- 6 変分自己符号化器**

6 変分自己符号化器

- 生成モデル
 - 変分自己符号化器 (VAE) の概要
 - 変分自己符号化器 (VAE) の理論

- 生成モデルの目的

- データ x に関する分布 $p(x)$ を推定する

- データ生成過程

- データ x は一般に高次元である
- 但し、実際にデータが分布しているのは、ごく限られた一部の低次元の領域であると考えられる (多様体仮説)
- データ x 自体は高次元だが、本質的には低次元の情報しか持たないと考えられる
- データ x を、より低次元なベクトル z を使って、表現することを考える
- データに関する分布 $p(x)$ を、潜在変数 z に関する分布と、うまく組み合わせる
- 潜在変数からデータが生成されるまでの過程を組み込んで、 $p(x)$ を記述する

- 6 変分自己符号化器
 - 生成モデル
 - 変分自己符号化器 (VAE) の概要
 - 変分自己符号化器 (VAE) の理論

変分自己符号化器 (VAE) の概要

- 深層学習における生成モデル
 - 主に以下の 2 つの手法が存在する
 - 敵対的生成ネットワーク (Generative Adversarial Networks, GAN)
 - 変分自己符号化器 (Variational Auto Encoders, VAE)
 - ここでは変分自己符号化器 (VAE) について扱う
 - VAE を、異常検知 (不良品の検出など) に使った例がある

変分自己符号化器 (VAE) の概要

- VAE におけるグラフィカルモデル

- 図 16 のような、潜在変数を含んだグラフィカルモデルを考える
- データ x について、ある一つの潜在変数 z が対応しているとする
- 各データ x は、分布 $p(x)$ から独立にサンプルされるとする
- 従って、データ $\{x_1, \dots, x_N\}$ は独立同分布標本とする
- θ は、潜在変数 z からデータ x を取得する際に使用されるパラメータ
- ϕ は、データ x から潜在変数 z を生成する際に使用されるパラメータ
- N は、データ数である

変分自己符号化器 (VAE) の概要

- データ x の生成過程

- データ x の生成過程は、次のように考える
 - 分布 $p(\mathbf{z}|\theta)$ から、潜在変数 z_i がサンプルされる
 - 分布 $p(\mathbf{x}|z_i, \theta)$ から、データ x_i がサンプルされる
- これより、データ x の分布を次のように表現できる

$$p(\mathbf{x}|\theta) = \int p(\mathbf{x}|\mathbf{z}, \theta)p(\mathbf{z}|\theta)d\mathbf{z} \quad (615)$$

- 潜在変数 z をデータ x から取得する過程

- 潜在変数 z_i をデータ x_i から得る過程は、次のように考える
 - 分布 $q(\mathbf{z}|\mathbf{x}_i, \phi)$ から、潜在変数 z_i がサンプルされる

変分自己符号化器 (VAE) の概要

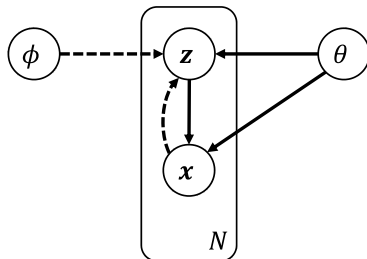


図 16: 変分自己符号化器 (VAE) におけるグラフィカルモデル

変分自己符号化器 (VAE) の概要

- 確率分布のニューラルネットワークによる表現
 - 潜在変数を含む確率モデルについて、パラメータの最尤解を求めるために、EM アルゴリズムを導出した
 - EM アルゴリズムでは、潜在変数に関する事後分布 $p(z|x, \theta)$ を計算する必要があった
 - この事後分布 $p(z|x, \theta)$ の計算が困難であるとき、 $p(z|x, \theta)$ を別の分布 $q(z|x, \phi)$ で近似し、変分推論によって $q(z|\phi)$ の最適解を求めた
- VAE は変分推論の変種であり、近似事後分布 $q(z|x, \phi)$ と、 $p(x|z, \theta)$ の2つをニューラルネットワークで表現する
- データ x を潜在変数 z に対応付けるニューラルネットワークを、Encoder という
- 潜在変数 z からデータ x を復元するニューラルネットワークを、Decoder という
- 分布 $q(z|x, \phi)$ は Encoder、分布 $p(x|z, \theta)$ は Decoder に相当する

- 6 変分自己符号化器
 - 生成モデル
 - 変分自己符号化器 (VAE) の概要
 - 変分自己符号化器 (VAE) の理論

変分自己符号化器 (VAE) の理論

- 変分自己符号化器 (VAE) の理論
 - 変分下界 $\mathcal{L}(q)$ は次のようであった

$$\mathcal{L}(q) = \int q(z|\mathbf{x}) \ln \frac{p(\mathbf{x}, z)}{q(z|\mathbf{x})} dz \quad (616)$$

$$= \int q(z|\mathbf{x}) \ln \frac{p(\mathbf{x}|z)p(z)}{q(z|\mathbf{x})} dz \quad (617)$$

$$= \int q(z|\mathbf{x}) \ln p(\mathbf{x}|z) dz + \int q(z|\mathbf{x}) \ln \frac{p(z)}{q(z|\mathbf{x})} dz \quad (618)$$

$$= \int q(z|\mathbf{x}) \ln p(\mathbf{x}|z) dz - \text{KL}(q(z|\mathbf{x})||p(z)) \quad (619)$$

$$= \mathbb{E}_{z \sim q(z|\mathbf{x})} [\ln p(\mathbf{x}|z)] - \text{KL}(q(z|\mathbf{x})||p(z)) \quad (620)$$

