KAUNO TECHNOLOGIJOS UNIVERSITETAS INFORMATIKOS FAKULTETAS

Intelektikos pagrindai (P176B101)

Antro laboratorinio darbo ataskaita

Atliko:

IFF-1/1 gr. Studentas

Vytenis Kriščiūnas

Priėmė:

lekt. Nečiūnas Audrius

lekt. Budnikas Germanas

TURINYS

1.	Pasirinktas duomenų rinkinys	3
2.	Pasirinkti atributai	3
3.	Turimo duomenų rinkinio suskaidymas	3
4.	Įvestys ir išvestys	4
5.	Sprendimų medžio sudarymas	4
6.	Grafinis sprendinių medžio atvaizdavimas	4
7.	Sprendinių medžio testavimas	6
8.	Sprendimų medžio gylio keitimas	6
9.	Atsitiktinio miško sudarymas	7
10.	Atsitiktinį mišką sudarančių medžių kiekio keitimas	7
11.	Sprendinių medžio ir atsitiktinio miško palyginimas	7

1. Pasirinktas duomenų rinkinys

Naudojau tą patį duomenų rinkinį, kaip ir pirmojo laboratorinio darbo metu su normalizuotais duomenimis.

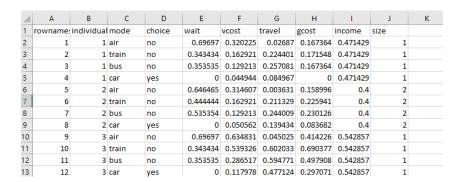
2. Pasirinkti atributai

Įvesties atributai: "wait" ir "travel". Jie šie atributai padeda rasti geriausius prognozuojamus rezultatus.

Išvesties prognozuojamas atributas: "mode", kurio kardinalumas lygus 4.

3. Turimo duomenų rinkinio suskaidymas

Duomenų rinkinį suskaidžiau į 70-30% dalis, atitinkamai: apmokymo imtis ir testavimo imtis.



1	Α	В	С	D	E	F	G	Н	1	J
1	rowname	individua	mode	choice	wait	vcost	travel	gcost	income	size
2	588	147	car	no	0	0.258427	0.567901	0.460251	0.028571	1
3	589	148	air	no	0.646465	0.567416	0.073348	0.410042	0.542857	2
4	590	148	train	yes	0.10101	0.219101	0.484386	0.51046	0.542857	2
5	591	148	bus	no	0.535354	0.230337	0.625999	0.640167	0.542857	2
5	592	148	car	no	0	0.095506	0.606391	0.523013	0.542857	2
7	593	149	air	no	0.646465	0.792135	0.0748	0.577406	0.542857	1
3	594	149	train	yes	0.10101	0.455056	0.553377	0.74477	0.542857	1
9	595	149	bus	no	0.535354	0.382022	0.604938	0.736402	0.542857	1
0	596	149	car	no	0	0.269663	0.615832	0.661088	0.542857	1
1	597	150	air	yes	0.454545	0.516854	0.150327	0.389121	0.828571	1
2	598	150	train	no	0.343434	0.5	0.619463	0.669456	0.828571	1
3	599	150	bus	no	0.353535	0.247191	0.6122	0.476987	0.828571	1
A	600	150			0	0 264046	0 507500	n 476007	0000571	-1

4. Įvestys ir išvestys

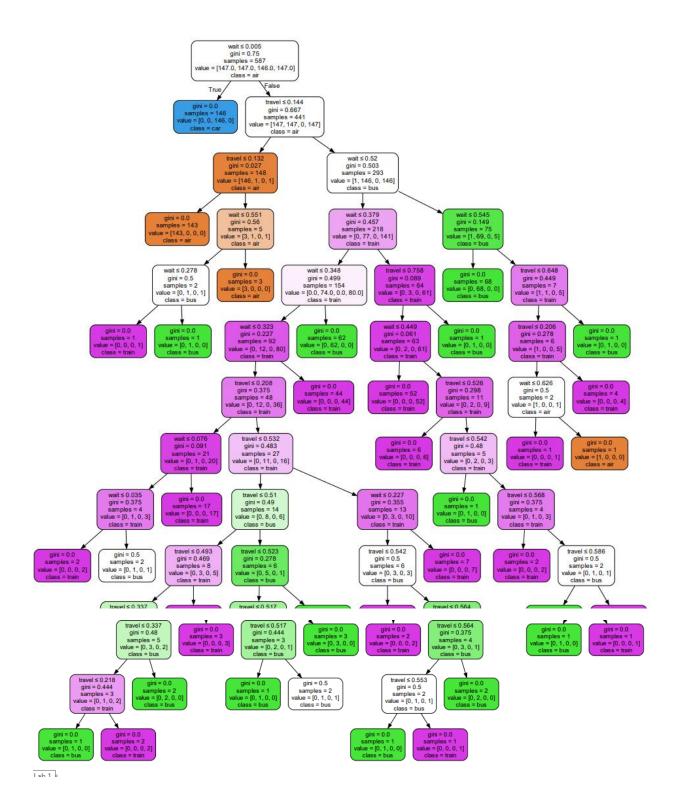
Atributo "mode" galima variantai: *car, air, bus* ir *train*. Juos bandau nuspėti naudodamas tolydinius duomenis iš "wait" ir "travel" stulpelių.

