

# Priprema podataka

NIKOLA MILIKIĆ

EMAIL: [nikola.milikic@fon.bg.ac.rs](mailto:nikola.milikic@fon.bg.ac.rs)

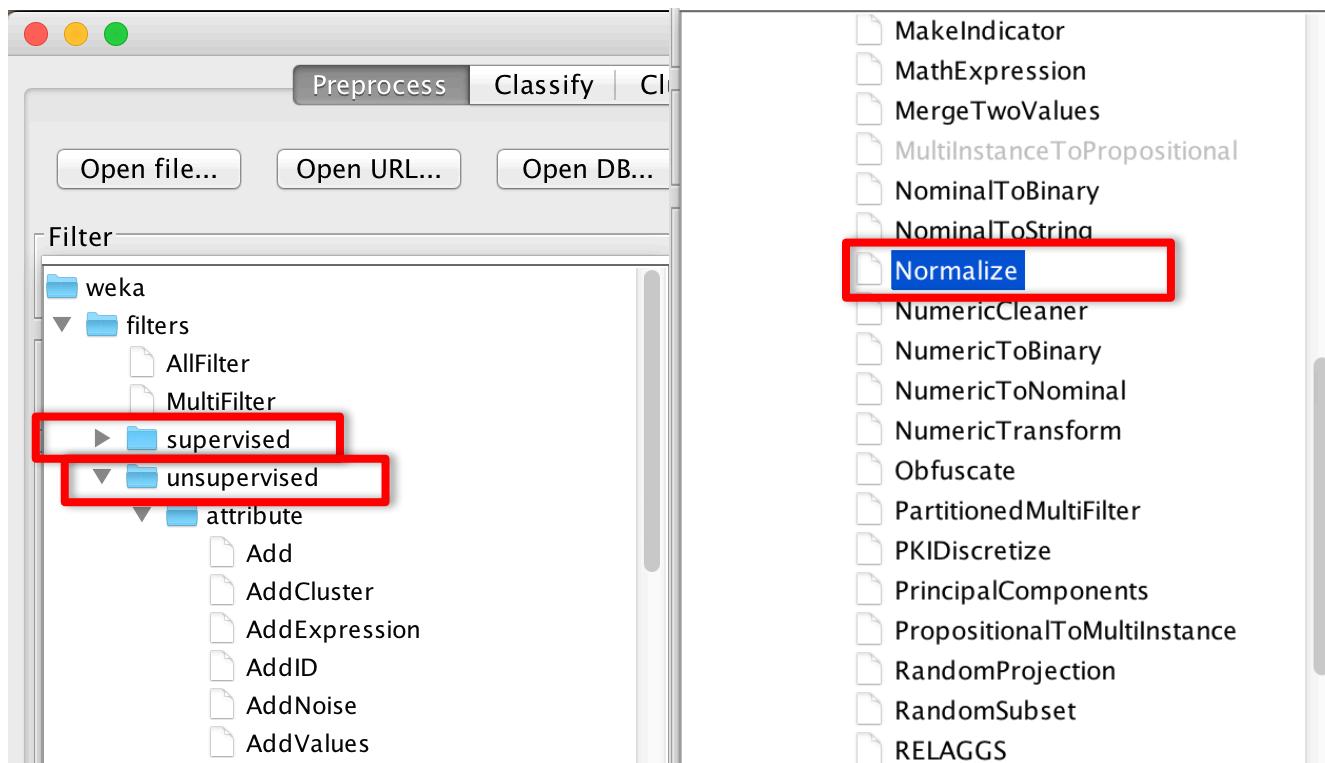
URL: <http://nikola.milikic.info>

# Normalizacija

**Normalizacija** je svođenje vrednosti na neki opseg (obično 0 - 1). Koristi se kada atributi imaju razlicite merne jedinice ili opsege vrednosti.

Nedostatak je ako postoje *outliers*.

