# ANALIZA SKUPA PODATAKA
# COFFEE QUALITY DATABASE FROM CQI

ANJA CVETKOVIĆ

ISTRAŽIVANJE PODATAKA 1

# ANALIZA PODATAKA I PRETPROCESIRANJE

**Mere kvaliteta**

- Aroma – miris kafe

- Flavor – ukus kafe

- Aftertaste – da li kafa ostavlja prijatnu senzaciju na paleti

- Acidity – prijatna oštrinu u ukusu kafe

- Body – izražava teksturu kafe, za koju je poželjno da ima težinu i da bude kremasta

- Balance – ukus se ne lokaliziju na jednom mestu palete

- Uniformity – konzistentnost ukusa

- Cup Cleanliness – "čista šoljica" označava da se ne javljaju arome koje nisu od kafe a koje su posledica defekata

- Sweetness – slatkoća koja se oseti na vrhu jezika, poželjna osobina

- Moisture – preporučena vlažnost je oko 11.5% potpuno procesuiranog zrna

- Defects – mogu da budu prve i druge kategorije, na osnovu uzorka od 350g, defekti prvi kategorije su skroz crna ili kisela zrna, tragovi rastinja ili kamenčići, dok su defekti druge kategorije insekti u uzorku, oštećenja od vode i slično.

- Total Cup Points – zbir ocen prethodnih mera, vrednost na [1,100] i ciljna promenljiva
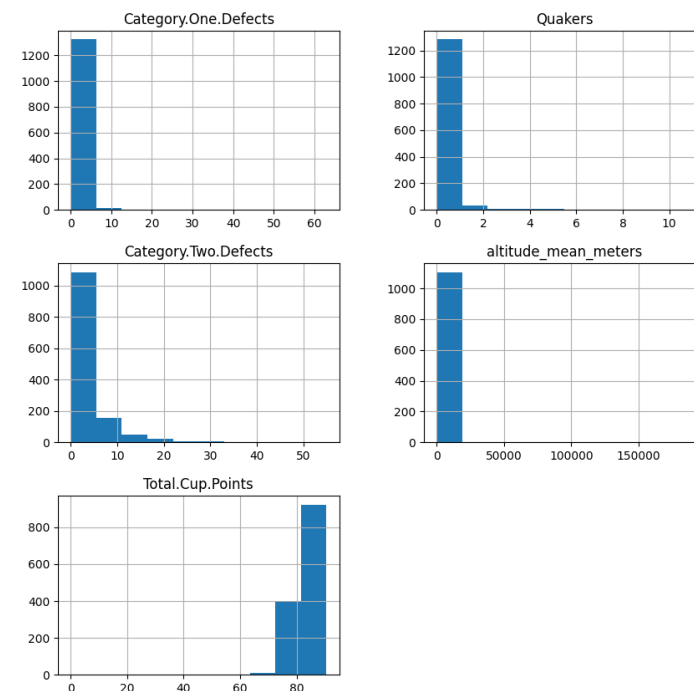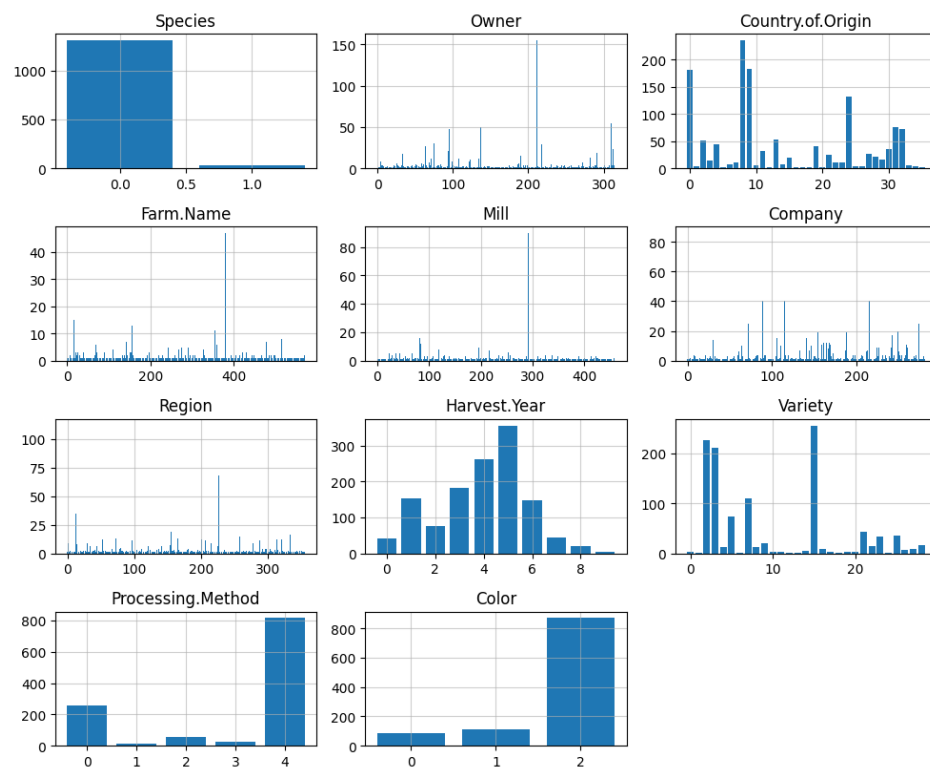
# ANALIZA PODATAKA I PRETPROCESIRANJE

