

# Paskaita 4

## Klasifikācija ir regresīvā medicīnā

Naprindotam grāvis algoritms (cit-tādumam, kas ir ļoti efektīvs (dažādi optimāli) algoritmi parasti saņemti ir reāli, nenoņemot netīšas datus struktūras.

Tarp jū vērā populāriem ir ja medzīo datus struktūra.

Šajā paskaitē ~~ir~~ klasifikācija uzdevumi spriedumi paritēnā medzīs - sudarīšnie klasifikācija medicīnā (algoritms)

Analogiski galā būs neapņemas ir regresīvā medicīnā

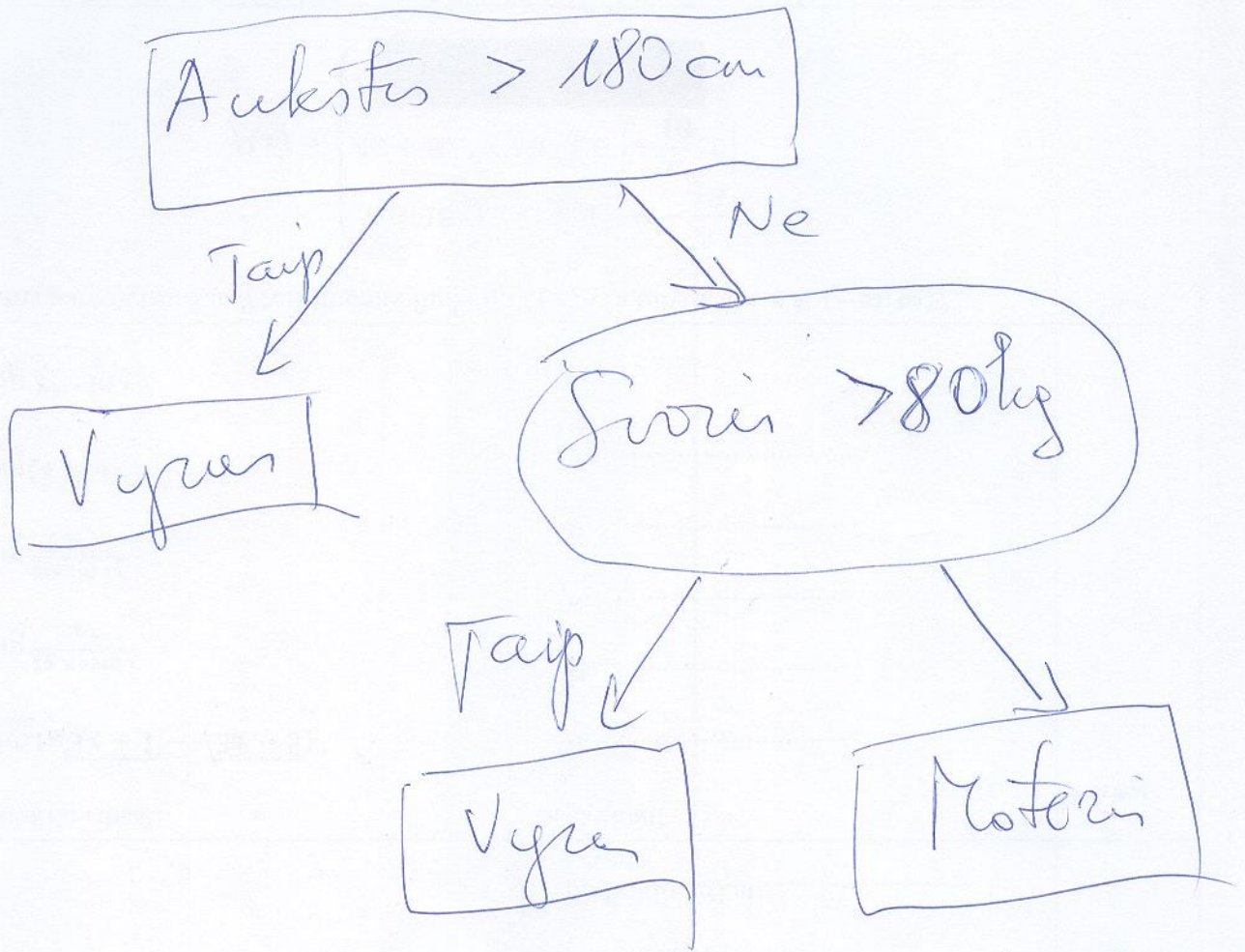
Jie bus naudojami kartu su  
mūsų jau išnagrinėtais regresiniais  
klasifikavimo algoritmais. Šios  
klasės algoritmus vadiname

CART (Classification And  
Regression Trees).

