

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 2, 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

目次

- 1 K-Means 法
- 2 混合ガウス分布
- 3 EM アルゴリズム

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

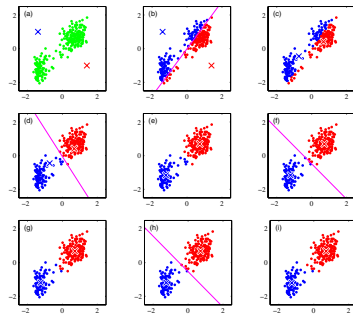


Figure 9.1 Illustration of the K -means algorithm using the re-scaled Old Faithful data set. (a) Green points denote the data set in a two-dimensional Euclidean space. The initial choices for centres μ_1 and μ_2 are shown by the red and blue crosses, respectively. (b) In the initial E step, each data 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.

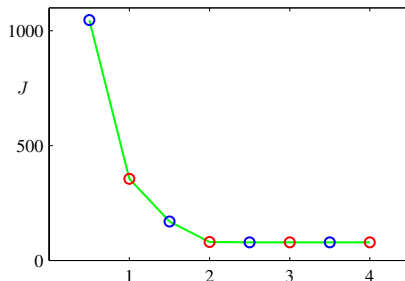


図 2: 目的関数 J の値の遷移

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.

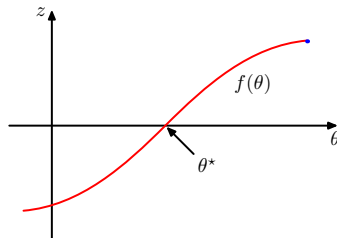


図 3: $f(\theta)$ のグラフ

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

目次

- 1 K-Means 法
- 2 混合ガウス分布
- 3 EM アルゴリズム

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
 - \mathbf{x} の N 個の観測点で構成されるデータ集合 $\mathcal{D} = \{\mathbf{x}_1, \dots, \mathbf{x}_N\}$ ($\mathbf{x}_i \in \mathbb{R}^D$)
 - データ集合 $\mathcal{D} = \{\mathbf{x}_1, \dots, \mathbf{x}_N\}$ 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 $\boldsymbol{\mu}_k$ 、共分散行列 $\boldsymbol{\Sigma}_k$ の正規分布 $\mathcal{N}(\mathbf{x}|\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)$ で表現できる**とする
 - 各クラスタの分布 $\mathcal{N}(\mathbf{x}|\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)$ 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

混合ガウス分布

- 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

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

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

- 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

- $\sum_{\mathbf{z}}$ 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、 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

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

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

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

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

混合ガウス分布

- 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

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

混合ガウス分布

- 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**

目次

- 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

- 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 X と、潜在変数 Z の同時分布 $p(X, Z|\theta)$ が与えられているとする
- 目的は、尤度関数 $p(X|\theta)$ 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 $p(Z|X, \theta^{\text{old}})$ 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

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

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

- これからの話の流れ

混合ガウス分布の再解釈

- 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

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

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

K-Means 法との関連

- 負担率は、以下のように変形できる

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

- 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

- 1 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

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
- 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 $\ln p(\mathbf{X}|\theta)$.

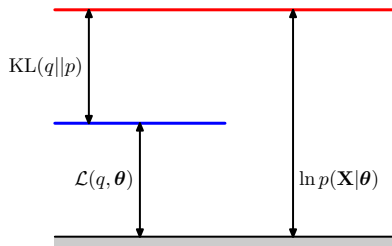


図 4: $\ln p(\mathbf{X}|\theta)$ の分解

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

$$\begin{aligned} & \ln p(\mathbf{X}|\boldsymbol{\theta}^{\text{old}}) \\ = & \mathcal{L}(q, \boldsymbol{\theta}^{\text{old}}) + \text{KL}(q||p) \end{aligned} \tag{175}$$

$$\begin{aligned} = & \mathcal{L}(q, \boldsymbol{\theta}^{\text{old}}) + \text{KL}(q(\mathbf{Z})||p(\mathbf{Z}|\mathbf{X}, \boldsymbol{\theta}^{\text{old}})) \\ & (\text{E ステップ前}) \end{aligned} \tag{176}$$

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

$$\ln p(\mathbf{X}|\boldsymbol{\theta}^{\text{old}}) \tag{177}$$

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

$$\begin{aligned} &= \mathcal{L}(q, \boldsymbol{\theta}^{\text{old}}) \\ &\quad (\text{E ステップ後}) \end{aligned} \tag{179}$$

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 参照）で示されている下界 $\mathcal{L}(q, \theta^{\text{old}})$ が、赤線（図 5 参照）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

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.

