

▼ Biblioteke

```
1 import pandas as pd
2 import matplotlib.pyplot as plt
3 import seaborn as sns
4 import numpy as np
5 from sklearn.preprocessing import StandardScaler
6 from sklearn.cluster import KMeans
7 import seaborn as sns
8 from sklearn.decomposition import PCA
9 from sklearn.cluster import MeanShift,BisectingKMeans
10 from sklearn.cluster import DBSCAN
11 import plotly as py
12 import plotly.graph_objs as go
13 from plotly.offline import download_plotlyjs, init_notebook_mode, plot, iplot
14 from sklearn.metrics import silhouette_score, calinski_harabasz_score,davies_bouldin_score
15 from sklearn.metrics import adjusted_rand_score
16 import warnings
17 warnings.filterwarnings('ignore')
```

▼ Ucitavanje

```
1 url = "https://raw.githubusercontent.com/aleksicmilica/ml-projekat2/main/bank_marketing_dataset.csv"
2 dataset= pd.read_csv(url)
```

```
1 dataset.head()
```

	age	job	marital	education	default	housing	loan	contact	month	day_of_w
0	56	housemaid	married	basic.4y	no	no	no	telephone	may	r
1	57	services	married	high.school	unknown	no	no	telephone	may	r
2	37	services	married	high.school	no	yes	no	telephone	may	r
3	40	admin.	married	basic.6y	no	no	no	telephone	may	r
4	56	services	married	high.school	no	no	yes	telephone	may	r

5 rows × 21 columns

```
1 dataset.columns
```

```
Index(['age', 'job', 'marital', 'education', 'default', 'housing', 'loan',
      'contact', 'month', 'day_of_week', 'duration', 'campaign', 'pdays',
      'previous', 'poutcome', 'emp.var.rate', 'cons.price.idx',
      'cons.conf.idx', 'euribor3m', 'nr.employed', 'subscribed'],
      dtype='object')
```

```
1 dataset.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 41188 entries, 0 to 41187
Data columns (total 21 columns):
#   Column                Non-Null Count  Dtype
---  -
0   age                    41188 non-null  int64
1   job                    41188 non-null  object
2   marital                41188 non-null  object
3   education              41188 non-null  object
4   default                41188 non-null  object
5   housing                41188 non-null  object
6   loan                   41188 non-null  object
7   contact                41188 non-null  object
8   month                  41188 non-null  object
9   day_of_week            41188 non-null  object
10  duration                41188 non-null  int64
11  campaign                41188 non-null  int64
12  pdays                  41188 non-null  int64
13  previous                41188 non-null  int64
14  poutcome               41188 non-null  object
15  emp.var.rate           41188 non-null  float64
```

```

16 cons.price.idx 41188 non-null float64
17 cons.conf.idx 41188 non-null float64
18 euribor3m      41188 non-null float64
19 nr.employed    41188 non-null float64
20 subscribed     41188 non-null object
dtypes: float64(5), int64(5), object(11)
memory usage: 6.6+ MB

```

✓ Predobrada

✓ Provera prisustva missing values

```

1 total_missing = dataset.isnull().sum()
2 print(total_missing)
3 if (total_missing == 0).all():
4     print("Sve vrednosti su prisutne u dataset-u")

```

```

age          0
job          0
marital      0
education    0
default      0
housing      0
loan         0
contact      0
month        0
day_of_week  0
duration     0
campaign     0
pdays      0
previous     0
poutcome    0
emp.var.rate 0
cons.price.idx 0
cons.conf.idx 0
euribor3m    0
nr.employed  0
subscribed   0
dtype: int64
Sve vrednosti su prisutne u dataset-u

```

✓ Provera da li su svi podaci odgovarajuceg tipa

```

1 irregular_values = (dataset["duration"] < 0).sum() + (dataset["campaign"] < 0).sum() + (dataset["pdays"] < 0).sum() + (dataset["previous
2 print(irregular_values)

```

```
0
```

```

1 numerical_type = pd.api.types.is_numeric_dtype(dataset["age"].dtype) and pd.api.types.is_numeric_dtype(dataset["duration"].dtype) and pd.
2 print(numerical_type)

```

```
True
```

```

1 if irregular_values == 0 and numerical_type :
2     print("Svi podaci su odgovarajuceg tipa")
3 else:
4     print("Podaci nisu odgovarajuceg tipa")

```

```
Svi podaci su odgovarajuceg tipa
```

✓ Vizuelizacija podataka

```

1 columns = dataset.columns
2 n = len(columns)

```

```
1 fig, axes = plt.subplots(n//4+1, 4, figsize=(20, 20))
2 axes = axes.flatten()
3 for i, column_name in enumerate(columns):
4     sns.histplot(x=column_name, data=dataset, ax=axes[i])
5     axes[i].set_title(column_name)
6     axes[i].set_xlabel('')
7     axes[i].set_ylabel('')
8
9 plt.tight_layout()
10 plt.show()
```



```
1 dataset_2 = dataset.copy()
```

