# "Modern and renaissance poetry"

# Analiza podataka i pretprocesiranje
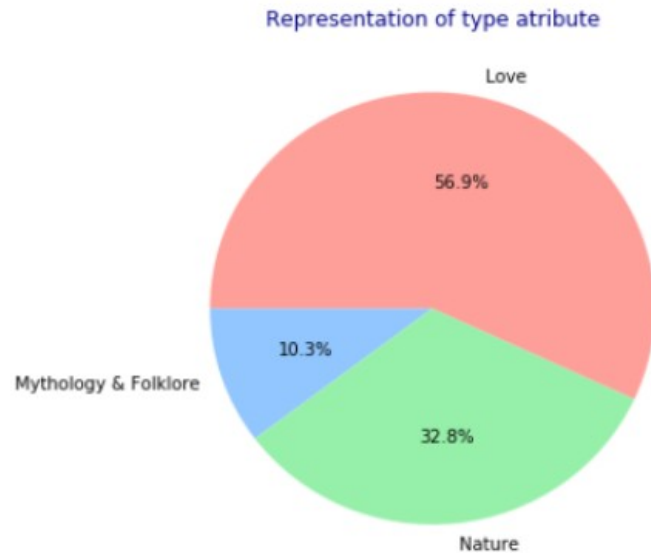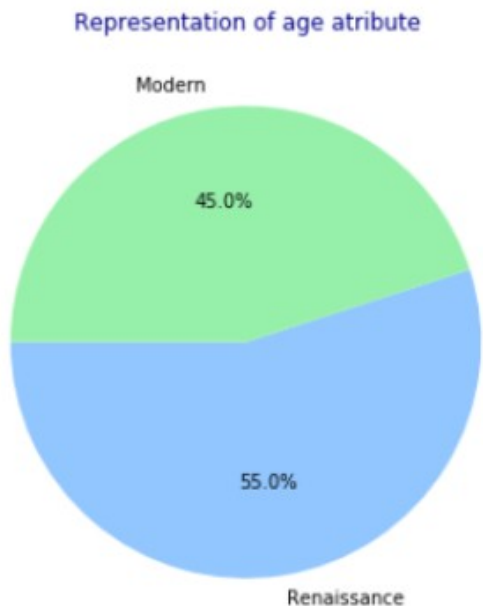
- 573 instance
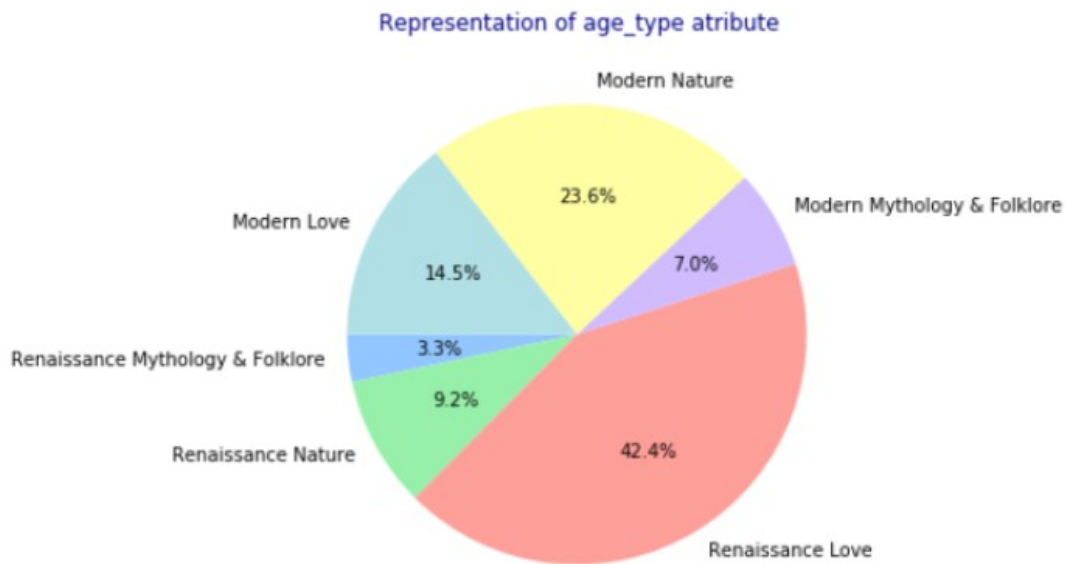- 5 atributa