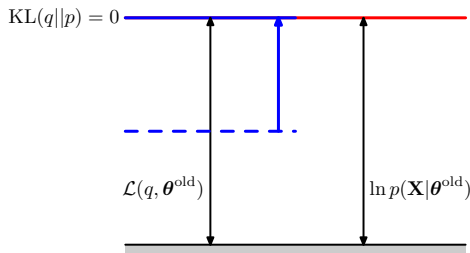


図 5: EM アルゴリズムの E ステップ

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

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

$$= \mathcal{L}(q, \theta) + \text{KL}(q(\mathbf{Z})||p(\mathbf{Z}|\mathbf{X}, \theta)) \tag{181}$$

$$\begin{aligned} = & \mathcal{L}(q, \theta) + \text{KL}(p(\mathbf{Z}|\mathbf{X}, \theta^{\text{old}})||p(\mathbf{Z}|\mathbf{X}, \theta)) \\ & (\text{M ステップ前}) \end{aligned} \tag{182}$$

- 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

- $\text{KL}(q||p)$ 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、 $\text{KL}(q||p) > 0$ となる

$$\begin{aligned} & \ln p(\mathbf{X}|\boldsymbol{\theta}^{\text{new}}) \\ = & \mathcal{L}(q, \boldsymbol{\theta}^{\text{new}}) + \text{KL}(p(\mathbf{Z}|\mathbf{X}, \boldsymbol{\theta}^{\text{old}})||p(\mathbf{Z}|\mathbf{X}, \boldsymbol{\theta}^{\text{new}})) \quad (183) \\ & (\text{M ステップ後}) \end{aligned}$$

- 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)

$$\begin{aligned} & \ln p(\mathbf{X}|\boldsymbol{\theta}^{\text{new}}) - \ln p(\mathbf{X}|\boldsymbol{\theta}^{\text{old}}) \\ = & \mathcal{L}(q, \boldsymbol{\theta}^{\text{new}}) + \text{KL}(p(\mathbf{Z}|\mathbf{X}, \boldsymbol{\theta}^{\text{old}})||p(\mathbf{Z}|\mathbf{X}, \boldsymbol{\theta}^{\text{new}})) - \\ & \mathcal{L}(q, \boldsymbol{\theta}^{\text{old}}) \quad (184) \end{aligned}$$

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

$$\geq \mathcal{L}(q, \boldsymbol{\theta}^{\text{new}}) - \mathcal{L}(q, \boldsymbol{\theta}^{\text{old}}) \quad (186)$$

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.

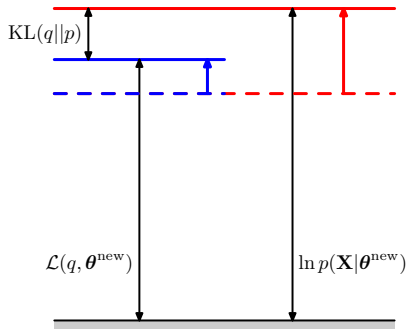


図 6: EM アルゴリズムの M ステップ

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)

$$\begin{aligned}\mathcal{L}(q, \theta) &= \sum_{\mathbf{Z}} q(\mathbf{Z}) \ln \frac{p(\mathbf{X}, \mathbf{Z}|\theta)}{q(\mathbf{Z})} \\ &= \sum_{\mathbf{Z}} p(\mathbf{X}, \mathbf{Z}|\theta^{\text{old}}) \ln \frac{p(\mathbf{X}, \mathbf{Z}|\theta)}{p(\mathbf{X}, \mathbf{Z}|\theta^{\text{old}})} \quad (187)\end{aligned}$$

$$\begin{aligned}&= \sum_{\mathbf{Z}} p(\mathbf{X}, \mathbf{Z}|\theta^{\text{old}}) \ln p(\mathbf{X}, \mathbf{Z}|\theta) - \\ &\quad \sum_{\mathbf{Z}} p(\mathbf{X}, \mathbf{Z}|\theta^{\text{old}}) \ln p(\mathbf{X}, \mathbf{Z}|\theta^{\text{old}}) \quad (188)\end{aligned}$$

