

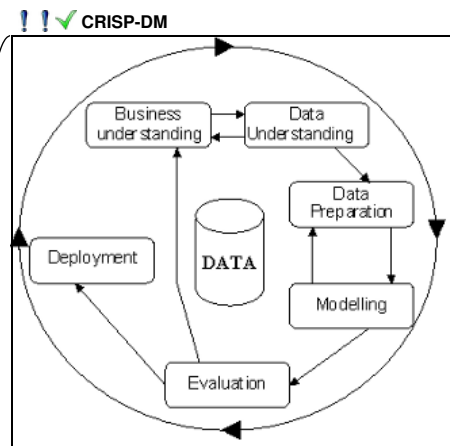
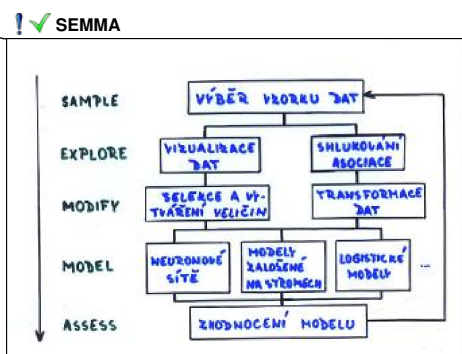
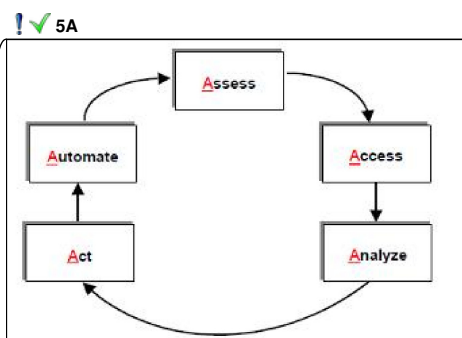
Data-Mining

1. Metodiky

poskytnout uživatelům jednotný rámec pro řešení různých úloh z oblasti dobývání znalostí

vyvinuté firmami (závislé na sw)

vyvinuté konsorciem (nezávislé na sw)



ASSESS - posouzení potřeb projektu Stanovení kontextu - cílů, strategií a procesů

ACCESS - shromáždění potřebných dat a jejich příprava

ANALYZE - provedení analýz
Data se přeměňují na informace a znalosti
→ použít vícero metod a porovnat jejich výsledky a efektivitu

ACT - přeměna znalostí na akční znalosti
Doporučení, dodatečné otázky a následná rozhodnutí
→ nalezené výsledky by měly být prezentovány jasně a srozumitelně

AUTOMATE - převedení výsledků analýzy do praxe
Může zahrnovat např. i vytvoření praktického rozhraní pro snadné použití
Umožnit aktualizaci modelů podle nových výsledků

vyvinula firma SPSS (nyní IBM)
používal ji systém SPSS Clementine
přestala se používat v roce 2001
nahradila ji CRISP-DM
podobá se hodně SEMMA

SAMPLE - výběr vhodných objektů

EXPLORE - vizuální explorace a redukce dat

MODIFY - seskupování objektů a hodnot atributů, datové transformace

MODEL - analýza dat

Neuronové sítě, rozhodovací stromy, statistické techniky, asociace a shlukování

ASSESS - porovnání modelů a interpretace

Srozumitelnost pro uživatele

vyvinula firma SAS
na svých stránkách tvrdí že SEMMA není metoda dataminingu, ale logická organizace klíčových úkolů data-miningu pro SAS Enterprise Miner :o)
používá ji systém SAS Enterprise Miner

CRoss-Industry Standard Process for Data Mining

pořadí fází není přesně určeno

Výsledky získané v jedné fázi ovlivňují volbu dalších kroků

Některé kroky a fáze je třeba provádět opakovaně

Business understanding

Pochopení cílů úlohy a požadavků na řešení (formulovaných z pohledu manažera)

Manažerskou formulaci je nutné převést na zadání úlohy pro dobývání znalostí z databází

"Revize" zdrojů (datových, výpočetních i lidských)

Hodnotí se možná rizika, náklady a přínos

Stanoví se předběžný plán prací

Data understanding

Prvotní sběr dat

Získání základní představy o datech

Posouzení kvality dat, vytipování zajímavých podmnožin záznamů v databázi, ...

Výpočet deskriptivních charakteristik dat

Četnost atributů, průměrné hodnoty, ...

Výhodou jsou vizualizační techniky

Data preparation

Vytvoření datového souboru, který bude zpracováván jednotlivými analytickými metodami

Data by měla obsahovat údaje podstatné pro danou úlohu a měla by být ve tvaru vyžadovaném algoritmy pro analýzu

Příprava dat zahrnuje:

Selekci dat, čištění dat, transformaci dat, vytváření dat, integrování dat, formátování dat, ...

Jednotlivé úkony se obvykle provádějí opakovaně a v nejrušnějším pořadí

Modeling

Použití analytických metod pro dobývání znalostí

Z možných metod vybrat ty nejhodnější a adekvátně nastavit jejich parametry

Iterativní činnost Opakovaná aplikace algoritmů s různými parametry

Může vést k potřebě modifikovat data

Ověření nalezených znalostí

Z pohledu manažerů

Evaluation

Byly splněny cíle formulované při zadání úlohy?

Rozhodnutí o způsobu využití výsledků

Upravit získané znalosti do podoby použitelné pro zákazníka (manažera, zadavatele)

Zákazník musí pochopit, co je třeba učinit pro efektivní využití dosažených výsledků!

