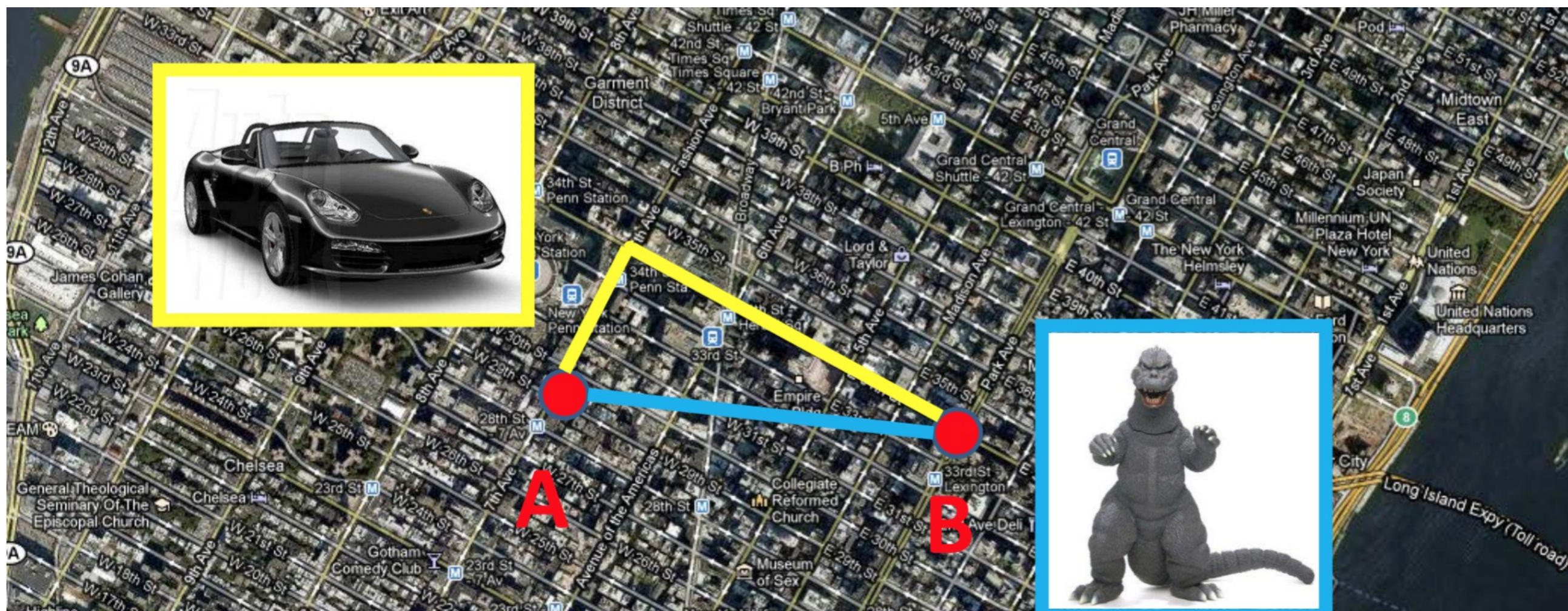


Pokročilá chemoinformatika

Neparametrické modely
únor 2017

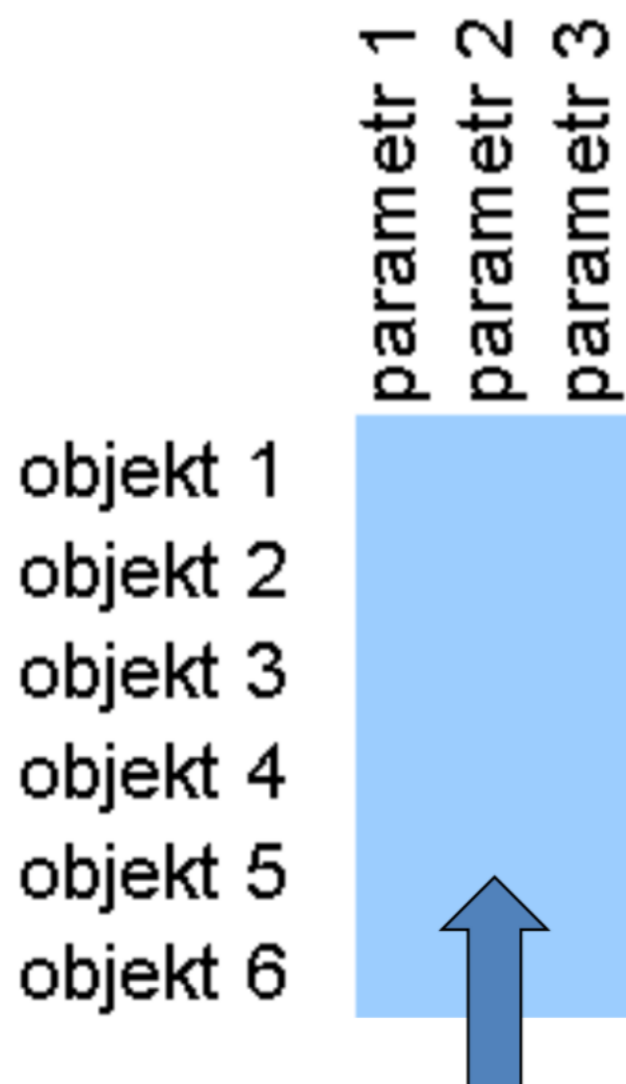
Vzdálenost

- Každá molekula je definována v N-dimenzionálním prostoru (N je počet deskriptorů).
- Mezi každou molekulou můžeme spočítat vzdálenost a pro celou sadu pak asociační matici.