**Metapodaci o zrnu**

- Processing Method
- Color
- Species (Arabica/Robusta)
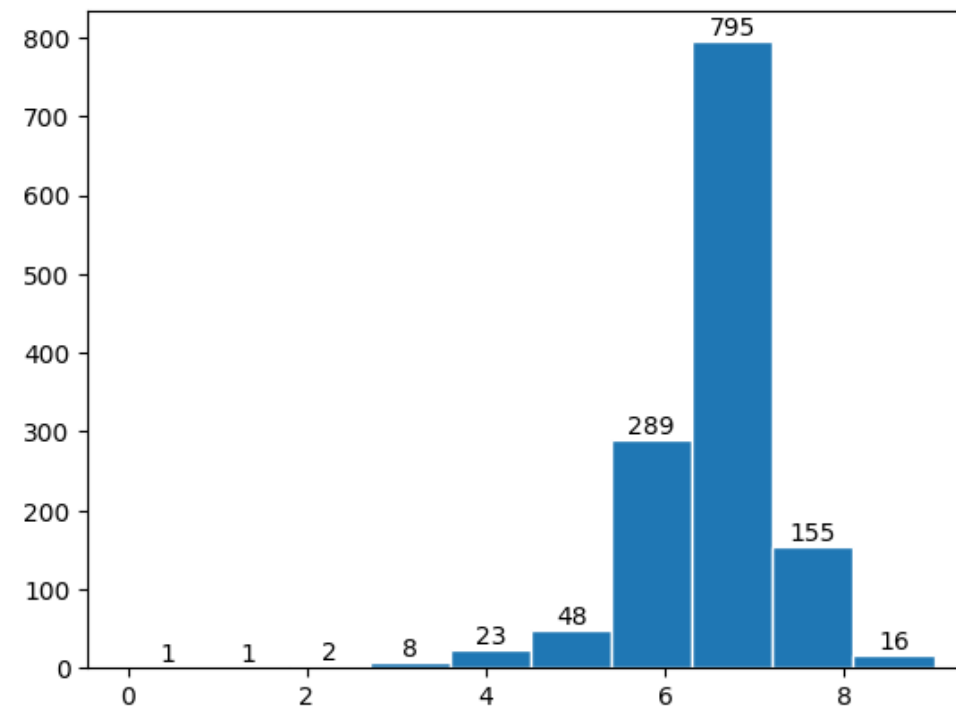- Variety – podvrsta
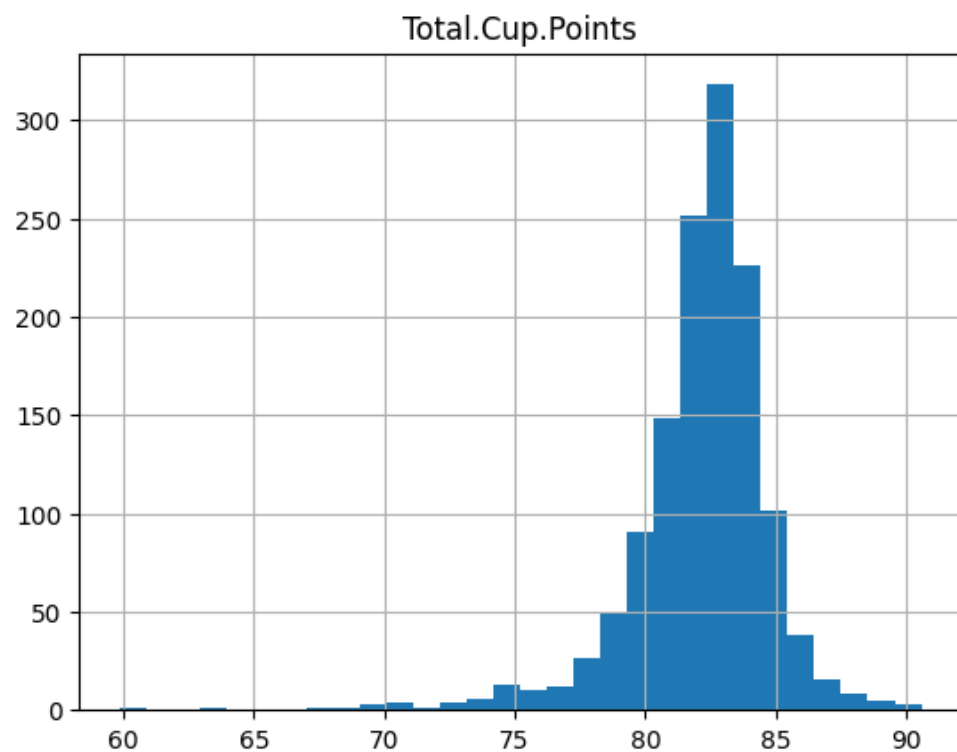- Quakers - broj nezrelih zrna

