

X ① K-means クラスタリングが適切なアルゴリズムである場合を選ぶ

① ユーザ情報のデータベースから抽出された時、異なる marketing segment に自動的にグループ分けしる ○

② スーパーマーケットのたしさんの品物の販売記録から、同時にかわれているものをグループ分けして同じ棚に置くようにしたい。 X 正しい X ○が正解であった。

K-means には任意の「2点間の距離」が定義される必要がある

同時に売れた回数（逆数）を距離にできるか？ 加法原理がないため「距離」ではない

③ 過去の天気記録から、明日の雨量を数値で predict する。

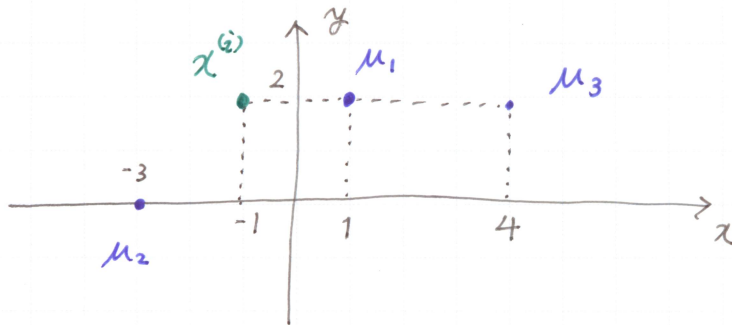
X K グループに分ける必要がない。

今日の天気からどのグループ化されるかわかっても雨量の数値はわかる
クラス内の平均？
最小？
最大？
最頻？

④ スーパーマーケット

○、各商品の今後の売れ高を予測する X

2 $\mu_1 = \begin{pmatrix} 1 \\ 2 \end{pmatrix}, \mu_2 = \begin{pmatrix} -3 \\ 0 \end{pmatrix}, \mu_3 = \begin{pmatrix} 4 \\ 2 \end{pmatrix}, x^{(i)} = \begin{pmatrix} -1 \\ 2 \end{pmatrix}$



最も $x^{(i)}$ に近いのは μ_1 $\therefore c^{(i)} = \mu_1$ の index = 1

3 K-means は 2つの step がくり返し実行される。2つの step を選べ。

① 7229 割り当て step で $c^{(i)}$ が更新される

② cluster centroid μ_j を動かす

③ μ_k を動かす (最も近い $x^{(i)}$ に合わせるように)

④ 7229 割り当て step で μ_j を最も近い $x^{(i)}$ に合わせるように μ_j を動かす

4 unlabeled dataset $\{x^{(1)}, \dots, x^{(m)}\}$, 50回の random initialization で 50 種の 7229 割り → どの分類結果を選ぶべき?

① ない X

② ラベルがない子母集合だけの子母集合、無理 X

③ $\frac{1}{m} \sum_{i=1}^m \|x^{(i)} - \mu_{c(i)}\|^2$ が最小のものを選ぶ

④ 50回の結果を選ぶ。収束しているものを 最後の

5

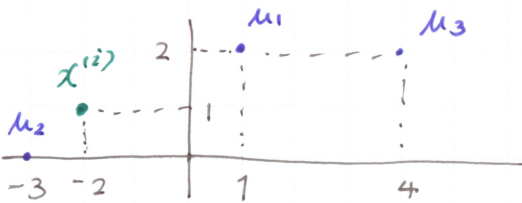
- ① 初期値として $\mu_1 = \dots = \mu_k = 0$ を選ぶ \times
- ② 「正しい」または「適切な」 k の値が あいまいなデータセットもある。
人間のエキスパートもわからない決定ではない。 \circ
- ③ local optimum に陥ることを心配するときは、複数回 random initialization を使う \circ
- ④ K-means は unsupervised なので、data に overfit することは無い。
計算可能な数のクラスター数は多い方がいい。 \times

Week 8 Quiz 1 2回目

1

- ① ユーザー情報から ○
- ② スーパーで同じ商品を売った人の ○
- ③ 過去の天気から 明日の雨量 ×
- ④ スーパー 今週の売れ残りを ×

2



$$c^{(i)} = 2 \quad (= \mu_2 \text{ の index})$$

3

K-meansに含まれる2つのstep

- ① Kを固定したときの elbow method ×
- ② 7329 列のデータを step 2- $c^{(i)}$ を更新 ○
- ③ cluster centroid 移動 step, μ_k を更新 ○
- ④ feature scaling ×

4

unlabeled $\{x^{(1)}, \dots, x^{(m)}\}$, 50回329~4 → 50回の結果

① $\frac{1}{m} \sum_i \|x^{(i)} - \mu_{c(i)}\|^2$ を最小にした ○

5

- ① $\mu_1 = \dots = \mu_k = 0$ を初期値 ×
- ② local optima を逃れようとする → 7329の初期値の乱数回数 ○
- ③ unsupervised なのに overfit は少ない, 7329 数は7-8を329より多く ×
- ④ どのkが正しいかわからないデータもある.