▼ Labeliranje kategorickih tipova

```

1 lista = ["job","marital", "education", "default", "housing", "loan", "contact", "month", "day_of_week", "poutcome", "subscribed"]
2 for el in lista:
3     possible_location = dataset[el].unique().tolist()
4     # possible_location.insert(0,'nan')
5     dict_help={}
6     for i in range(len(possible_location)):
7         dict_help[possible_location[i]] = i
8     print(dict_help)
9     dataset.replace({el:dict_help},inplace = True)
10 dataset

{'housemaid': 0, 'services': 1, 'admin.': 2, 'blue-collar': 3, 'technician': 4, 'retired': 5, 'management': 6, 'unemployed': 7, 'self-emp': 8, 'married': 0, 'single': 1, 'divorced': 2, 'unknown': 3}
{'basic.4y': 0, 'high.school': 1, 'basic.6y': 2, 'basic.9y': 3, 'professional.course': 4, 'unknown': 5, 'university.degree': 6, 'illiterate': 7}
{'no': 0, 'unknown': 1, 'yes': 2}
{'no': 0, 'yes': 1, 'unknown': 2}
{'no': 0, 'yes': 1, 'unknown': 2}
{'telephone': 0, 'cellular': 1}
{'may': 0, 'jun': 1, 'jul': 2, 'aug': 3, 'oct': 4, 'nov': 5, 'dec': 6, 'mar': 7, 'apr': 8, 'sep': 9}
{'mon': 0, 'tue': 1, 'wed': 2, 'thu': 3, 'fri': 4}
{'nonexistent': 0, 'failure': 1, 'success': 2}
{'no': 0, 'yes': 1}

```

	age	job	marital	education	default	housing	loan	contact	month	day_of_week	...	campaign	pdays	previous	poutcome	emp.var
0	56	0	0	0	0	0	0	0	0	0	...	1	999	0	0	
1	57	1	0	1	1	0	0	0	0	0	...	1	999	0	0	
2	37	1	0	1	0	1	0	0	0	0	...	1	999	0	0	
3	40	2	0	2	0	0	0	0	0	0	...	1	999	0	0	
4	56	1	0	1	0	0	1	0	0	0	...	1	999	0	0	
...	
41183	73	5	0	4	0	1	0	1	5	4	...	1	999	0	0	
41184	46	3	0	4	0	0	0	1	5	4	...	1	999	0	0	
41185	56	5	0	6	0	1	0	1	5	4	...	2	999	0	0	
41186	44	4	0	4	0	0	0	1	5	4	...	1	999	0	0	
41187	74	5	0	4	0	1	0	1	5	4	...	3	999	1	1	

41188 rows × 21 columns

```
1 dataset_bckup = dataset.copy()
```

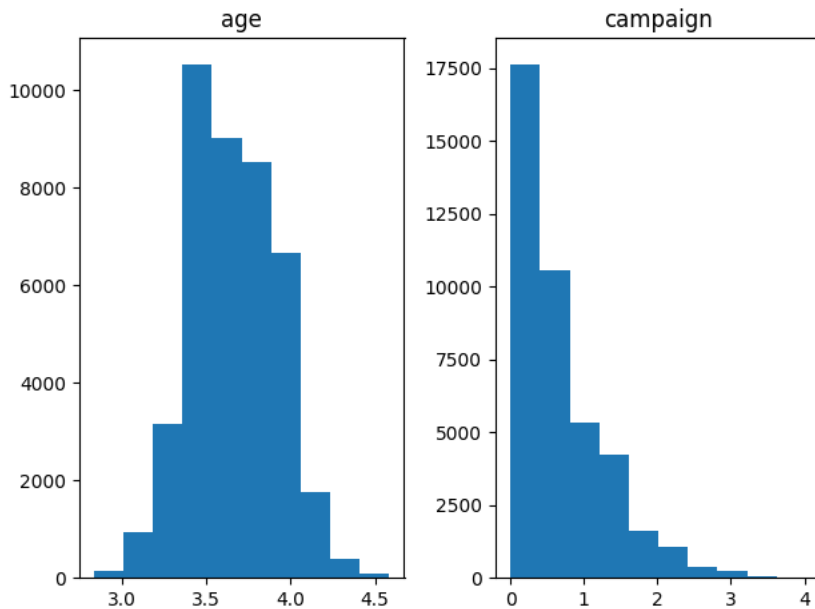
```
1 dataset = dataset_bckup.copy()
```