- ここでは、単一のデータ \mathbf{x} と、それに対応する潜在変数 z を考えている
- また、パラメータ θ, ϕ は省略している

変分自己符号化器 (VAE) の理論

- 第1項 $\mathbb{E}_{\mathbf{z} \sim q(\mathbf{z}|\mathbf{x})} [\ln p(\mathbf{x}|\mathbf{z})]$ を大きく、また第2項 $\text{KL}(q(\mathbf{z}|\mathbf{x})||p(\mathbf{z}))$ を小さくすることで、変分下界 $\mathcal{L}(q)$ を大きくできる
- VAE では、変分下界 $\mathcal{L}(q)$ を最大化するパラメータ θ, ϕ を求めるために、ニューラルネットを使用する (変分推論にニューラルネットをねじ込んだもの)
- KL ダイバージェンスの項は、後ほど求めることにする (解析的に求められる)
- 第1項は、分布 $q(\mathbf{z}|\mathbf{x})$ に関する期待値であり、VAE ではサンプリングで近似する

$$\mathbb{E}_{\mathbf{z} \sim q(\mathbf{z}|\mathbf{x})} [\ln p(\mathbf{x}|\mathbf{z})] \simeq \frac{1}{L} \sum_{l=1}^L \ln p(\mathbf{x}_i | \mathbf{z}_{i,l}) \quad (621)$$

変分自己符号化器 (VAE) の理論

- これより、変分下界 $\mathcal{L}(q)$ は以下のように書ける

$$\mathcal{L}(q) \simeq -\text{KL}(q(\mathbf{z}|\mathbf{x})||p(\mathbf{z})) + \frac{1}{L} \sum_{i=1}^L \ln p(\mathbf{x}_i|\mathbf{z}_{i,l}) \quad (622)$$

- パラメータ θ, ϕ を含めれば、次のように書ける

$$\mathcal{L}(q) \simeq -\text{KL}(q(\mathbf{z}|\mathbf{x}, \phi)||p(\mathbf{z}|\theta)) + \frac{1}{L} \sum_{i=1}^L \ln p(\mathbf{x}_i|\mathbf{z}_{i,l}, \theta) \quad (623)$$

変分自己符号化器 (VAE) の理論

- Encoder のニューラルネットの入出力
 - Encoder は、分布 $q(z|x, \phi)$ を表現するニューラルネット
 - 入力 x に対応する潜在変数 z を得る
- Encoder の入力、明らかにデータ x である
- 変分下界 $\mathcal{L}(q)$ において、項 $\mathbb{E}_z [\ln p(x|z)]$ は近似する必要があった
- z を、分布 $q(z|x, \phi)$ から L 回サンプリングした
- ニューラルネットで、 x から z を直接サンプリングするのは困難
- そこで、Encoder では、サンプルされたデータは出力しないことにする
- その代わりに、サンプルする分布のパラメータを出力する
- 例えば、サンプルする分布がガウス分布であれば、平均と分散の2つのパラメータを出力する

変分自己符号化器 (VAE) の理論

- 後述のように、分布 $q(z|\mathbf{x}, \phi)$ はガウス分布になるので、Encoder は平均ベクトル μ と共分散行列 Σ を出力する
- 共分散行列は、実際には対角行列であるため、実際には行列ではなく、行列の対角成分を要素にもつベクトルを出力する

変分自己符号化器 (VAE) の理論

● Encoder の損失関数

- VAE では、事前分布として、平均ベクトル $\mathbf{0}$ 、共分散行列 \mathbf{I} のガウス分布 $\mathcal{N}(\mathbf{z}|\mathbf{0}, \mathbf{I})$ を仮定する
- データ \mathbf{x} は高次元だが、実際にはそのうちの低次元な領域にまともって存在する (多様体仮説)
- 従って、データ \mathbf{x} の構造を、より低次元な潜在変数 \mathbf{z} の空間 $\mathbf{z} \sim \mathcal{N}(\mathbf{z}|\mathbf{0}, \mathbf{I})$ に押し込めることができる
- Encoder の損失関数は、KL ダイバージェンス $\text{KL}(q(\mathbf{z}|\mathbf{x}, \phi) || p(\mathbf{z}|\theta))$ で定義できる
- この KL ダイバージェンスを最小化することは、Encoder の分布 $q(\mathbf{z}|\mathbf{x}, \phi)$ を、 $p(\mathbf{z}|\theta) \sim \mathcal{N}(\mathbf{z}|\mathbf{0}, \mathbf{I})$ に近づける制約に相当する
- $p(\mathbf{z}|\theta)$ がガウス分布であれば、事後分布 $p(\mathbf{z}|\mathbf{x}, \theta)$ もガウス分布であり、従って $q(\mathbf{z}|\mathbf{x}, \phi)$ もガウス分布となる