**Metapodaci od poreklu**

- Owner
- Country of Origin
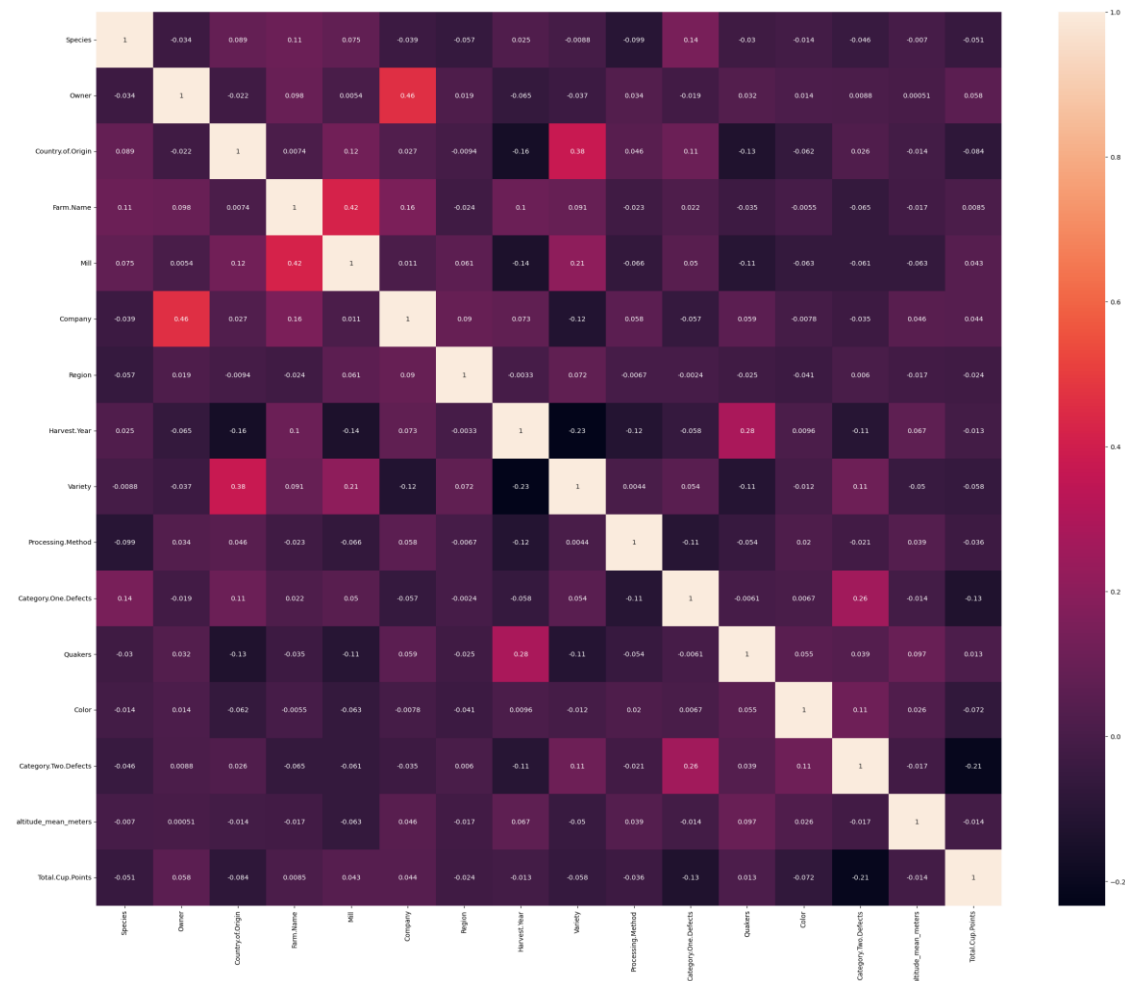- Farm Name
- Lot Number
- Mill
- Company
- Altitude
- Region