Iskoseni podaci normalizuju se primenom logaritma

```

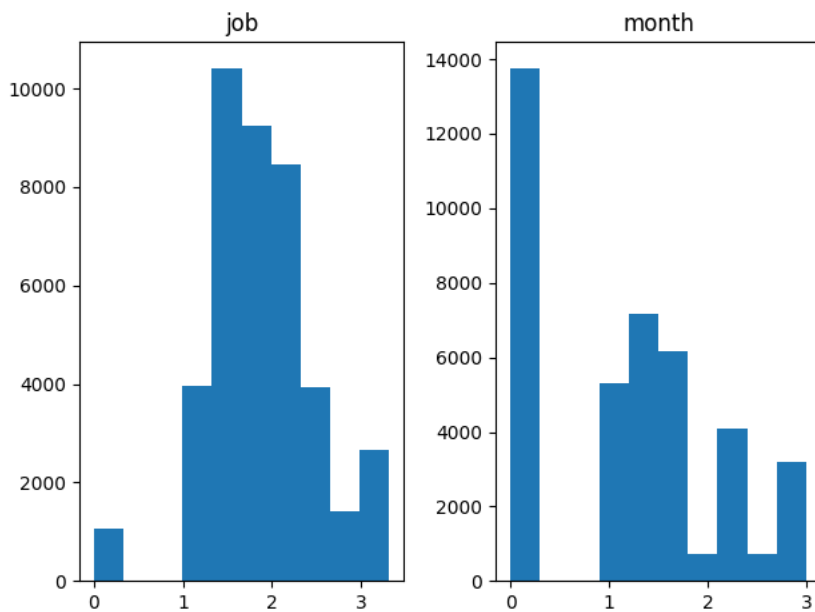
1 skewed_columns = ["age","campaign"]
2 fig, axes = plt.subplots(1,2)
3
4 for ind,col in enumerate(skewed_columns):
5     dataset[col] = np.log(dataset[col])
6
7     axes[ind].hist(dataset[col])
8     axes[ind].set_title(col)
9     axes[ind].set_xlabel('')
10    axes[ind].set_ylabel('')
11
12 plt.tight_layout()
13 plt.show()
14

```



Umereno iskoseni podaci normalizuju se primenom korena

```
1 moderate_skewed_columns = ["job", "month"]
2 fig, axes = plt.subplots(1, 2)
3
4 for ind, col in enumerate(moderate_skewed_columns):
5     dataset[col] = np.sqrt(dataset[col])
6
7     axes[ind].hist(dataset[col])
8     axes[ind].set_title(col)
9     axes[ind].set_xlabel('')
10    axes[ind].set_ylabel('')
11
12 plt.tight_layout()
13 plt.show()
14
```



Ostali numericki podaci standardizuju se

```

1 scaler = StandardScaler()
2 normalize_columns = ["education", "duration", "pdays", "previous", "poutcome", "emp.var.rate", "cons.price.idx", "cons.conf.idx", "euribor3m"]
3 # normalize_columns = ["education", "duration"]
4 normalized_data= scaler.fit_transform(dataset[normalize_columns])
5 dataset[normalize_columns] = normalized_data
6

