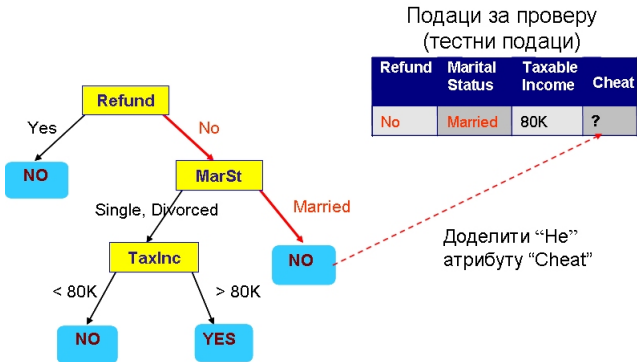


Класификација

Ненад Митић

Математички факултет
nenad@matf.bg.ac.rs

Примена модела - пример (наставак)



Конструкција дрвета одлучивања

Изазови

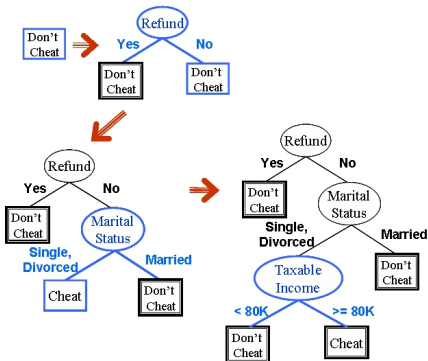
- Како одабрати атрибут(е) по коме се врши подела?
- Како формирати упит за различите типове атрибута?
- Како одредити најбољу поделу?
- На који начин индуктивно применити претходне критеријуме на дрво по дубини?
- Када стати са конструкцијом дрвета?

Конструкција дрвета одлучивања

Изазови

- Како одабрати атрибут(е) по коме се врши подела?
- Како формирати упит за различите типове атрибута?
- Како одредити најбољу поделу?
- На који начин индуктивно применити претходне критеријуме на дрво по дубини?
- Када стати са конструкцијом дрвета?
- Начин рада са недостајућим вредностима?
- Који је критеријум за процену грешке у генерализацији?
- Цена и перформансе модела?

Хантов алгоритам - илустрација



<i>Tid</i>	Refund	Marital Status	Taxable Income	Cheat
1	Yes	Single	125K	No
2	No	Married	100K	No
3	No	Single	70K	No
4	Yes	Married	120K	No
5	No	Divorced	95K	Yes
6	No	Married	60K	No
7	Yes	Divorced	220K	No
8	No	Single	85K	Yes
9	No	Married	75K	No
10	No	Single	90K	Yes

Избор атрибута за поделу

- Атрибут за поделу материјала се бира стратегијом похлепе (енг. *greedy*)
 - Поделити слоге према атрибуту који оптимизује одређени критеријум
- Одлуке које треба донети
 - Како поделити слоге и која је најбоља подела?
 - Како навести услове за тестирање атрибута
 - Када стати са деобом

Мере нечистоће - примери израчунавања

Чвор H_1	Број
Класа 1	0
Класа 2	6

$$\text{Гини} = 1 - (0/6)^2 - (6/6)^2 = 0$$

$$\text{Ентропија} = -(0/6)\log_2(0/6) - (6/6)\log_2(6/6) = 0$$

$$\text{Грешка класификације} = 1 - \max[0/6, 6/6] = 0$$

Чвор H_2	Број
Класа 1	1
Класа 2	5

$$\text{Гини} = 1 - (1/6)^2 - (5/6)^2 = 0.278$$

$$\text{Ентропија} = -(1/6)\log_2(1/6) - (5/6)\log_2(5/6) = 0.650$$

$$\text{Грешка класификације} = 1 - \max[1/6, 5/6] = 0.167$$

Чвор H_3	Број
Класа 1	3
Класа 2	3

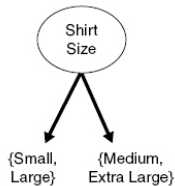
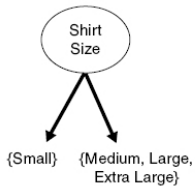
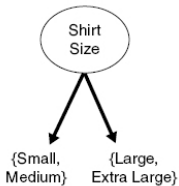
$$\text{Гини} = 1 - (3/6)^2 - (3/6)^2 = 0.5$$

$$\text{Ентропија} = -(3/6)\log_2(3/6) - (3/6)\log_2(3/6) = 1$$

$$\text{Грешка класификације} = 1 - \max[3/6, 3/6] = 0.5$$

Подела заснована на редним атрибутима

- Бинарна подела или подела на више грана
- Да ли су све поделе коректне? Очување редоследа!