Deployment

Implementace klasifikačního algoritmu v user-friendly podobě

Příprava uživatelského manuálu

Instalace programu na pobočkách banky a zaškolení uživatelů

Změna metodiky poskytování úvěrů a příslušná změna vnitřních předpisů banky

dnes praktický standard

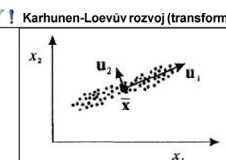
používá např. systém SPSS Modeller (dříve Clementine)

vyvinulo konsorcium: SPSS, NCR (Teradata), Daimler AG, OHRA (pojistovna)

- ✓ 2. Předzpracování dat
- ✓ 4. Rozhodovací stromy
- ✓ 5. Asociacní pravidla
- ✓ 6. Bayesovské modely
- ✓ 7. Neuronové sítě
- ✓ 8. ELM-sítě
- ✓ 9. Genetické algoritmy
- ✗ 10. Pokročilé předzpracování dat

1. Metodiky

Výběr příznaků

- selekcce**
 - Cíl: máme zbytečnou moc příznaků => chceme vybrat jenom ty s největším přínosem
 - z původní množiny ponecháme jen n nejvíce informativních příznaků
 - může snížit cenu měření příznaků
 - transformace příznaků na menší počet jiných příznaků
 - od výsledku předpokládáme stejnou informační přínos
 - musíme měřit všechny příznaky
- analýza hlavních komponent (PCA)**
 - něco jako neperiodický Fourier, vliv prvků rozvoje se s indexem monotónně snižuje
 - cílem metody je transformace $V: X^m \rightarrow Y^p$ vektorů m -dimenzionálního příznakového prostoru vektorů X^m do p -dimenzionálního prostoru vektorů Y^p , kde $p \leq m$, který minimalizuje střední kvadratickou odchylku
 - máme: $\{x_i \in R^m, i=1..m\}$
 - hledáme $\{e_i \in R^m, i=1..p, p < m\}$ pro aproximaci vektorů x_i : $y_i = \sum_{j=1}^p c_{ij} e_j$
 - vtzorce je lin. kombinace vektorů e_i tak aby kvadrát odchylky ϵ byl minimální: $\epsilon^2 = \|x_i - y_i\|^2$
 - při daném počtu členů rozvoje poskytuje ze všech možných aproximací **nejmenší střední kvadratickou odchylku** od pův. vzorů
 - při použití disperzní matice jsou nové transformované příznakové proměnné **nekorelované**
 - členy rozvoje nepřispívají k aproximaci rovnoměrně vliv každého z členů uspořádané posloupnosti aproximace se zmenšuje s jeho pořadím určeným velikostí odpovídajících vlastních čísel
 - velikost chyby aproximace neovlivňuje strukturu rozvoje** změna požadavků na velikost střední kvadratické odchylky nevyžaduje přepočítat celý rozvoj, je třeba pouze změnit počet jeho členů
- princip**
 - ! **Karhunen-Loeuvův rozvoj (transformace)**
- vlastnosti**
 - extrakce 

Kontingenční tabulka - pro dva atributy matice, pro každou kombinaci hodnot četnost, marginály se sumary přes řádky/sloupce

	Y_1	Y_2	...	Y_k	Σ
X_1	a_{11}	a_{12}	...	a_{1k}	$r_{1.}$
X_2	a_{21}	a_{22}	...	a_{2k}	$r_{2.}$
...
X_k	a_{k1}	a_{k2}	...	a_{kk}	$r_{k.}$
Σ	$s_{.1}$	$s_{.2}$...	$s_{.k}$	n

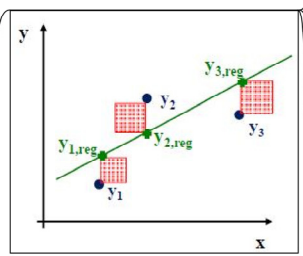
a_{ij} ... četnost (frekvence) kombinace $(X = X_i) \wedge (Y = Y_j)$
 $r_{k.}, s_{.j}$... řádkové, sloupcové součty (tzv. marginální hodnoty)

Statistické testy

- χ^2 -test** - jsou atributy nezávislé? stupně volnosti (skoro počet kombinací).
 Hladina významnosti (pravděpodobnost závěru)
 $\chi^2 = \sum_{i,j} \frac{(a_{ij} - c_{ij})^2}{c_{ij}}$
- Fisherův test** - máme málo dat a jen dva booleovské atributy (= čtyřpolní kontingenční tabulka)
 $P = \frac{r_1! r_2! s_1! s_2!}{n! (a_{11} - i)! (a_{12} + i)! (a_{21} + i)! (a_{22} - i)!}$
 n je celkový součet
 Je-li $P \leq \alpha$, zamítne se nulová hypotéza o nezávislosti na hladině významnosti α
- počet stupňů volnosti** = (počet řádků - 1) * (počet sloupců - 1)
 Laicky: Měříme, jak daleko je očekávaná a reálná četnost, vzhledem k očekávané četnosti
 Je-li χ^2 hodnotavětší než "tabulková" hodnota χ^2 -distribuce, atributy jsou závislé
 χ^2 -test lze použít jen v případě dostatečně velkých četností: pro $(r_{k.}, s_{.j}) / n \geq 5 \forall k, j$

Regresní analýza - určit, jaký vztah má proměnná Y k jedné anebo vícero jiným proměnným X_1, \dots, X_n

- typy závislosti 2 proměnných**
 - funkční vztah - 2 závislé proměnné, určité x odpovídá jediná hodnota y
 - korelace - 2 náhodné (nezávislé) proměnné, síla vztahu popisovaná kor. koef. R
 - regrese - vztah náhodné (nezávislé, vysvětlující) proměnné x a závislé (vysvětlované) proměnné y, má pravděpodobnostní rozdělení
 - pro x a y hledáme přímku (její parametry) vystihující průběh jejich závislosti = je nejbližší všem bodům
 - závislost y na x můžeme popsat (aproximovat) pomocí **regresní rovnice (přímky)**: $y = q_1 x + q_0 + \epsilon$
 - ϵ je regresní odchylka
 - alternativně používaný zápis: $y = bx + a$
- Jednorozměrná regrese** - jaké parametry má lin. závislost mezi 2 veličinami
- Lineární regrese** - jaké parametry má lin. závislost mezi 2 veličinami
- metoda nejmenších čtverců** - minimalizace rozdílů mezi skutečnou a očekávanou hodnotou
- regresní koeficienty**
 - hledáme: $\min \sum_{i=1}^n (y_i - f(x_i))^2 = \min \sum_{i=1}^n (y_i - q_1 x_i - q_0)^2$
 - sumu zderivujeme podle všech q_i a získáme rovnice rovné 0 (jako při hledání minima)
 - odvození**
 - q_1 - směrnice
 - q_0 - absolutní člen
 - \tan úhlu s osou x