```

```

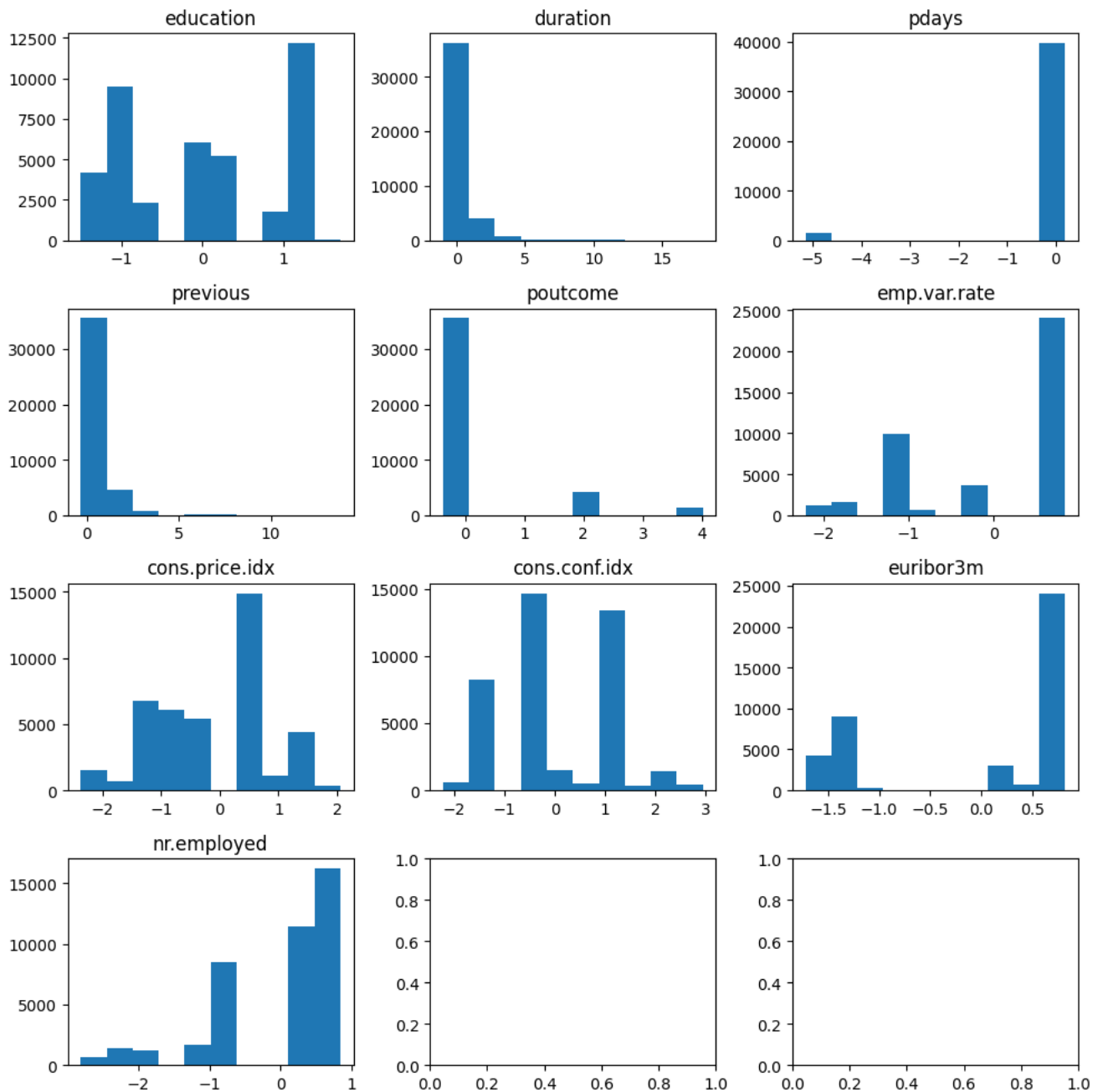
1

```

```

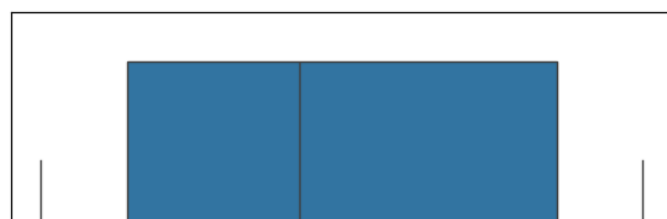
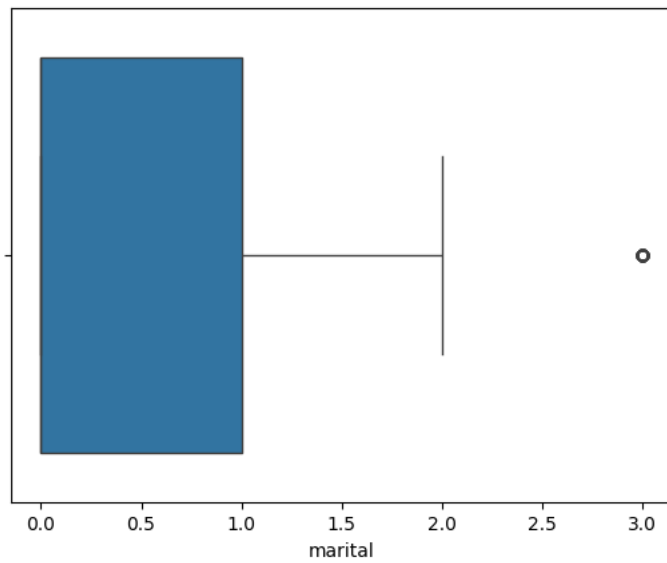
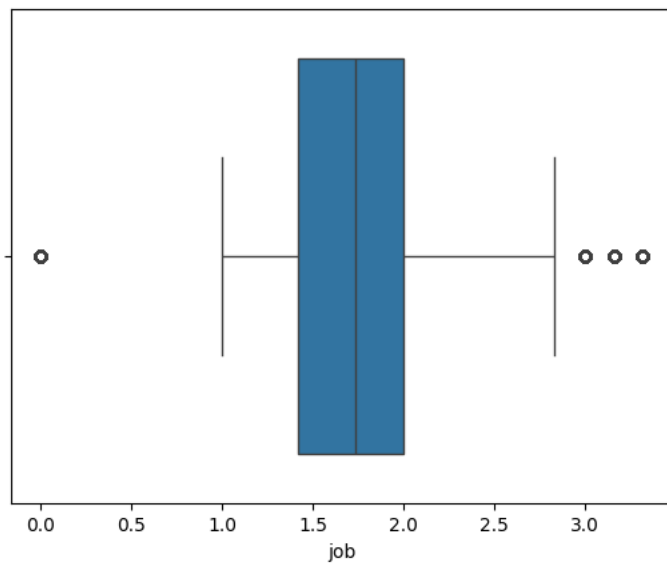
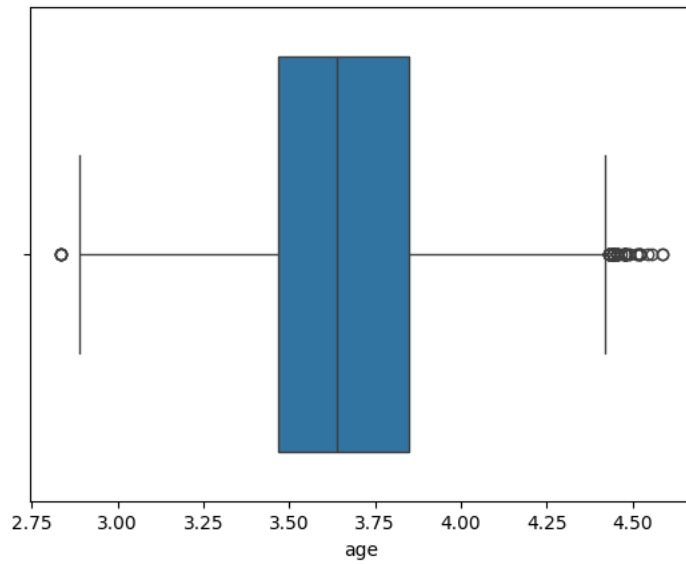
1 fig, axes = plt.subplots(4,3,figsize=(10,10))
2 axes = axes.flatten()
3 for ind,col in enumerate(normalize_columns):
