

# GMM の学習方法について。

1. GMMとは
2. 学習方法

## 1. GMMとは

混合ガウスモデル, Gaussian Mixture Model.

データが混合ガウス分布で生成したと考える。

与えられたデータに適合する混合ガウス分布のパラメータを求めよう。

### 1.1 混合ガウスモデル

混合分布のひとつ。

各分布も、ガウス分布の線形結合として表現する。

#### 1.1.1 混合分布

データ、分布の形状によっては、それを1つの分布でうまく表現できないことがある。

例えば、2つの分布を複数の分布を重ね合わせる(線形結合)ことで表現するもの。

十分な数のガウス分布を用い、各分布の重みやパラメータ(平均, 共分散)を調節すれば、ほぼ任意の連続密度関数も、任意の精度で近似できる。

$$p(x) = \sum_{k=1}^K \pi_k \mathcal{N}(x | \mu_k, \Sigma_k)$$

ここで、

$$\mathcal{N}(x | \mu_k, \Sigma_k) = \frac{1}{(2\pi)^{D/2}} \frac{1}{|\Sigma_k|^{D/2}} \exp\left\{-\frac{1}{2}(x - \mu_k)^T \Sigma_k^{-1} (x - \mu_k)\right\}$$

は混合要素(mixture component)といい、それぞれ平均  $\mu_k$ 、共分散  $\Sigma_k$  をパラメータとして持つ。

また、パラメータ  $\pi_k$  は混合係数 といふ。

$$\sum_{k=1}^K \pi_k = 1, \quad 0 \leq \pi_k \leq 1.$$

である。

分布のパラメータ

$$\pi = \{\pi_1, \pi_2, \dots, \pi_K\},$$

$$\mu = \{\mu_1, \mu_2, \dots, \mu_K\},$$

$$\Sigma = \{\Sigma_1, \Sigma_2, \dots, \Sigma_K\}$$

のなかでは、最尤推定, EMアルゴリズムが

よく  
→ ch2.

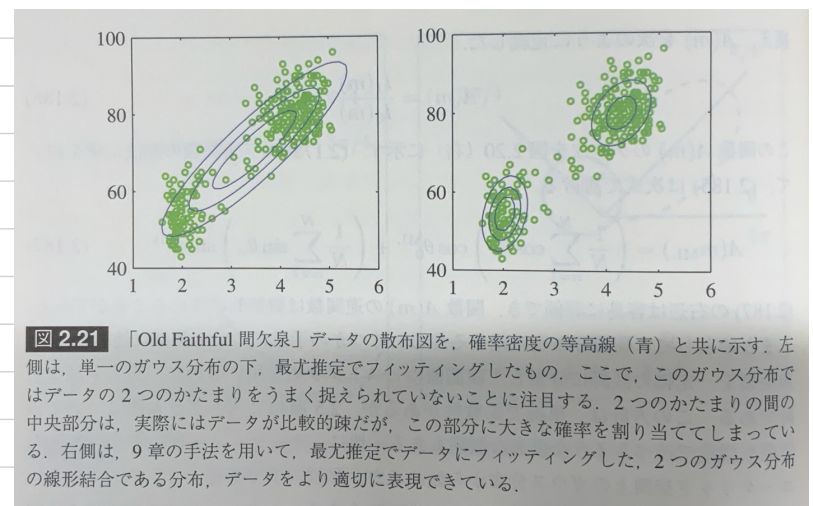
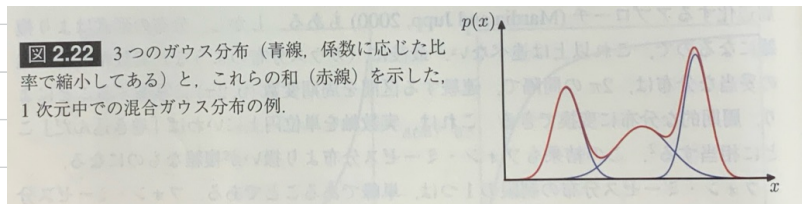


図 2.21 「Old Faithful 間欠泉」データの散布図を、確率密度の等高線(青)と共に示す。左側は、単一のガウス分布の下、最尤推定でフィッティングしたもの。ここで、このガウス分布ではデータの2つのかたまりをうまく捉えられていないことに注目する。2つのかたまりの間の中央部分は、実際にはデータが比較的疎だが、この部分に大きな確率を割り当ててしまっている。右側は、9章の手法を用いて、最尤推定でデータにフィッティングした、2つのガウス分布の線形結合である分布。データをより適切に表現できている。

## 2. 学習方法

### 2.1 分布のパラメータを推定する

#### 2.1.1 最尤推定

#### 2.1.2 ガウス分布の最尤推定

#### 2.1.3 混合ガウス分布の最尤推定

### 2.2 EMアルゴリズム

## 2.1 分布のパラメータを推定する。

分布の形が分かっているとき、入力のデータに最もフィットするようには、分布のパラメータを決めない  
パラメータの決め方は、最尤推定と EM アルゴリズムがある。

### 2.1.1 最尤推定

尤度も最大化する。尤度は積の形であるため、最大値を求めにくい(微分が面倒...)