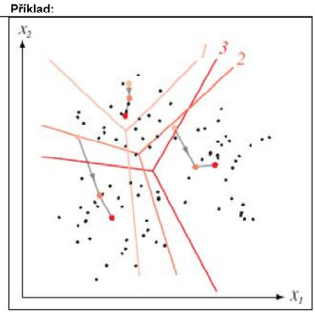


Mnohorozměrná regrese

- Lineární**
 - řešení metodou nejmenších čtverců $\vec{q} = (X^T X)^{-1} X^T \vec{y}$
- Nelineární**
 - Logistická regrese**
 - předpokládáme složitější funkční závislost mezi y a vektorem \vec{x} (kvadratickou, exponenciální)
 - předpokládáme, že závislá veličina y je kategoriální, např. dvouhodnotová
 - modelujeme pravděpodobnost, že y má konkrétní hodnotu v závislosti na kombinaci hodnot nezávislých veličin vektoru \vec{x}
 - podmíněná šance: $P(y|x) / (1 - P(y|x))$

Diskriminační analýza - v podstatě zkoušíme vyrobit perceptron, měříme chybu klasifikace příkladů do zadaných tříd

- Lze pozorované vzory rozdělit do skupin (shluků) vzájemně si blízkých vzorů?
- Volba míry vzdálenosti závislá na měřítku veličin -> veličiny normovat
- Centroid - střed shluku
- Shlukování metodou k-středů (k-means)**
 - Inicializace:** Náhodně zvol rozklad do k shluků
 - Update-step:** Určí centroidy pro všechny shluky v aktuálním rozkladu
 - Assignment-step:** Pro každý vzor x Určí vzdálenosti $d(x, c_k)$
 - Učítá vzdálenosti $d(x, c_k)$
 - c_k - centroid k-tého shluku
 - Nechť $d(x, c_l) = \min_k d(x, c_k)$
 - z této zvláštní konstrukce získáme "index" nejbližšího shluku l
 - není-li x součástí shluku l, přesuň do shluku l
 - nedošlo-li k žádnému přesunu -> KONEC
 - ! U ušní me jen nechala spočítat si nějaký centroid, resp. ukázat jak to udelam.

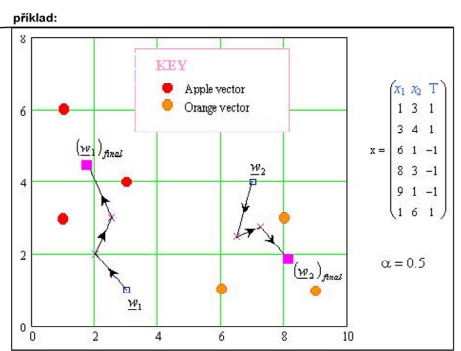


Algoritmus hierarchického shlukování (metodou „zdola nahoru“)

- na začátku je každý bod cluster
- dokud je více než 1 shluk
- Určí vzájemné vzdálenosti mezi všemi clustery
- vždy shluknu dva nejbližší body do nového clusteru

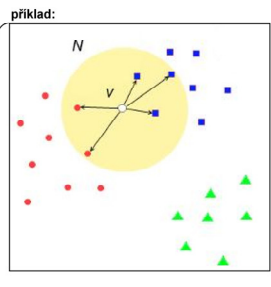
Shluková analýza (klastrování)

- Learning Vector Quantization**
 - reprzentace shluků vektory
 - máme klasifikátor příznakového prostoru, chceme hranice klasifikátoru promítnout i do mapy
 - všech vektorů w_k
 - váha vektoru v čase 0
 - Inicializace:** (0) parametru učení (0)
 - najdi index q tak že $w_{qj}(k)$ je váhový vektor s nejmenší vzdáleností $\sqrt{\sum_{j=1}^m (x_j - w_{qj}(k))^2}$
 - posuň váhový w_k k trénovacímu x
 - pozn.: $(x_i - w_{qi}) = ((x_{i1} - w_{q1}), \dots, (x_{im} - w_{qm}))$
 - C_{xi} je definovaná třída vektoru x_i
 - posuň w od x
 - Hlavní cyklus:** \forall trénovací vstupní vektor x_i
 - if $C_{xi} \neq C_{xi}$ $w_{qj}(k+1) = w_{qj}(k) + \mu(k) [x_{ij} - w_{qj}(k)]$
 - if $C_{xi} = C_{xi}$ $w_{qj}(k+1) = w_{qj}(k) - \mu(k) [x_{ij} - w_{qj}(k)]$
 - posuň parametru učení (k+1)
 - Ukončovací podmínka: nějaký počet iterací nebo nízná hodnota μ



Klasifikace k-Nearest Neighbour

- trénovací množina klasifikovaných vzorů M
- Vstup:** vzor v, který chceme klasifikovat a celé kladné číslo k
- Algoritmus:** M
 - v množině M najdeme množinu N obsahující k vzorů, které jsou k vzorům nejbližší ze všech vzorů z M
 - vzoru v dáme klasifikaci, která je nejčastější v množině N
 - lze použít i na regresi a klastrování
 - probírá se jenom na cvičení



- 4. Rozhodovací stromy
- 5. Asociční pravidla
- 6. Bayesovské modely
- 7. Neuronové sítě
- 8. ELM-sítě
- 9. Genetické algoritmy
- 10. Pokročilé předzpracování dat

Data-Mining

- 1. Metodiky
- 2. Předzpracování dat

Vstup: Vektory příznaků; Výstup: Třídy klasifikace

Vnitřní uzel se větví podle hodnoty svého atributu (s atributy rodičů již zafixovanými)

- klasifikace množiny na základě atributu každého prvku
- posloupnost dotazu na jednotlivé atributy (na rovnost nebo větší/menší)
- chceme mít co nejkratší cestu
- složitost rozhodovacích stromů:** $h \log(n)$... kde h je počet atributu