Asociační matice

NxP MATICE

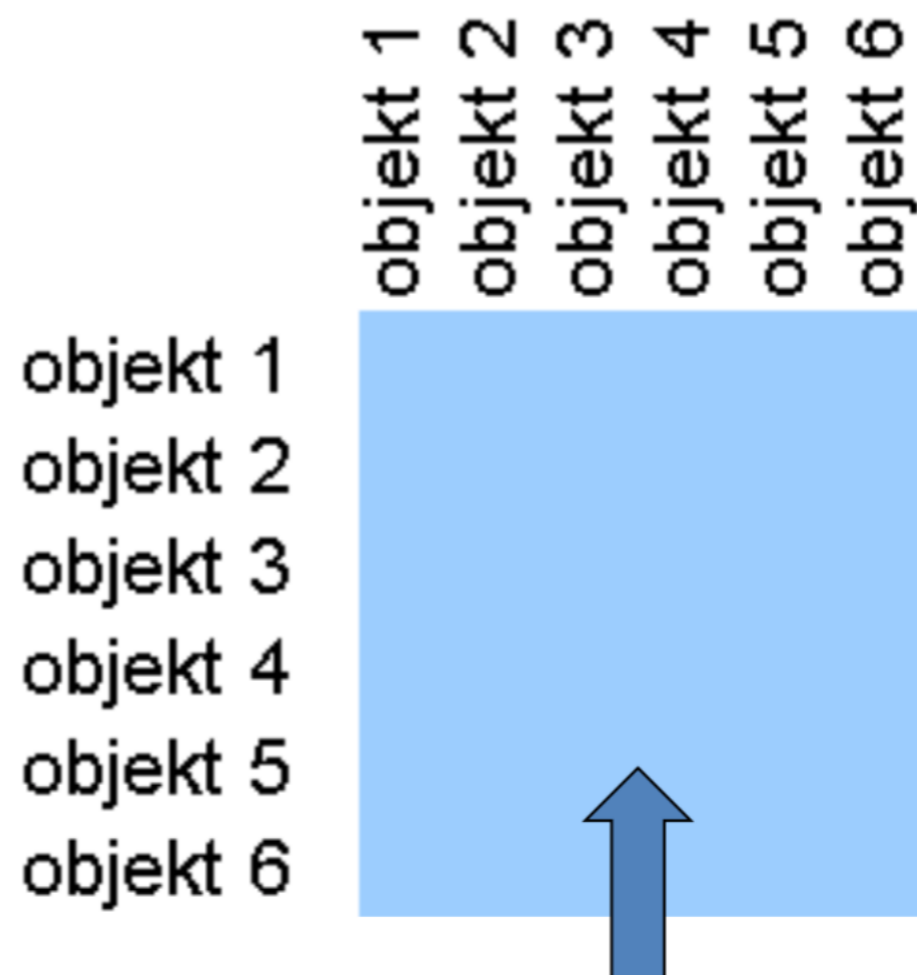


Hodnoty parametrů pro jednotlivé objekty

Výpočet metriky
podobností/
vzdáleností



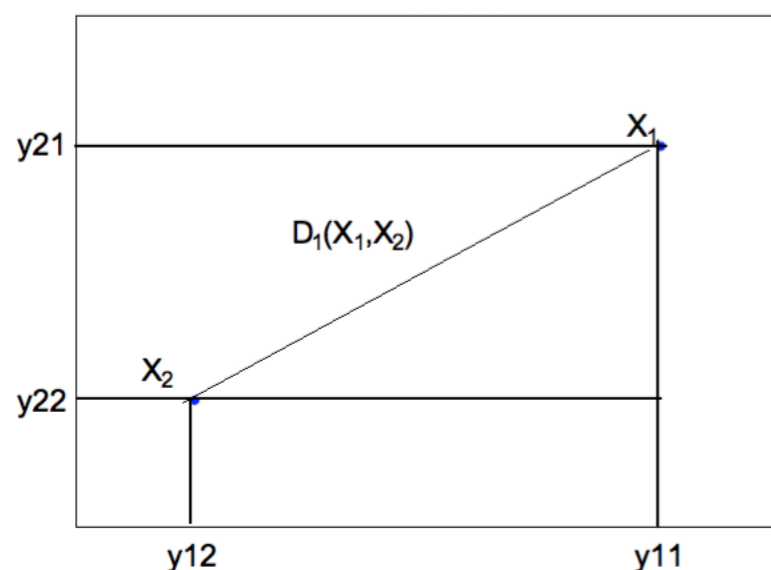
ASOCIAČNÍ MATICE



Korelace, kovariance, vzdálenost, podobnost

Základní vzdálenosti

