

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

Imię i nazwisko.....

Ocena.....

K O L O K W I U M 2 12 kwietnia 2019 r.

Statystyczna Eksploracja Danych

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

Zad 1. W przypadku idealnie separowalnym metody SVM hiperpłaszczyzna marginesu:

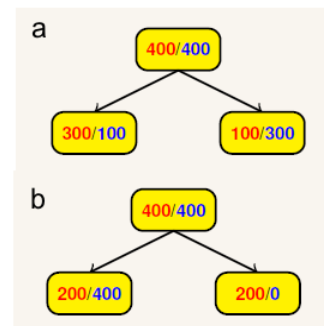
- (a) przechodzi zawsze przez jeden punkt danych,
- (b) nie przechodzi przez punkt danych,
- (c) może przechodzić przez wiele punktów danych.

Zad 2. Wybierz **prawdziwą** informację:

- (a) aby zastosować jądro w metodzie SVM, należy znać przestrzeń, do której transformuje ono obserwacje
- (b) obliczenia za pomocą jądra dotyczą jedynie iloczynu skalarnego,
- (c) można stosować dowolne jądra.

Zad 3. Obejrzyj przypadki **a** i **b** podziału rozdziału elementów do węzłów-dzieci i zaznacz prawidłową odpowiedź:

- (a) różnica różnorodności klas pomiędzy węzłem rodzicem i węzłami-dzieci jest zawsze taka sama w obu przypadkach,
- (b) różnica różnorodności klas pomiędzy węzłem rodzicem i węzłami-dzieci jest taka sama w obu przypadkach dla indeksu Gini'ego,
- (c) różnica różnorodności klas pomiędzy węzłem rodzicem i węzłami-



Zad 4. Zaznacz **nieprawdziwą** informację, dotyczącą drzew klasyfikujących:

- (a) wszystkie podziały w drzewie mogą zostać wykonane na jednym atrybucie,
- (b) drzewo może składać się z jednego węzła,
- (c) w każdym drzewie muszą zostać wykorzystane podziały na każdym atrybucie.

Zad 5. W przypadku zespołu $M < 100$ klasyfikatorów, każdy o skuteczności $p > 0.5$:

- (a) klasyfikatory jednomyślnie dokonują decyzji,
- (b) większość klasyfikatorów nigdy się nie myli,
- (c) jest duże prawdopodobieństwo, że większość klasyfikatorów dokona poprawnej decyzji.

Case study (5 pkt)

Zakładając, że stosowaną miarą różnorodności jest indeks Gini'ego [$Q=2p(1-p)$], skonstruuj dla poniższych danych drzewo klasyfikujące (obserwacje leżą na prostej $y=x$, przy czym współrzędne są kolejnymi liczbami naturalnymi od 1 do 45). Narysuj to drzewo. Czy istnieje inne rozwiązanie? Czy sens ma zastosowanie algorytmu kosztu-łożoności?