Šio duomenų struktūra padės mums  
atlikti /pildyti sprendimą, kolektai  
klasių paskirsto duomenys, todėl  
tokius algoritmus dar vadiname  
sprendimo medžiais (decision trees)

Nagrinėdami pavyzdį, kai  
turime du klases Vyrai / Moterys.  
keičiu kintamųjų parametrus (Aukštis,  
Svoris).

# Klasifikatorius

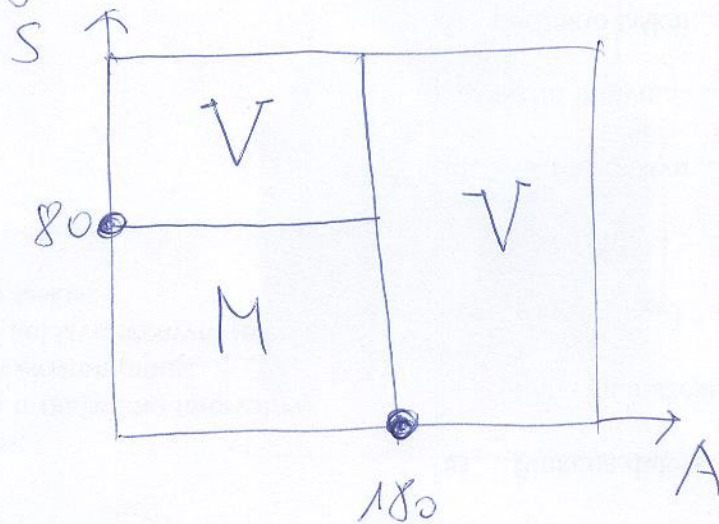


Teigu svorine šley klasifikatorius  
(dvejetauis medis) - je svorine  
is auksitu nelaukustu, sau per  
galvime įvertitu šleyis pradedu  
duomenis (164, 62)  
Ats.: ???

Skaidymo žingsniuose (aprašymo metu)

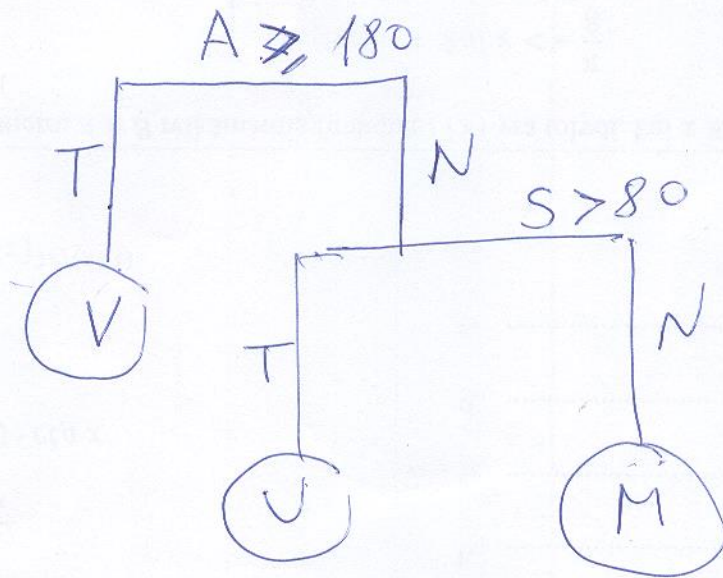
- teorine parinkti kiti atlikime skaidymų
- Kokios yra skaidymo parametro reikšmės