4     axes[ind].hist(dataset[col])
5     axes[ind].set_title(col)
6     axes[ind].set_xlabel('')
7     axes[ind].set_ylabel('')
8
9 plt.tight_layout()
10 plt.show()
11

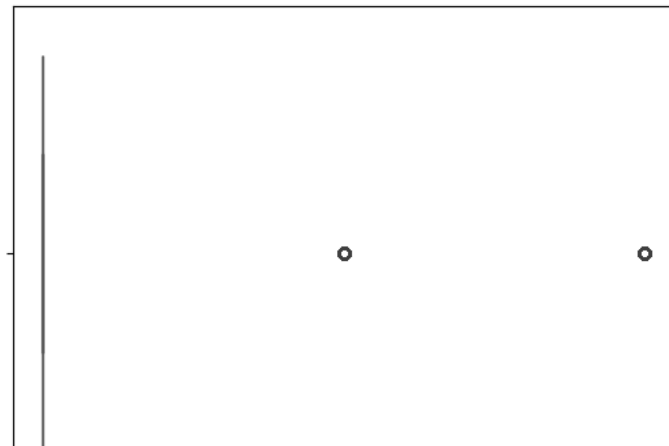
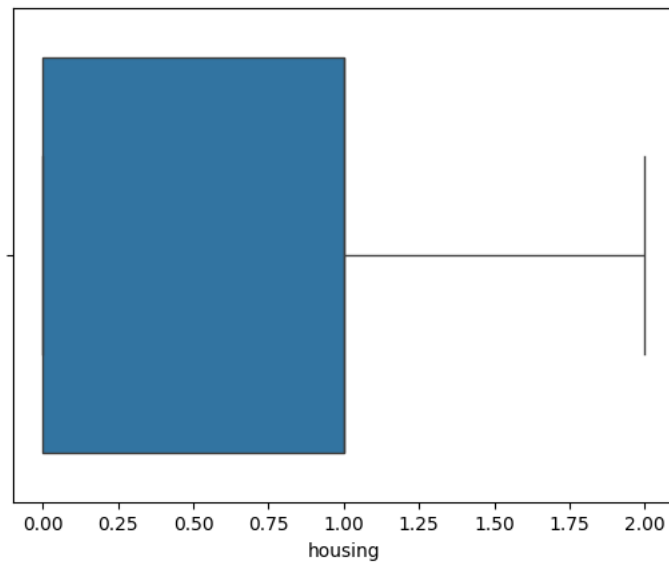
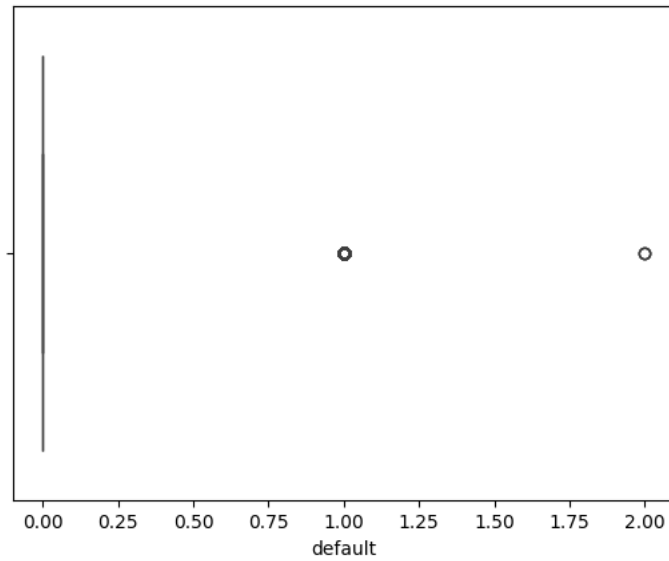
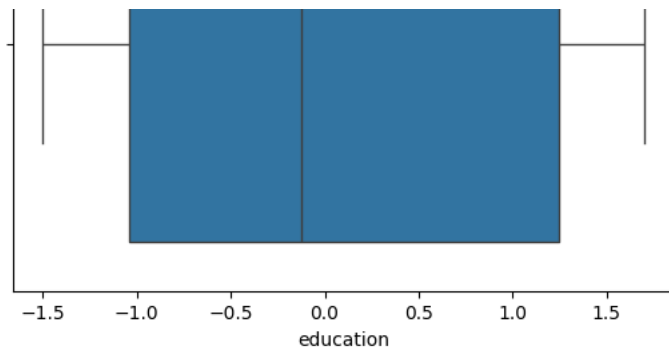
```