- Zastupljenost klasa

# Analiza podataka i pretprocesiranje

- Zastupljenost klasa

# Analiza podataka i pretprocesiranje

- Nedostajuce vrenosti, duplirani podaci

- 506 instanci

- Transformacija podataka:

```
'let bird loudest layon sole arabian treeherald sad trumpet whose sound chaste wings obey thou shrieking harbinger f
oul precurrer fiend augur fevers end troop come thou near session interdictevery fowl tyrant wing save eagle feather
d king keep obsequy strict let priest surplice white defunctive music death divining swan lest requiem lack right th
ou treble dated crow thy sable gender makstwith breath thou givst takst mongst mourners shalt thou go anthem doth co
mmence love constancy dead phoenix turtle fledin mutual flame hence lovd love twainhad essence one two distincts div
ision none number love slain hearts remote yet asunder distance space seentwixt turtle queen wonder love shinethat t
urtle saw rightflaming phoenix sight either others mine property thus appalledthat self single natures double namene
ither two one called reason confounded saw division grow together yet either neither simple well compounded cried tr
ue twainseemeth concordant one love reason reason none parts remain whereupon made threneto phoenix dove co supremes
stars love chorus tragic scene threnosbeauty truth rarity grace simplicity enclosd cinders lie death phoenix nest tu
rtles loyal breastto eternity doth rest leaving posterity twas infirmity married chastity truth may seem cannot beau
ty brag tis truth beauty buried urn let repairthat either true fair dead birds sigh prayer'
```

- TF – IDF matrica :

# Klasifikacija - KNN

- Na neredukovanim podacima

  k = 23, kosinusna mera

- Na redukovanim podacima – prokletstvo dimenzionalnosti

  k = 23, euklidska mera
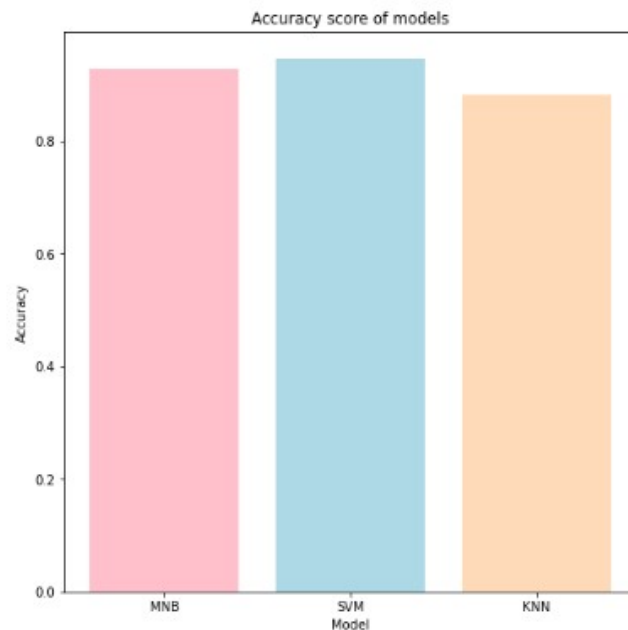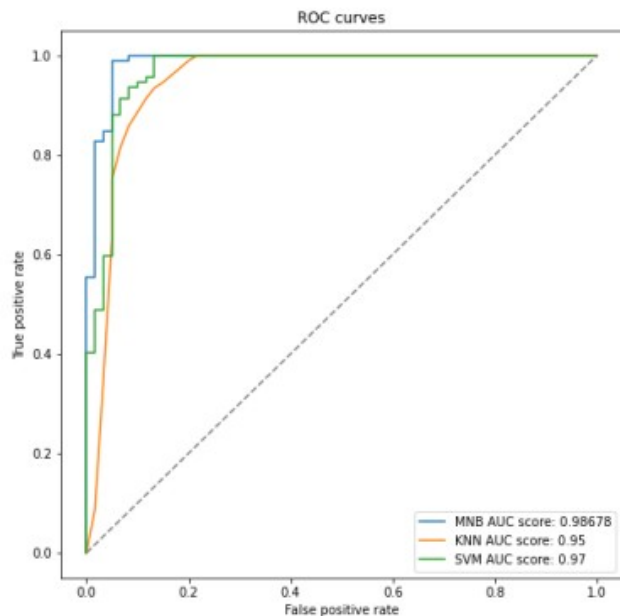
- Multinominalni Naivni Bajes
- Primenjen na neredukovanom skupu
- Određivanje hiperparametara pohlepnim algoritmom

  alpha = 0.1

- Primena nad neredukovanim podacima
- Određivanje hiperparametara pohlepnim algoritmom

