

Úkol B – Hledání podskupin v datech – jaterní testy ZS 2018

Cíl

Cílem druhého úkolu je seznámit se s praktickým použitím shlukovacích metod. Shlukovací metody lze využít k odhalení vnitřní struktury dat. Poznání vnitřní struktury dat může pomoci k lepšímu porozumění datům a jevům, které zachycují, stejně tak ke konstrukci lepších klasifikátorů.

Data

Dataset Indian Liver Patient Records obsahuje 416 záznamů pacientů s cirhózou a 167 pacientů se zdravými játry naměřených v Andhra Pradesh v Indii. Data obsahují jak základní demografické údaje, tak i biochemické markery jako je hladina bilirubinu nebo albuminu v krvi.

Hledání podskupin v datech

Pokuste se v datech nalézt vnitřní strukturu (podskupiny) na základě hodnot příznaků. Zvolte vhodnou metodu. Zamyslete se nad tím, zda je vhodné data normalizovat. Výsledky vizualizujte a pokuste se je interpretovat, např. pomocí typických (průměrných) reprezentantů. Prozkoumejte, které z příznaků se nejvíce podílí na rozdělení pozorování do jednotlivých shluků. Uvažte, jestli je možné dát nalezenou strukturu do souvislosti s výskytem cirhózy jater.

Požadované kroky analýzy [10 bodů]

- upravte data pro shlukování [1 b]
- zvolte vhodnou metodu shlukování [1 b]
- zvolte vhodné parametry shlukovací metody [1 b]
- prezentujte výsledky shlukování pomocí typických reprezentantů shluků [3 b]
- zjistěte v jakých příznacích a jak se jednotlivé shluky liší [3 b]
- porovnejte výsledky shlukování s výskytem cirhózy jater [1 b]

Výsledky upravte do formy zprávy, která bude obsahovat stručný **úvod**, popis metod, které jste použili, v sekci **metody**, výsledky jejich aplikace na data v sekci **výsledky** a závěry, které jste zjistili interpretací výsledků v sekci **závěr**. Maximální délka zprávy je 3 stránky. Zprávu ve formátu *pdf* a stručně okomentovaný, přehledný zdrojový *R skript* odevzdejte pomocí *UploadSystemu*.

Název parametru	Popis parametru
Age	Věk (pacientům nad 89 let je přiřazen věk 90)
Gender	Pohlaví
Total_bilirubin	Celkový bilirubin
Direct_bilirubin	Konjugovaný bilirubin
Alkaline_phosphotase	Alkalická fosfatáza
Alamine_Aminotransferase	Alaninaminotransferáza
Aspartate_Aminotransferase	Aspartátaminotransferáza
Total_protiens	Celkový sérový protein
Albumin	Albumin
Albumin_Globulin_Ration	Poměr albuminu a globulinu
Dataset	Diagnoza (1 = cirhóza jater, 2 = zdravá játra)

Tabulka 1: Výčet parametrů

Bonusová úloha – simulovaná data

Na rozdíl od hlavního úkolu je bonusová úloha zaměřená na simulovaná data. Data k bonusové úloze představují závislosti, která jsou těžká pro standardní algoritmy shlukování jako je k-means. Vaším úkolem je v tomto případě vyzkoušet na data metody shlukování prezentované v rámci DVZ a pokusit se najít takovou kombinaci parametrů metod (metriky, způsob linkování pozorování) a transformace příznaků (polynomiální příznaky, kernelová metoda), které by umožnili správně klasifikovat data v souborech *jain.csv*, *spiral.csv* a *pathbased.csv*. Všechna data obsahují správné řešení v proměnné *class*. Volbu metody zdůvodněte a vysvětlíte, proč Váš přístup data správně shlukuje. Za bonusovou úlohu je možné získat až 3 body.