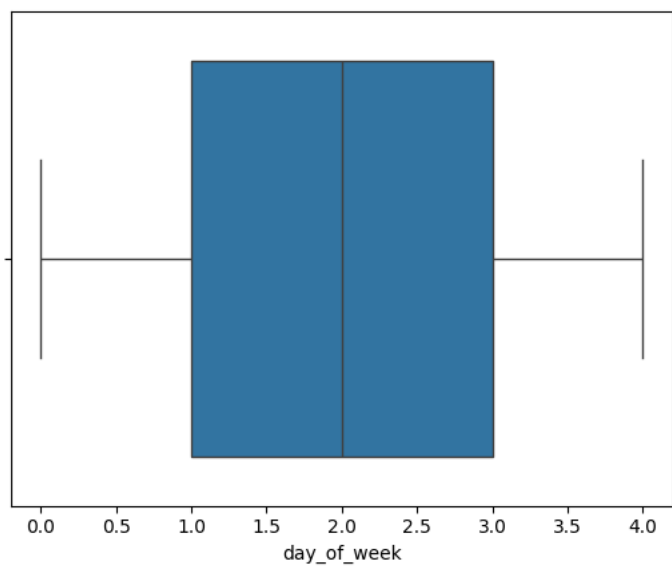
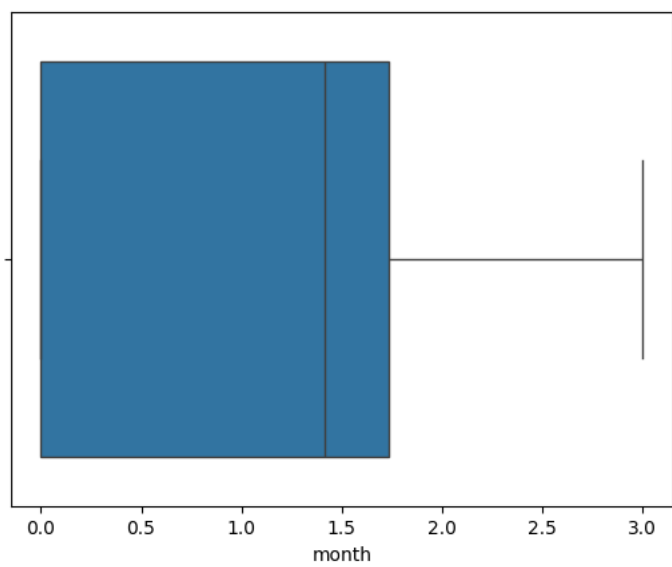
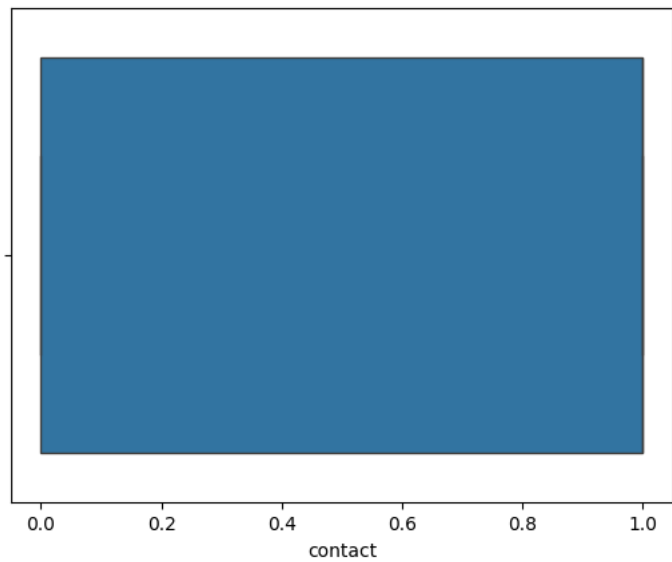
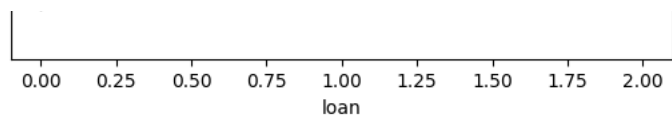