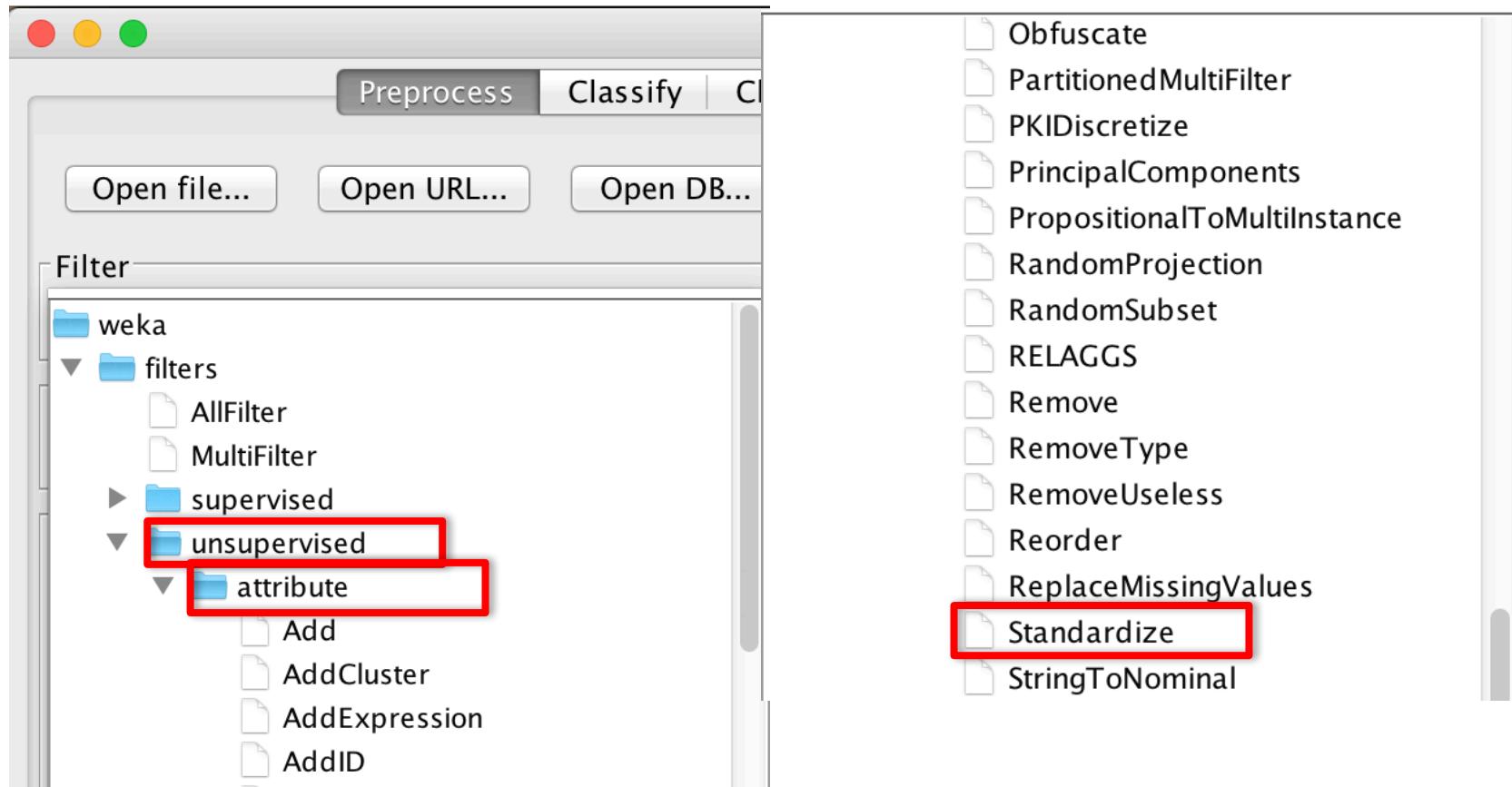
FishersIrisDataset.arff



# Standardizacija

**Standardizacija** je svođenje srednje vrednosti na 0, a standardne devijacije na vrednost 1

FishersIrisDataset.arff



# Diskretizacija atributa

**Diskretizacija** je proces transformacije numeričkih podataka u nominalne tako što se numeričke vrednosti smeštaju u odgovarajuće grupe kojih ima konačan broj.

Najčešći pristupi diskretizacije su:

- Nenadgledani pristupi:
  - Jednake širine opsega (Equal-width binning)
  - Jednaka pojavljivanja u opsezima (Equal-frequency binning)
- Nadgledani pristup – uzima u obzir klase

# Jednake širine opsega

*Jednake širine opsega* (eng. *Equal-width binning*) deli opseg mogućih vrednosti na **N** podopsega iste širine.

$$\text{Širina} = (\text{maks. vrednost} - \text{min. vrednost}) / N$$

Primer: Ako je opseg posmatranih vrednosti između 0 – 100, možemo kreirati 5 podopsega na sledeći način:

$$\text{Širina} = (100 - 0) / 5 = 20$$

Opsezi su: [0-20], (20-40], (40-60], (60-80], (80-100]

Obično se prvi i poslednji opsezi proširuju kako bi uključili vrednosti van opsega.

# Jednaka pojavljivanja u opsezima

*Jednaka pojavljivanja u opsezima* (eng. *Equal-frequency ili equal-height binning*) deli opseg mogućih vrednosti na  $N$  podopsega gde svaki podopseg sadrži isti broj instanci.

Primer: Prepostavimo da želimo da smestimo u 5 podopsega vrednosti:

5, 7, 12, 35, 65, 82, 84, 88, 90, 95

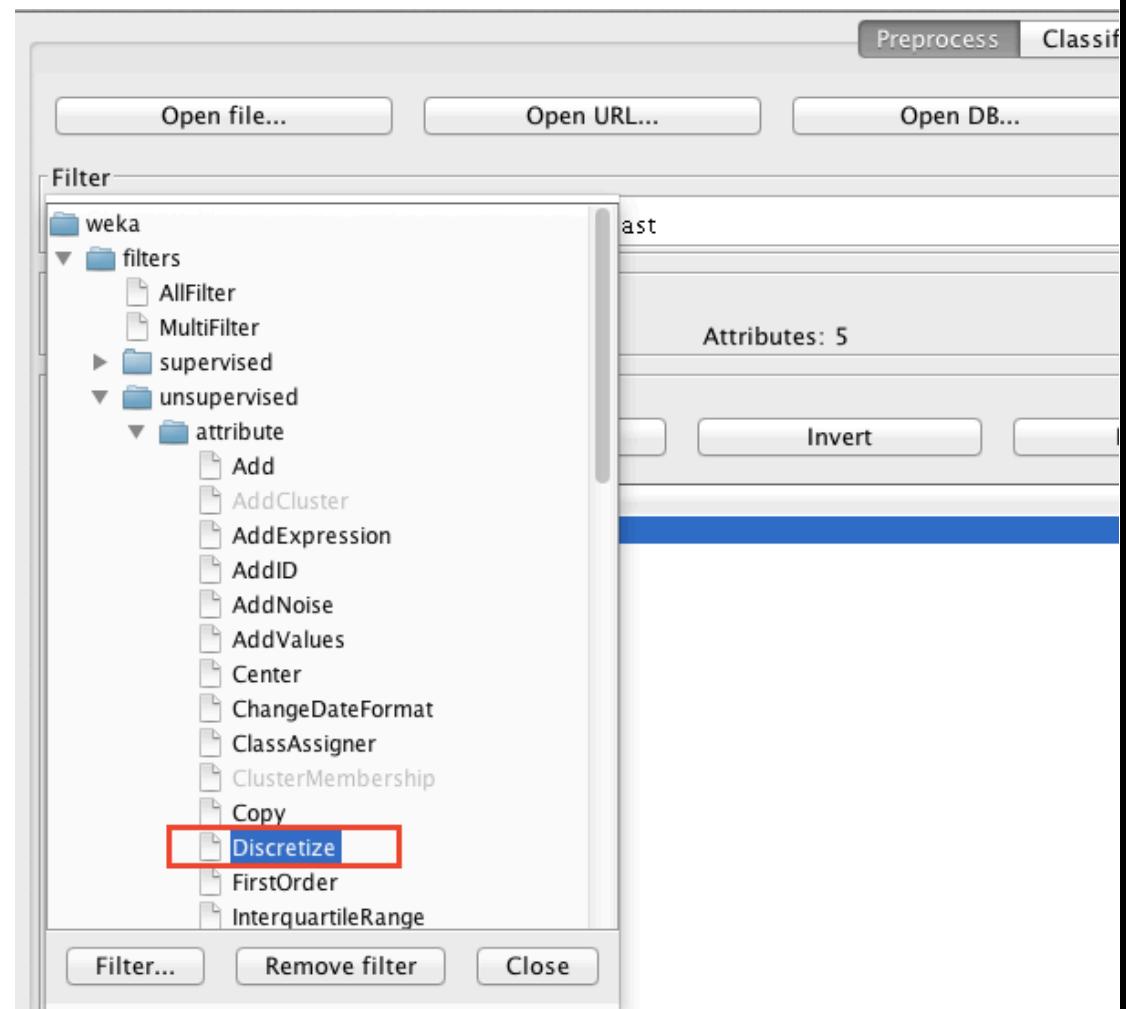
Podopsege ćemo podeliti tako što će svaki sadržati po dve instance:

5, 7, | 12, 35, | 65, 82, | 84, 88, | 90, 95

# Diskretizacija u Weka-i

Atributi se diskretizuju tako što se nad njihovim vrednostima primeni odgovarajući *Filter*.