変分自己符号化器 (VAE) の理論

- よって $\text{KL}(q(z|\mathbf{x}, \phi) || p(z|\theta))$ は、2つのガウス分布間の KL ダイバージェンスである
- 一般に、2つのガウス分布 $p(z) = \mathcal{N}(z|\boldsymbol{\mu}_0, \boldsymbol{\Sigma}_0)$ 、 $q(z) = \mathcal{N}(z|\boldsymbol{\mu}_1, \boldsymbol{\Sigma}_1)$ 間の KL ダイバージェンスは、解析的に計算できる
- KL ダイバージェンス $\text{KL}(p(z) || q(z))$ を順番に求めてみよう

$$\begin{aligned} & \text{KL}(p(z) || q(z)) \\ &= \int p(z) \ln \frac{p(z)}{q(z)} dz \end{aligned} \quad (624)$$

$$= \int p(z) (\ln p(z) - \ln q(z)) dz \quad (625)$$

$$= \int p(z) (\ln \mathcal{N}(z|\boldsymbol{\mu}_0, \boldsymbol{\Sigma}_0) - \ln \mathcal{N}(z|\boldsymbol{\mu}_1, \boldsymbol{\Sigma}_1)) dz \quad (626)$$

$$= \mathbb{E}_{z \sim p(z)} [\ln \mathcal{N}(z|\boldsymbol{\mu}_0, \boldsymbol{\Sigma}_0) - \ln \mathcal{N}(z|\boldsymbol{\mu}_1, \boldsymbol{\Sigma}_1)] \quad (627)$$

- データを D 次元、潜在変数を K 次元とする

変分自己符号化器 (VAE) の理論

- ここで

$$\begin{aligned} & \ln \mathcal{N}(z|\mu_0, \Sigma_0) \\ = & \ln \left(\frac{1}{(2\pi)^{\frac{K}{2}}} \frac{1}{|\Sigma_0|^{\frac{1}{2}}} \exp \left(-\frac{1}{2} (z - \mu_0)^T \Sigma_0^{-1} (z - \mu_0) \right) \right) \\ = & -\frac{K}{2} \ln 2\pi - \frac{1}{2} \ln |\Sigma_0| - \frac{1}{2} (z - \mu_0)^T \Sigma_0^{-1} (z - \mu_0) \quad (628) \end{aligned}$$

$$\begin{aligned} & \ln \mathcal{N}(z|\mu_1, \Sigma_1) \\ = & -\frac{K}{2} \ln 2\pi - \frac{1}{2} \ln |\Sigma_1| - \frac{1}{2} (z - \mu_1)^T \Sigma_1^{-1} (z - \mu_1) \quad (629) \end{aligned}$$

であるので

$$\begin{aligned} & \ln \mathcal{N}(z|\mu_0, \Sigma_0) - \ln \mathcal{N}(z|\mu_1, \Sigma_1) \\ = & -\frac{K}{2} \ln 2\pi - \frac{1}{2} \ln |\Sigma_0| - \frac{1}{2} (z - \mu_0)^T \Sigma_0^{-1} (z - \mu_0) - \end{aligned}$$

変分自己符号化器 (VAE) の理論

$$\begin{aligned} & \left(-\frac{K}{2} \ln 2\pi - \frac{1}{2} \ln |\Sigma_1| - \frac{1}{2} (z - \mu_1)^T \Sigma_1^{-1} (z - \mu_1) \right) \\ = & \frac{1}{2} \ln \frac{|\Sigma_1|}{|\Sigma_0|} - \frac{1}{2} (z - \mu_0)^T \Sigma_0^{-1} (z - \mu_0) + \\ & \frac{1}{2} (z - \mu_1)^T \Sigma_1^{-1} (z - \mu_1) \end{aligned} \quad (630)$$

• よって

$$\begin{aligned} & \text{KL}(p(z) || q(z)) \\ = & \mathbb{E}_{z \sim p(z)} [\ln \mathcal{N}(z | \mu_0, \Sigma_0) - \ln \mathcal{N}(z | \mu_1, \Sigma_1)] \\ = & \mathbb{E}_{p(z)} \left[\frac{1}{2} \ln \frac{|\Sigma_1|}{|\Sigma_0|} - \frac{1}{2} (z - \mu_0)^T \Sigma_0^{-1} (z - \mu_0) + \right. \\ & \left. \frac{1}{2} (z - \mu_1)^T \Sigma_1^{-1} (z - \mu_1) \right] \quad (631) \\ = & \frac{1}{2} \ln \frac{|\Sigma_1|}{|\Sigma_0|} - \frac{1}{2} \mathbb{E}_{p(z)} \left[(z - \mu_0)^T \Sigma_0^{-1} (z - \mu_0) \right] + \end{aligned}$$

変分自己符号化器 (VAE) の理論

$$\frac{1}{2} \mathbb{E}_{p(z)} \left[(z - \mu_1)^T \Sigma_1^{-1} (z - \mu_1) \right] \quad (632)$$