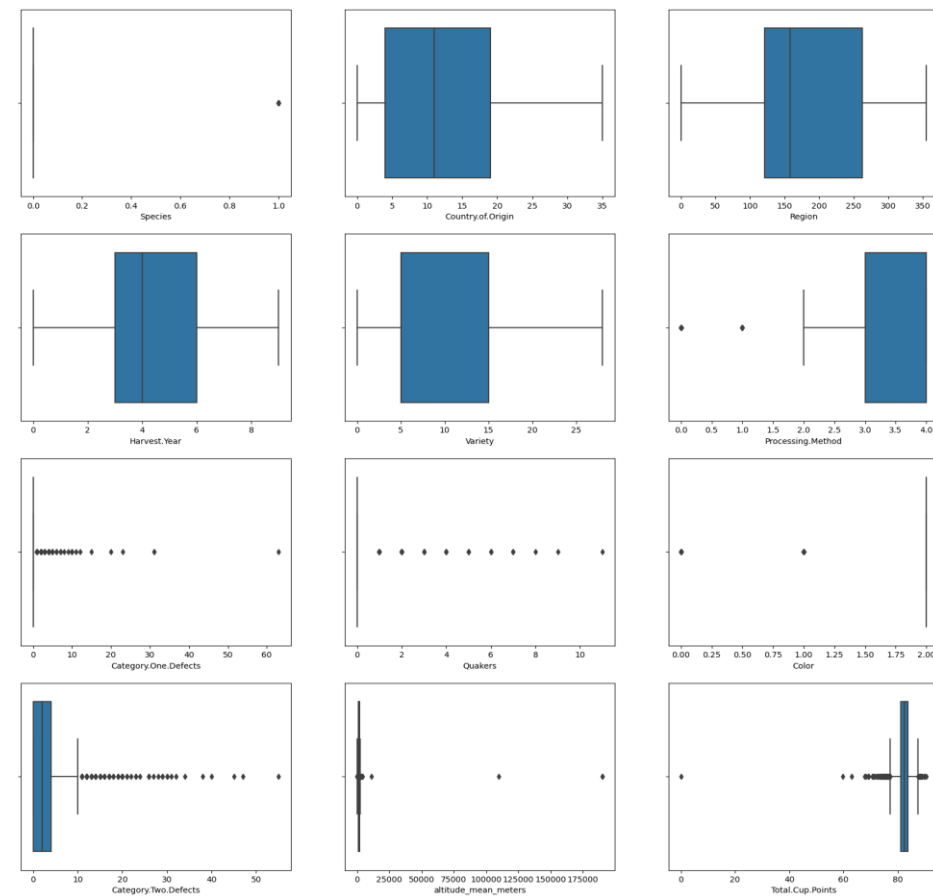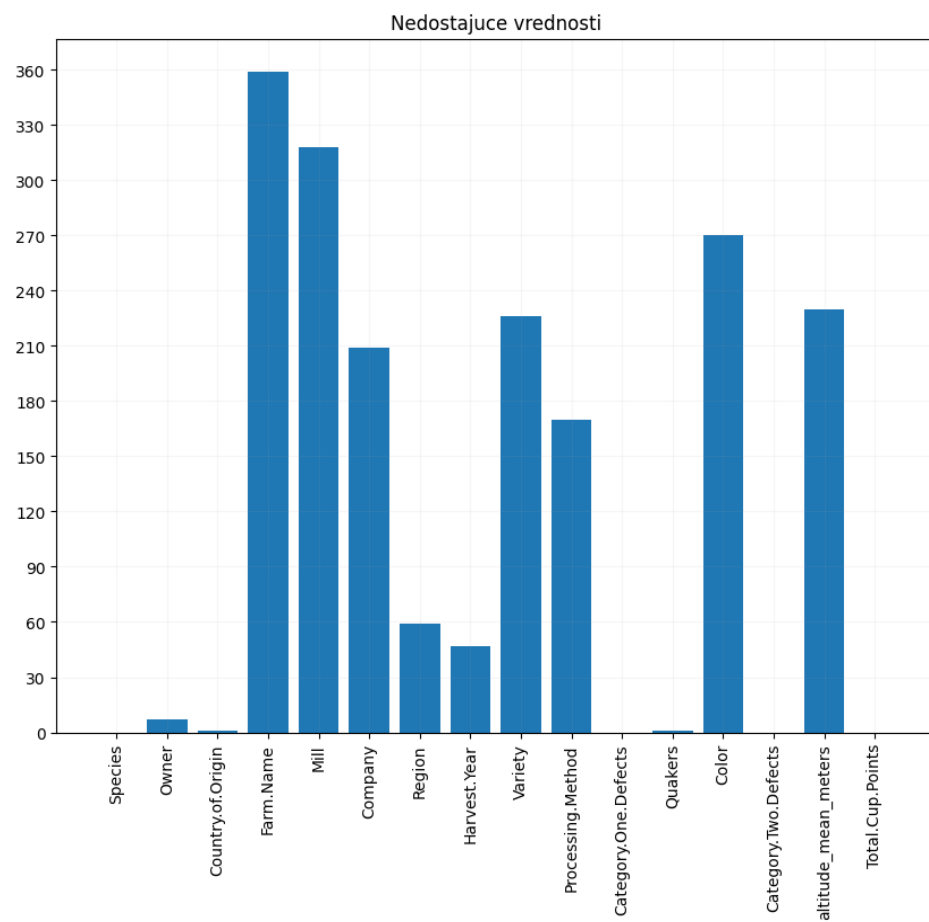# ANALIZA PODATAKA I PRETPROCESIRANJE



Total.Cup.Points

# ANALIZA PODATAKA I PRETPROCESIRANJE

# ANALIZA PODATAKA I PRETPROCESIRANJE

# KLASIFIKACIJA-STABLA ODLUČIVANJA



Stablo odlučivanja sa nasumičnim hiper-parametrima i nebalansiranim podacima

# KLASIFIKACIJA-STABLA ODLUČIVANJA



Stablo odlučivanja sa podešenim hiper-parametrima i balansiranim podacima

# KLASIFIKACIJA-RANDOM FOREST

Classification report for model RandomForestClassifier on training data
Accuracy: 0.9974886991461577

---

|      | pre  | rec  | spe  | f1   | geo  | iba  | sup  |
|------|------|------|------|------|------|------|------|
| 4.0  | 1.00 | 1.00 | 1.00 | 1.00 | 1.00 | 1.00 | 467  |
| 5.0  | 1.00 | 1.00 | 1.00 | 1.00 | 1.00 | 1.00 | 425  |
| 6.0  | 1.00 | 1.00 | 1.00 | 1.00 | 1.00 | 1.00 | 192  |
| 7.0  | 1.00 | 0.98 | 1.00 | 0.99 | 0.99 | 0.98 | 123  |
| 8.0  | 0.99 | 0.99 | 1.00 | 0.99 | 1.00 | 0.99 | 296  |
| 9.0  | 0.99 | 1.00 | 1.00 | 1.00 | 1.00 | 1.00 | 488  |
| avg / total | 1.00 | 1.00 | 1.00 | 1.00 | 1.00 | 1.00 | 1991 |

---

Confusion matrix for model RandomForestClassifier on training data

---

|   | 4   | 5   | 6   | 7   | 8   | 9   |
|---|-----|-----|-----|-----|-----|-----|
| 4 | 466 | 0   | 0   | 0   | 1   | 0   |
| 5 | 0   | 425 | 0   | 0   | 0   | 0   |
| 6 | 0   | 0   | 192 | 0   | 0   | 0   |
| 7 | 0   | 0   | 0   | 121 | 1   | 1   |
| 8 | 0   | 0   | 0   | 0   | 294 | 2   |
| 9 | 0   | 0   | 0   | 0   | 0   | 488 |

Classification report for model RandomForestClassifier on training data
Accuracy: 0.3401826484018265