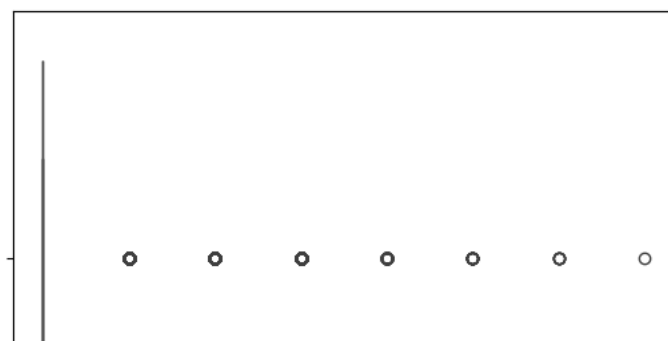
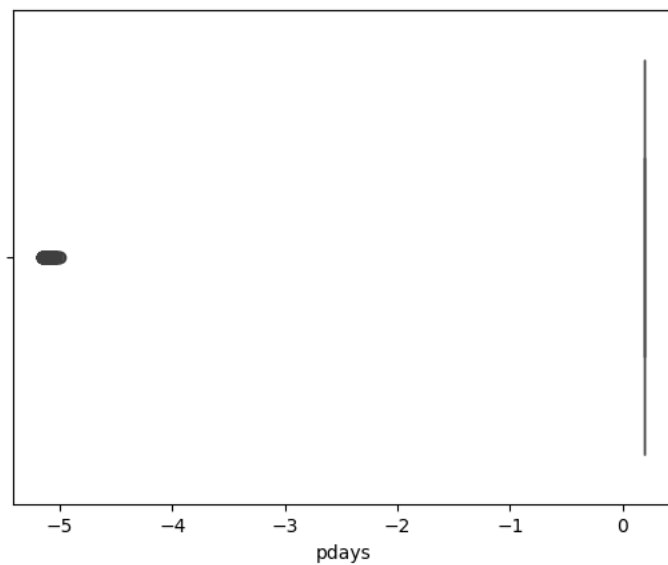
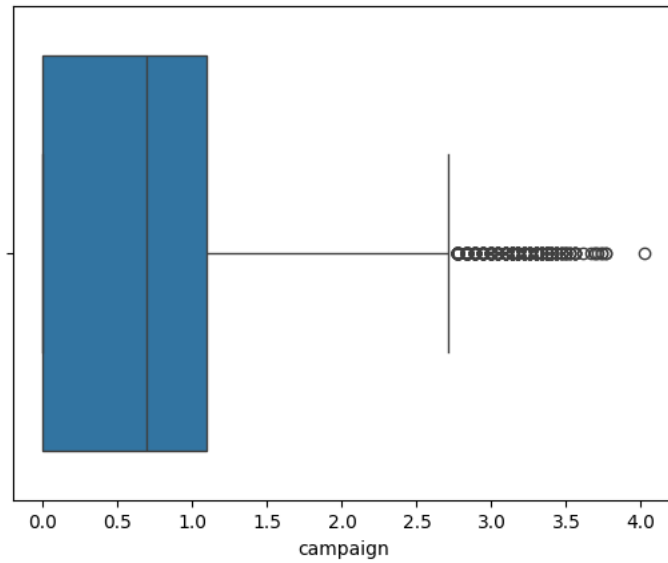
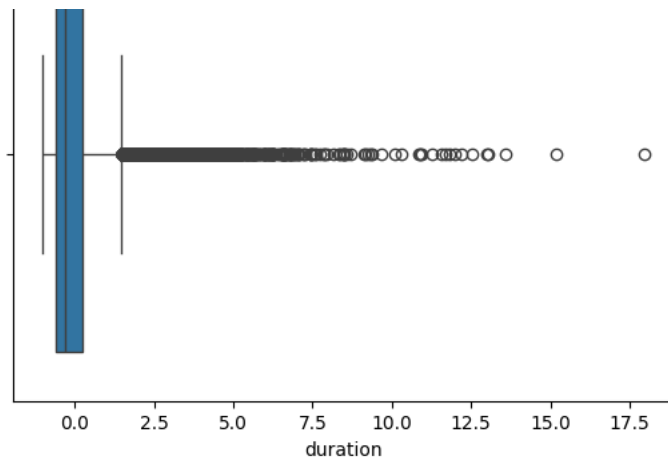
✓ Detekcija i otklanjanje outlier-a, normalizacija podataka