- Euklidovská vzdálenost



$$D_1^2(x_1, x_2) = \sum_{j=1}^p (y_{1j} - y_{2j})^2$$

- Průměrná vzdálenost

$$D_2^2(x_1, x_2) = \frac{1}{p} \sum_{j=1}^p (y_{1j} - y_{2j})^2$$

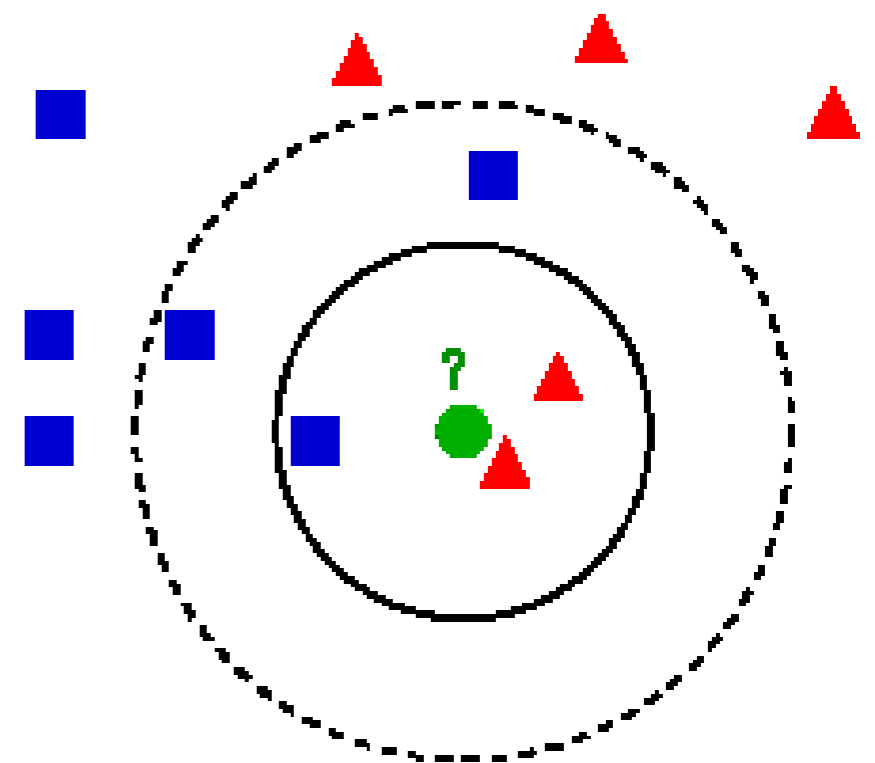
- Chord distance

- Manhattanska vzdálenost

KNN (K-Nearest Neighbors)