- ここで、期待値についての式を導出しておく
- $\mathbb{E}[z] = \mu$ 、 $\mathbb{E}[(z - \mu)(z - \mu)^T] = \Sigma$ とする
- z の i 成分を z_i 、 μ の i 成分を μ_i 、 Σ の i, j 成分を Σ_{ij} とする
- このとき

$$\begin{aligned} \Sigma_{ij} &= \mathbb{E}[(z_i - \mu_i)(z_j - \mu_j)] \\ &= \mathbb{E}[z_i z_j - z_i \mu_j - z_j \mu_i + \mu_i \mu_j] \\ &= \mathbb{E}[z_i z_j] - \mu_j \mathbb{E}[z_i] - \mu_i \mathbb{E}[z_j] + \mu_i \mu_j \\ &= \mathbb{E}[z_i z_j] - \mu_j \mu_i - \mu_i \mu_j + \mu_i \mu_j \\ &= \mathbb{E}[z_i z_j] - \mu_i \mu_j \end{aligned} \quad (633)$$

であるから

$$\mathbb{E}[z_i z_j] = \Sigma_{ij} + \mu_i \mu_j \quad (634)$$

変分自己符号化器 (VAE) の理論

- そして、行列 A の i, j 成分を A_{ij} とすれば、以下を得る

$$\begin{aligned}\mathbb{E}[z^T A z] &= \mathbb{E}\left[\sum_i \sum_j z_i A_{ij} z_j\right] \\&= \sum_i \sum_j A_{ij} \mathbb{E}[z_i z_j] \\&= \sum_i \sum_j A_{ij} (\Sigma_{ij} + \mu_i \mu_j) \\&= \sum_i \sum_j A_{ij} \Sigma_{ij} + \sum_i \sum_j A_{ij} \mu_i \mu_j \\&= \sum_i \sum_j A_{ij} \Sigma_{ji} + \sum_i \sum_j \mu_i A_{ij} \mu_j \\&= \sum_i (A \Sigma)_{ii} + \mu^T A \mu \\&= \text{Tr}(A \Sigma) + \mu^T A \mu\end{aligned}\tag{635}$$

変分自己符号化器 (VAE) の理論

- 上式の変形では、共分散行列 Σ が対称行列ゆえ、 $\Sigma_{ij} = \Sigma_{ji}$ が成立することを用いた
- また、ベクトル \mathbf{a} の i 成分を a_i とすれば、以下を得る

$$\begin{aligned}\mathbb{E}[\mathbf{a}^T \mathbf{z}] &= \mathbb{E}[\mathbf{z}^T \mathbf{a}] \\ &= \mathbb{E}\left[\sum_i z_i a_i\right] \\ &= \sum_i a_i \mathbb{E}[z_i] \\ &= \sum_i a_i \mu_i \\ &= \mathbf{a}^T \boldsymbol{\mu} = \boldsymbol{\mu}^T \mathbf{a}\end{aligned}\tag{636}$$

変分自己符号化器 (VAE) の理論

- これより、 \mathbf{a}, \mathbf{B} をそれぞれ適当なベクトル、行列とすれば、以下を得る

$$\begin{aligned} & \mathbb{E} \left[(\mathbf{z} - \mathbf{a})^T \mathbf{B} (\mathbf{z} - \mathbf{a}) \right] \\ = & \mathbb{E} \left[\mathbf{z}^T \mathbf{B} \mathbf{z} - \mathbf{z}^T \mathbf{B} \mathbf{a} - \mathbf{a}^T \mathbf{B} \mathbf{z} + \mathbf{a}^T \mathbf{B} \mathbf{a} \right] \\ = & \mathbb{E} \left[\mathbf{z}^T \mathbf{B} \mathbf{z} \right] - \mathbb{E} \left[\mathbf{z}^T \mathbf{B} \mathbf{a} \right] - \mathbb{E} \left[\mathbf{a}^T \mathbf{B} \mathbf{z} \right] + \mathbf{a}^T \mathbf{B} \mathbf{a} \\ = & \left(\text{Tr}(\mathbf{B} \boldsymbol{\Sigma}) + \boldsymbol{\mu}^T \mathbf{B} \boldsymbol{\mu} \right) - \boldsymbol{\mu}^T \mathbf{B} \mathbf{a} - \mathbf{a}^T \mathbf{B} \boldsymbol{\mu} + \mathbf{a}^T \mathbf{B} \mathbf{a} \\ = & \text{Tr}(\mathbf{B} \boldsymbol{\Sigma}) + \boldsymbol{\mu}^T \mathbf{B} \boldsymbol{\mu} - 2\boldsymbol{\mu}^T \mathbf{B} \mathbf{a} + \mathbf{a}^T \mathbf{B} \mathbf{a} \end{aligned} \quad (637)$$

特に、 $\mathbf{a} = \boldsymbol{\mu}, \mathbf{B} = \boldsymbol{\Sigma}^{-1}$ とすれば

$$\begin{aligned} & \mathbb{E} \left[(\mathbf{z} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1} (\mathbf{z} - \boldsymbol{\mu}) \right] \\ = & \text{Tr}(\boldsymbol{\Sigma}^{-1} \boldsymbol{\Sigma}) + \boldsymbol{\mu}^T \boldsymbol{\Sigma}^{-1} \boldsymbol{\mu} - 2\boldsymbol{\mu}^T \boldsymbol{\Sigma}^{-1} \boldsymbol{\mu} + \boldsymbol{\mu}^T \boldsymbol{\Sigma}^{-1} \boldsymbol{\mu} \\ = & \text{Tr}(\boldsymbol{\Sigma}^{-1} \boldsymbol{\Sigma}) \\ = & \text{Tr}(\mathbf{I}) = K \end{aligned} \quad (638)$$