---

|      | pre  | rec  | spe  | f1   | geo  | iba  | sup  |
|------|------|------|------|------|------|------|------|
| 4.0  | 0.00 | 0.00 | 0.92 | 1.00 | 0.00 | 0.00 | 8    |
| 5.0  | 0.08 | 0.31 | 0.87 | 0.13 | 0.52 | 0.26 | 16   |
| 6.0  | 0.28 | 0.26 | 0.81 | 0.27 | 0.46 | 0.20 | 95   |
| 7.0  | 0.70 | 0.32 | 0.79 | 0.44 | 0.51 | 0.24 | 263  |
| 8.0  | 0.29 | 0.61 | 0.80 | 0.39 | 0.70 | 0.48 | 51   |
| 9.0  | 0.12 | 0.60 | 0.95 | 0.19 | 0.75 | 0.55 | 5    |
| avg / total | 0.52 | 0.34 | 0.81 | 0.40 | 0.51 | 0.26 | 438  |

---

Confusion matrix for model RandomForestClassifier on training data

---

|   | 4  | 5  | 6  | 7  | 8  | 9  |
|---|----|----|----|----|----|----|
| 4 | 0  | 6  | 1  | 1  | 0  | 0  |
| 5 | 4  | 5  | 4  | 2  | 1  | 0  |
| 6 | 18 | 13 | 25 | 23 | 15 | 1  |
| 7 | 12 | 32 | 57 | 85 | 59 | 18 |
| 8 | 1  | 3  | 3  | 9  | 31 | 4  |
| 9 | 0  | 0  | 0  | 1  | 1  | 3  |