  C = 1, linear

  Preprilagođavanje rešeno C = 0.5

- Primena na redukovanim podacima u cilju vizuelizacije

  Promena optimalnih hiperparametara

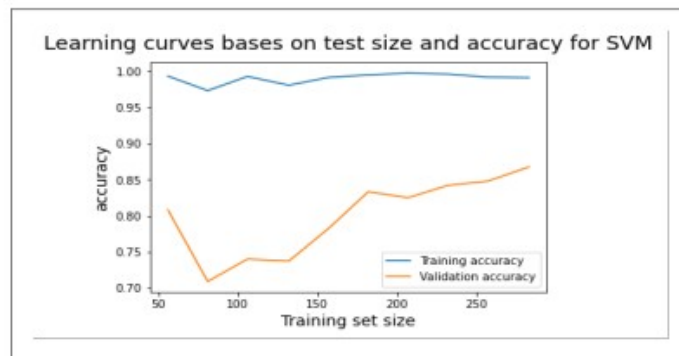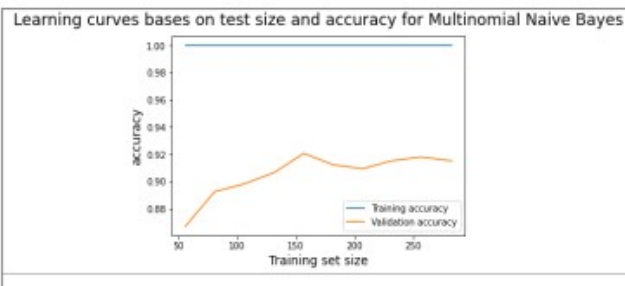  C = 1, rbf jezgro

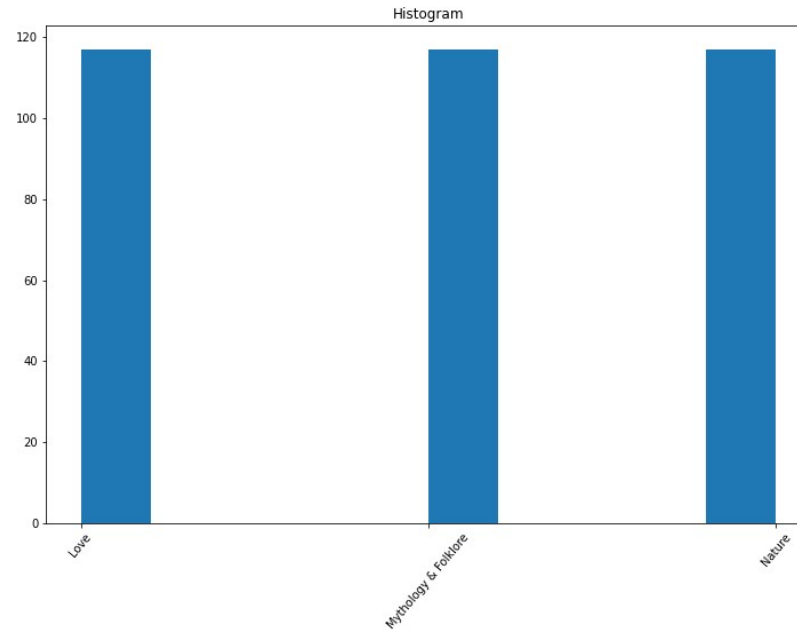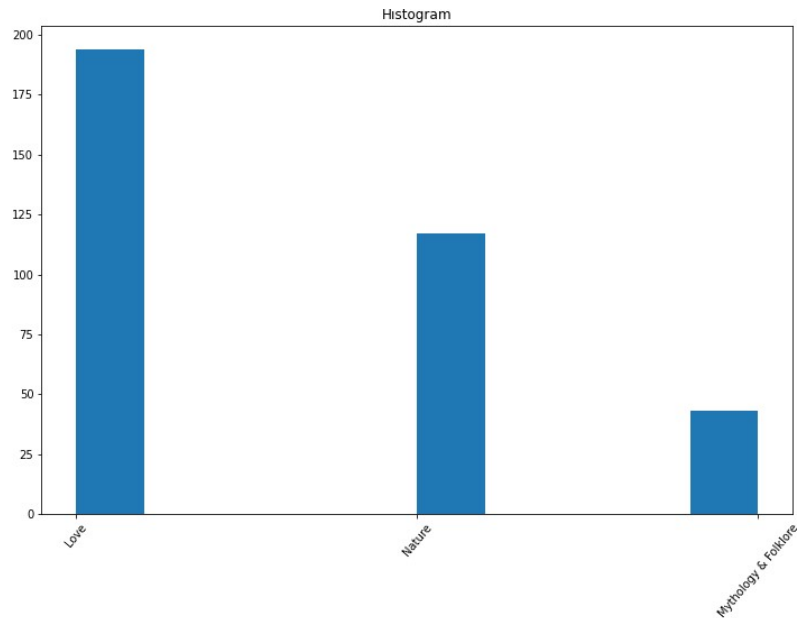# Klasifikacija - poređenje