Na *Preprocess* tabu se bira opcija  
*Choose -> Filter* i u folderu *filters/unsupervised/attrib ute* se odabira filter *Discretize.*

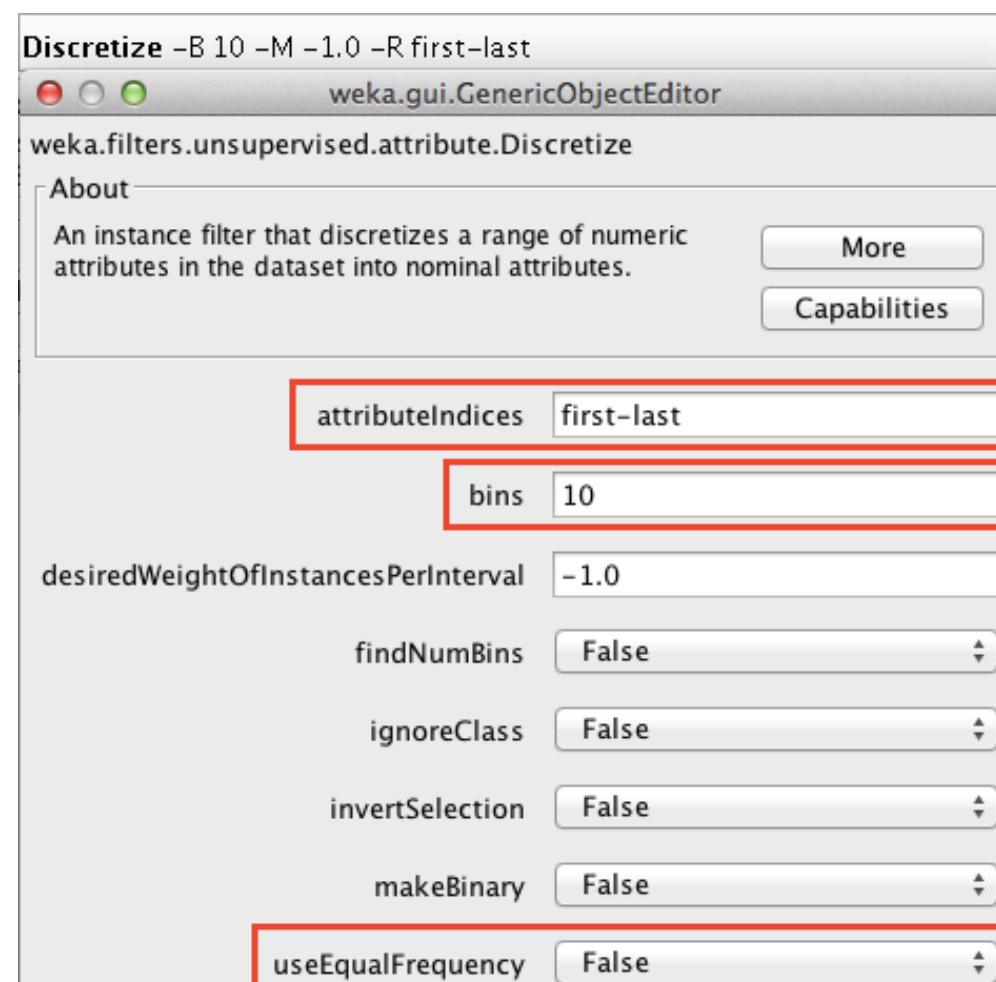


FishersIrisDataset.arff

# Diskretizacija u Weka-i

Po defaultu se primenjuje Diskretizacija sa jednakim širinama opsega.

- *attributeIndices* - vrednost *first-last* označava da diskretizjemo sve attribute. Mogu se navesti i redni brojevi atributa
- *bins* - željeni broj opsega
- *useEqualFrequency* – *true* ako se koristi diskretizacija sa jednakim pojavljivanjima u opsezima, *false* ako se koristi Diskretizacija sa jednakim širinama opsega



# Diskretizacija u Weka-i

Pritisom na **Apply** se primjenjuje odabrani filter

Preprocess Classify Cluster Associate Select attributes Visualize

Open file... Open URL... Open DB... Generate... Undo Edit...

Choose Discretize -B 10 -M -1.0 -R first-last Apply

Current relation Relation: FishersIrisDataset-weka.filters.unsupervised.attribute.Remove... Instances: 150 Attributes: 5

Attributes All None Invert Pattern

No.	Name
1	<input checked="" type="checkbox"/> Sepal Length
2	<input type="checkbox"/> Sepal Width
3	<input type="checkbox"/> Petal Length
4	<input type="checkbox"/> Petal Width
5	<input type="checkbox"/> Species

Dobijeni podopsezi vrednosti

Selected attribute Name: Sepal Length Type: Nominal Missing: 0 (0%) Distinct: 10 Unique: 0 (0%)

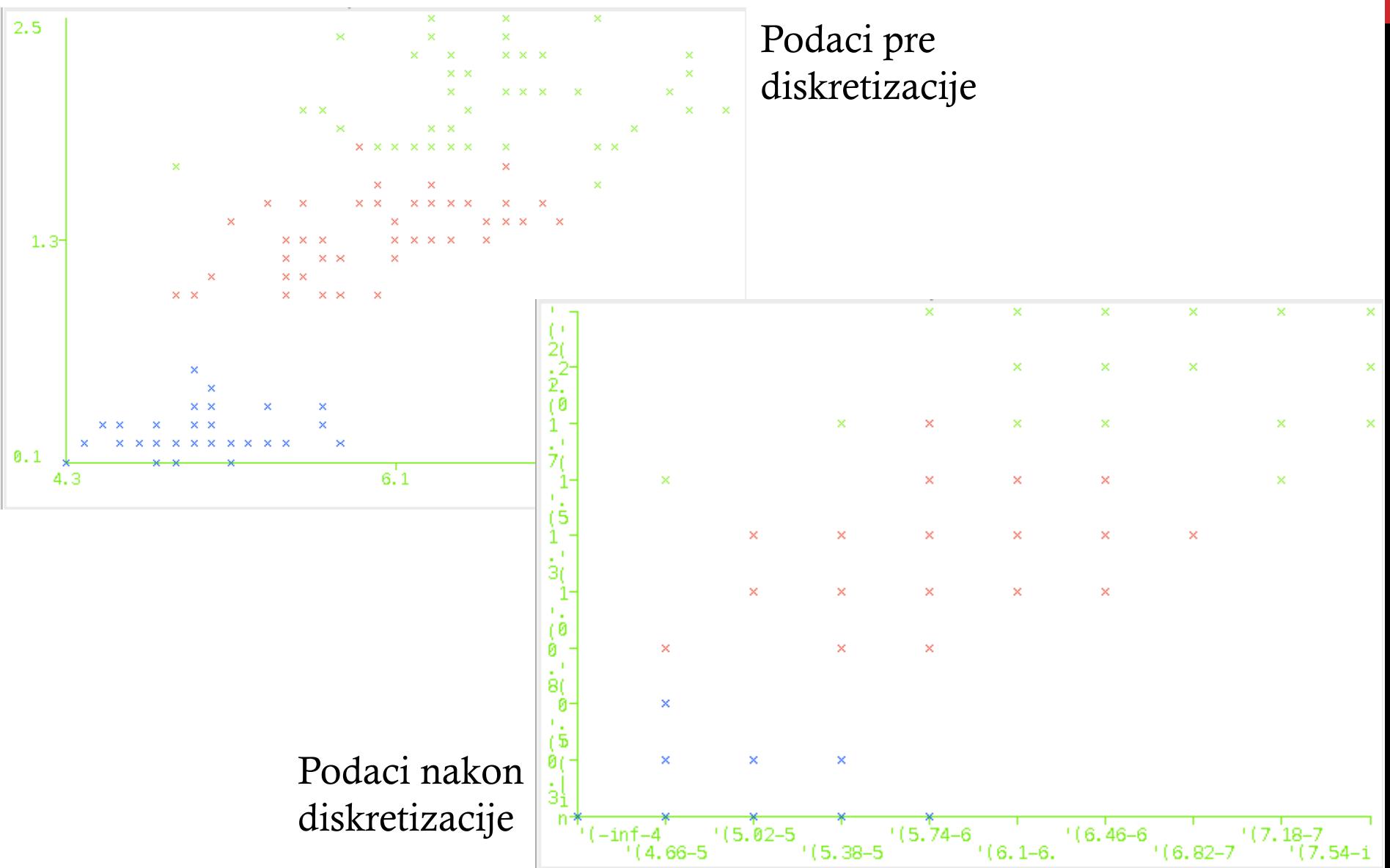
No.	Label	Count
1	'(-inf-4.66]'	9
2	'(4.66-5.02]'	23
3	'(5.02-5.38]'	14
4	'(5.38-5.74]'	27
5	'(5.74-6.1]'	22
6	'(6.1-6.46]'	20
7	'(6.46-6.82]'	18
8	'(6.82-7.18]'	6
9	'(7.18-7.54]'	5
10	'(7.54-inf)'	6