---



ROC multiclass one vs rest for random forest

# KLASIFIKACIJA-KNN

# KLASIFIKACIJA-KNN

# KLASIFIKACIJA-SVM

```
Classification report for model SVC on training data
-----------------------------------------------------------

              precision    recall  f1-score   support

         4.0       1.00      1.00      1.00       467
         5.0       1.00      1.00      1.00       425
         6.0       1.00      0.99      1.00       192
         7.0       0.99      0.99      0.99       123
         8.0       0.99      0.98      0.98       296
         9.0       0.99      1.00      0.99       488

    accuracy                           0.99      1991
   macro avg       1.00      0.99      0.99      1991
weighted avg       0.99      0.99      0.99      1991
```

```
Classification report for model SVC on test data
-----------------------------------------------------------

              precision    recall  f1-score   support

         4.0       0.03      0.12      0.05         8
         5.0       0.15      0.50      0.23        16
         6.0       0.23      0.20      0.21        95
         7.0       0.68      0.38      0.49       263
         8.0       0.22      0.39      0.28        51
         9.0       0.04      0.20      0.06         5

    accuracy                           0.34       438
   macro avg       0.22      0.30      0.22       438
weighted avg       0.49      0.34      0.38       438
```

```
-----------------------------------------------------------
Confusion matrix for model SVC on training data
-----------------------------------------------------------

      4    5    6    7    8    9
4   467    0    0    0    0    0
5     1  424    0    0    0    0
6     0    0  191    0    1    0
7     0    0    0  122    1    0
8     0    0    0    1  289    6
9     0    0    0    0    0  488

-----------------------------------------------------------
```
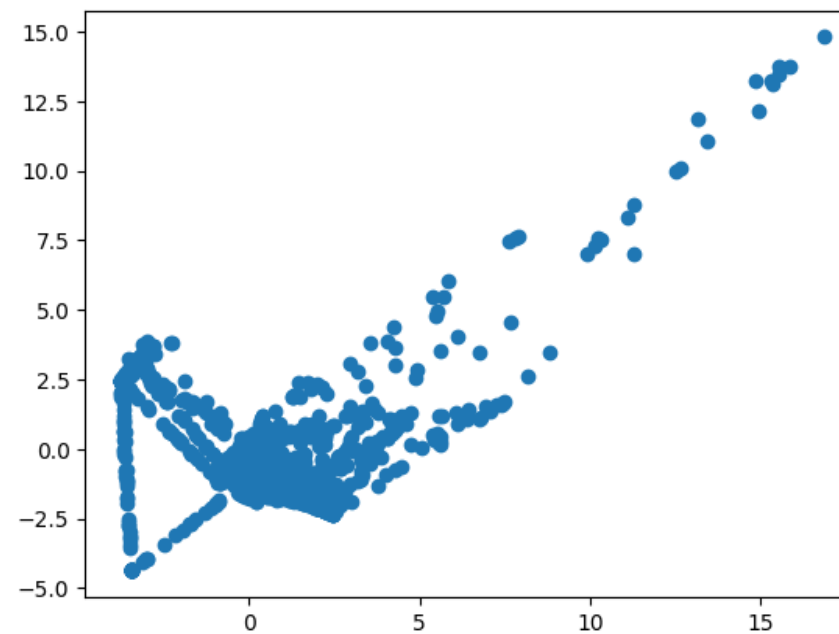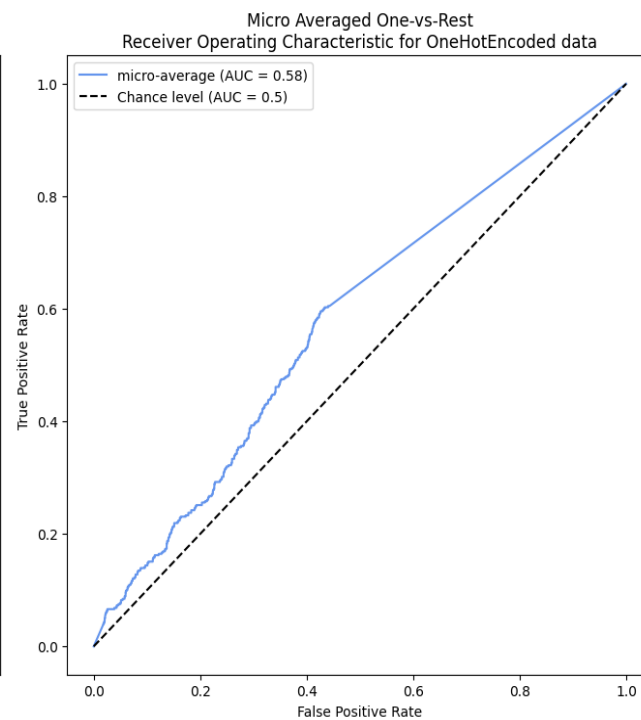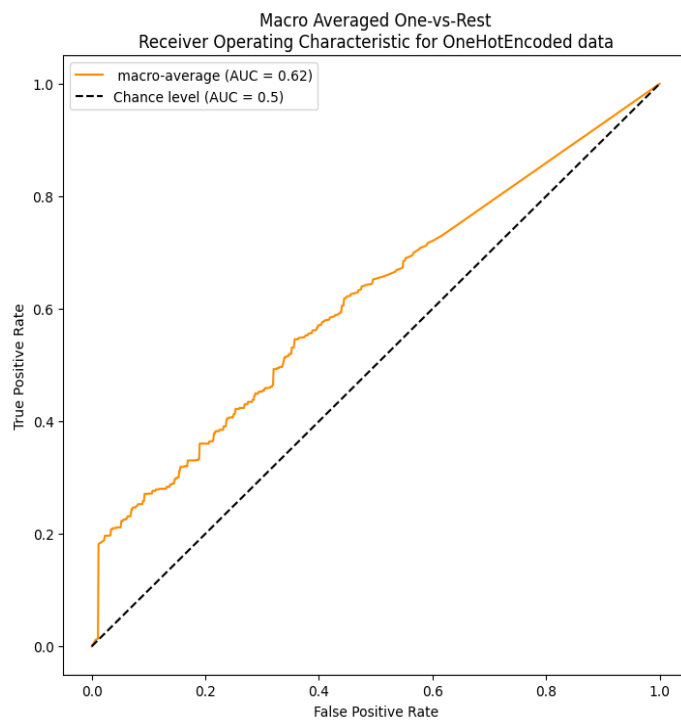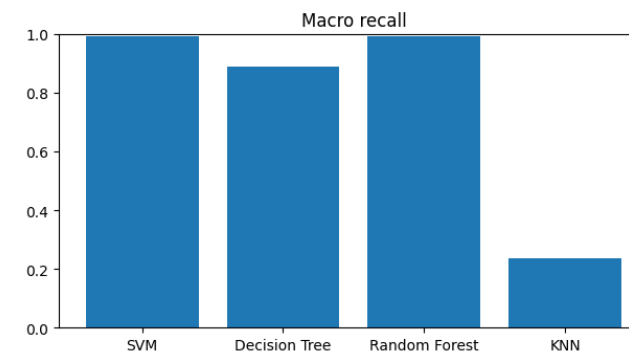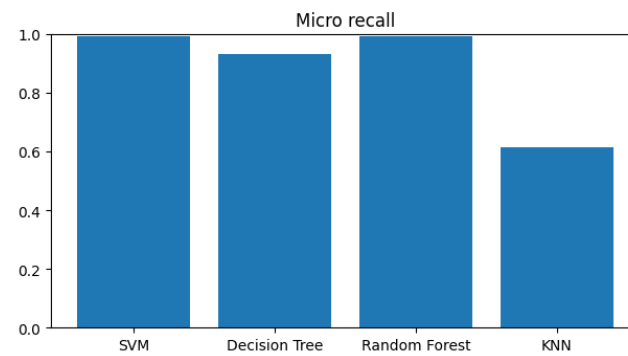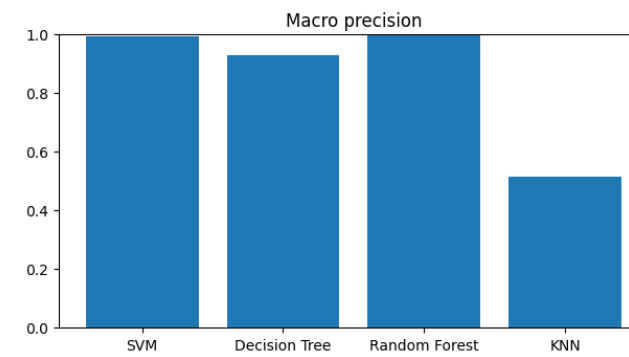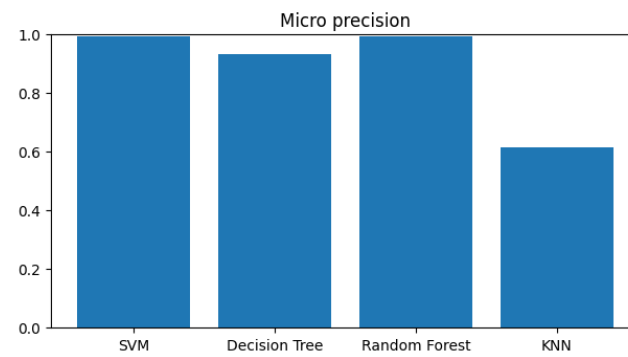
```
-----------------------------------------------------------
Confusion matrix for model SVC on test data
-----------------------------------------------------------

      4    5    6    7    8    9
4     1    4    1    2    0    0
5     1    8    1    4    2    0
6    15   13   19   28   17    3
7    17   26   56  100   48   16
8     2    2    5   14   20    8
9     0    0    1    0    3    1

-----------------------------------------------------------
```
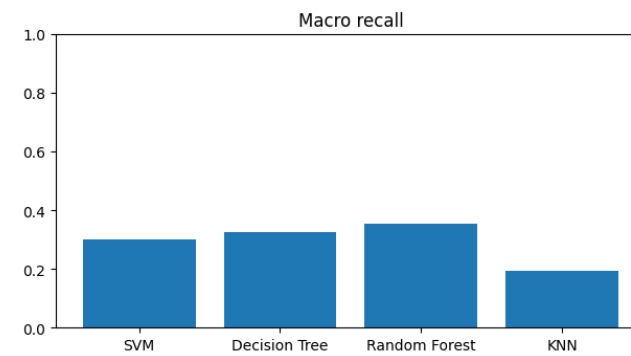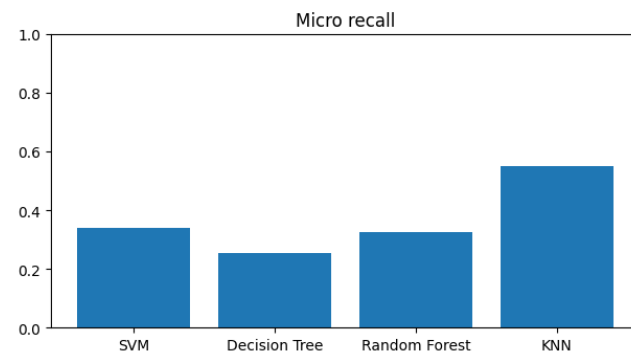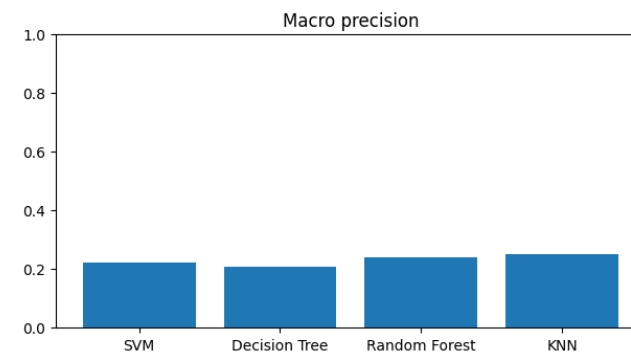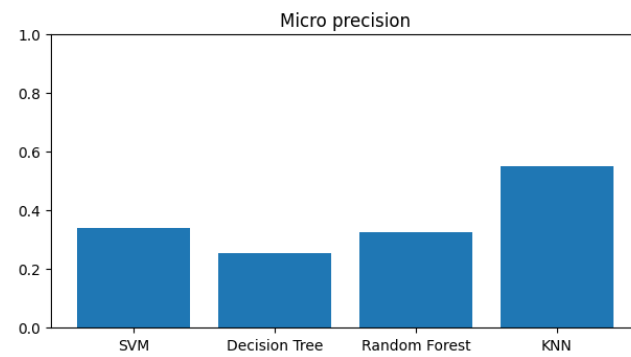
# KLASIFIKACIJA