$$= \mathcal{Q}(\theta, \theta^{\text{old}}) + \text{Const.} \quad (189)$$

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

$$\begin{aligned} &= \sum_{\mathbf{Z}} p(\mathbf{Z}|\mathbf{X}, \boldsymbol{\theta}^{\text{old}}) \ln \frac{p(\mathbf{X}, \mathbf{Z}|\boldsymbol{\theta}^{\text{new}})}{p(\mathbf{X}, \mathbf{Z}|\boldsymbol{\theta}^{\text{old}})} + \\ &\quad \sum_{\mathbf{Z}} p(\mathbf{Z}|\mathbf{X}, \boldsymbol{\theta}^{\text{old}}) \ln \frac{p(\mathbf{Z}|\mathbf{X}, \boldsymbol{\theta}^{\text{old}})}{p(\mathbf{Z}|\mathbf{X}, \boldsymbol{\theta}^{\text{new}})} \\ &= \sum_{\mathbf{Z}} p(\mathbf{Z}|\mathbf{X}, \boldsymbol{\theta}^{\text{old}}) \ln p(\mathbf{X}, \mathbf{Z}|\boldsymbol{\theta}^{\text{new}}) - \\ &\quad \sum_{\mathbf{Z}} p(\mathbf{Z}|\mathbf{X}, \boldsymbol{\theta}^{\text{old}}) \ln p(\mathbf{X}, \mathbf{Z}|\boldsymbol{\theta}^{\text{old}}) - \\ &\quad \sum_{\mathbf{Z}} p(\mathbf{Z}|\mathbf{X}, \boldsymbol{\theta}^{\text{old}}) \ln \frac{p(\mathbf{Z}|\mathbf{X}, \boldsymbol{\theta}^{\text{old}})}{p(\mathbf{Z}|\mathbf{X}, \boldsymbol{\theta}^{\text{new}})} \\ &= \mathcal{Q}(\boldsymbol{\theta}^{\text{old}}, \boldsymbol{\theta}^{\text{new}}) - \mathcal{Q}(\boldsymbol{\theta}^{\text{old}}, \boldsymbol{\theta}^{\text{old}}) + \\ &\quad \text{KL} (p(\mathbf{Z}|\mathbf{X}, \boldsymbol{\theta}^{\text{old}}) || p(\mathbf{Z}|\mathbf{X}, \boldsymbol{\theta}^{\text{new}})) \end{aligned} \tag{190}$$

$$\geq \mathcal{Q}(\boldsymbol{\theta}^{\text{old}}, \boldsymbol{\theta}^{\text{new}}) - \mathcal{Q}(\boldsymbol{\theta}^{\text{old}}, \boldsymbol{\theta}^{\text{old}}) \tag{191}$$

$$\geq 0$$

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

$$\begin{aligned}\ln p(\mathbf{X}|\boldsymbol{\theta}) &= \ln \sum_{\mathbf{Z}} p(\mathbf{X}, \mathbf{Z}|\boldsymbol{\theta}) \\ &= \ln \sum_{\mathbf{Z}} q(\mathbf{Z}) \frac{p(\mathbf{X}, \mathbf{Z}|\boldsymbol{\theta})}{q(\mathbf{Z})} \\ &\geq \sum_{\mathbf{Z}} q(\mathbf{Z}) \ln \frac{p(\mathbf{X}, \mathbf{Z}|\boldsymbol{\theta})}{q(\mathbf{Z})} \\ &= \mathcal{L}(q, \boldsymbol{\theta})\end{aligned}\tag{192}$$

- 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

$$\begin{aligned} &= - \sum_{\mathbf{Z}} q(\mathbf{Z}) \ln \frac{p(\mathbf{Z}|\mathbf{X}, \boldsymbol{\theta})}{q(\mathbf{Z})} \\ &= \text{KL}(q||p) \end{aligned} \tag{193}$$

ゆえ、 $\text{KL}(q||p)$ となることが分かったので

$$\ln p(\mathbf{X}|\boldsymbol{\theta}) = \mathcal{L}(q, \boldsymbol{\theta}) + \text{KL}(q||p) \tag{194}$$

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

$$\begin{aligned} & \left. \frac{\partial}{\partial \theta} \ln p(\mathbf{X}|\theta) \right|_{\theta^{\text{old}}} \\ &= \left. \frac{\partial}{\partial \theta} \mathcal{L}(q, \theta) \right|_{\theta^{\text{old}}} + \left. \frac{\partial}{\partial \theta} \text{KL}(q||p) \right|_{\theta^{\text{old}}} \\ &= \left. \frac{\partial}{\partial \theta} \mathcal{L}(q, \theta) \right|_{\theta^{\text{old}}} \end{aligned} \tag{195}$$

- 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.