Class: Species (Nom) Visualize All

Bin Range	Setosa	Versicolor	Virginica	Total
'(-inf-4.66]'	9	0	0	9
'(4.66-5.02]'	23	0	0	23
'(5.02-5.38]'	14	0	0	14
'(5.38-5.74]'	27	0	0	27
'(5.74-6.1]'	22	0	0	22
'(6.1-6.46]'	20	0	0	20
'(6.46-6.82]'	18	0	0	18
'(6.82-7.18]'	6	0	0	6
'(7.18-7.54]'	5	0	0	5
'(7.54-inf)'	6	0	0	6

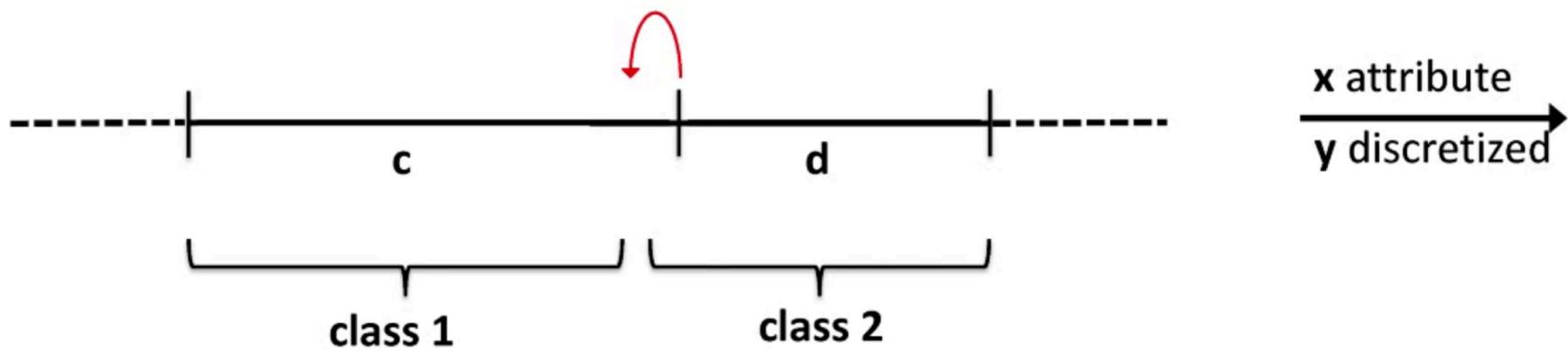
Status OK Log x 0

# Podaci pre i posle diskretizacije



# Nadgledana diskretizacija

- Šta ako sve instance u jednom binu pripadaju jednoj klasi, a sve instance drugog bina pripadaju drugoj klasi osim prvog koji pripada prvoj klasi?



- Nadgledana diskretizacija uzima u obzir i klasu

# Nadgledana diskretizacija

- Jedan od pristupa je koristiti entropiju
  - U primeru *weather.numeric.arff*, kod atributa *temperature*

64	65	68	69	70		71	72	75	80	81	83	85
yes	no	yes	yes	yes		no	no	yes	no	yes	yes	no
						yes	yes					

4 yes, 1 no      5 yes, 4 no  
entropy = 0.934 bits

- Uzima se razdelenik sa najmanjom entropijom (najvećom informacionom dobiti)