TDIDT indukce

algoritmus

- Top Down Induction of Decision Trees
- Zvol jeden atribut jako kořen dílčího stromu (nikdy 2x stejný v 1 větvi)
- Rozděl data v tomto uzlu na podmnožiny podle hodnot zvoleného atributu a přidej uzel pro každou podmnožinu
- existuje-li uzel, pro který nepatří všechna data do téže třídy, opakuj pro tento uzel postup od bodu 1; jinak skonči

!!! Entropie - míra neuspořádanosti v systému

$$H = - \sum_i p_i \log p_i$$

pi - relativní četnost třídy i, pro 0 doplníme výsledek limitou = 0

!!! ID3 indukce - Jako TDIDT, ale pracuje na výstupu místo třídy jen s bool hodnotami

algoritmus

- Iterative Dichotomiser 3
- spocítej entropii pro každý atribut množiny S
- rozděl S do podmnožin použitím atributu nejlépe klasifikující množinu (nejnižší součet entropií podmnožin)
- vytvoř uzel stromu obsahující atribut (podmnožin je stejně jako hodnot atributu)
- rekurze na podmnožiny pomocí zbývajících atributů (pokud mají všechny stejnou klasifikaci, nebo dojdou atributy -> list)

!!! C4.5 indukce

- modifikace ID3
- chybějící data: při konstrukci stromu se ignorují; při klasifikaci se odhadnou
- spojitá data: kategorizace podle trénovací množiny
- pruning:
 - validace - rozdělení trénovací množiny: 2/3trénovací, 1/3validační
 - reduced-error pruning: náhrada podstromu listem ohodnocených nejčastěji třídou (na příslušných trénovacích vzorech); nový strom zkontrolujeme na validační množině aby nebyl horší
 - roubování: subtree raising; přenesení nejvíce používanějších podstromů výše
 - rulepost-pruning: konverze stromu (cest do listů) na ekvivalentní sady pravidel; prořezání odstraněním předpokladů pokud to nezhorší očekávanou správnost klasifikace; usporadání pravidel podle přesnosti
 - dolní odhad správnosti: vypočet správnosti na trénovacích datech; vypočet smerodatné odchylky
 - kriterium pro dělení: použijeme atribut s nejvyšším poměrným informačním ziskem (gain ratio)

✗ C5.0/See5: boosting

hlasování několika paralelních klasifikátorů

✗ CART

Classification And Regression Trees

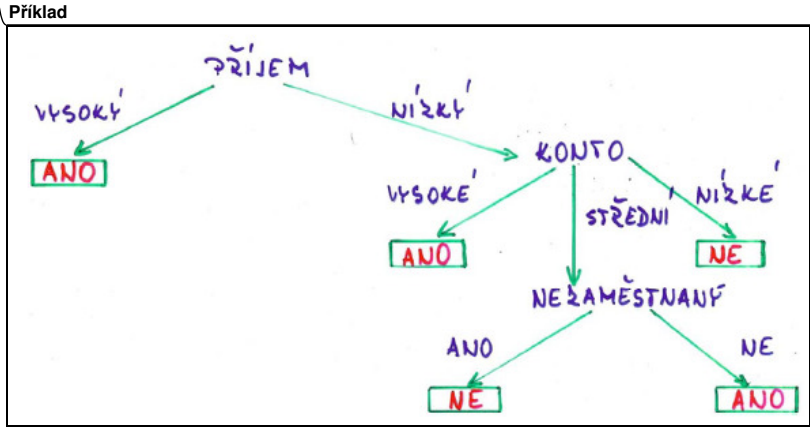
generuje binární rozhodovací stromy

prostor hodnot atributu rozpívotován dělicí hodnotou

Scalable PaRallelizable INduction of decision Trees

SPRINT

- kriterium pro dělení: Gini-index
 - $GINI(D) = 1 - \sum p_i^2$ (D - databáze, pi - frekvence Ci v D)
 - $GINI_{SPLIT}(D) = n1/n GINI(D1) + n2/n GINI(D2)$ (n - počet prvků)
- volba atributu pro dělení: Rainforest (AVC tabulka pro každý uzel)
- indukce rozhodovacího stromu:
 - projdou se trénovací data
 - vytvoříme AVC-tabulku
 - určí se nejlepší atribut pro dělení
 - rozdělení trénovacích dat
 - konstrukce AVC pro následující uzel



- 5. Asociční pravidla
- 6. Bayesovské modely
- 7. Neuronové sítě
- 8. ELM-sítě
- 9. Genetické algoritmy
- ✗ 10. Pokročile předzpracování dat

- ✓ 1. Metodiky
- ✓ 2. Předzpracování dat
- ✓ 4. Rozhodovací stromy

Množina produktů v transakci, které se spolu vyskytují nejčastěji? Virtuální položky pro "metadata".

!! MBA

3 míry [0; 1]

- podpora(pravidlo) = $\frac{\#i\&j}{\#vsech}$
- spolehlivost(pravidlo) = $\frac{\#i\&j}{\#i}$
- zlepseni(pravidlo) = $\frac{p(i\&j)}{p(i)p(j)} = \frac{\#i\&j\#vsech}{\#i\#j}$ 💡 pokud je <1 je pravidlo horší než náhodná volba

hlavní kroky

- zvolte položky na adekvátní úrovni
- vytvořte if-then pravidla z tabulky četností
- určete nejlepší pravidla (zlepšení by mělo být > 1)

Prořezávání, zobecňování (značky, výrobci), disociační pravidla (NOT; pochybná užitečnost), časové řady (rozšíření do okénka)

✓ 5. Asociacní pravidla

!! **Algoritmus APRIORI** - generování asociacních pravidel (frequent itemsets)

Generování kombinací do šířky - kombinace délky k z kombinací délky k-1

Každou kombinaci délky k-1 porovnáme s každou jinou, shodují-li se v k-2 ítelech, spojíme

Spojené zpětně profilujeme, chybí-li nám nějaká jejich podkombinace délky k-1, zahodíme je