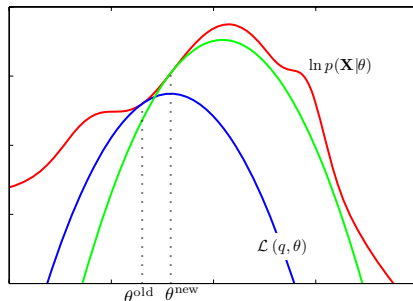


図 7: EM アルゴリズムの手続き

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

$$p(\mathbf{X}, \mathbf{Z}) = \prod_i p(x_i, z_i) \quad (196)$$

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

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

$$\begin{aligned} &= \frac{\prod_i p(\mathbf{x}_i, \mathbf{z}_i | \boldsymbol{\theta})}{\prod_i \sum_{\mathbf{z}} p(\mathbf{x}_i, \mathbf{z}_i | \boldsymbol{\theta})} \\ &= \frac{\prod_i p(\mathbf{x}_i, \mathbf{z}_i | \boldsymbol{\theta})}{\prod_i p(\mathbf{x}_i | \boldsymbol{\theta})} \\ &= \prod_i \frac{p(\mathbf{x}_i, \mathbf{z}_i | \boldsymbol{\theta})}{p(\mathbf{x}_i | \boldsymbol{\theta})} \\ &= \prod_i p(\mathbf{z}_i | \mathbf{x}_i, \boldsymbol{\theta}) \end{aligned} \tag{197}$$

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

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

$$\boldsymbol{\theta}^* = \arg \max_{\boldsymbol{\theta}} p(\boldsymbol{\theta}|\mathbf{X}) \quad (198)$$

$$= \arg \max_{\boldsymbol{\theta}} \frac{p(\mathbf{X}, \boldsymbol{\theta})}{p(\mathbf{X})} \quad (199)$$

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

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

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

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

$$\begin{aligned} &= \ln \frac{p(\mathbf{X}|\boldsymbol{\theta})p(\boldsymbol{\theta})}{p(\mathbf{X})} \\ &= \ln p(\mathbf{X}|\boldsymbol{\theta}) + \ln p(\boldsymbol{\theta}) - p(\mathbf{X}) \end{aligned} \quad (203)$$

$$= \mathcal{L}(q, \boldsymbol{\theta}) + \text{KL}(q||p) + \ln p(\boldsymbol{\theta}) - p(\mathbf{X}) \quad (204)$$

$$\geq \mathcal{L}(q, \boldsymbol{\theta}) + \ln p(\boldsymbol{\theta}) - p(\mathbf{X}) \quad (205)$$

$$= \mathcal{L}(q, \boldsymbol{\theta}) + \ln p(\boldsymbol{\theta}) + \text{Const.} \quad (206)$$

- $\ln p(\mathbf{X})$ 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

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

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

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

$$p(Z|X, \theta) = \prod_i p(z_i|x_i, \theta) \quad (207)$$

- 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)

$$d = \gamma^{\text{new}}(z_{mk}) - \gamma^{\text{old}}(z_{mk}) \quad (208)$$

- 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)

$$N_k^{\text{new}} = N_k^{\text{old}} + \gamma^{\text{new}}(z_{mk}) - \gamma^{\text{old}}(z_{mk}) = N_k^{\text{old}} + d \quad (209)$$

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

$$\boldsymbol{\mu}_k^{\text{old}} = \frac{1}{N_k^{\text{old}}} \sum_i \gamma^{\text{old}}(z_{ik}) \mathbf{x}_i \quad (210)$$

$$\boldsymbol{\Sigma}_k^{\text{old}} = \frac{1}{N_k^{\text{old}}} \sum_i \gamma^{\text{old}}(z_{ik}) (\mathbf{x}_i - \boldsymbol{\mu}_k^{\text{old}})(\mathbf{x}_i - \boldsymbol{\mu}_k^{\text{old}})^T \quad (211)$$

$$\pi_k^{\text{old}} = \frac{N_k^{\text{old}}}{N} \quad (212)$$

但し

$$N_k^{\text{old}} = \sum_i \gamma^{\text{old}}(z_{ik}) \quad (213)$$

逐次型の EM アルゴリズムの例

- 平均の更新式は

$$\boldsymbol{\mu}_k^{\text{new}} = \frac{1}{N_k^{\text{new}}} \left(\sum_{i \neq m} \gamma^{\text{old}}(z_{ik}) \mathbf{x}_i + \gamma^{\text{new}}(z_{mk}) \mathbf{x}_m \right) \quad (214)$$