64	65	68	69	70	71	72	75	80	81	83	85
yes	no	yes	yes	yes	no	no	yes	no	yes	yes	no

# Nadgledana diskretizacija u Weka-i

weather.numeric.arff

Preprocess Classify Cluster

Open file... Open URL... Open DB... C

Filter

Choose None

Current relation  
Relation: weather  
Instances: 14 Attributes: 5

Attributes

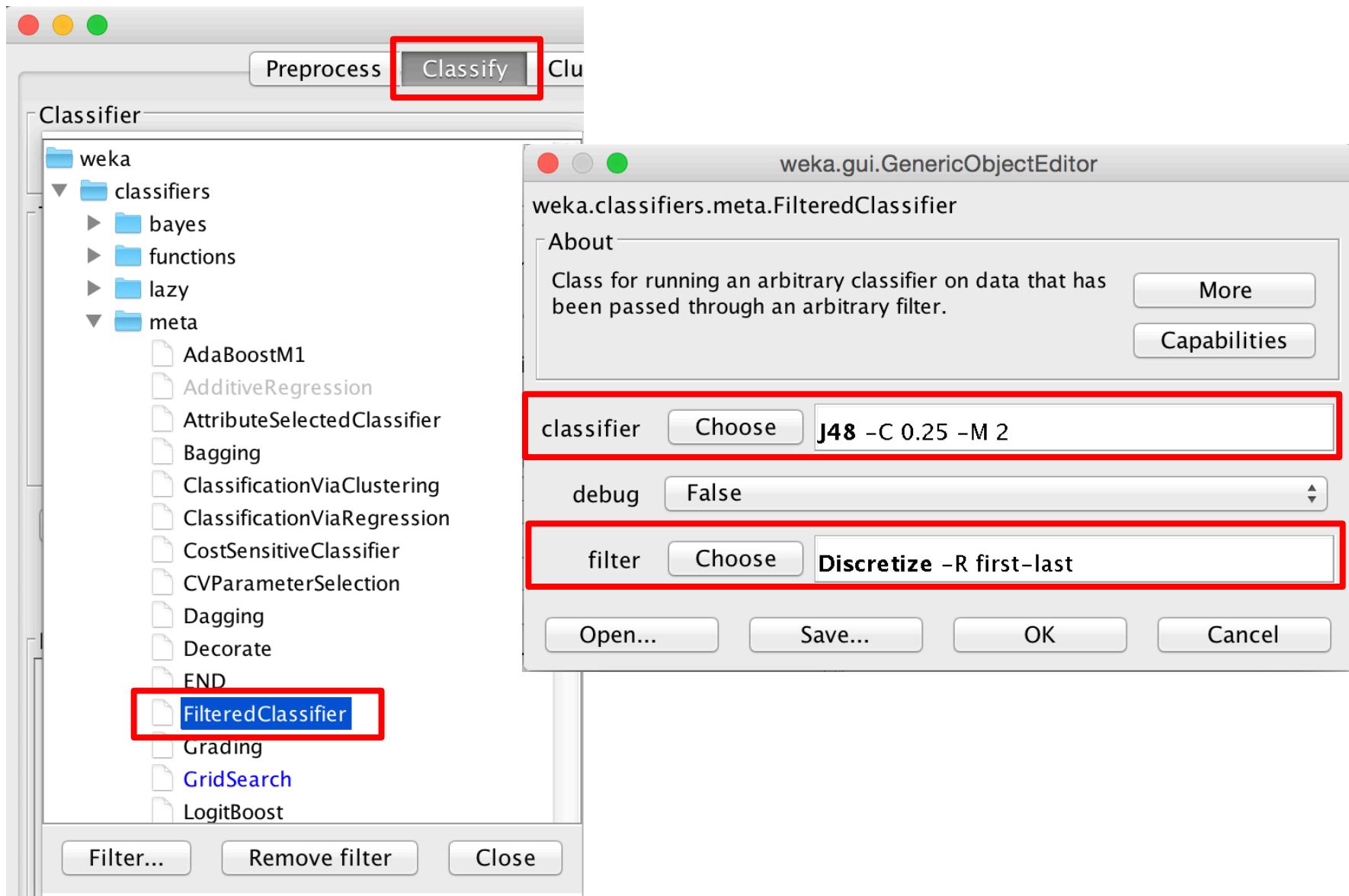
Problem je što prilikom nadgledane diskretizacije koristimo podatke iz celog dataset-a, pa samim tim i test podatke nad kojima ćemo posle vršiti testiranje performansi klasifikatora

weka

filters

- AllFilter
- MultiFilter
- supervised
- attribute
  - AddClassification
  - AttributeSelection
  - ClassOrder
  - Discretize
  - NominalToBinary
  - PLSFilter
- instance
- unsupervised

# meta>FilteredClassifier



# Selekcija atributa

**Selekcija atributa** (eng. Attribute Selection ili Feature Selection) je proces odabira podskupa relevantnih atributa koji će se koristiti.

Primenjuje se u slučajevima kada se u datasetu nalaze atributi koji su redundantni ili nerelevantni.

- Redundanti atributi su oni koji ne pružaju nikakve dodatne informacije u odnosu na već selektovane attribute.
- Nerelevantni atributi su oni koji ne pružaju nikakve informacije u datom kontekstu.

# Prednosti primene selekcije atributa

Suvišni atributi mogu degradirati performanse modela.

Prednosti selekcije atributa:

- Poboljšava čitljivost modela time što se model sastoji samo iz relevantnih atributa
- Kraće vreme treniranja
- Povećana generalizacija time što smanjuje mogućnosti za overfitting

Najbolji način za selekciju atributa je ručno ukoliko se dobro poznaje problem koji se rešava. I automatizovani pristupi selekcije daju dobre rezultate.

# Pristupi selekcije atributa

Postoje dva pristupa:

- *Filter* metoda – koriste se procene atributa na osnovu generalnih svojstava podataka. Svakom atributu se dodeli neki skor. Na osnovu skora se atributi rangiraju i zadržavaju ili odbacuju. Npr. mogu se koristiti Hi kvadrat test, informaciona dobit, stepen korelacijske, itd..
- *Wrapper* metoda – podskupovi atributa se evaluiraju primenom algoritma mašinskog učenja koji će se koristiti nad skupom podataka. Naziv Wrapper se koristi iz razloga što je algoritam učenja “zapakovan” u samom procesu selekcije. Biće odabran onaj podskup atributa za koje dati algoritam učenja daje najbolje rezultate.

# Primer selekcije atributa

census90-income.arff

Preprocess Classify Cluster Associate Select attributes Visualize

Open file... Open URL... Open DB... Generate... Undo Edit... Save...

Filter Choose None Apply

Current relation  
Relation: 1990census  
Instances: 32561 Attributes: 15

Attributes  
All None Invert Pattern

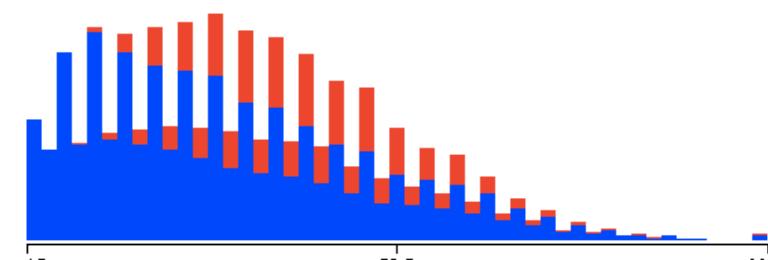
No.	Name
1	<input checked="" type="checkbox"/> age
2	<input type="checkbox"/> workclass
3	<input type="checkbox"/> fnlwgt
4	<input type="checkbox"/> education
5	<input type="checkbox"/> education-num
6	<input type="checkbox"/> marital-status
7	<input type="checkbox"/> occupation
8	<input type="checkbox"/> relationship
9	<input type="checkbox"/> race
10	<input type="checkbox"/> sex
11	<input type="checkbox"/> capital-gain
12	<input type="checkbox"/> capital-loss
13	<input type="checkbox"/> hours-per-week
14	<input type="checkbox"/> native-country
15	<input type="checkbox"/> income

Remove

Status OK Log x 0

Selected attribute  
Name: age Type: Numeric  
Missing: 0 (0%) Distinct: 73 Unique: 2 (0%)  
Statistic Value  
Minimum 17  
Maximum 90  
Mean 38.582  
StdDev 13.64

Class: income (Nom) Visualize All



# Primer selekcije atributa

Preprocess Classify Cluster Associate Select attributes Visualize

Open file... Open URL... Open DB... Generate... Undo Edit...