したがって、尤度の対数をとって、対数尤度を最大化する問題を解く。

解き方はいろいろ。ラグランジュとか。

#### 2.1.1.1 尤度

確率変数  $X$ 。サンプルデータ  $D = \{x_1, x_2, \dots, x_n\}$  について、 $x_i$  が独立に同一の確率分布に従うとき、 $D$  が生成される確率  $P(D)$  は、

$$P(D) = \prod_{x_i \in D} p(x_i)$$

と書ける。この  $P(D)$  を、尤度という。

尤度の対数をとったものは、対数尤度という。

$$\ln P(D) = \ln \prod_{x_i \in D} p(x_i)$$

$$= \sum_{x_i \in D} \ln p(x_i)$$

である。

尤度の最大化と、対数尤度の最大化は等価である。

### 2.1.2 ガウス分布の最尤推定

ある多変量ガウス分布から、観測値  $\{x_n\}$  が独立に得られたと仮定したデータ集合  $X = (x_1, x_2, \dots, x_n)^T$  があるとき、この分布のパラメータは最尤推定で求められる。

対数尤度は

$$\ln p(X | \mu, \Sigma) = -\frac{ND}{2} \ln(2\pi) - \frac{N}{2} \ln |\Sigma| - \frac{1}{2} \sum_{n=1}^N (x_n - \mu)^T \Sigma^{-1} (x_n - \mu)$$

を与えられる。整理すると、尤度関数はデータに対して

$$\sum_{n=1}^N x_n, \sum_{n=1}^N x_n x_n^T$$

のみに依存していることがわかる。この2つを、ガウス分布の十分統計量という。

対数尤度の  $\mu$  についての導関数は、

$$\frac{\partial}{\partial \mu} \ln p(X | \mu, \Sigma) = \sum_{n=1}^N \Sigma^{-1} (x_n - \mu)$$

であり、これを0とすると、平均の最尤値が得られる。

$$\mu_{ML} = \frac{1}{N} \sum_{n=1}^N x_n \quad \text{--- ①}$$

ML: Maximum Likelihood.

共分散については、Magnus and Neudecker (1999) より、

$$\Sigma_{ML} = \frac{1}{N} \sum_{n=1}^N (x_n - \mu_{ML})(x_n - \mu_{ML})^T \quad \text{--- ②}$$

である。 $\mu_{ML}$  が含まれているのは、この式が  $\mu, \Sigma$  について同時に最大化したものであるから。

①より、 $\mu_{ML}$  は  $\Sigma_{ML}$  に依存しないので、①で  $\mu_{ML}$  を求めてから②で  $\Sigma_{ML}$  を求める。

真の分布下で  $\mu_{ML}, \Sigma_{ML}$  の値を評価すると、

$$\mathbb{E}[\mu_{ml}] = \mu$$

$$\mathbb{E}[\Sigma_{ml}] = \frac{N-1}{N} \Sigma$$

となる、 $\Sigma_{ml} < \Sigma$  であるが、これは②式を次のように補正することで、真の値が得られる。

$$\hat{\Sigma} = \frac{1}{N-1} \sum_{n=1}^N (X_n - \mu_{ml})(X_n - \mu_{ml})^T$$

### 2.1.3 混合ガウス分布と最尤推定

混合ガウス分布は、 $k$  個のガウス分布の線形結合で、1次元から

$$\pi = \{\pi_1, \pi_2, \dots, \pi_k\}$$

$$\mu = \{\mu_1, \mu_2, \dots, \mu_k\}$$

$$\Sigma = \{\Sigma_1, \Sigma_2, \dots, \Sigma_k\}$$

であり、定義は

$$P(X) = \sum_{k=1}^k \pi_k \mathcal{N}(X | \mu_k, \Sigma_k)$$

である。  $X = \{x_1, x_2, \dots, x_N\}$  についての対数尤度関数は

$$\ln P(X | \pi, \mu, \Sigma) = \sum_{n=1}^N \ln \left\{ \sum_{k=1}^k \pi_k \mathcal{N}(x_n | \mu_k, \Sigma_k) \right\}$$