変分自己符号化器 (VAE) の理論

- これをを用いれば、KL ダイバージェンスは次のようになる

$$\begin{aligned} & \text{KL}(p(\mathbf{z})||q(\mathbf{z})) \\ = & \frac{1}{2} \ln \frac{|\Sigma_1|}{|\Sigma_0|} - \frac{1}{2} \mathbb{E}_{p(\mathbf{z})} \left[(\mathbf{z} - \boldsymbol{\mu}_0)^T \Sigma_0^{-1} (\mathbf{z} - \boldsymbol{\mu}_0) \right] + \\ & \frac{1}{2} \mathbb{E}_{p(\mathbf{z})} \left[(\mathbf{z} - \boldsymbol{\mu}_1)^T \Sigma_1^{-1} (\mathbf{z} - \boldsymbol{\mu}_1) \right] \\ = & \frac{1}{2} \ln \frac{|\Sigma_1|}{|\Sigma_0|} - \frac{1}{2} K + \frac{1}{2} \left(\text{Tr}(\Sigma_1^{-1} \Sigma_0) + \boldsymbol{\mu}_0^T \Sigma_1^{-1} \boldsymbol{\mu}_0 - \right. \\ & \left. 2\boldsymbol{\mu}_0^T \Sigma_1^{-1} \boldsymbol{\mu}_1 + \boldsymbol{\mu}_1^T \Sigma_1^{-1} \boldsymbol{\mu}_1 \right) \\ = & \frac{1}{2} \ln \frac{|\Sigma_1|}{|\Sigma_0|} - \frac{1}{2} K + \frac{1}{2} \text{Tr}(\Sigma_1^{-1} \Sigma_0) + \\ & \frac{1}{2} (\boldsymbol{\mu}_0 - \boldsymbol{\mu}_1)^T \Sigma_1^{-1} (\boldsymbol{\mu}_0 - \boldsymbol{\mu}_1) \\ = & \frac{1}{2} \left(\ln \frac{|\Sigma_1|}{|\Sigma_0|} - K + \text{Tr}(\Sigma_1^{-1} \Sigma_0) + \right. \end{aligned}$$

変分自己符号化器 (VAE) の理論

$$(\boldsymbol{\mu}_0 - \boldsymbol{\mu}_1)^T \boldsymbol{\Sigma}_1^{-1} (\boldsymbol{\mu}_0 - \boldsymbol{\mu}_1) \Big) \quad (639)$$

- これで、2つのガウス分布 $p(\mathbf{z}) = \mathcal{N}(\mathbf{z}|\boldsymbol{\mu}_0, \boldsymbol{\Sigma}_0)$ 、 $q(\mathbf{z}) = \mathcal{N}(\mathbf{z}|\boldsymbol{\mu}_1, \boldsymbol{\Sigma}_1)$ 間の KL ダイバージェンスが、次のようになることが分かった

$$\begin{aligned} & \text{KL}(p(\mathbf{z})||q(\mathbf{z})) \\ &= \frac{1}{2} \left(\ln \frac{|\boldsymbol{\Sigma}_1|}{|\boldsymbol{\Sigma}_0|} - K + \text{Tr}(\boldsymbol{\Sigma}_1^{-1} \boldsymbol{\Sigma}_0) + \right. \\ & \quad \left. (\boldsymbol{\mu}_0 - \boldsymbol{\mu}_1)^T \boldsymbol{\Sigma}_1^{-1} (\boldsymbol{\mu}_0 - \boldsymbol{\mu}_1) \right) \end{aligned} \quad (640)$$

- ここでは、 $p(\mathbf{z}) = \mathcal{N}(\mathbf{z}|\boldsymbol{\mu}_0, \boldsymbol{\Sigma}_0)$ は、Encoder のニューラルネットが表現する分布 $q(\mathbf{z}|\mathbf{x}, \phi)$ である
- そして、分布 $q(\mathbf{z}|\mathbf{x}, \phi)$ はガウス分布であったので、ここでは平均 $\boldsymbol{\mu}_0$ と共分散行列 $\boldsymbol{\Sigma}_0$ を使って、 $q(\mathbf{z}|\mathbf{x}, \phi) = \mathcal{N}(\mathbf{z}|\boldsymbol{\mu}_0, \boldsymbol{\Sigma}_0)$ と表すことにする

変分自己符号化器 (VAE) の理論

- また $q(z) = \mathcal{N}(z|\mu_1, \Sigma_1)$ は、潜在変数に関する事後分布 $p(z|\theta) = \mathcal{N}(z|\mathbf{0}, I)$ である
- 結局、KL ダイバージェンス $\text{KL}(q(z|x, \phi)||p(z|\theta))$ は、先程の式に $\mu_1 = \mathbf{0}$ と $\Sigma_1 = I$ を代入すれば得られる

$$\begin{aligned} & \text{KL}(q(z|x, \phi)||p(z|\theta)) \\ &= \text{KL}(\mathcal{N}(z|\mu_0, \Sigma_0)||\mathcal{N}(z|\mathbf{0}, I)) \\ &= \frac{1}{2} (-\ln |\Sigma_0| - K + \text{Tr}(\Sigma_0) + \mu_0^T \mu_0) \end{aligned} \quad (641)$$