Подела по категоричким атрибутима - израчунавање Гини индекса

- За сваку од различитих вредности атрибута израчуна се број слогова у свакој класи података
- Одреди се Гинијев индекс за сваки чвор и пронађе минимална вредност

Подела у више грана

	CarType		
	Family	Sports	Luxury
C1	1	2	1
C2	4	1	1
Gini	0.393		

Подела у две групе
(налази се најбоља подела по вредности)

	CarType	
	{Sports, Luxury}	{Family}
C1	3	1
C2	2	4
Gini	0.400	

	CarType	
	{Sports}	{Family, Luxury}
C1	2	2
C2	1	5
Gini	0.419	

Подела по непрекидним атрибутима - израчунавање Гини индекса

За ефикасно израчунавање за сваки атрибут се

- врши сортирање по вредностима
- добијене вредности линеарно скенирају уз ажурирање броја и слогова и рачунање Гинијевог индекса
- бира се позиција за поделу са најмањим Гинијевим индексом

Cheat	No		No		No		Yes		Yes		Yes		No		No		No					
	Taxable Income																					
	60		70		75		85		90		95		100		120		125		220			
Сортиране вредности →	55	65	72	80	87	92	97	110	122	172	230											
Позиције поделе →	<=	>	<=	>	<=	>	<=	>	<=	>	<=	>	<=	>	<=	>	<=	>	<=	>		
Yes	0	3	0	3	0	3	0	3	1	2	2	1	3	0	3	0	3	0	3	0	3	0
No	0	7	1	6	2	5	3	4	3	4	3	4	3	4	4	3	5	2	6	1	7	0
Gini	0.420		0.400		0.375		0.343		0.417		0.400		0.300		0.343		0.375		0.400		0.420	

Алтернативни критеријум поделе

- Мере као што су Гинијев индекс и ентропија фаворизују атрибуте са великим бројем различитих вредности
- У неким случајевима се добијају погрешни резултати (нпр. број индекса, матични број, ...)
- Стратегије превазилажења проблема
 - употреба искључиво бинарне поделе
 - укључивање броја могућих вредности

$$\text{Однос добити за поделу} = \frac{\Delta_{podela}}{Split_info}$$

$$Split_info = - \sum_{i=1}^k P(v_i) \log_2 P(v_i)$$

где је k број подела атрибута

Модел алгоритма за формирање дрвета одлучивања

```
/* Skup atributa F; skup trening slogova E */
Formiranje drveta(E, F )
if uslov_zaustavljanja(E,F) then
  list=formiraj_cvor()
  labela lista=Klasifikuj(E)
  return list
else
  koren=formiraj_cvor()
  koren.test_uslov = najdi_najbolju_podelu(E,F)
  let V={v|v je moguci ishod za koren.test_uslov}
  for each v in V do
    E_v = {e | koren.test_uslov(e) = v i e pripada E}
    dete = Formiranje_drveta(E_v,F)
    dodaj dete cvor kao potomak korenog cvora i
      oznaci granu (koren-dete) sa v
  end for
end if
return koren

/* Posle formiranj drvo treba da se potkrese */
```

Критеријум заустављања

Формирање нових чвора деобом слога текућег чвора се зауставља (функција `uslov_zaustavljanja` из претходног алгоритма) када

- сви слогови у текуће, чвору припадају истој класи
- сви слогови имају исте вредности (свих) атрибута
- испуњен критеријум ранијег заустављања (достигнута дубина дрвета,)

Карактеристике дрвета одлучивања

- Јевтина за конструисање
- Јако брза у класификацији непознатог материјала
- Лака за интерпретацију за дрвета мале величине
- Прецизност је упоредива са осталим техникама класификације за једноставне типове података
- Избор мере нечистоће нема велики утицај на перформансе

Карактеристике дрвета одлучивања

- Применљивост (на све типове података)
- Изражајност (могу да представе сваку функцију дискретних атрибута)
- Ефикасност израчунавања
- Рад са недостајућим вредностима
- Рад са ирелевантним атрибутима и редундантним атрибутима