Filter

weka  
filters  
supervised  
attribute  
AttributeSelection

15

Selected attribute  
Name: age  
Missing: 0 (0%)  
Distinct: 73  
Statistic Value  
Minimum 17  
Maximum 90  
Mean 38.582  
StdDev 13.64

Želimo da primenimo selekciju atributa

Class: income (Nom)

Filter... Remove filter Close

# Primer selekcije atributa

Filter

Choose AttributeSelection -E "weka.attributeSelection.ClassifierSubsetEval -B weka.filters.supervised.attribute.AttributeSelection

Current relation weka.gui.GenericObjectEditor

Relation: 1 Instances: 3

Attributes

All

No. | N  
1 |  
2 |  
3 |  
4 |  
5 |  
6 |  
7 |  
8 |  
9 |  
10 |  
11 |  
12 |

evaluator  
search  
Open...

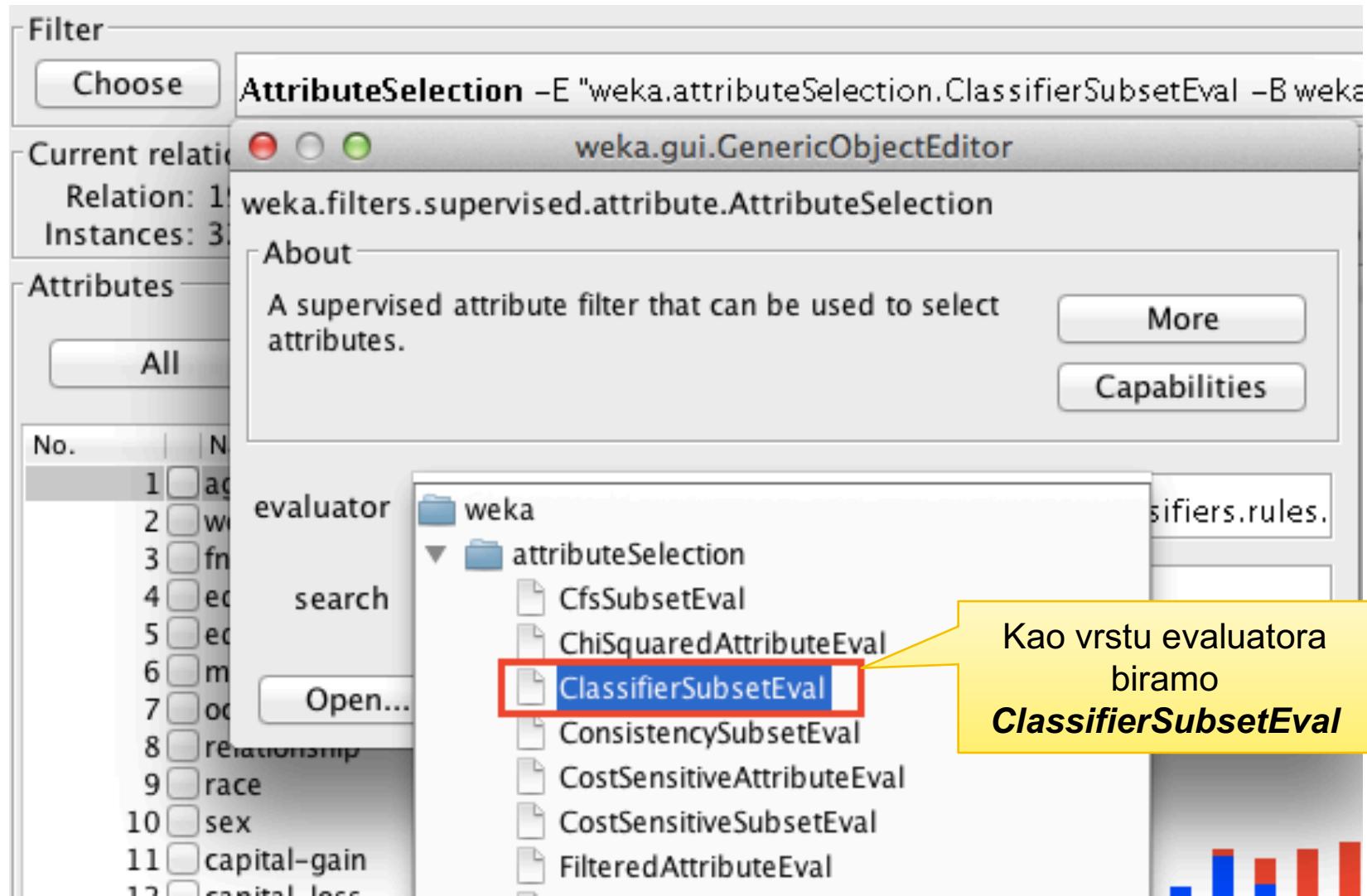
About

A supervised attribute filter that can be used to select attributes.