となる。対数の内部は  $k$  個の和があり、複雑であり、閉形式の解析解では最尤推定解が得られない。

尤度関数を最大化するアプローチは、繰り返しの数値最適化、すなわち EM アルゴリズムがある。

## 2.2 EM アルゴリズム

EM アルゴリズムは、有名な繰り返しの1次元推定手法。

GMM のおおよそ潜在変数をもつモデルの最尤1次元を求めるエレガントな方法。

### 2.2.1 潜在変数

これまで、入力されたデータは

$$X = \{x_n : n=1, \dots, N\}$$

としていた。GMM では各データ点  $x_n$  がどのガウス分布から生成されたものかわからない。

そこで、潜在変数  $z$  を導入する。

$$Z = \{z_{nk} : n=1, \dots, N \text{ and } k=1, \dots, k\}$$

$$z_{nk} = \begin{cases} 1 & (\text{if } x_n \text{ is generated by } k\text{th Gaussian}) \\ 0 & (\text{else}) \end{cases}$$

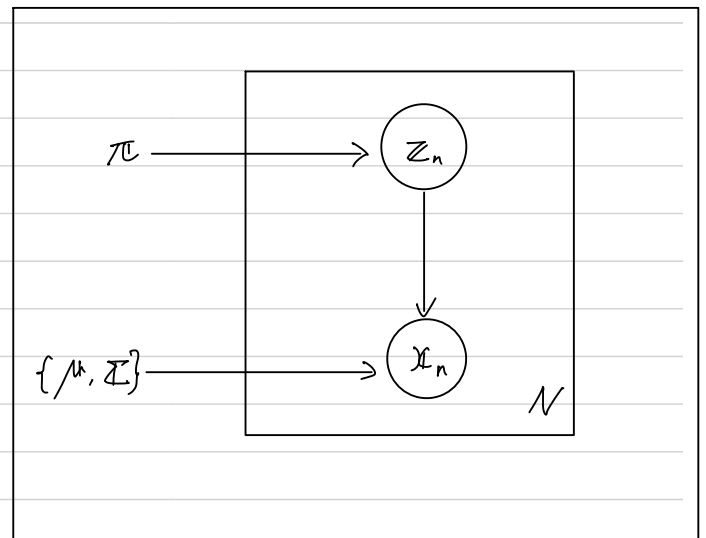
このとき、混合分布は右図のように表せる。

図からわかるように、 $x_n$  は  $z_n$  に依存している。

したがって、周辺尤度  $P(x_n)$  は、潜在変数  $z_n$

marginalizing out することで得られる。

$$\begin{aligned} P(x_n) &= \sum_{z_n} P(z_n) P(x_n | z_n) \\ &= \sum_{k=1}^k P(z_{nk}=1) P(x_n | z_{nk}=1) \\ &= \sum_{k=1}^k \pi_k \mathcal{N}(x_n | \mu_k, \Sigma_k) \end{aligned}$$



$n=1, \dots, N$  について、

$$Z_n = [z_{n1}, \dots, z_{nk}]^T \text{ が潜在変数、}$$

$\mathcal{O}_n$  が観測ベクトル

であることを示す。

これはすなわち、混合ガウス分布は、その潜在変数を陽に含た形で記述できることを示す。

EMにかかるとき、 $\{\mathcal{X}, \mathcal{Z}\}$  は完全データセット、 $\mathcal{X}$  は不完全データセットとよばれる。

EMアルゴリズムでは  $\log p(\mathcal{X}|\Lambda)$  を最大化の代わりに、 $\log p(\mathcal{X}, \mathcal{Z}|\Lambda)$  を最大化する。

重要なのは、 $\log p(\mathcal{X}, \mathcal{Z}|\Lambda)$  の最大化が簡単で、各イテレーションで閉形式の解が得られること。

(各ガウスのパラメータ  $(\mu_k, \Sigma_k)$  と  $\Lambda = \{\mu_k, \Sigma_k | k=1, \dots, K\}$  とする)

(しかし、実際は  $\mathcal{Z}$  も知り得ないため、 $\log p(\mathcal{X}, \mathcal{Z}|\Lambda)$  を計算できない。幸い、この事後分布  $(P(\mathcal{Z}|\mathcal{X}, \Lambda))$  はベイズの定理から導ける)

## 2.2.2 負担率

事後分布  $P(\mathcal{Z}|\mathcal{X}, \Lambda)$ 、特に  $\mathcal{X}_n$  についての事後確率は、

$$\begin{aligned} \gamma(z_{nk}) &= P(z_{nk}=1 | \mathcal{X}_n, \Lambda) \\ &= \frac{P(z_{nk}=1 | \Lambda) p(\mathcal{X}_n | z_{nk}=1, \Lambda)}{P(\mathcal{X}_n | \Lambda)} \\ &= \frac{\pi_k \mathcal{N}(\mathcal{X}_n | \mu_k, \Sigma_k)}{\sum_{j=1}^K \pi_j \mathcal{N}(\mathcal{X}_n | \mu_j, \Sigma_j)} \end{aligned}$$