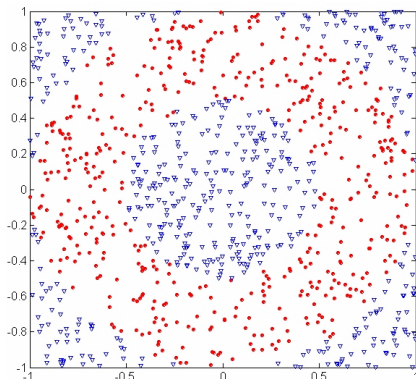
Карактеристике дрвета одлучивања

- Слаба могућност рада са повезаним атрибутима (спајање два атрибута даје значајну информацију)
- Критеријум поделе само један атрибут → интервал је правоугаоног облика
- Критеријум поделе комбинација више атрибута → могуће је обрадити и нпр. *косе* податке
- Избор начина поткресивања значајно утиче на резултате
- Проблем преприлагођавања и потприлагођавања

Преприлагођеност и потприлагођеност

- Модел који исувише добро класификује податке за тренинг може да има лошије карактеристике при генерализацији од модела који има већу грешку у процесу прављења модела - превише прилагођен модел (преприлагођен модел, енг. *model overfitting*)
- Ако је модел исувише једноставан грешке при тренирању и уопштавању могу да буду јако високе - премало прилагођен модел (потприлагођен модел, енг. *model underfitting*)

Преприлагођеност и потприлагођеност



500 кружних и 500
троугаоних тачака

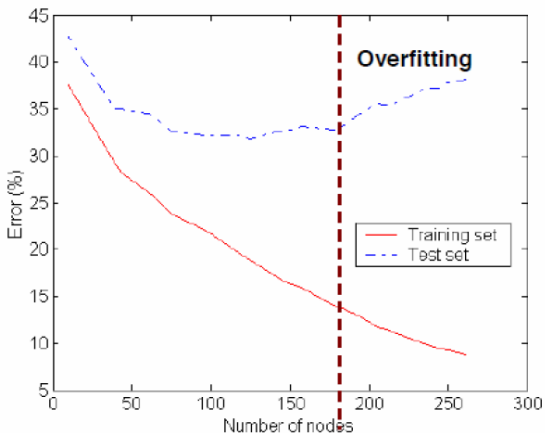
Кружне тачке:
 $0.5 \leq \sqrt{x_1^2 + x_2^2} \leq 1$

Троугаоне тачке:
 $\sqrt{x_1^2 + x_2^2} < 0.5$ или
 $\sqrt{x_1^2 + x_2^2} > 1$

30% тачака се бира за
тренинг, остале за
тест

Дрво са Гинијевим
индексом као мером
се примењује на
тренинг податке

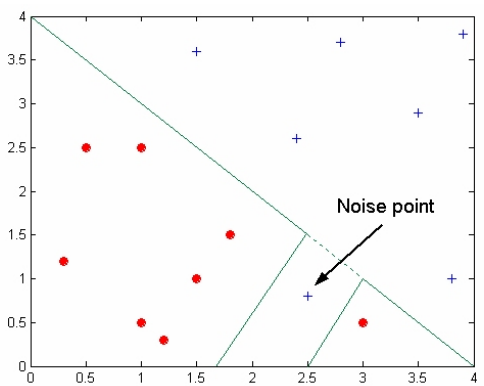
Преприлагођеност и потприлагођеност



Потприлагођавање: ако је модел исувише једноставан и тренинг и тест грешка су велике.

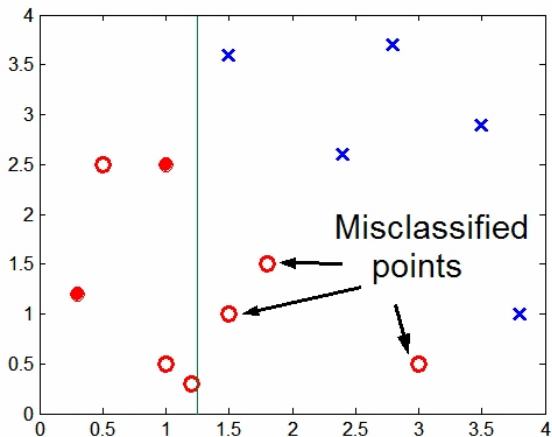
Поткресивањем дрвета на различитим нивоима

Преприлагођеност и потприлагођеност



Границе поделе могу да се искриве и због постојања шума

Преприлагођеност и потприлагођеност



Преприлагођеност може да се јави због непостојања репрезентативних примера

Процена грешке у генерализацији

- Нека је T дрво, t чвор, N број листова у дрвету T , $e(t)$ број погрешно класификованих слогова у t , и $e(T)$ укупан број грешака у класификацији по дрвету T
- Грешка поновне замене (грешка при тренирању)
 $\min(\sum e(t))$ - бира се модел са најмањом грешком при тренирању
- Грешка при генерализацији: грешка при тестирању
 $\sum e'(t)$