Algoritmus k-nejbližších sousedů

- patří mezi nejjednodušší “machine learning” algoritmy
- Je to jedna z mnoha hierarchickým metod shlukování
- může sloužit ke klasifikaci nebo predikci
- pokud budeme klasifikovat zelené kolečko podle 3-KK (bereme v úvahu 3 sousedy) bude patřit do trojúhelníků, v 5-KK do čtverců
- v případě regrese se hodnota bude počítat na základě nejbližších sousedů a vzdálenosti

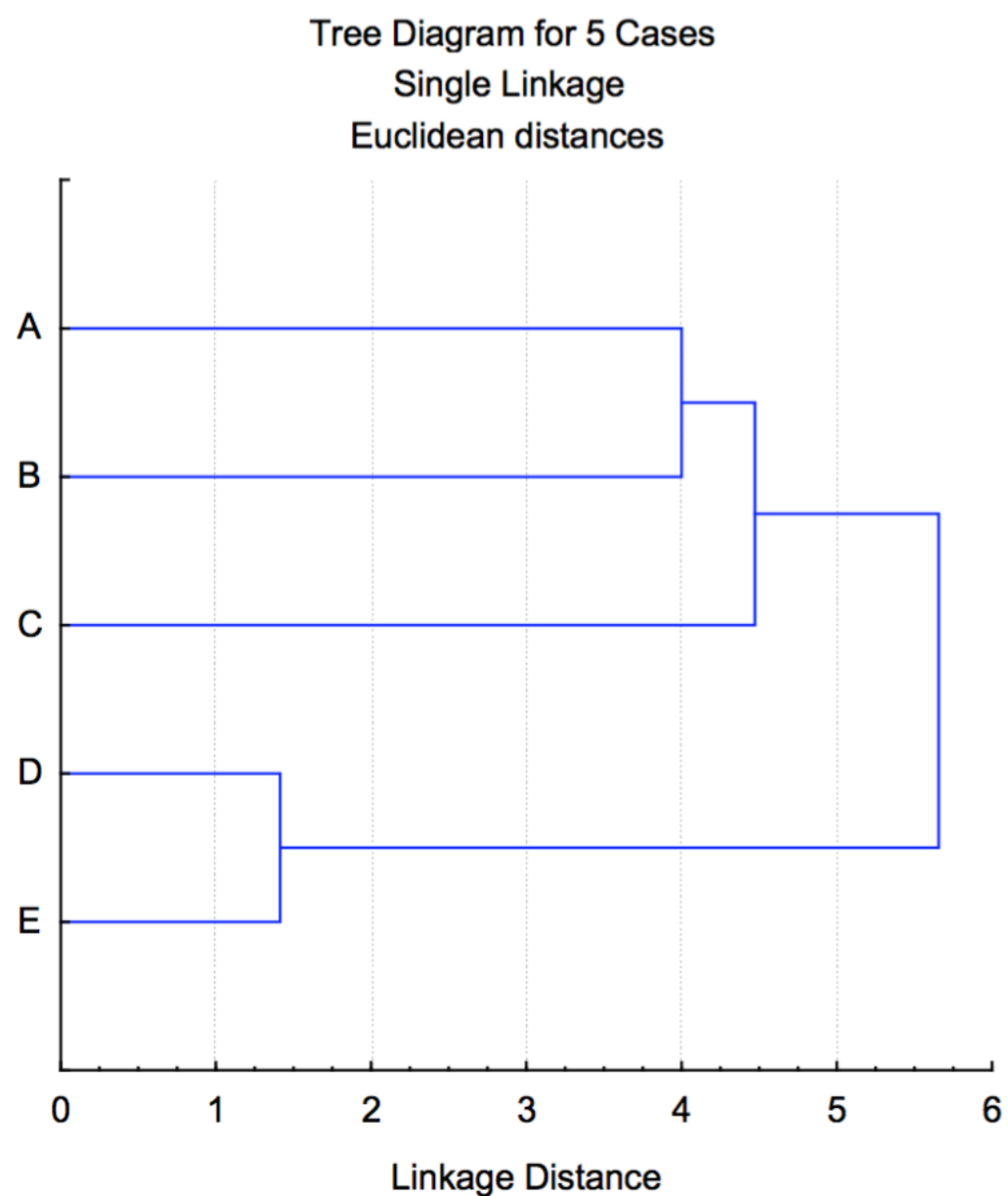


Tvorba KNN

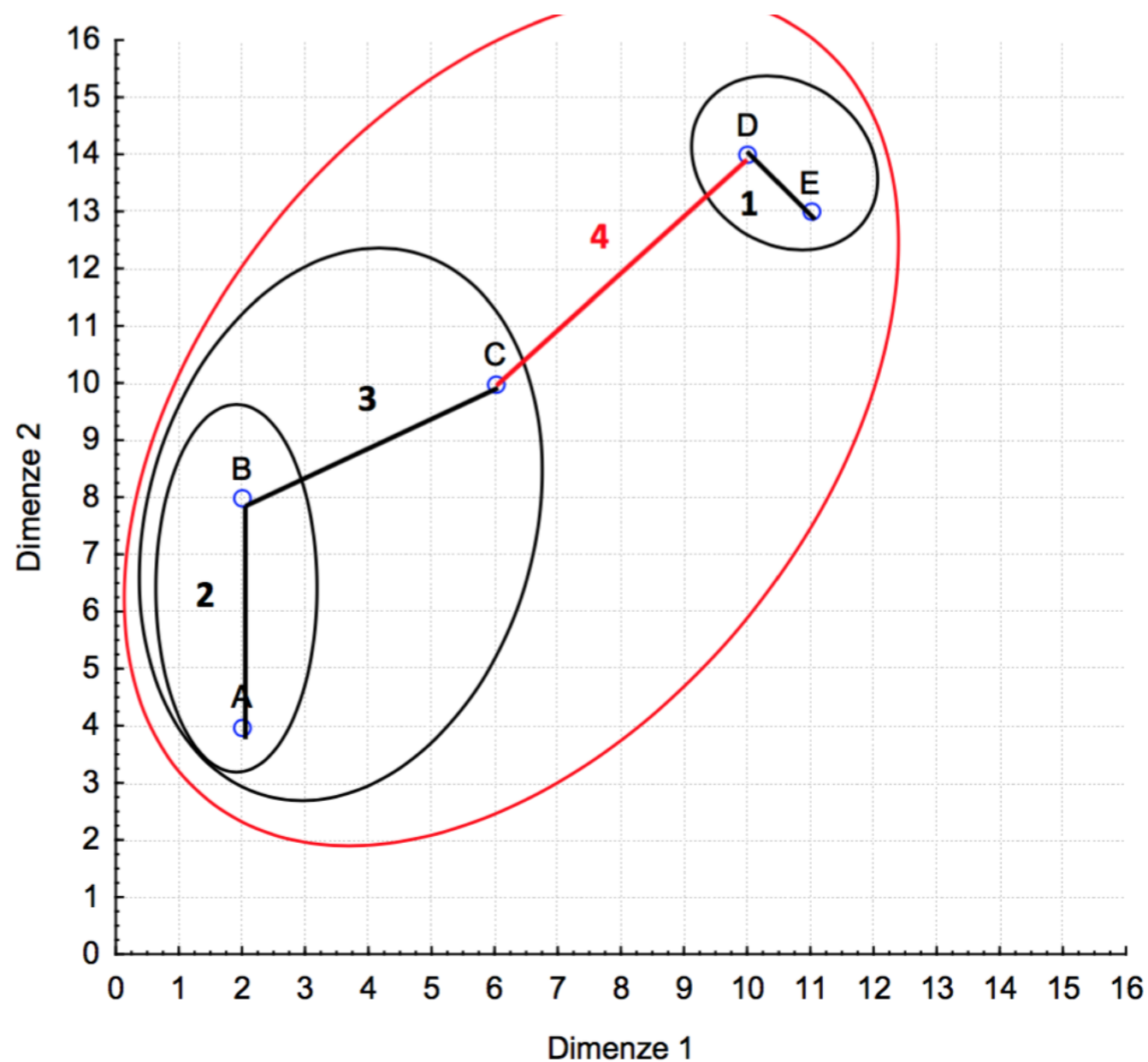
1. Výpočet asociační matice
2. Definice prvního shluku (nebližší prvky v matici)
3. Přepočítání matice, nový shluk je brán jako nový prvek
4. Opakujeme kroky 2 - 3, dokud není pouze jeden "prvek"

