

Case study.....pkt,
Pytania zamknięte.....pkt
Razem..... pkt

Imię i nazwisko.....

Ocena.....

K O L O K W I U M 1

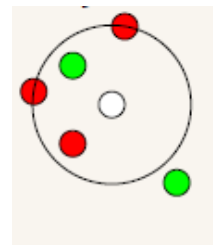
15 marca 2019 r.

Statystyczna Eksploracja Danych

Pytania zamknięte (20 pytań po 0.5 pkt każde, jedna odpowiedź prawidłowa)

Zad 1. Ile punktów zostanie wziętych pod uwagę podczas wyznaczania przynależności do klasy punktu środkowego metodą 3-nn (rysunek po prawej)?

- (a) 3
- (b) 4
- (c) 2



Zad 2. Ten sam rysunek, co w poprzednim pytaniu. Czy w przypadku metody (4,3)-nn punkt centralny zostanie

- (a) zaliczony do klasy czerwonej,
- (b) zaliczony do klasy zielonej,
- (c) decyzja nie zostanie podjęta.

Zad 3. Jakie założenia prowadzą do różnych postaci reguły klasyfikacyjnej w metodach LDA i QDA?

- (a) jednowymiarowy rozkład Gaussa (LDA), dwuwymiarowy rozkład Gaussa (QDA),
- (b) różne macierze kowariancji (QDA), równe macierze kowariancji (LDA),
- (c) stosowalność dla dwóch klas (LDA), stosowalność dla trzech klas (QDA).

Zad 4. Wybierz **prawdziwą** cechę naiwnego klasyfikatora Bayesowskiego:

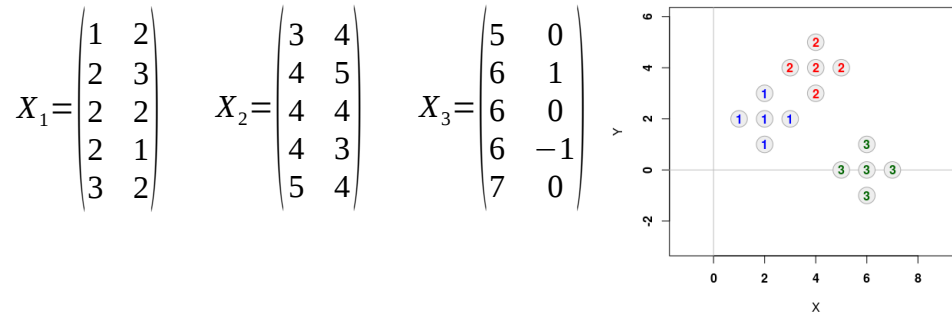
- (a) multiplikatywny wpływ kolejnych składowych wektora obserwacji,
- (b) Gaussowski rozkład obserwacji $p(k|\mathbf{x})$,
- (c) niezależność składowych wektora obserwacji.

Zad 5. W metodzie bootstrap:

- (a) losujemy bez zwracania
- (b) losujemy ze zwracaniem
- (c) tworzymy tylko jedną pseudopróbe.

Case study (5 pkt)

Dane są trzy serie danych (x,y) : \mathbf{X}_1 , \mathbf{X}_2 oraz \mathbf{X}_3 , każda pochodząca z innej klasy (czyli $g=3$), w każdej z nich jest po 5 punktów:



Oblicz optymalny kierunek rozdzielający \mathbf{a} za pomocą metody Fishera. Oceń, czy metoda umożliwia rozdzielenie klas. Dlaczego drugi kierunek (odpowiadający drugiej wartości własnej) jest gorszy?

Wzory:

Wartość oczekiwana w każdej klasie k (n_k to liczba punktów w danej klasie):

$$\mathbf{m}_k = \frac{1}{n_k} \sum_{i=1}^{i=n_k} \mathbf{x}_i$$

Macierz kowariancji w każdej klasie k :

$$\mathbf{S}_k = \frac{1}{n_k - 1} \sum_{i=1}^{i=n_k} (\mathbf{x}_i - \mathbf{x}_k)(\mathbf{x}_i - \mathbf{x}_k)^T$$

Macierz kowariancji wewnątrzgrupowej \mathbf{W} (\mathbf{n} to całkowita suma punktów $\mathbf{n} = \mathbf{n}_1 + \mathbf{n}_2 + \mathbf{n}_3$):

$$\mathbf{W} = \frac{1}{n - g} \sum_{k=1}^{k=g} (n_k - 1) \mathbf{S}_k$$

Macierz kowariancji międzygrupowej \mathbf{B} (\mathbf{m} to średnia ze wszystkich punktów):

$$\mathbf{B} = \frac{1}{g - 1} \sum_{k=1}^{k=g} n_k (\mathbf{m}_k - \mathbf{m})(\mathbf{m}_k - \mathbf{m})^T$$

Optymalny kierunek \mathbf{a} otrzymujemy maksymalizując wyrażenie

$$J = \frac{\mathbf{a}^T \mathbf{B} \mathbf{a}}{\mathbf{a}^T \mathbf{W} \mathbf{a}}$$

wyznaczając wektor własny odpowiadający największej wartości własnej równania

$$\mathbf{A} = \mathbf{W}^{-1} \mathbf{B}$$

Wartości własne λ można otrzymać wyznaczając pierwiastki wielomianu charakterystycznego

$$\det(\mathbf{A} - \lambda \mathbf{I}) ,$$

gdzie \mathbf{I} to macierz jednostkowa (same zera oprócz jedynek na diagonalu), a odpowiadające poszczególnym wartościom własnym λ_i wektory własne \mathbf{x}_i rozwiązują równania

$$(\mathbf{A} - \lambda_i \mathbf{I}) \cdot \mathbf{x}_i = 0$$