Процена грешке у генерализацији

На пример, за $\Omega = 0.5$ укупан број грешака је $e'(T) = e(T) + N \times 0.5$. За дрво са 30 листова и 10 грешака на тренинг скупу од 1000 ставки:

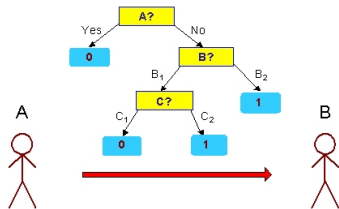
- Грешка тренирања = $10/1000$ (1%)
- Грешка уопштавања = $(10 + 30 \times 0.5)/1000$ (2.5%)
- Смањење грешке поткресивањем (дрвета)
 - Користи се посебан скуп података за процену грешке уопштавања

Процена грешке у генерализацији

- Што је модел сложенији веће су шансе за преприлагођеност
- Параметар процене је и сложеност модела
- Модел не треба непотребно усложњавати
- Од два модела са сличном грешком генерализације изабрати онај који је једноставнији - Окамов (*Осст*) принцип (жилет)
- Код сложених модела већа је шанса да случајно дође до уклапања због грешаак у подацима

Принцип најмање дужине описа (*MDL*)

X	y
X ₁	1
X ₂	0
X ₃	0
X ₄	1
...	...
X _n	1



X	y
X ₁	?
X ₂	?
X ₃	?
X ₄	?
...	...
X _n	?

- Начин да се укључи сложеност модела
- Дат је скуп слогова са познатим атрибутом X
- Особа A познаје све вредности класа атрибута y
- Особа B нема ту информацију и може да пошаље захтев A да проследи све ознаке класа - $\Theta(n)$ битова где је n број слогова
- Алтернативно, A може да направи модел, кодира и пошаље B

Баратање преприлагођавањем у индукцији по дрвету

Поткресивање по завршетку

- Дрво одлучивања расте до крајњих граница
- Исеку се чворови у дрвету од дна ка врху
- Ако се грешка генерализације побољша после отсецања поддрво се замени са чвором који је лист
- Ознака класе листа се одређују према већини класа инстанци поддрвета
- За поткресивање по завршетку се може користити и *MDL*

Пример поткресивања по завршетку

Class = Yes	20
Class = No	10
Error = 10/30	

Тренинг грешка (пре деобе) = 10/30

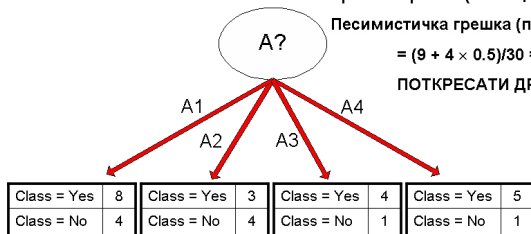
Песимистичка грешка = $(10 + 0.5)/30 = 10.5/30$ за $\Omega=0.5$

Тренинг грешка (после деобе) = 9/30

Песимистичка грешка (после деобе)

$$= (9 + 4 \times 0.5)/30 = 11/30$$

ПОТКРЕСАТИ ДРВО!



Руковање атрибутима са недостајућим вредностима

Утицај НВ на дрво одлучивања

- Шта ако атрибут по коме се дели чвор има НВ (како рачунати меру нечистоће)
- Како дистрибуирати инстанце са НВ на децу чворове
- Како класификовати тест инстанцу са НВ

Мерење преформанси израчунавања

Утицај НВ на дрво одлучивања

- Нагласак је на предвиђачким особинама модела
 - предност у односу на брзину класификације, изградње модела, скалабилност, ...
- Матрица конфузије

		Предвиђена класа	
		Класа=Да	Класа=Не
Стварна класа	Класа=Да	a	b
	Класа=Не	c	d

- a: TP (true positive)
- b: FN (false negative)
- c: FP (false positive)
- d: TN (true negative)

Мерење преформанси израчунавања

	Предвиђена класа		
	Класа=Да	Класа=Не	
Стварна класа	Класа=Да	a TP	b FN
	Класа=Не	c FP	d TN

Најчешће коришћена метрика

$$\text{Тачност} = \frac{a + d}{a + b + c + d} = \frac{TP + TN}{TP + FN + FP + TN}$$

Ограничења прецизности

- Размотримо проблем 2 класе
 - Број примерака класе 0 = 9990
 - Број примерака класе 1 = 10

- Ако модел предвиђа да ће сви слогови бити класе 0
тада је прецизност $9990/10000 = 99.9 \%$
 - Прецизност нема значаја (ни смисла) јер модел не
открива ни један слог класе