Vizualizace shlukové analýzy

dendrogram



klastry



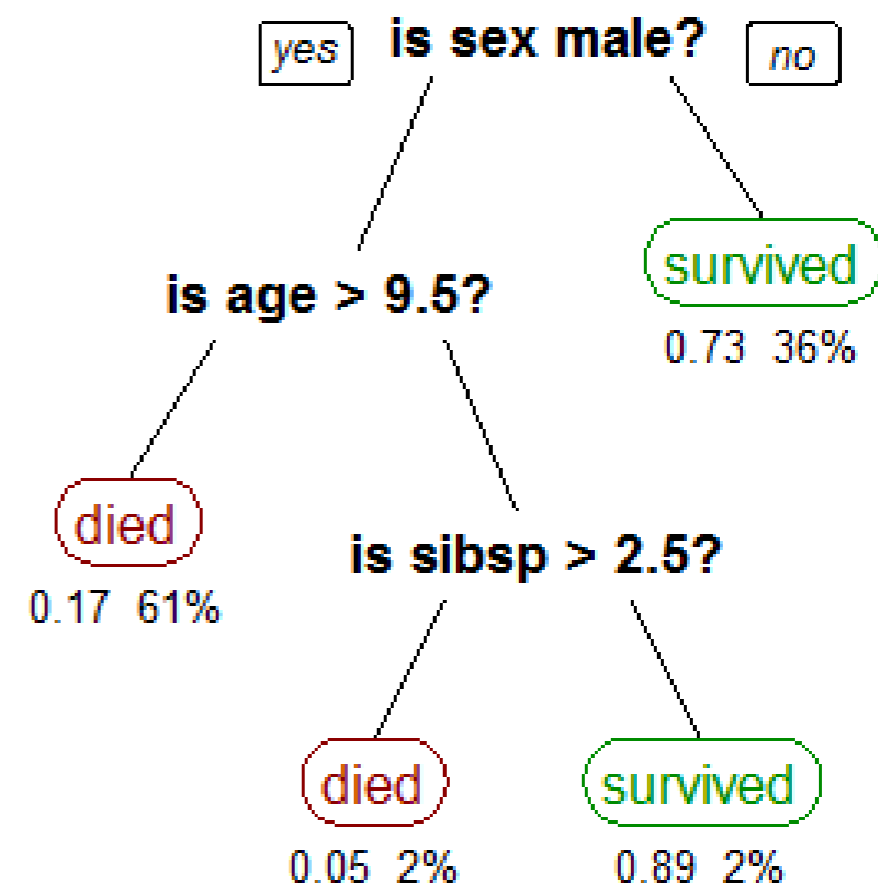
Výhody a nevýhody KNN

- Výhody:
 - Robustní na nesourodé data
 - Efektivní při velkém množství dat
- Nevýhody:
 - Potřeba definovat parametr K
 - Potřeba výběru metody výpočtu vzdálenosti
 - Výpočetně náročné – nutnost výpočet všech vzdáleností