Do L1 přiřadíme všechny kategorie, které dosahují alespoň požadované četnosti

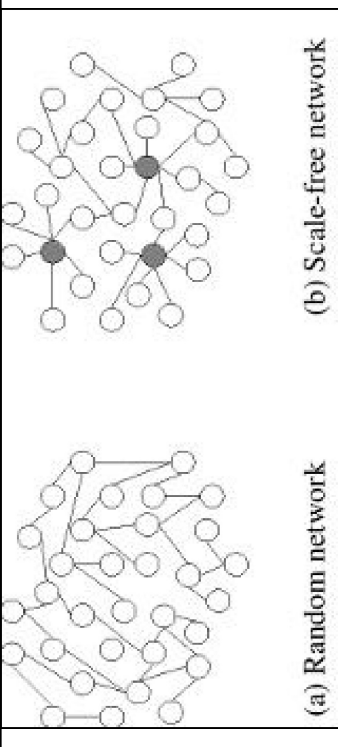
for(k=2; L_{k-1} != {}; k++)
 vygeneruj na základě L_{k-1} kandidáty C_k

kartézský součin L_{k-1} x L_{k-1} a odstraň pokud její podkombinace délky k-1 není obsažena v L_{k-1}

Do L_k zařaď ty kombinace z C_k, které dosáhly alespoň požadované četnosti

return U L_k;

! **Scale-Free síť:** Síť, kde některé uzly mají mnohem více hran (huby) než ostatní



💡 nalezení vztahů mezi údaji, vizualizace linků a vztahů

odolné proti náhodným poruchám, zranitelné při koordinovaném útoku

využití:

- počítačové sítě
- sociální sítě
- monitorování vztahů mezi společnostmi
- marketing

- ✓ 6. Bayesovské modely
- ✓ 7. Neuronové sítě
- ✓ 8. ELM-sítě
- ✓ 9. Genetické algoritmy
- ✗ 10. Pokročile předzpracování dat

Data-Mining

- ✓ 1. Metodiky
- ✓ 2. Předzpracování dat
- ✓ 4. Rozhodovací stromy
- ✓ 5. Asociční pravidla

Data-Mining

6. Bayesovské modely

Naivní Bayesovský klasifikátor: mám jev a klasifikaci, při tréninku počítám četnosti a P(jev|klasifikace)

Připoužití podle Bayesova pravidla počítám podmíněnou P(klasifikace|jev)

$$P(H|E) = \frac{P(H \cap E)}{P(E)} = \frac{P(E|H)P(H)}{P(E)}$$

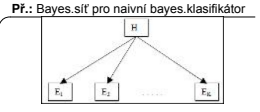
podmíněná=aposteriori

klasifikace pomocí naivního bayesovského klasifikátoru budeme hledat hypotézu s největší aposteriori pravděpodobností

$$H_i = P(H_i) \times \prod P(E_k|H_i) = \frac{\text{Počet vzorů z t}}{\text{Počet všech vzorů}} \times \prod \frac{\text{Počet vzorů z t splňujících } E_k}{\text{Počet vzorů z t}}$$

Naivní, protože předpokládám, že jevy jsou nezávislé; v praxi obvykle nejsou, přesto funguje překvapivě dobře

na půl cesty mezi pravidly a neuronkami



Podmíněná nezávislost: A je podm. nezávislá na B zname-li C, pokud P(A|B,C)=P(A|C)P(B|C)

Bayesovská síť je DAG graf

definice: každému uzlu u (~náhodné veličině) je přiřazeno pst rozdělení tvaru P(u|rodiče(u))

Sdružená pst distribuce celé sítě - P(u1,...,un) = Π P(uj|rodiče(uj))

Inference (pravděpodobnostní odvozování): diagnostická (bottom-up) která příčina má max. aposteriori pst; kauzální (top-bottom) pst jevu za předpokladu nějaké příčiny

Učení

Známa struktura

veličiny plně pozorovatelné (vše vidíme)

veličiny částečně pozorovatelné (některé veličiny=uzly sítě nelze pozorovat nebo šum)

Neznáma struktura

veličiny plně pozorovatelné (prohledávání prostoru modelů)

veličiny částečně pozorovatelné (průsvih)

metoda max. věrohodného odhadu

spočítat z dat odhady podmíněných pravděpodobnostních distribucí pro jednotlivé uzly sítě

gradientní metoda: výpočet gradientu funkce ln P(D|h)

EM algoritmus (expectation maximisation): Inicializace: náhodně zvol parametry distribuce

Maximalizace: z nich určí max. věrohodný odhad parametrů distribuce

Sdružená distribuce: P(Z,K,D,M) = P(Z)P(K|Z)P(D|Z)P(M|K,D)

tabulka podm. pstí uzlů:

Z	P(K=0)	P(K=1)
0	0.5	0.5
1	0.9	0.1

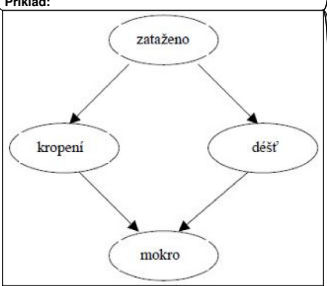
P(K|Z)

Z	P(D=0)	P(D=1)
0	0.8	0.2
1	0.2	0.8

P(D|Z)

K	D	P(M=0)	P(M=1)
0	0	1.0	0.0
1	0	0.1	0.9
0	1	0.1	0.9
1	1	0.01	0.99

P(M|K,D)



diagnostická (bottom-up) inference:

$$P(K=1|M=1) = \frac{P(K=1, M=1)}{P(M=1)} = \frac{\sum_{z,d} P(Z=z, K=1, D=d, M=1)}{P(M=1)} = \frac{\sum_{z,d} (P(Z=z)P(K=1|Z=z)P(D=d|Z=z)P(M=1|K=1, D=d))}{P(M=1)} = \frac{0.2781}{0.4581}$$

Maximální aposteriori pravděpodobnost má tedy příčina déšť: P(D=1|M=1) > P(K=1|M=1).