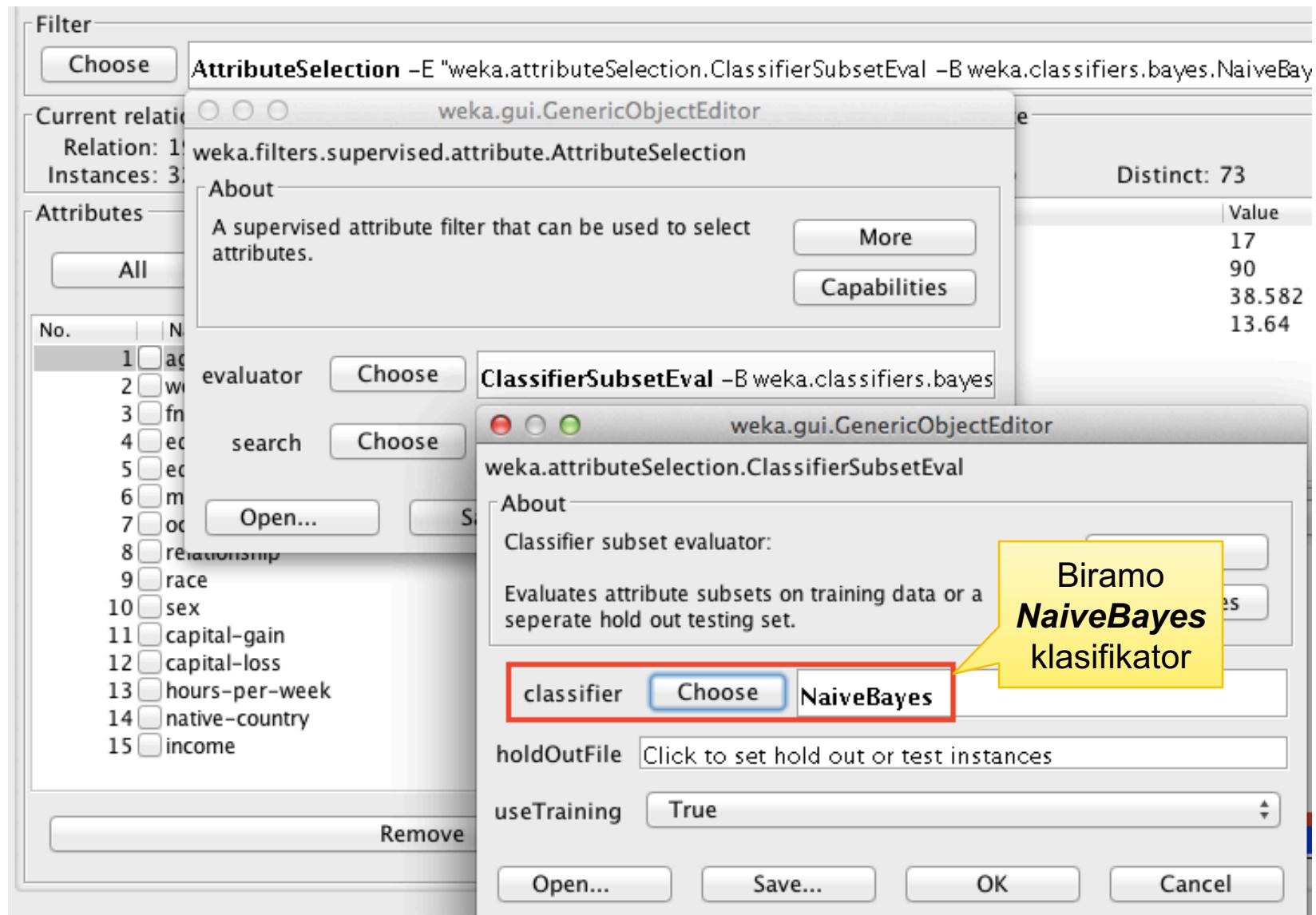
More  
Capabilities

weka  
attributeSelection  
CfsSubsetEval  
ChiSquaredAttributeEval  
**ClassifierSubsetEval**  
ConsistencySubsetEval  
CostSensitiveAttributeEval  
CostSensitiveSubsetEval  
FilteredAttributeEval

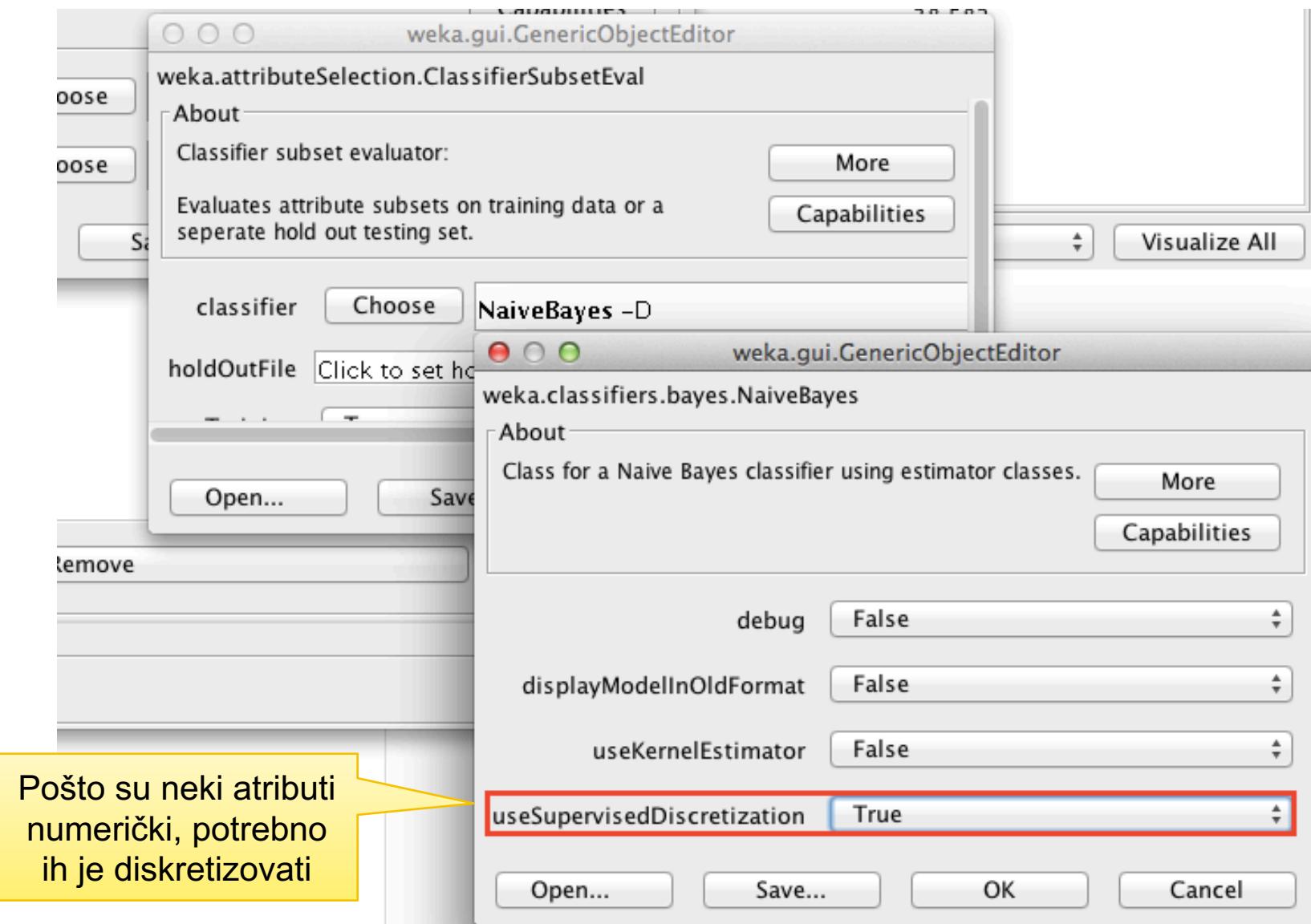
Kao vrstu evaluatorsa biramo **ClassifierSubsetEval**