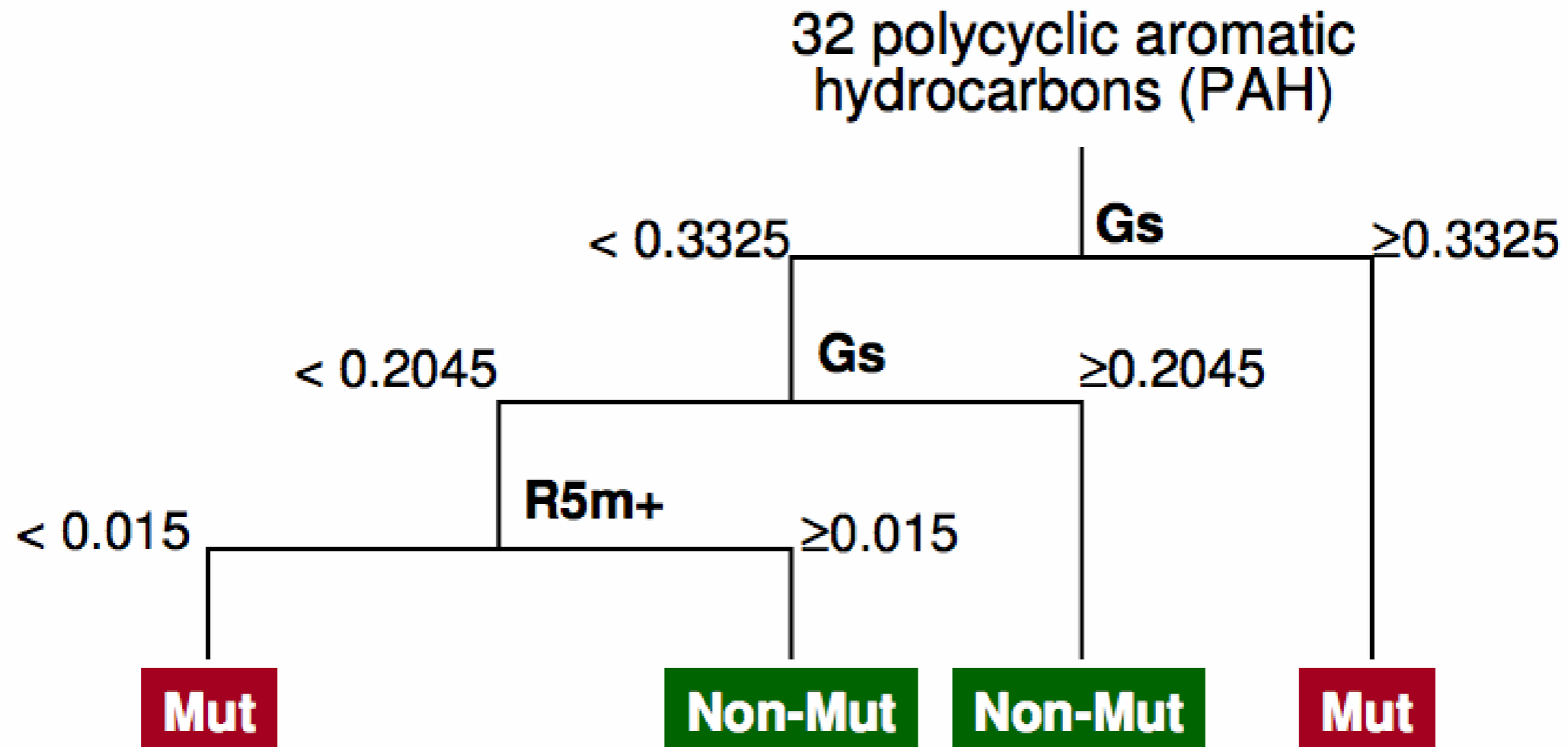
Decision Tree

Rozhodovací stromy

- může sloužit ke klasifikaci (classification trees) nebo predikci (regression trees)
- V těchto stromových strukturách představují listy (leafs) třídy a větve představují spojky mezi třídami.