Skelbime, keičiant fiksuotą aprašytą nedidelių duomenų aibę



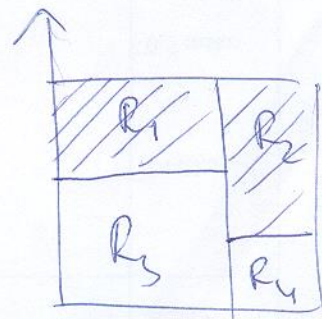
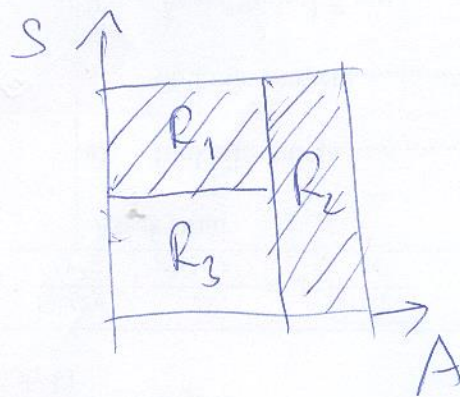
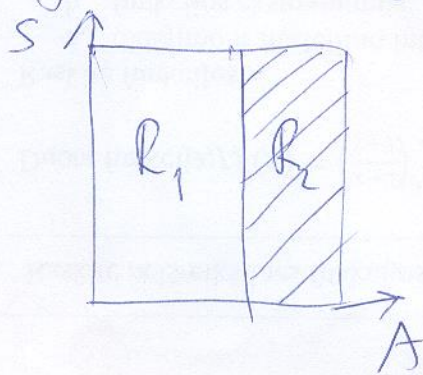
Atvairdavome prognozių erdvės skaidymų 2 3 sritys (2 skaidymo žingsniai)



Regressiivis meotseis sprianta  
 zymeta (vurduots) faips



Geometrisis aistulimors Grābime

regrueta dermate, suti (aekstis, sooris)  
 u jo rekursivā skaidome y suti R<sub>j</sub>



 - vye.  - mot

Skaidrums zingsnyje pasvintu svitu R  
skaidroune i divi dalis (slygdzotai ja  
vadulimive kairija (Left) ir desimaja (Right)  
Turime K skirtingu klasiu (daivai  
K=2 - vyrai ir moterys, kaip ir ne).

Pazymejume

$n_{k,L}$  - deomenys, apibedizanti k-klase  
ir priklausanti L arbei, skaidis

$n_{k,R}$  - — " — priklausanti R arbei,  
skaidis

$$n_L = \sum_{k=1}^K n_{k,L}$$

$$n_R = \sum_{k=1}^K n_{k,R}$$

$$p_{k,L} = \frac{n_{k,L}}{n_L}$$

$$p_{k,R} = \frac{n_{k,R}}{n_R}$$

(santykinė dalis) k. to sus  
klasi deomenys).

Apmokymo metu klasifikuojamo fiksuojamo  
vertiname naudojant Gini funkciją

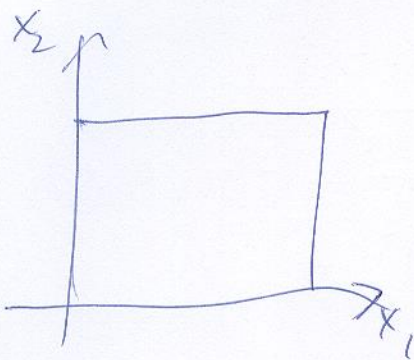
$$G(s_k) = \sum_{j=L,R} \sum_{k=1}^K p_{k,j} (1 - p_{k,j})$$

Tegu kažkurioje dalyje  $j$  (kairėje ar dešinėje)  
visos dalelės priklauso tai pačiai klasei,  
tai šios dalelės  $G_j = 0$  - tobulas  
išskaidymas

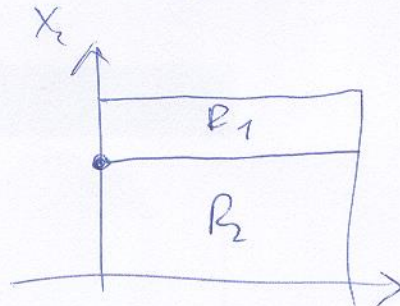
Tegu, pvz.  $K=2$  ir turime mišinį  
50:50 procentų, tai tada gauname  
blogiausią klasifikatorių  $G_j = \frac{1}{2}$ .

