

1 K-Means 法

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

$$\begin{aligned} & \mu_k^T x_i \\ = & (\mu_k^T x_i)^T \quad (\because \text{スカラーであるため転置してもよい}) \\ = & x_i^T (\mu_k^T)^T \quad (\because (ab)^T = b^T a^T) \\ = & x_i^T \mu_k \end{aligned}$$

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

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

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

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

K-Means 法によるクラスタリング

- これより、 μ_k について解くと次のようになる

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

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

K-Means 法によるクラスタリング

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

K-Means 法によるクラスタリング

- $\{\mu_k\}$ と $\{r_{ik}\}$ についての最適化

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

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

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

K-Means 法によるクラスタリング

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

1 $\mathcal{M} = \{\mu_k\}$ の初期値を設定

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

K-Means 法によるクラスタリング

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

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

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

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

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

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

K-Means 法によるクラスタリング

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

K-Means 法によるクラスタリング

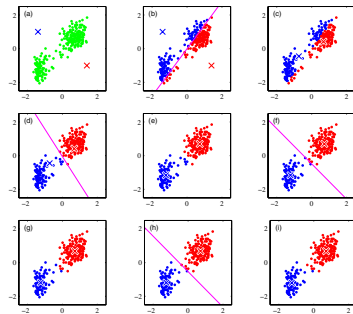


Figure 9.1 Illustration of the K -means algorithm using the re-scaled Old Faithful data set. (a) Green points denote the data set in a two-dimensional Euclidean space. The initial choices for centres μ_1 and μ_2 are shown by the red and blue crosses, respectively. (b) In the initial E step, each data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

図 1: K-Means アルゴリズムの動作

K-Means 法によるクラスタリング

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

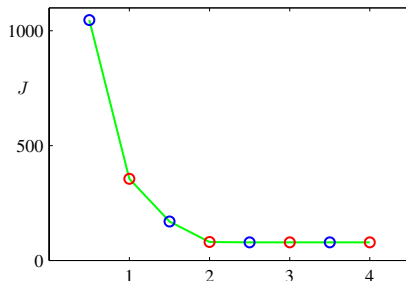


図 2: 目的関数 J の値の遷移

K-Means 法によるクラスタリング

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

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

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

逐次版の K-Means 法

- 逐次版の K-Means 法の導出

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

- Robbins-Monro 法

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

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

- このとき、根 θ^* の逐次計算式は、次のように記述される

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

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

- Robbins-Monro 法を適用するための条件

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

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

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

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

3 学習率の系列 $\{\eta_N\}$ は次の条件を満たす

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

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

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

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

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

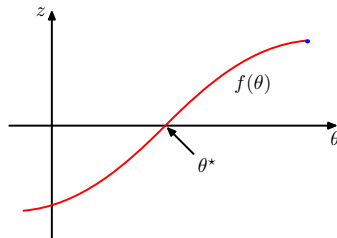


図 3: $f(\theta)$ のグラフ

逐次版の K-Means 法

● 逐次版の K-Means 法の導出

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

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

これは以下の式と等価である ($N_k = \sum_i r_{ik}$)

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

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

逐次版の K-Means 法

- このとき、上式から r_{ik} を省略でき、次のようになる

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

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

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

逐次版の K-Means 法

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

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

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

K-Means アルゴリズムの特徴

- K-Means アルゴリズムの特徴

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

K-Means アルゴリズムのまとめ

- K-Means アルゴリズムの目的

- 各データが属するクラスタを決定する (ハード割り当て)

- K-Means アルゴリズムで行っていること

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

K-Means アルゴリズムのまとめ

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