と書ける。これは、コンポーネント  $k$  が  $\mathcal{X}$  の観測値を「説明する」割合を表す「負担率」として解釈できる。

## 2.2.3 Q関数

$\gamma(z_{nk})$  とこれまで推定したパラメータ  $\Lambda^{old}$  から、

$\mathcal{Z}$  の事後分布の下で  $\log p(\mathcal{X}, \mathcal{Z}|\Lambda)$  の期待値を計算することによって新しいパラメータの推定値  $\Lambda$  を

計算できる。

$$\begin{aligned} Q(\Lambda | \Lambda^{old}) &= \mathbb{E}_{\mathcal{Z}} \{ \log p(\mathcal{X}, \mathcal{Z} | \Lambda) | \mathcal{X}, \Lambda^{old} \} \\ &= \sum_{n=1}^N \sum_{k=1}^K P(z_{nk}=1 | \mathcal{X}_n, \Lambda^{old}) \log p(\mathcal{X}_n, z_{nk}=1 | \Lambda) \\ &= \sum_{n=1}^N \sum_{k=1}^K \gamma(z_{nk}) \log p(\mathcal{X}_n, z_{nk}=1 | \Lambda) \\ &= \sum_{n=1}^N \sum_{k=1}^K \gamma(z_{nk}) \log p(\mathcal{X}_n | z_{nk}=1, \Lambda) P(z_{nk}=1 | \Lambda) \\ &= \sum_{n=1}^N \sum_{k=1}^K \gamma(z_{nk}) \log [\mathcal{N}(\mathcal{X}_n | \mu_k, \Sigma_k) \pi_k] \end{aligned}$$

この関数は、補助関数、あるいはQ関数とよばれる。

$\Lambda$  について、

$$\frac{\partial Q(\Lambda | \Lambda^{old})}{\partial \Lambda} = 0$$

とし、 $Q(\Lambda | \Lambda^{old})$  を最大化すればよい。この条件から、

$$\begin{cases} \mu_k = \frac{\sum_{n=1}^N \gamma(z_{nk}) \mathcal{X}_n}{\sum_{n=1}^N \gamma(z_{nk})} \\ \Sigma_k = \frac{\sum_{n=1}^N \gamma(z_{nk}) (\mathcal{X}_n - \mu_k)(\mathcal{X}_n - \mu_k)^T}{\sum_{n=1}^N \gamma(z_{nk})} \\ \pi_k = \frac{1}{K} \sum_{n=1}^N \gamma(z_{nk}) \end{cases} \quad \text{--- (1)}$$

を得る。ここで、

$$\begin{cases} n_k = \sum_{n=1}^N \gamma(z_{nk}) \\ \mathbf{t}_k = \sum_{n=1}^N \gamma(z_{nk}) \mathcal{X}_n \end{cases}$$

語彙認識では、コンポーネントの事後分布を計算することも、アライメントとよぶ。これは、音声フレーム  $\mathcal{X}_n$  が、どのコンポーネントに近いかを表すため。

$$S_k = \sum_{n=1}^N r(z_{nk}) x_n x_n^T$$

よって、①は

$$\begin{cases} M_k = \frac{1}{N_k} f_k \\ \Sigma_k = \frac{1}{N_k} S_k - M_k M_k^T \\ \pi_k = \frac{1}{N} N_k \end{cases}$$

よける。

この値を用いて、 $Q(\Lambda | \Lambda^{old})$  の値を求めよことが出来る。

1 テレ-ジョンを通じてパラメータを更新し、 $Q$  関数の値の変化量が小さい値を下の回、終了。

## 2.2.4 E-M

Init.

- ・  $X$  からランダムに  $k$  個のデータをえらぶ。
- ・ えらんだデータを  $\{M_k | k=1, \dots, k\}$  にあてがう。
- ・  $\pi_k = \frac{1}{k}$ ,  $\Sigma_k = I$  にそれぞれ設定する。 ( $k=1, \dots, k$ )。

E-STEP

- ・  $Q$  関数が  $M_k, \Sigma_k, \pi_k$  (for  $k=1, \dots, k$ ) で表れるように、 $\tau(z_{nk})$  を、全学習サンプルについて求める。

M-STEP

- ・ E-step で求めた  $\tau(z_{nk})$  を用いて、新しいパラメータを求める。
- ・  $\tau(z_{nk})$  を求めたパラメータから、 $Q$  関数の値を求める。
- ・ 前式行時の  $Q$  関数の値と、今回求めた  $Q$  関数の値を比較する。
  - ・ 値の増加が止まっていれば、アルゴリズム終了。
  - ・ 止まなければ、goto E-step.