Decision tree – příklad predikce mutagenity



Gs: G total symmetry index/weighted by atomic electrotopological states (3D-WHIM descriptor)

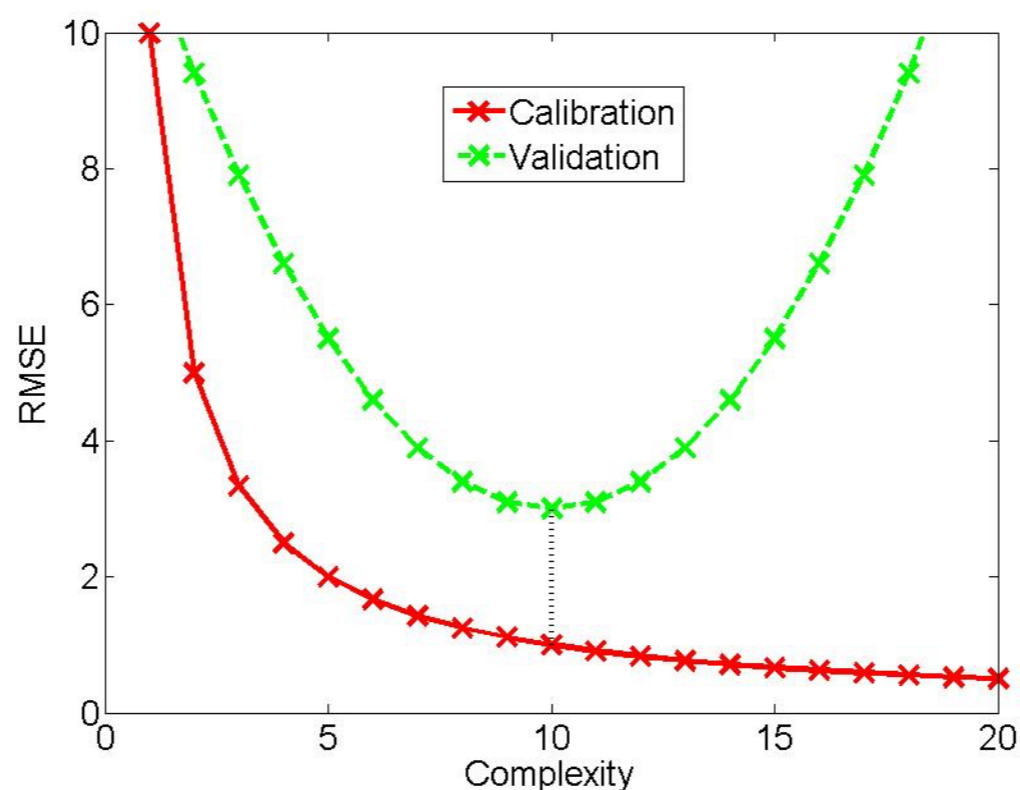
R5m+: R maximal autocorrelation of lag 5/weighted by atomic masses (3D-GETAWAY descriptor)

Tvorba rozhodovacího stromu CART

1. Regresní stromy
Kriteriální statistika: Least square deviation
$$Q = 1/n \sum (y_i - \bar{y})^2$$
 2. Klasifikační stromy
Kriteriální statistika: Gini index
- Podle kriteriální statistiky je provedene nejvhodnější rozdělení na dceřinné uzly
 - Růst je zastaven podle zvolených kriterií (nebo vyčerpání možností)

Testování výsledného stromu a jeho prořezání

- Každý strom je nutné testovat pomocí externí validace nebo křížové validace!



- Výsledná komplexita stromu je redukována dle výsledku testování.

Výhody a nevýhody rozhodovacích stromů

- Výhody:
 - Snadné grafické znázornění
 - Odolný vůči odlehlým hodnotám
 - Vhodné použití pro regresi i klasifikaci
- Nevýhody:
 - Nestabilní
 - Stromy jsou nevhodné pro malý počet vzorků
 - Vyžaduje znalosti s nastavením stromů

Náhodné lesy

- Tato metoda používá pro klasifikaci nebo regresi velké množství stromů
- Vhodné metoda pro použití v případě velkého množství deskriptorů a dobře funguje i na malých datových sadách
- Možnosti dalších analýz:
 - Měření významnosti proměnných
 - Efekt proměnných na predikci
 - Prototyp kategorií a možnost jejich překryvu
 - Detekce odlehlých hodnot
 - Doplnění chybějících hodnot

Genetický algoritmus

- TODO