

はじめてのパターン認識

第5章 k最近傍法

11t4062a 広原裕亮

はじめに

- 最近傍法
 - 入力データとすべての学習データとの距離からクラスを識別
- K最近傍法
 - 最も近いk個の学習データを選ぶ方法
- 近似最近傍探索

5. 1 最近傍法とボロノイ境界 (1)

- K 個のクラス $\Omega = \{C_1, \dots, C_k\}$
- i 番目のクラスの学習データ数を $N(i)$
- 学習データの集合 $S_i = \{X_1^{(i)}, \dots, X_{N(i)}^{(i)}\}$
- 入力データと学習データの類似度
 $d(x, x_j^{(i)}) = ||x - x_j^{(i)}||$

5. 1 最近傍法とボロノイ境界 (2)

• 識別クラス =

$\{\arg \min d(x, x_j^{(i)}) \mid \min d(x, x_j^{(i)}) < t \text{ の時}$

$\{\text{reject} \mid \min d(x, x_j^{(i)}) \geq t \text{ の時}$

tは、どの学習データとも距離が大きい場合にリジェクトするためのしきい値

5. 1 最近傍法とボロノイ境界 (3)

- ボロノイ領域

- 各学習データが隣接する学習データと等距離にある境界で囲まれた領域

- ボロノイ境界

- 各学習データから等距離の点の集合

$$(\bar{x} - x)^T n = 0, \quad \bar{x} = (x_i + x_j) / 2, \quad n = (x_i - x_j)$$

- 上式を満たす x が表す超平面が半空間に2分割する

5.1.2 学習データの数と識別性能

- 学習データを増やすと識別性能は上昇する
- しかし計算量が増加してしまう

5. 2 K最近傍法(1)

- 識別クラス =

$$\{j \mid \{k_j\} = \max\{k_1, \dots, k_K\}\}$$

$$\{\text{リジェクト} \mid \{k_i, \dots, k_j\} = \max\{k_1, \dots, k_K\}\}$$

- K個のデータのうち最も多いクラスjと識別
- タイがある場合リジェクト

5. 2 K最近傍法(2)

- 最近傍法の場合($k=1$)
 - 孤立した識別境界が見られる
 - 計算量が多い
- 11最近傍法($k=11$)
 - 滑らかな境界となる
 - 識別精度が下がる

5.3 kNN法の計算量とその低減法

- データと全てのクラスの学習データとのユークリッド距離
- 距離を昇順にソート



データが増えると多くの計算時間と記憶容量が必要。
これらを緩和する方法がある。

誤り削除型k最近傍法

- 正しくないクラスの領域に存在している学習データを削除する



不可能

- 誤ってるのに削除されない
- 正しいのに削除される

圧縮型k最近傍法

- 識別境界のデータが重要

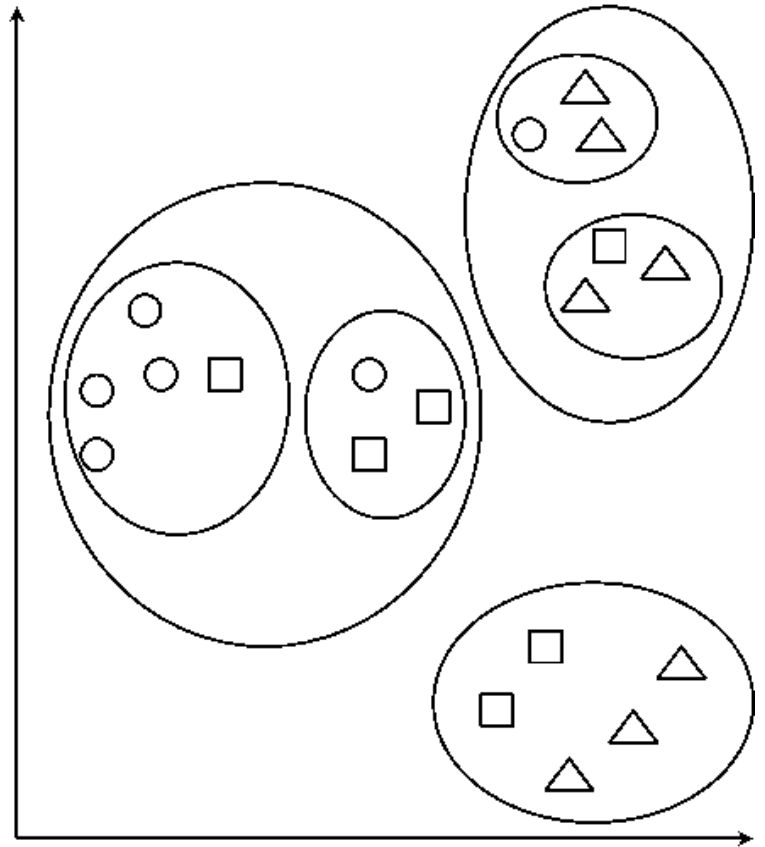
$$\varepsilon(x) = \min[P(C_1|x), P(C_2|x)]$$

— ベイズ誤り率を小さくする

- 削除前と誤り率が同じになるようにする

分岐限定法(1)

- ・クラスタリングで学習データを組織化
- ・平均ベクトル m_i
- ・最遠のデータまでの距離 d_i を算出



分岐限定法(2)

・一番上の階層のクラスタから順に平均ベクトルの最も近いクラスタを探す

→ほかのクラスタからは？

$$d(x, m_n) > d(x, x_i) + d_n$$

が成り立てば、そのクラスタを探索する必要はない。

近似最近傍探索(1)

- 最近傍より少し遠くても許容する
— $d(q, x) \leq (1 + \varepsilon)d(q, x^*)$
の距離にあるデータまで許容する。
 x^* とは最適解である。

近似最近傍探索(2)

Step1:2分木を作成する。

Step2:入力データと同じ領域内のデータとの距離を算出

Step3:次に近い領域との距離を算出し比較

Step4:次にデータとの距離を算出し比較。以降それを繰り返す

近似最近傍探索 (3)

- $d(q,x) / (1+\varepsilon) \leq d(q,X)$ を満たすと近似解

↓ なぜ言えるのか

上記の式を満たした領域までの距離を r'

$$d(q,x) / (1+\varepsilon) \leq r' \leq d(q, x')$$

↓

$$d(q,x) \leq (1+\varepsilon)r' \leq (1+\varepsilon)d(q, x')$$