5. Sprendimų medžio sudarymas

Naudojau CART algoritmą, nes ID3 algoritmas gali dirbti tik su kategoriniais duomenimis, o man reikia ir tolydinių duomenimų. Šis algoritmas naudoja Gini indeksą, kuris nurodo, kaip dažnai atsitiktinai pasirinkti elementai būtų neteisingai išdėstyti sprendimų medyje.

6. Grafinis sprendinių medžio atvaizdavimas

Naudojau graphviz biblioteką sprendimų medžio formavimui.



7. Sprendinių medžio testavimas

Prognozavimo tikslumui nustatyti naudojau tikslumo metiką:

```
tikslumas = \frac{teisingų spėjimų skaičius}{visas spėjimų skaičius}
```

Prognozavimo tikslumas:

```
Prognozavimo tikslumas: 0.9565217391304348
```

Sumaišymo matrica (stulpelis nurodo tikrąją atributo reikšmę, o eilutė spėjamas reikšmes):

```
Sumaišymo matrica:
[[61 0 0 2]
[ 3 55 0 5]
[ 0 0 64 0]
[ 0 1 0 62]]
```

8. Sprendimų medžio gylio keitimas

Galima pastebėti, kad gyliui augant auga ir jo formavimo trukmė. Tikslumas pagerėja pasiekus pakanamai dedidelį gylį ir praktiškai nusistovi, taigi labai didelis medžio gylis nesuteikia jokio pranašumo prognozavimo tikslumo prasme.

```
Rezultatai:
Gylis
       Formavimo trukmė (s)\Tikslumas
       0.004044532775878906
                                        0.7272727272727273
       0.003522634506225586
                                        0.8458498023715415
       0.0035212039947509766
                                        0.9604743083003953
                                        0.9525691699604744
       0.003018617630004883
32
       0.003017425537109375
                                        0.9604743083003953
64
       0.005030393600463867
                                        0.9565217391304348
128
       0.007094860076904297
                                        0.9565217391304348
                                        0.9525691699604744
       0.009632110595703125
```

9. Atsitiktinio miško sudarymas

Galima pastebėti, kad atsitiktinio miško prognozės tikslumas yra šiek tiek mažesnis nei vieno iš sprendimų medžio, taip yra todėl, kad medžių kiekis bei mažesnio tikslumo medžiai gali nulemti neteisingų rezultatų priėmimą (vote). Vis dėlto gautas tikslumas panaudojus atisitiktinį mišką dažnai yra didesnis už pavienius medžius.

```
Prognozavimo tikslumas atsitiktiniui miškui: 0.9407114624505929

Sprendimų medis 0 Prognozavimo tikslumas: 0.9486166007905138, Gylis: 13

Sprendimų medis 1 Prognozavimo tikslumas: 0.857707509881423, Gylis: 14

Sprendimų medis 2 Prognozavimo tikslumas: 0.9209486166007905, Gylis: 14

Sprendimų medis 3 Prognozavimo tikslumas: 0.9130434782608695, Gylis: 13

Sprendimų medis 4 Prognozavimo tikslumas: 0.9090909090909091, Gylis: 16

Geriausias gylis: 13 su prognozavimo tikslumu: 0.9486166007905138
```

10. Atsitiktinį mišką sudarančių medžių kiekio keitimas

Mažiausiai medžių turintis atsitiktinis miškas – 3, gavo geriausią prognozės rezultatą. Galima daryti išvadą, kad atsitiktinai suformuotų medžių kiekis ir jų rūšis lemia atsitiktinio miško prognozės tikslumą.

```
Medžių kiekis: 3, Tikslumas: 0.9488048674489352
Medžių kiekis: 4, Tikslumas: 0.9266261045922063
Medžių kiekis: 5, Tikslumas: 0.9266695639576996
Medžių kiekis: 6, Tikslumas: 0.9369114877589455
Medžių kiekis: 7, Tikslumas: 0.9488048674489352
Medžių kiekis: 8, Tikslumas: 0.941981747066493
Medžių kiekis: 9, Tikslumas: 0.9470664928292047
Geriausias medžių kiekis: {'n_estimators': 3}
Geriausias tikslumas: 0.9488048674489352
```

11. Sprendinių medžio ir atsitiktinio miško palyginimas

Sprendinių medžio tikslumas šiuo atvėju gavosi geresnis nei atsitiktinio miško, tai galėjo nulemti medžių skirtumai ir atsitiktiniame miške naudotų sprendimų medžių duomenų pasikartojimas – kuo yra daugiau vienodų duomenų tuo sunkiau priimti korektiškus sprendimus.

Suformavus naują atsitiktinį mišką yra gaunamas geresnis tikslumas už pavienį sprendimų medį:

Prognozavimo tikslumas: 0.9565217391304348

Prognozavimo tikslumas atsitiktiniui miškui: 0.9644268774703557