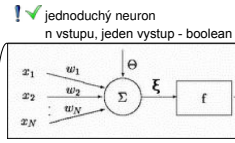
kauzální (top-down) inference:

$$P(M=1|Z=1) = \frac{P(M=1, Z=1)}{P(Z=1)} = \frac{\sum_{k,d} P(Z=1, K=k, D=d, M=1)}{P(Z=1)} = \frac{P(Z=1) \sum_{k,d} (P(K=k|Z=1)P(D=d|Z=1)P(M=1|K=k, D=d))}{P(Z=1)} = 0.7452$$

- ✓ 7. Neuronové sítě
- ✓ 8. ELM-sítě
- ✓ 9. Genetické algoritmy
- ✗ 10. Pokročile předzpracování dat

- 1. Metodiky
- 2. Předzpracování dat
- 4. Rozhodovací stromy
- 5. Asociční pravidla
- 6. Bayesovské modely

Jednoduchý perceptron je výpočetní jednotka s prahem θ , která pro vstupy $x_1, \dots, x_n \in \mathbb{R}$ a váhy w_1, \dots, w_n dává výstup 1, jestliže platí $\sum w_i x_i \geq \theta$ (tzn. $w \cdot x \geq \theta$) a jinak 0.



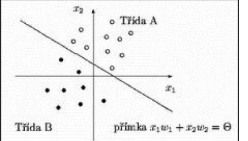
práh neuronu: $\theta \in \mathbb{R}$ zpravidla se v literatuře označuje jako theta (θ, Θ, Θ)
potenciál neuronu: $\xi = \sum w_i x_i + \theta$ označuje se jako $\xi(x)$ nebo jako u
přenosová (aktivační) funkce: $f(\xi)$
 typy výstupů: $f = 1$ neuron je aktivní, $f = 1/2$ tichý, $f = 0$ pasivní
 příklady: sgn, jednotkový skok, sigmoidea: $\frac{1}{1 + e^{-\xi}}$
výstup neuronu: $y = f(\xi = \sum w_i x_i + \theta)$

Příklad: množina co není lin. separabilní

Perceptrony

(absolutně) lineárně separabilní. A a B - v n-rozměrném prostoru, pokud $\exists w_1, \dots, w_n, \theta \in \mathbb{R}$ že $\forall (x_1, \dots, x_n) \in A$ splňuje $\sum w_i x_i (\geq \theta)$ a $\forall (x_1, \dots, x_n) \in B$ splňuje $\sum w_i x_i < \theta$

Algoritmus učení (Hebb)
 hledáme dělící nadrovinu a na ní kolmý vektor w



Inicializace: náhodně zvolíme vstupní váhy $w_i(0) (1 \leq i \leq n+1)$
 předložíme požadovaný výstup pro trénovací vzor x
 spočteme skut. výstup $y(t) = \text{sgn}(\sum w_i(t)x_i(t))$
Učení: **adaptační pravidlo:** $w_i(t+1) = w_i(t) + \eta \cdot x_i(t) \cdot (d(t) - y(t))$
 $0 \leq \eta \leq 1$ je koeficient učení (η je eta)
 Chybně klasifikovaný vektor uvnitř (ve směru váhového vektoru) odtáhneme (otáčíme w od něj)
 Chybně klasifikovaný vektor vně přitáhneme (otáčíme w k němu)
 pokud t nedosáhli požadované hodnoty, opakuj (pro další trénovací dvojice x a $d(t)$)
 zkráceně: postupně beru prvky z trénovací množiny, v případě nesprávné klasifikace otáčím vektor (přitahám k němu prvek) kódu prvku (vynásobím učícím parametrem $\eta < 1$)

Konvergence učení - P a N jsou konečné a lin. separabilní množiny => perc.alg učení provede konečný počet aktualizací vektoru w_i

Přihrádkový (ratchet) algoritmus - vždy si pamatují váhový vektor, který zatím prošel nejvíce úspěšnými klasifikacemi za sebou, až zastavím učení, vydám ten

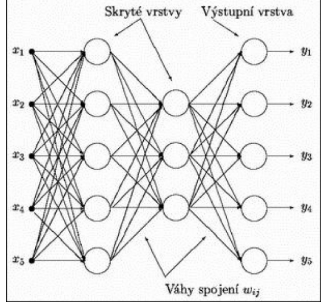
Vícevrstvá neuronová síť, potřebují propagovat updaty váhových vektorů

Použití: klasifikace obrazů, aproximace funkcí, řízení, predikce časových řad

neexistují zpětné vazby ani vazby v rámci 1 vrstvy
 vnitřní (skryté) neurony
 výstupní neurony

Chybová fce: Minimalizují na trénovací množině odchylku přes výstupní neurony
 $E = 1/2 \sum_p \sum_j (y_{jp} - d_{jp})^2$
 chceme minimalizovat chybu $E = y - d$
 p - vzory
 j - výstupní neurony
 y_{jp} - skutečná odezva
 d_{jp} - požadovaná odezva
Váhy update: $w_{ij}(t+1) = w_{ij}(t) + \Delta w_{ij}(t)$
 proti směru gradientu chybové funkce
 od výstupní vrstvy směrem ke vstupní
 proto se metoda jmenuje **back-prop**
 $\Delta w_{ij}(t) = \eta \cdot \delta_j(t) \cdot y_i(t)$ - přírůstek váhy w_{ij} k minimalizaci E

Obrázek

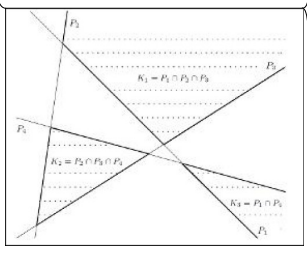


Back-propagation

Adaptační pravidla (odvození) aktualizace synaptických vah pro:

Výstupní neuron (odvození):
 $\Delta w_{ij} = -\frac{\partial E}{\partial w_{ij}} = -\frac{\partial E}{\partial \xi_j} \cdot \frac{\partial \xi_j}{\partial w_{ij}} = -\frac{\partial E}{\partial \xi_j} \cdot y_i = -\delta_j \cdot y_i$
 parc. derivace - $\partial E / \partial w_{ij}$ určuje směr největšího úbytku fce
 spec. případem je hledání minima fce pomocí nulových bodů její derivace
Vnitřní neuron (odvození):
 $\Delta w_{ij} = -\frac{\partial E}{\partial w_{ij}} = -\sum_k \frac{\partial E}{\partial \xi_k} \frac{\partial \xi_k}{\partial w_{ij}} = -\sum_k \delta_k \cdot w_{jk} \cdot \frac{\partial \xi_k}{\partial w_{ij}} = -\sum_k \delta_k \cdot w_{jk} \cdot f'(\xi_k) \cdot y_i = \delta_j \cdot y_i$

Algoritmus trénování hledáme průniky nadrovin

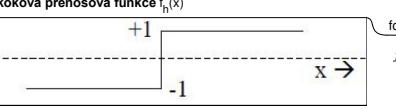


Inicializace: náhodně zvolíme váhy $w_i(0) (1 \leq i \leq n+1)$
 předložíme požadovaný výstup pro trénovací vzor x
 spočteme skut. výstup $y_j = f(\xi_j) = \frac{1}{1 + e^{-\xi_j}}$ kde $\xi_j = \sum w_{ij} x_i$
 výstup pak tvoří vstup na další vrstvě
Učení: **adaptační pravidlo:** $w_{ij}(t+1) = w_{ij}(t) + \eta \delta_j y_i + \alpha_m (w_{ij}(t) - w_{ij}(t-1))$
 $0 \leq \eta \leq 1$ je koeficient učení (η je eta)
 $\alpha_m \leq 1$ - moment učení
 δ_j je chyba učení
 $\delta_j = \begin{cases} \lambda y_j (1 - y_j) (d_j - y_j) & \text{pro výstupní uzly} \\ \lambda y_j (1 - y_j) \sum_k \delta_k w_{jk} & \text{pro skryté uzly} \end{cases}$
 k - index pro neuron následující za neuronem j
 λ - je strmost přenosové funkce
 máme opět trénovací množinu, pro x spočteme p , $f(p) = y$, porovnáme s d
 opakuj dokud chyba není menší než stanovená hodnota

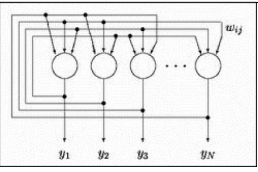
jednovrstvá rekurentní síť s pevnými vahami

Bipolární vstupy i výstupy $\{+1, -1\}$ mohou být binární

Použití: asociativní paměť, optimalizační úlohy, rekonstrukce neúplných a šumem poškozených obrazů



Skoková přenosová funkce $f_H(x)$
 formálně: $f_H(x) = \begin{cases} 1 & \text{pro } x > 0 \\ -1 & \text{pro } x \leq 0 \end{cases}$



každé 2 neurony jsou spojené symetrickými vahami $w_{ij} = w_{ji}$, každý neuron je vstupní i výstupní

Hopfieldovy sítě

Učení (Hebb): "what fires together, wires together"
 $w_{ij} = \begin{cases} \sum_k x_i^k x_j^k & \text{pro } i \neq j \\ 0 & \text{pro } i = j \end{cases}$
 w_{ij} je váha synapse mezi neurony i a j
 x_i^k je i -tá složka s -tého vektoru

Algoritmus výpočtu
 Inicializace: $y_i(0) = x_i$ předložíme neznámý vstupní vzor
 $y_i(t)$ výstup neuronu i v čase t
 Adaptační pravidlo: neurony se iterativně updatují
 $y_i(t+1) = f_H(\sum_j w_{ij} y_j(t))$ $1 \leq i \leq n$
 opakujeme dokud se výstupy neuronů neustálí
 výstupy pak reprezentují trén. vzor nejlépe odpovídající předloženému vzoru
 taky může zůstat oscilující mezi dvěma stavy

Konvergence při vybavení - energetická funkce neustále klesá

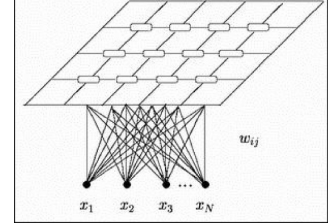
Energetická funkce vyjadřuje energii sítě ve stavu
 $E(x) = -\frac{1}{2} \sum_{i,j} w_{ij} x_i x_j$
 Věta: síť dosáhne stabilního stavu v lokálním minimum energetické funkce

Kohonenovy mapy

Self-organizing maps (SOM)

Algoritmus učení
 Inicializace: náhodně zvolím váhy w_{ij} je váha mezi vstupním neuronem i a výstupním j
 předložíme trénovací vzor
 vypočtu vzdálenost mezi vst. a váhovým vektorem každého neuronu a předám kompetiční vrstvě
 $d_i = \sum_j (x_j(t) - w_{ij}(t))^2$
 vyberu výstupní neuron c s minimální d_i a označím ho jako vítěze
Učení: **Adaptační pravidlo v okolí c**
 $\Delta w_{ij} = \alpha(t) \Phi(c, j) (x_i(t) - w_{ij}(t))$
 vigilační koef. α se s časem snižuje
 vigilance ~ bdělost
 opakuj pro všechny vstupní vektory dokud se váhy mění
 Kompetice neuronů o vzorek, vítěz inhibuje ostatní
 Neurony topologicky rozprostřeny, soutěží o váhový vektor nejbližší vstupu
 chceme aby neurony s váhovým vektorem blízko v prázdném prostoru byly vedle sebe i ve výstupním prostoru

Obrázek



Použití: shlukování, analýza dat, vytváření sémantických map

Pro učení s učitelem
 learning vector quantization
 máme klasifikátor prázdného prostoru, chceme hranice klasifikátoru promítnout i do mapy
 LVQ1: adaptují pouze nejbližší neuron, ale v případě špatné třídy v neuronu opacným směrem!
 LVQ2.1: adaptují, jsou-li dva nejbližší neurony v různých třídách (jeden v mojí, druhý v cizí) a jsem v rámci okenka dělicí nadplochy
 LVQ3: trosku adaptují i když jsou oba sousedé v me třídě