- Roc krive i tačnost:

- **Krive učenja:**

- Atribut type

- Atribut age

- Na neredukovanim podacima:

  Prokletstvo dimenzionalnosti



Inertias and Silhouettes

- Na redukovanim podacima:

  k=3

- ## Spoljašnji kriterijumi:



Cluster vs Age

```
[[ 84 138]
 [  0 284]]
              precision    recall  f1-score   support

      Modern       1.00      0.38      0.55       222
  Renaissance      0.67      1.00      0.80       284

    accuracy                           0.73       506
   macro avg       0.84      0.69      0.68       506
weighted avg       0.82      0.73      0.69       506
```

- ## Klasterovanje kao dobar načina analize pre nadgledanog učenja

- Sakupljajuće hijerarhijsko klasterovanje
- Dendogram, Vardov metod

- Unutrašnji kriterijumi

  1. Siluate skor

  2. Kriterijum odnosa

  varijnse

  3. Kofenetski koeficijent

  korelacije

- Potrošačka korpa – pretprocesiranje

# Pravila pridruživanja

- Potrošačka korpa – pravila

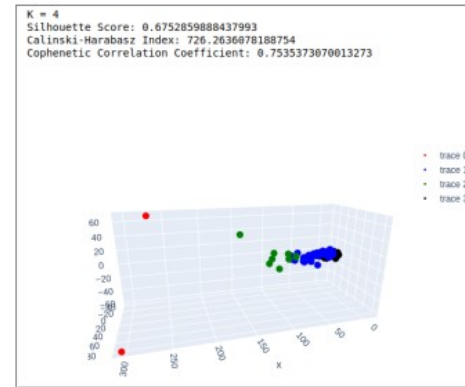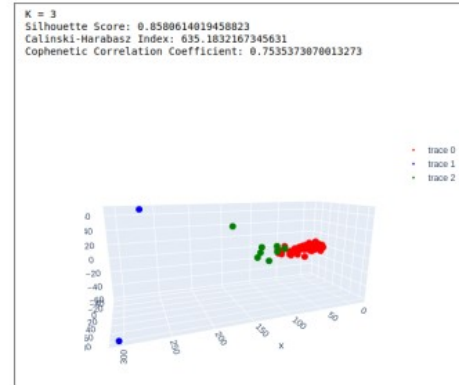| | antecedents | consequents | antecedent support | consequent support | support | confidence | lift | leverage | conviction | zhangs_metric |
|---|---|---|---|---|---|---|---|---|---|---|
| 0 | (poems) | (permission) | 0.114625 | 0.154150 | 0.112648 | 0.982759 | 6.375332 | 0.094979 | 49.059289 | 0.952303 |
| 1 | (publishing) | (permission) | 0.059289 | 0.154150 | 0.059289 | 1.000000 | 6.487179 | 0.050149 | inf | 0.899160 |
| 2 | (reprinted) | (permission) | 0.086957 | 0.154150 | 0.086957 | 1.000000 | 6.487179 | 0.073552 | inf | 0.926407 |
| 3 | (used) | (permission) | 0.071146 | 0.154150 | 0.065217 | 0.916667 | 5.946581 | 0.054250 | 10.150198 | 0.895551 |
| 4 | (used) | (poems) | 0.071146 | 0.114625 | 0.061265 | 0.861111 | 7.512452 | 0.053110 | 6.374704 | 0.933288 |
| 5 | (thee) | (thy) | 0.154150 | 0.209486 | 0.124506 | 0.807692 | 3.855588 | 0.092214 | 4.110672 | 0.875612 |