- この KL ダイバージェンスが、Encoder の損失関数として定義される
- Encoder のニューラルネットは、入力としてデータ x を取り、平均 μ_0 と共分散行列 Σ_0 を出力する
- 従って、Encoder の出力と、潜在変数の次元 K を上式に代入すれば、損失関数を容易に計算できる

変分自己符号化器 (VAE) の理論

- Σ_0 は対称行列であるため、実際に出力されるのは、 Σ_0 の対角成分を並べたベクトルである
- 一般的な VAE の Encoder は次の図 17 のように表せる

変分自己符号化器 (VAE) の理論

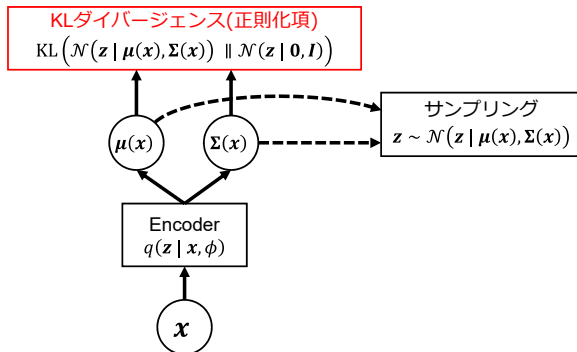


図 17: VAE の Encoder の概要

変分自己符号化器 (VAE) の理論

● Encoder のニューラルネットの構造

- VAE が最初に提案された論文では、隠れ層は 1 層となっている
- 最終層は、 μ_0 と $\sqrt{\Sigma_0}$ を出力する 2 つのユニットから成る
- $\sqrt{\Sigma_0}$ とは、行列 Σ_0 の各要素の平方根を取った行列である
- 隠れ層の重みを W_h 、バイアスを b_h 、活性化関数を $f(\cdot)$ 、層の出力を h とする
- 隠れ層で行う処理は、次の式で表される

$$h = f(W_h x + b_h) \quad (642)$$

- μ_0 は、重み W_m とバイアス b_m を使って、以下のように計算される

$$\mu_0 = W_m h + b_m \quad (643)$$

- $\sqrt{\Sigma_0}$ は、重み W_s とバイアス b_s から、以下のように計算される

$$\sqrt{\Sigma_0} = W_s h + b_s \quad (644)$$

変分自己符号化器 (VAE) の理論

- Decoder のニューラルネットの入出力
 - Decoder は、分布 $p(x|z, \theta)$ を表現するニューラルネット
 - 潜在変数 z から元のデータ x を復元する
- Decoder の入力は、Encoder によってサンプリングされた z となる
- もう少し正確に表現すると、Encoder から出力されるのは、分布のパラメータ μ, Σ である
- そして、そのパラメータを使って分布 $\mathcal{N}(z|\mu, \Sigma)$ を構成し、分布から z をサンプリングする
- Decoder の出力は、復元されたデータ y である

変分自己符号化器 (VAE) の理論

● Decoder の損失関数

- 画像データでは通常、各ピクセルの値が 0 から 1 までになるようにスケーリングされている
- このとき、 $p(x|z)$ はベルヌーイ分布と仮定していることになる
- 出力層のユニット j は、 $y_j = p(x_j = 1|z)$ を出力しているとみなせる
- y_j は、再構成 y の j 番目の要素であり、元のデータ x の j 番目の要素 x_j と対応する
- VAE が最初に提案された論文では、隠れ層は 1 層のみである
- 再構成 y は、次のように計算される

$$y = f_{\sigma}(W_o \tanh(W_h z + b_h) + b_o) \quad (645)$$

- $f_{\sigma}(\cdot)$ は、行列の各要素にシグモイド関数 $\sigma(\cdot)$ を適用する活性化関数
- W_h, b_h は隠れ層の重みとバイアス、 W_o, b_o は出力層の重みとバイアス

変分自己符号化器 (VAE) の理論

- このとき $\ln p(\mathbf{x}|\mathbf{z})$ は次のように記述できる

$$\begin{aligned}\ln p(\mathbf{x}|\mathbf{z}) &= \ln \prod_{j=1}^D (p(x_j = 1|\mathbf{z}))^{x_j} (p(x_j = 0|\mathbf{z}))^{1-x_j} \\ &= \ln \prod_j (p(x_j = 1|\mathbf{z}))^{x_j} (1 - p(x_j = 1|\mathbf{z}))^{1-x_j} \\ &= \ln \prod_j y_j^{x_j} (1 - y_j)^{1-x_j} \\ &= \sum_j (x_j \ln y_j + (1 - x_j) \ln (1 - y_j))\end{aligned}\tag{646}$$

- \mathbf{z} は、分布 $q(\mathbf{z}|\mathbf{x}, \phi)$ からサンプリングされている
- 従って、上記を最大化することは、 $\mathbb{E}_{q(\mathbf{z})} [\ln p(\mathbf{x}|\mathbf{z})]$ を最大化することに等しい