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

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

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

$$= \boldsymbol{\mu}_k^{\text{old}} + \frac{\gamma^{\text{new}}(z_{mk}) - \gamma^{\text{old}}(z_{mk})}{N_k^{\text{new}}} (\mathbf{x}_m - \boldsymbol{\mu}_k^{\text{old}}) \quad (215)$$

逐次型の EM アルゴリズムの例

$$= \boldsymbol{\mu}_k^{\text{old}} + \frac{d}{N_k^{\text{new}}} (\mathbf{x}_m - \boldsymbol{\mu}_k^{\text{old}}) \quad (216)$$

- 分散の更新式は

$$\boldsymbol{\Sigma}_k^{\text{new}} = \frac{1}{N_k^{\text{new}}} \sum_i \gamma^{\text{new}}(z_{ik}) \mathbf{x}_i \mathbf{x}_i^T - \boldsymbol{\mu}_k^{\text{new}} (\boldsymbol{\mu}_k^{\text{new}})^T \quad (217)$$

$$= \frac{1}{N_k^{\text{new}}} \left(\sum_{i \neq m} \gamma^{\text{old}}(z_{ik}) \mathbf{x}_i \mathbf{x}_i^T + \gamma^{\text{new}}(z_{mk}) \mathbf{x}_m \mathbf{x}_m^T \right) - \boldsymbol{\mu}_k^{\text{new}} (\boldsymbol{\mu}_k^{\text{new}})^T \quad (218)$$

$$= \frac{1}{N_k^{\text{new}}} \left(\left(\sum_i \gamma^{\text{old}}(z_{ik}) \mathbf{x}_i \mathbf{x}_i^T - \gamma^{\text{old}}(z_{mk}) \mathbf{x}_m \mathbf{x}_m^T \right) + \gamma^{\text{new}}(z_{mk}) \mathbf{x}_m \mathbf{x}_m^T \right) - \boldsymbol{\mu}_k^{\text{new}} (\boldsymbol{\mu}_k^{\text{new}})^T \quad (219)$$

逐次型の EM アルゴリズムの例

$$\begin{aligned} &= \frac{1}{N_k^{\text{new}}} \left(N_k^{\text{old}} \left(\boldsymbol{\Sigma}_k^{\text{old}} + \boldsymbol{\mu}_k^{\text{old}} (\boldsymbol{\mu}_k^{\text{old}})^T \right) - \right. \\ &\quad \left. \gamma^{\text{old}}(z_{mk}) \mathbf{x}_m \mathbf{x}_m^T + \gamma^{\text{new}}(z_{mk}) \mathbf{x}_m \mathbf{x}_m^T \right) - \\ &\quad \boldsymbol{\mu}_k^{\text{new}} (\boldsymbol{\mu}_k^{\text{new}})^T \end{aligned} \quad (220)$$

$$\begin{aligned} &= \frac{1}{N_k^{\text{new}}} \left(N_k^{\text{old}} \boldsymbol{\Sigma}_k^{\text{old}} + N_k^{\text{old}} \boldsymbol{\mu}_k^{\text{old}} (\boldsymbol{\mu}_k^{\text{old}})^T + \right. \\ &\quad \left(\gamma^{\text{new}}(z_{mk}) - \gamma^{\text{old}}(z_{mk}) \right) \mathbf{x}_m \mathbf{x}_m^T - \\ &\quad \left. N_k^{\text{new}} \boldsymbol{\mu}_k^{\text{new}} (\boldsymbol{\mu}_k^{\text{new}})^T \right) \\ &= \frac{1}{N_k^{\text{new}}} \left(N_k^{\text{old}} \boldsymbol{\Sigma}_k^{\text{old}} + N_k^{\text{old}} \boldsymbol{\mu}_k^{\text{old}} (\boldsymbol{\mu}_k^{\text{old}})^T + \right. \\ &\quad \left. d\mathbf{x}_m \mathbf{x}_m^T - N_k^{\text{new}} \boldsymbol{\mu}_k^{\text{new}} (\boldsymbol{\mu}_k^{\text{new}})^T \right) \end{aligned} \quad (221)$$