- Kategorički atributi – pretprocesiranje

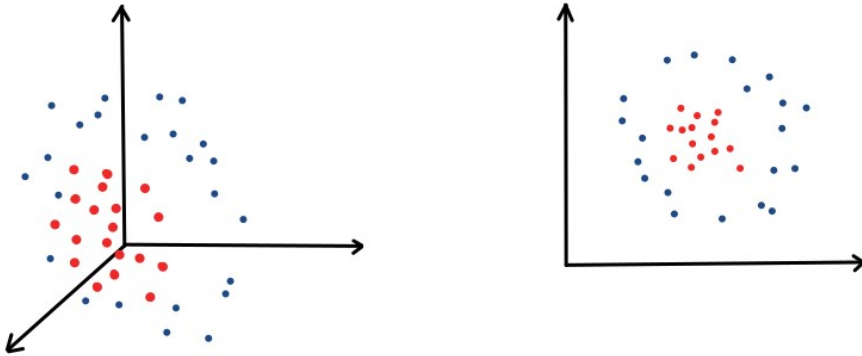| | Love | Modern | Mythology & Folklore | Nature | Renaissance | archibald macleish | asil bunting | carl sandburg | christopher marlowe | conrad aiken | ... | t s eliot | thomas bastard | thomas campion | thomas heywood | thomas lodge | thomas nashe |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 0 | False | False | True | False | True | False | False | False | False | False | ... | False | False | False | False | False | False |
| 1 | False | False | True | False | True | False | False | False | False | False | ... | False | False | False | False | False | False |
| 2 | False | False | True | False | True | False | False | False | False | False | ... | False | True | False | False | False | False |
| 3 | False | False | True | False | True | False | False | False | False | False | ... | False | False | False | False | False | False |
| 4 | False | False | True | False | True | False | False | False | False | False | ... | False | False | False | False | False | False |

- Kategorički atributi – pravila

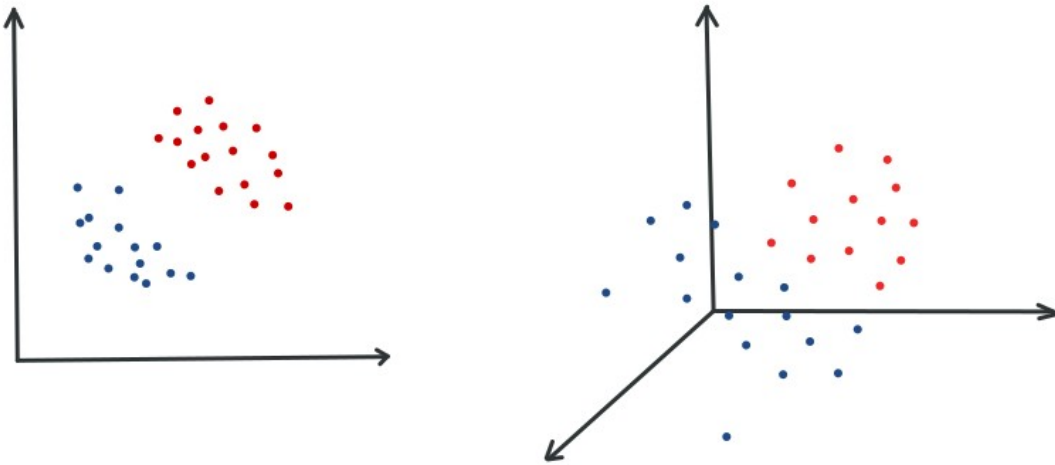| | antecedents | consequents | antecedent support | consequent support | support | confidence | lift | leverage | conviction | zhangs_metric |
|---|---|---|---|---|---|---|---|---|---|---|
| 0 | (Love) | (Renaissance) | 0.568935 | 0.549738 | 0.424084 | 0.745399 | 1.355916 | 0.111318 | 1.768498 | 0.608937 |
| 1 | (Renaissance) | (Love) | 0.549738 | 0.568935 | 0.424084 | 0.771429 | 1.355916 | 0.111318 | 1.885908 | 0.582974 |
| 2 | (Nature) | (Modern) | 0.328098 | 0.450262 | 0.235602 | 0.718085 | 1.594817 | 0.087872 | 1.950015 | 0.555094 |
| 3 | (william shakespeare) | (Renaissance) | 0.123909 | 0.549738 | 0.123909 | 1.000000 | 1.819048 | 0.055792 | inf | 0.513944 |

- Nije potrebno domensko predznanje

- Klasterovanje kao mocan alat za analizu i pretprocesiranja

- Problemi nebalansiranosti klasa

- Prokletstvo dimenzionalnosti

- Prokletstvo dimenzionalnosti



- Uvek treba imati na umu da su podaci, iako transormisani u numeričke, ustvari tekstualni!