変分自己符号化器 (VAE) の理論

- 上記は、 x_j と y_j のいずれもベルヌーイ分布に従う (二値変数) ときの、**負の交差エントロピー**となっていることが分かる
- 従って、 $\mathbb{E}_{q(z)} [\ln p(x|z)]$ を最大化することは、**交差エントロピーを最小化**することに相当
- Decoder の損失関数は、以下の**交差エントロピー誤差**として定義できる

$$E = - \sum_j (x_j \ln y_j + (1 - x_j) \ln (1 - y_j)) \quad (647)$$

- 元データ x_j と、その再構成 y_j との**差が大きければ大きいほど、上記の誤差は増大**する
- これより、上記の誤差は**再構成誤差** (Reconstruction Error) とよばれる
- これより、VAE の Encoder と Decoder は次の図 18 のように表せる

変分自己符号化器 (VAE) の理論

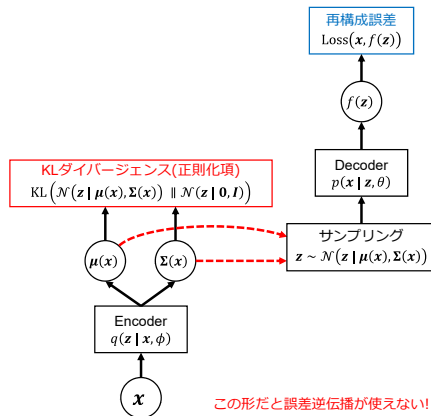


図 18: VAE の Encoder と Decoder の概要

変分自己符号化器 (VAE) の理論

● Reparameterization Trick

- VAE が先程の図 18 のようであるとき、**大きな問題が生じる**
- サンプルを行うと、計算グラフが途中で途切れるため、**誤差逆伝播法を実行できない**
- そこで、次の図 19 のように構成する
- $z \sim (z|\mu, \Sigma)$ として、 z を分布から直接サンプリングするのではない
- z を、**決定論的な関数** $g(\epsilon, x|\phi)$ から決定する
- 但し、 ϵ は、分布 $p(\epsilon)$ からサンプリングされる
- ニューラルネットの最適化には無関係な項 ϵ と、Encoder のパラメータ ϕ で z を表現することで、誤差逆伝播法を実行可能にする
- このテクニックを、**Reparameterization Trick** という
- $\epsilon \sim (\epsilon|0, I)$ とすれば、 z は次のように計算できる

$$z = g(\epsilon, x|\phi) = \mu + \Sigma^{\frac{1}{2}} \epsilon \quad (648)$$

変分自己符号化器 (VAE) の理論

- 共分散行列 Σ が対角行列であれば、上記の $\Sigma^{\frac{1}{2}}\epsilon$ は、単なる要素ごとの積 (対角行列の各要素を並べたベクトルと、 ϵ の要素ごとの積) として書ける
- z の式は以下のように導出できる
- 確率変数 z, ϵ 間の関係が、次のようになっているとする

$$z = \mu + U\Lambda^{\frac{1}{2}}\epsilon \quad (649)$$

- 但し、正定値対称行列 Σ が、固有値分解によって $\Sigma = U\Lambda U^T = U\Lambda^{\frac{1}{2}}(U\Lambda)^T$ と表せるとする
- U は固有ベクトルを並べた行列、 Λ は対角成分に固有値をもつ対角行列とする

変分自己符号化器 (VAE) の理論

- 確率分布 $p(\mathbf{z})$ と $p(\boldsymbol{\epsilon})$ との関係は、ヤコビ行列 $\mathbf{J} = \partial \boldsymbol{\epsilon} / \partial \mathbf{z}$ により次のように記述できる

$$p(\mathbf{z}) = p(\boldsymbol{\epsilon}) |\det(\mathbf{J})| = p(\boldsymbol{\epsilon}) \left| \det \left(\frac{\partial \mathbf{z}}{\partial \boldsymbol{\epsilon}} \right) \right| \quad (650)$$

- ヤコビ行列 \mathbf{J} を計算すると次のようになる

$$\mathbf{J} = \frac{\partial \boldsymbol{\epsilon}}{\partial \mathbf{y}} = \frac{\partial}{\partial \mathbf{z}} \left(\left(\mathbf{U} \boldsymbol{\Lambda}^{\frac{1}{2}} \right)^{-1} (\mathbf{z} - \boldsymbol{\mu}) \right) = \left(\left(\mathbf{U} \boldsymbol{\Lambda}^{\frac{1}{2}} \right)^{-1} \right)^T \quad (651)$$

変分自己符号化器 (VAE) の理論

- ヤコビ行列 J の行列式は次のようになる

$$\begin{aligned}\det(J) &= \left| \left((U\Lambda^{\frac{1}{2}})^{-1} \right)^T \right| \\ &= \left| (U\Lambda^{\frac{1}{2}})^{-1} \right| \quad (\because |A^T| = |A|) \\ &= \frac{1}{|U\Lambda^{\frac{1}{2}}|} \quad \left(\because |A^{-1}| = \frac{1}{|A|} \right)\end{aligned}\tag{652}$$