逐次型の EM アルゴリズムの例

ここで

$$\begin{aligned} & N_k^{\text{new}} \boldsymbol{\mu}_k^{\text{new}} (\boldsymbol{\mu}_k^{\text{new}})^T \\ = & (N_k^{\text{new}} \boldsymbol{\mu}_k^{\text{old}} + d(\mathbf{x}_m - \boldsymbol{\mu}_k^{\text{old}})) (\boldsymbol{\mu}_k^{\text{new}})^T \\ = & \frac{1}{N_k^{\text{new}}} (N_k^{\text{new}} \boldsymbol{\mu}_k^{\text{old}} + d(\mathbf{x}_m - \boldsymbol{\mu}_k^{\text{old}})) \\ & (N_k^{\text{new}} \boldsymbol{\mu}_k^{\text{old}} + d(\mathbf{x}_m - \boldsymbol{\mu}_k^{\text{old}}))^T \\ = & N_k^{\text{new}} \boldsymbol{\mu}_k^{\text{old}} (\boldsymbol{\mu}_k^{\text{old}})^T + 2\boldsymbol{\mu}_k^{\text{old}} d(\mathbf{x}_m - \boldsymbol{\mu}_k^{\text{old}})^T + \\ & \frac{1}{N_k^{\text{new}}} d^2 (\mathbf{x}_m - \boldsymbol{\mu}_k^{\text{old}}) (\mathbf{x}_m - \boldsymbol{\mu}_k^{\text{old}})^T \end{aligned} \quad (222)$$

$$\begin{aligned} = & (N_k^{\text{old}} + d) \boldsymbol{\mu}_k^{\text{old}} (\boldsymbol{\mu}_k^{\text{old}})^T + \\ & 2d \boldsymbol{\mu}_k^{\text{old}} \mathbf{x}_m^T - 2d \boldsymbol{\mu}_k^{\text{old}} (\boldsymbol{\mu}_k^{\text{old}})^T + \\ & \frac{1}{N_k^{\text{new}}} d^2 (\mathbf{x}_m - \boldsymbol{\mu}_k^{\text{old}}) (\mathbf{x}_m - \boldsymbol{\mu}_k^{\text{old}})^T \end{aligned} \quad (223)$$

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

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

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

$$\begin{aligned} &= \frac{N_k^{\text{old}}}{N_k^{\text{new}}} \left(\boldsymbol{\Sigma}_k^{\text{old}} + \frac{d}{N_k^{\text{new}}} \right. \\ &\quad \left. (\mathbf{x}_m - \boldsymbol{\mu}_k^{\text{old}}) (\mathbf{x}_m - \boldsymbol{\mu}_k^{\text{old}})^T \right) \\ &= \frac{N_k^{\text{old}}}{N_k^{\text{new}}} \left(\boldsymbol{\Sigma}_k^{\text{old}} + \frac{\gamma^{\text{new}}(z_{ik}) - \gamma^{\text{old}}(z_{ik})}{N_k^{\text{new}}} \right. \\ &\quad \left. (\mathbf{x}_m - \boldsymbol{\mu}_k^{\text{old}}) (\mathbf{x}_m - \boldsymbol{\mu}_k^{\text{old}})^T \right) \end{aligned} \quad (225)$$

- 混合係数の更新式は

$$\pi_k^{\text{new}} = \frac{N_k^{\text{new}}}{N} = \frac{N_k^{\text{old}} + \gamma^{\text{new}}(z_{mk}) - \gamma^{\text{old}}(z_{mk})}{N} \quad (226)$$

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

$$\mu_k^{\text{new}} = \mu_k^{\text{old}} + \frac{\gamma^{\text{new}}(z_{mk}) - \gamma^{\text{old}}(z_{mk})}{N_k^{\text{new}}} (\mathbf{x}_m - \mu_k^{\text{old}}) \quad (227)$$

$$\Sigma_k^{\text{new}} = \frac{N_k^{\text{old}}}{N_k^{\text{new}}} \left(\Sigma_k^{\text{old}} + \frac{\gamma^{\text{new}}(z_{ik}) - \gamma^{\text{old}}(z_{ik})}{N_k^{\text{new}}} (\mathbf{x}_m - \mu_k^{\text{old}}) (\mathbf{x}_m - \mu_k^{\text{old}})^T \right) \quad (228)$$

$$\pi_k^{\text{new}} = \frac{N_k^{\text{new}}}{N} = \frac{N_k^{\text{old}} + \gamma^{\text{new}}(z_{mk}) - \gamma^{\text{old}}(z_{mk})}{N} \quad (229)$$

但し

$$N_k^{\text{new}} = N_k^{\text{old}} + \gamma^{\text{new}}(z_{mk}) - \gamma^{\text{old}}(z_{mk}) \quad (230)$$

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