- 8. ELM-sítě
- 9. Genetické algoritmy
- 10. Pokročile předzpracování dat

Data-Mining

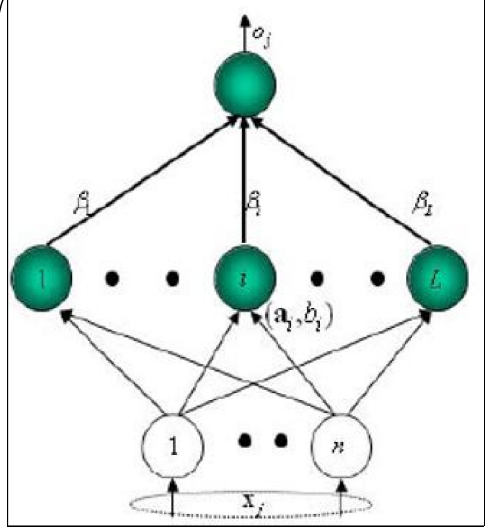
Učení s učitelem

Učení bez učitele

Data-Mining

- ✓ 1. Modely
- ✓ 2. Předzpracování dat
- ✓ 4. Rozhodovací stromy
- ✓ 5. Asociční pravidla
- ✓ 6. Baysovské modely
- ✓ 7. Neuronové sítě

Model



- 1 výstupní neuron
- L skrytých neuronů
- n výstupních neuronů
- a_i váhový vektor skrytého neuronu
- b_i práh i-tého skrytého neuronu
- dopředná síť s jednou skrytou vrstvou neuronů s libovolnou počtem částech
- spojitou přenosovou funkcí $G(a_i, b_i, x)$
- β_j váha mezi i-tým skrytým a výstupním neuronem

Výstup sítě: $f_L(x_j) = \sum_{i=1}^L \beta_i G(a_i, b_i, x_j) = t_j$

můžeme zapsat jako $H\beta = T$

$$H = \begin{bmatrix} G(a_1, b_1, x_1) & \dots & G(a_L, b_L, x_1) \\ \dots & \dots & \dots \\ G(a_1, b_1, x_N) & \dots & G(a_L, b_L, x_N) \end{bmatrix}$$

N je počet trénovacích vzorů
t jsou požadované výstupy

✓ 8. ELM-sítě

! ✓ algoritmus učení

- Trénovací množina $\mathcal{X} = \{(x_i, t_i) \mid x_i \in \mathbb{R}^n, t_i \in \mathbb{R}^m, i=1, \dots, N\}$
- Vstup:** výstupní funkce skrytých neuronů $G(a, b, x)$
počet skrytých neuronů L
 - Inicializace:** náhodně zvol parametry skrytých neuronů (a_i, b_i)
Spočítej matici výstupů skryté vrstvy H
 - Učení:** Spočítej váhy pro výstupní neurony $\beta: \beta = H^\dagger T$
- H[†] označuje Moore-Penroseovu pseudoinverzní matici pro matici výstupů skryté vrstvy neuronové sítě H
protože inverzní matice ne vždy existuje (čtvercová singulární, nebo obdélníková) tak MP matice je prostě nejbližší matice k matici inverzní kdyby existovala
jednoduchý neiterativní výpočet/učení

výhody

- jednoduchý neiterativní výpočet, který se skládá ze 3 kroků (mnohé z tradičních algoritmů učení jsou výrazně složitější)
- velmi rychlý a přímočarý proces učení (není třeba řešit otázky lokálních minim, vhodné inicializace parametrů, přeučení a další)
- parametry skrytých neuronů a_i a b_i jsou nezávislé na trénovacích datech na sobě navzájem
- mohou generovat parametry skrytých neuronů již před předložením trénovacích dat (konvenční metody učení musí trénovací data „vidět“ předtím, než nastaví parametry skrytých neuronů)
- lze použít s libovolnou omezenou po částech spojitou nekonzstantní přenosovou funkcí (gradientní algoritmy učení pracují jen s hladkými přenosovými funkcemi)

✓ 9. Genetické algoritmy

! Jednoduchý genet. algoritmus

Chromozom (genotyp) - řetězec symbolů kódující vlastnosti jedince (fenotyp)

Fitness funkce $\Phi: \text{genotyp} \rightarrow \mathbb{R}$

Populace P(t) umělých chromozomů

! jak dobře řeší problém, čím vyšší tím lepší

Inicializace: vygeneruj populaci $P(0) = \{x_1, \dots, x_N\}$ náhodných chromozomů

spočítej fitness všech

vytvoř populaci $P(t+1)$

selekcí **reprodukce ruletou** - naškálujeme jedince na ruletu podle fitness a otočíme ruletu N-krát
elitismus - S nejlepších jedinců kopírujeme vždy

Iterace:

- proved **křížení** ! náhodně se spárují a zkříží podle zadané pravděpodobnosti
- mutuj** chromozómy - náhodná změna bitů podle zadané pravděpodobnosti

jednobodové

vícibodové

skončí pokud: najdeme hledaného jedince nebo fitness nejlepšího neroste

! může popisovat část chromozomu zaručující vysokou fitness

Schema - abstraktní "šablona" genomu, ternární řetězec délky m, kromě 1/0 jeste * (wildcard)

- řád schématu $o(H)$** - počet pevných pozic
- délka schématu $\delta(H)$** - vzdálenost mezi první a poslední pevnou pozicí
- fitness schématu $f(H)$** - průměrná fitness odp. řetězců v populaci

! **Věta o schématech** - v populaci během GA (exponenciálně) roste počet schémat, která jsou krátká, mají nízký řád a mají nadprůměrnou fitness

! **Idea důkazu:** Díváme se postupně pro všechny genetické operátory, na čem závisí očekávaná četnost schématu v další populaci - ukazuje se, že právě na délce atd.

✗ 10. Pokročile předzpracování dat