Pašalinkus skaidymo kintamąjį,  
ūškome geriausias skaidymo vietas (godžiogi  
strategija), nors toks pasirinkimas  
gali pabloginti galimybes rasti geresnį spren-  
dimą kituose žingsniuose

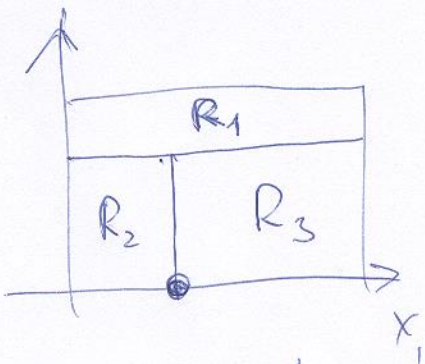
# Rekursinis skaidymas 2D atveju



a) prad. situ.



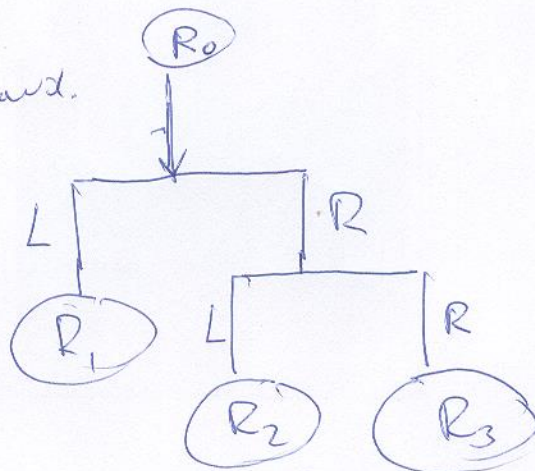
b) 2 sūtyjs po pirmos skaid.



c) trys sūtyjs po antros skaid.

- Likę lokų skaidymo žingsnių reikia atlikti?

?



Skaidymo procesą stabdomė, kai vienoje iš naujų sūčių yra mažiau nei  $N$  elementų ( $N \approx 3$  ar  $8$ ). Išvengti persimokymo!



## Sakņu karpjuma algoritmai

Sielūdzame, kad gautieji klasifikāciju  
metodes būtu kļuvis vienkāršāki: jā  
aizņem i lapu skaitis būtu kļuvis  
mazāks. Tai papildina ne tik klasifi-  
kācijas efektivitāte (skaidrības apmērs  
samazinās), bet algoritms ~~stabilizē~~ stabilizē.

Todiel paskatīsimies etape fiksācijā  
is eitis visus lapas / skaidrojumus  
metris turpinājumā ir patīk-  
nams

- ar pašreizējo / turpmāko skaidrojumu  
nepablogējamā klasifikācija fiksācija  
testamais datums ar būs atbilstošs (  
kas kada leidzotāme šādā tādā fiksācijā  
netāp pablogēti)
- ja jau šāds, tad šāds skaidrojums pašā-  
līnams.

## Modifikuojamas regresinis metodas sudarymui.

Regresijos metodus naudojamas, kai turime prognozuoti funkcijos reikšmę (o tam fiksuojame srities reikšmų intervalą) naujiems duomenims.

Tada modifikuojame klasifikatoriaus algoritmus:

- kiekvienam suskaidytos srities stačialaipsniui priskiriame reikšmę lygiai visų tikime stačialaipsnių esančių duomenų vidurkiui

$$\bar{y}_j = \frac{1}{n_j} \sum_{i=1}^{n_j} y_i, \quad \begin{array}{l} n_j - \text{duomenų} \\ \text{skaičius } j\text{-toje} \\ \text{sirtyje } R_j \end{array}$$

Paklaup vertiname imtami:

$$E = \sum_{j=1}^g \sum_{i=1}^{n_j} (y_i - \bar{y}_j)^2$$