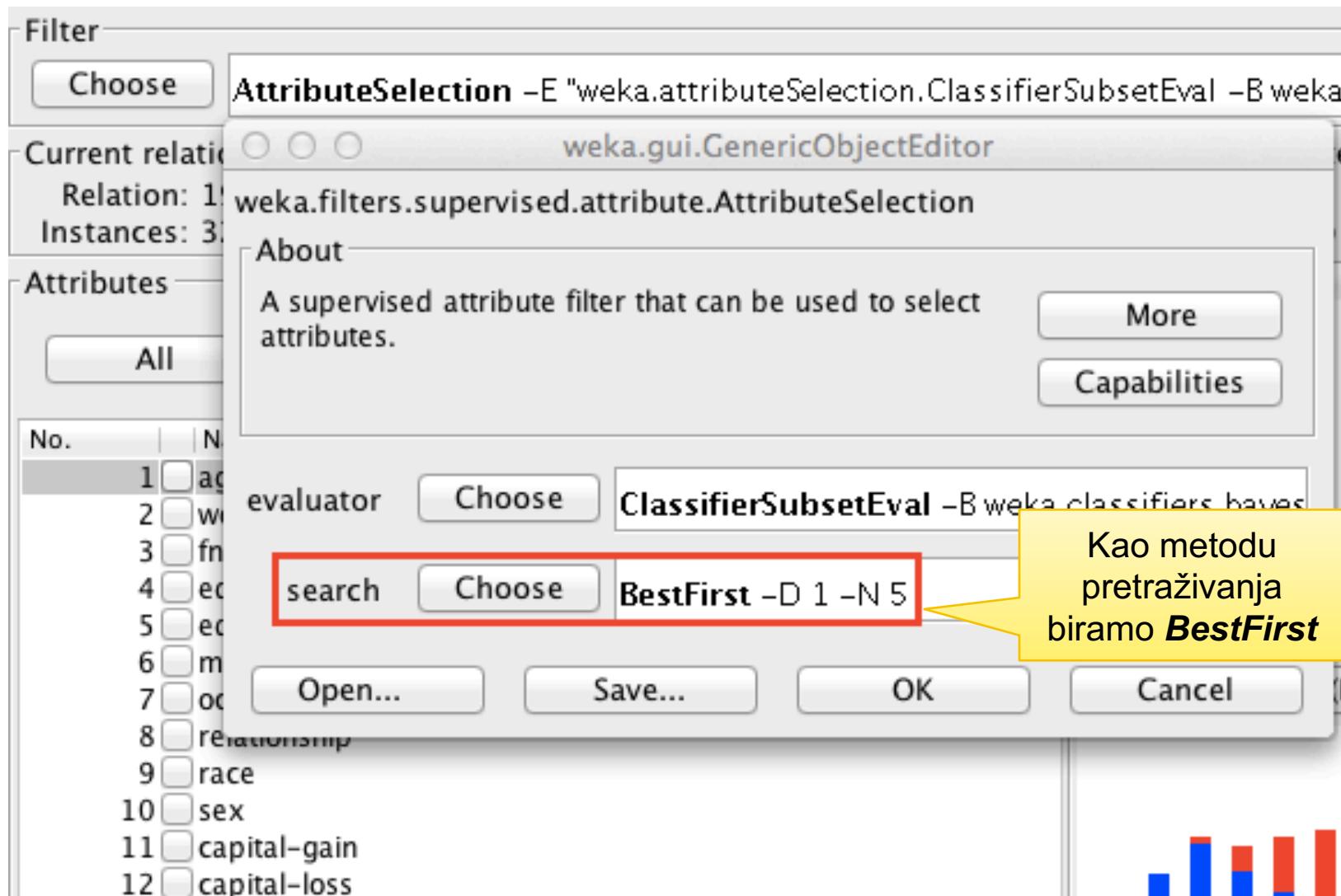
# Primer selekcije atributa



# Primer selekcije atributa



# Primer selekcije atributa



# Primer selekcije atributa

Screenshot of the Weka interface showing the "Select attributes" tab.

The "Selected attribute" panel shows "age" selected, with statistics: Name: age, Missing: 0 (0%), Distinct: 73, Type: Nominal, and Unique: 73.

A yellow callout box points to the "Apply" button, which is highlighted with a red border, indicating that the filter can be applied to other attributes.

**Filter je podešen i može biti primjenjen nad atributima**

The "Attributes" panel lists 15 attributes: age, workclass, fnlwgt, education, education-num, marital-status, occupation, relationship, race, sex, capital-gain, capital-loss, hours-per-week, native-country, and income. The "age" checkbox is selected.

The "Preprocess" tab is active at the top.

The status bar at the bottom left says "Status OK".

The status bar at the bottom right shows "Log" and "x 0".

# Primer selekcije atributa

Preprocess Classify Cluster Associate Select attributes Visualize

Open file... Open URL... Open DB... Generate... Undo Edit... Save...

Filter Choose AttributeSelection -E "weka.attributeSelection.ClassifierSubsetEval -B weka.classifiers.bayes.NaiveBayes -T -H \Click Apply

Current relation  
Relation: 1990census-weka.filters.supervised.attribute.Attribute...  
Instances: 32561 Attributes: 7

Attributes

All None Invert Pattern

No.	Name
1	<input checked="" type="checkbox"/> age
2	<input type="checkbox"/> education
3	<input type="checkbox"/> relationship
4	<input type="checkbox"/> race
5	<input type="checkbox"/> capital-gain
6	<input type="checkbox"/> capital-loss
7	<input type="checkbox"/> income

Selected attribute  
Name: age Type: Numeric  
Missing: 0 (0%) Distinct: 73 Unique: 2 (0%)

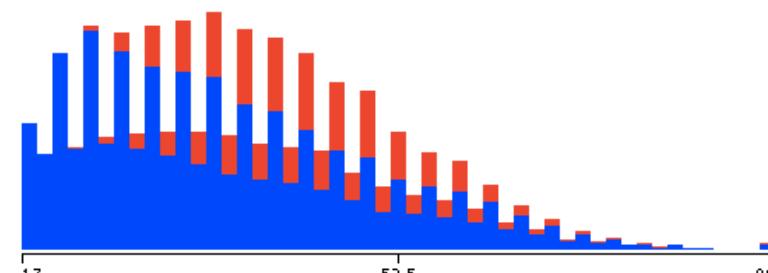
Statistic	Value
Minimum	17
Maximum	90
Mean	38.582
StdDev	13.64

Class: income (Nom) Visualize All

Broj atributa je redukovana na 7

Remove

Status OK Log x 0



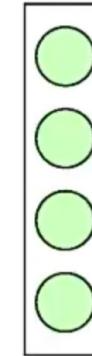
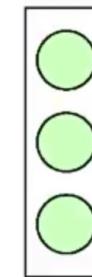
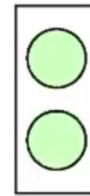
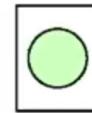
# Metod pretrage kod selekcije atributa

- Exhaustive search (512 podskupova atributa)

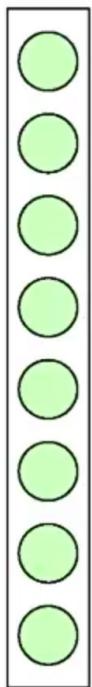
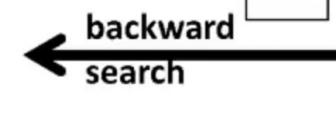
- Best First: *Forward, Backward, Bi-directional*

- *searchTermination* atribut određuje koliko podskupova koji ne poboljšavaju performanse testirati pre nego što prekine pretragu

0 attributes  
(ZeroR)



...



all 9  
attributes

forward  
search →

bidirectional  
search ↘

backward  
search ←

# Preporuke i zahvalnice

"Data Mining with Weka" and "More Data Mining with Weka": MOOCs from the University of Waikato. A self-paced session of "Data Mining with Weka" runs until 23 October June 2015.

- Link: <https://www.youtube.com/user/WekaMOOC/>

(Anonimni) upitnik za vaše kritike,  
komentare, predloge:

<http://goo.gl/cqdp3I>

# Pitanja?

NIKOLA MILIKIĆ

EMAIL: [nikola.milikic@fon.bg.ac.rs](mailto:nikola.milikic@fon.bg.ac.rs)

URL: <http://nikola.milikic.info>