- Σ の行列式は次のように表せる

$$|\Sigma| = \left| U\Lambda^{\frac{1}{2}} (U\Lambda^{\frac{1}{2}})^T \right| = |U\Lambda^{\frac{1}{2}}| \left| (U\Lambda^{\frac{1}{2}})^T \right| = |U\Lambda^{\frac{1}{2}}|^2 \tag{653}$$

- 従って $|U\Lambda^{\frac{1}{2}}| = |\Sigma|^{\frac{1}{2}}$ である

変分自己符号化器 (VAE) の理論

- Σ は正定値である ($|\Sigma| > 0$) から、 $|\Sigma|^{\frac{1}{2}} = |U\Lambda^{\frac{1}{2}}| > 0$ が成立し、 $U\Lambda^{\frac{1}{2}}$ も正定値行列となる
- これより、逆行列 $(U\Lambda^{\frac{1}{2}})^{-1}$ が存在するので、ヤコビ行列 J は計算できることが確認される
- ヤコビ行列 J の行列式の絶対値は、次のようになる

$$|\det(J)| = \left| \frac{1}{|U\Lambda^{\frac{1}{2}}|} \right| = \left| \frac{1}{|\Sigma|^{\frac{1}{2}}} \right| = \frac{1}{|\Sigma|^{\frac{1}{2}}} \quad (654)$$

- $p(\epsilon)$ がガウス分布 $\mathcal{N}(\epsilon|\mathbf{0}, I)$ であるとする、 $p(z)$ は次のようになる

$$\begin{aligned} p(z) &= p(\epsilon) \left| \det \left(\frac{\partial z}{\partial \epsilon} \right) \right| \\ &= \frac{1}{(2\pi)^{\frac{K}{2}}} \exp \left(-\frac{1}{2} \epsilon^T \epsilon \right) \frac{1}{|\Sigma|^{\frac{1}{2}}} \end{aligned} \quad (655)$$

変分自己符号化器 (VAE) の理論

- ここで、指数関数の中身は次のように書ける

$$\begin{aligned} & (z - \mu)^T \Sigma^{-1} (z - \mu) \\ = & \left(U \Lambda^{\frac{1}{2}} \epsilon \right)^T (U \Lambda U^T)^{-1} \left(U \Lambda^{\frac{1}{2}} \epsilon \right) \\ = & \epsilon^T \left(\Lambda^{\frac{1}{2}} \right)^T U^T (U^T)^{-1} \Lambda^{-1} U^{-1} U \Lambda^{\frac{1}{2}} \epsilon \\ = & \epsilon^T \left(\Lambda^{\frac{1}{2}} \right)^T \Lambda^{-1} \Lambda^{\frac{1}{2}} \epsilon \\ = & \epsilon^T \Lambda^{\frac{1}{2}} \Lambda^{-1} \Lambda^{\frac{1}{2}} \epsilon \\ = & \epsilon^T \epsilon \end{aligned} \tag{656}$$

変分自己符号化器 (VAE) の理論

- これより、 $p(\mathbf{z})$ は平均 $\boldsymbol{\mu}$ 、共分散行列 $\boldsymbol{\Sigma}$ のガウス分布である

$$\begin{aligned} p(\mathbf{z}) &= \frac{1}{(2\pi)^{\frac{K}{2}}} \exp\left(-\frac{1}{2}\boldsymbol{\epsilon}^T \boldsymbol{\epsilon}\right) \frac{1}{|\boldsymbol{\Sigma}|^{\frac{1}{2}}} \\ &= \frac{1}{(2\pi)^{\frac{K}{2}}} \frac{1}{|\boldsymbol{\Sigma}|^{\frac{1}{2}}} \exp\left(-\frac{1}{2}(\mathbf{z} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1}(\mathbf{z} - \boldsymbol{\mu})\right) \quad (657) \end{aligned}$$

$$= \mathcal{N}(\mathbf{z}|\boldsymbol{\mu}, \boldsymbol{\Sigma}) \quad (658)$$

- $\mathbf{z} \sim \mathcal{N}(\mathbf{z}|\boldsymbol{\mu}, \boldsymbol{\Sigma})$ について、共分散行列 $\boldsymbol{\Sigma}$ が既に対角行列であれば、 $\mathbf{U} = \mathbf{I}$ 、 $\boldsymbol{\Lambda} = \boldsymbol{\Sigma}$ であるので、結局以下が言える

$$\mathbf{z} = \boldsymbol{\mu} + \mathbf{U}\boldsymbol{\Lambda}^{\frac{1}{2}}\boldsymbol{\epsilon} = \boldsymbol{\mu} + \boldsymbol{\Sigma}^{\frac{1}{2}}\boldsymbol{\epsilon} \quad (659)$$

変分自己符号化器 (VAE) の理論

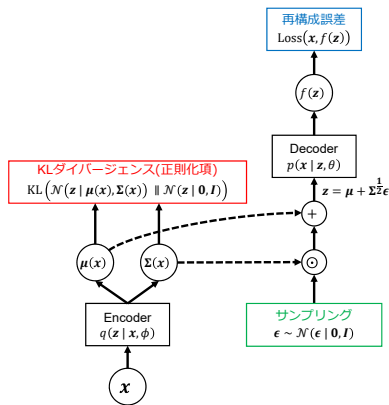


図 19: VAE の構造