

## Lista 4 - Selekcja cech

### Selekcja cech z wykorzystaniem GUI Weki

Należy otworzyć skonstruowany w ramach poprzedniej listy zbiór danych *XXXXXXL3\_1.arff* poprzez udostępnione przez *Wekę* GUI. Należy zapoznać się działaniem modułu selekcji cech (Zakładka *Select Attributes*).

### Metody selekcji cech wykorzystujące entropię

Typowymi metodami stosowanymi do selekcji cech są algorytmy wykorzystujące pojęcie entropii:

$$H(X) = - \sum_{x \in \mathbb{X}} p(x) \log p(x) \quad (1)$$

oraz entropii warunkowej:

$$H(X|Y) = \sum_{y \in \mathbb{Y}} p(y) H(X|y) \quad (2)$$

W *Wece* wyróżnić można dwie metody które wykorzystują entropię do oceny istotności atrybutów: **GainRatioAttributeEval**, oraz **InfoGainAttributeEval**. Pierwszy z nich bada istotność atrybutów ze względu współczynnik *GainRatio* definiowany w następujący sposób:

$$GainRatio(Class, Attribute) = \frac{H(Class) - H(Class|Attribute)}{H(Attribute)} \quad (3)$$

natomiast drugi z nich wykorzystuje tzn. *InfoGain*:

$$InfoGain(Class, Attribute) = H(Class) - H(Class|Attribute) \quad (4)$$

Opisane metody wykonują ocenę każdego atrybutu ze względu na przyjęte kryterium niezależnie, dla każdego z atrybutów osobno.

## Zadania

Wszystkie zadania zostaną wykonane na pliku *XXXXXXL3\_1.arff*.

1. Należy dokonać dyskretyzacji zmiennych numerycznych z wykorzystaniem filtra pracującego w trybie nadzorowanym. W dalszej kolejności należy zapoznać się z działaniem filtrów do selekcji cech **GainRatioAttributeEval**, oraz **InfoGainAttributeEval**. Należy wybrać cechy dla których zarówno *GainRatio*, jak i *InfoGain* przyjmują wartości wyższe niż 0.001. Należy uszeregować atrybuty rosnąco względem *GainRatio* i zbiór po procesie selekcji i uszeregowaniu zapisać jako *XXXXXXL4\_1.arff* (3 pkt).
2. Należy własnoręcznie (bez wykorzystywania klas **GainRatioAttributeEval**, **InfoGainAttributeEval**) zaimplementować metodę **GainRatioAttributeEval** i zweryfikować jej działanie na zbiorze *XXXXXXL3\_1.arff*. Należy zidentyfikować podstawę logarytmu, jaką wykorzystuje implementacja **GainRatioAttributeEval** w *Wece* zadając jej wartość jako parametr programu (5 pkt).