# KLASIFIKACIJA

# KLASIFIKACIJA



Micro Averaged One-vs-Rest Receiver Operating Characteristic

- micro-average AUC SVM (AUC = 0.67)
- micro-average AUC Decision Tree (AUC = 0.62)
- micro-average AUC Random Forest (AUC = 0.71)
- micro-average AUC KNN (AUC = 0.86)
- Chance level (AUC = 0.5)

Macro Averaged One-vs-Rest Receiver Operating Characteristic

- macro-average SVM (AUC = 0.64)
- macro-average Decision Tree (AUC = 0.59)
- macro-average Random Forest (AUC = 0.65)
- macro-average KNN (AUC = 0.60)
- Chance level (AUC = 0.5)

# KLASTEROVANJE-KMEANS

# KLASTEROVANJE-KMEANS

# KLASTEROVANJE-FUZZY CMEANS

# KLASTEROVANJE-DBSCAN

# PRAVILA PRIDRUŽIVANJA-APRIORI

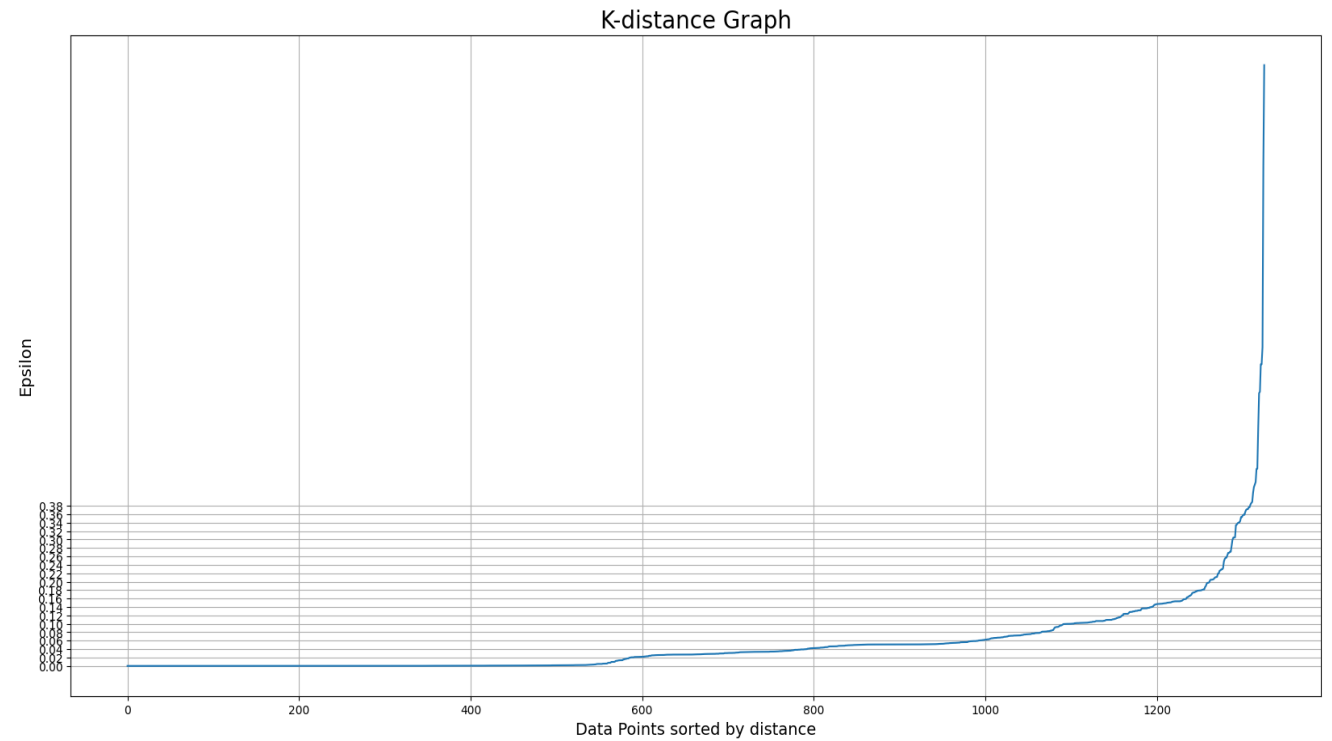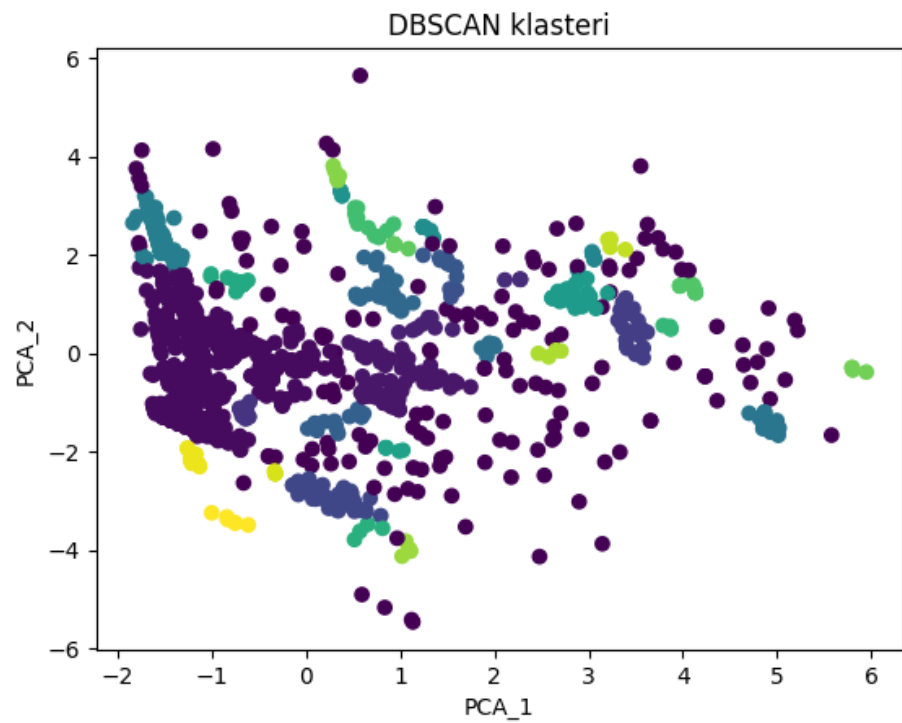| Consequent | Antecedent | Support % | Confidence % | Lift |
|---|---|---|---|---|
| Harvest.Year = 2012 | Variety = Typica<br>Country.of.Origin = Mexico | 10.256 | 82.353 | 2.676 |
| Harvest.Year = 2012 | Variety = Typica<br>Country.of.Origin = Mexico<br>Quakers = 0.000000 | 10.256 | 82.353 | 2.676 |
| Harvest.Year = 2012 | Variety = Typica<br>Country.of.Origin = Mexico<br>Species = Arabica | 10.256 | 82.353 | 2.676 |
| Harvest.Year = 2012 | Variety = Typica<br>Country.of.Origin = Mexico<br>Quakers = 0.000000<br>Species = Arabica | 10.256 | 82.353 | 2.676 |
| Variety = Caturra | Country.of.Origin = Colombia<br>Category.One.Defects = 0 | 10.935 | 95.172 | 2.64 |
| Variety = Caturra | Country.of.Origin = Colombia<br>Category.One.Defects = 0<br>Species = Arabica | 10.935 | 95.172 | 2.64 |
| Variety = Caturra | Country.of.Origin = Colombia<br>Category.One.Defects = 0<br>Quakers = 0.000000 | 10.256 | 94.853 | 2.631 |
| Variety = Caturra | Country.of.Origin = Colombia<br>Category.One.Defects = 0<br>Quakers = 0.000000<br>Species = Arabica | 10.256 | 94.853 | 2.631 |
| Variety = Caturra | Country.of.Origin = Colombia<br>Processing.Method = Washed / Wet<br>Color = Green | 10.709 | 93.662 | 2.598 |
| Variety = Caturra | Country.of.Origin = Colombia<br>Processing.Method = Washed / Wet<br>Color = Green | 10.709 | 93.662 | 2.598 |

# ZAKLJUČAK

- Poteškoće i izazovi

- Skup podataka ili tehnika?

- Pitanja?