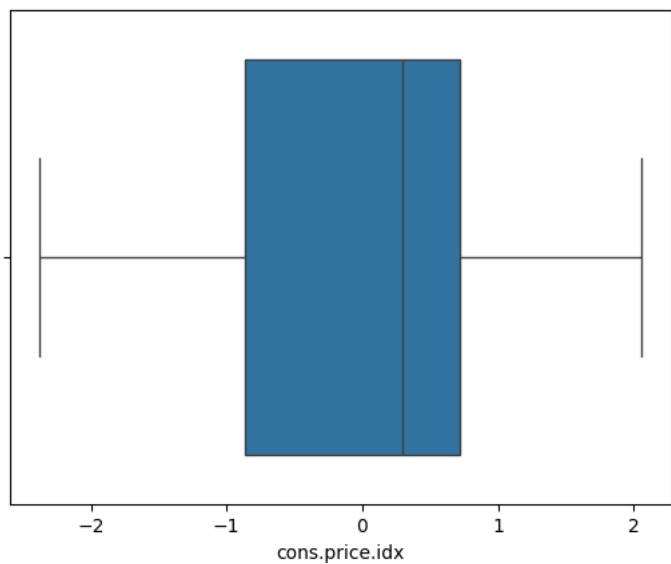
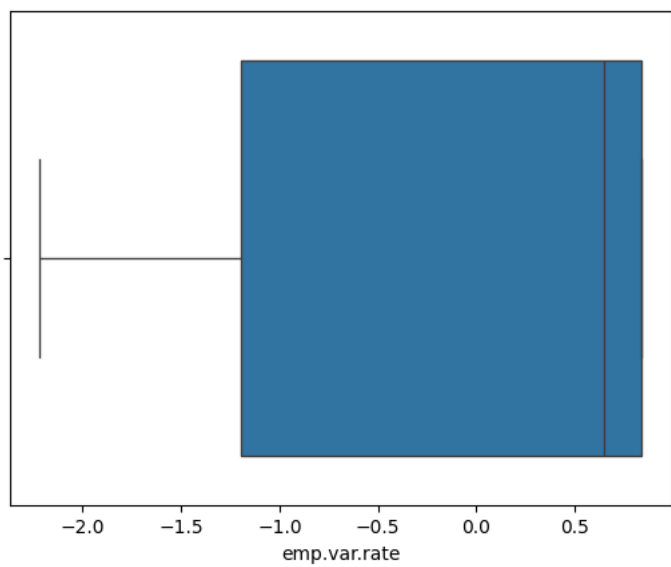
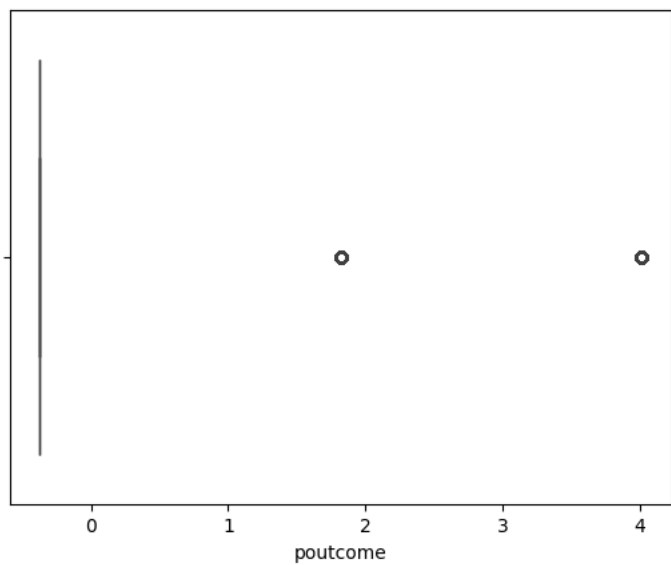
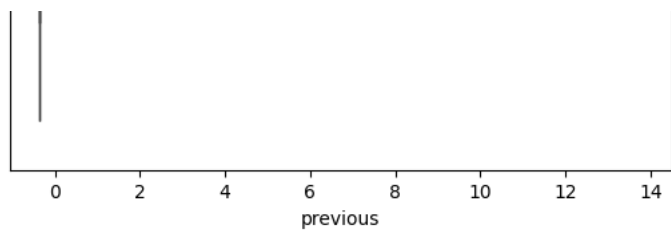
```
1 lista = ["age", "job", "marital", "education", "default", "housing", "loan", "contact", "month", "day_of_week", "duration", "campaign", "p  
2 for el in lista:  
3     sns.boxplot(x = dataset[el])  
4     plt.show()
```

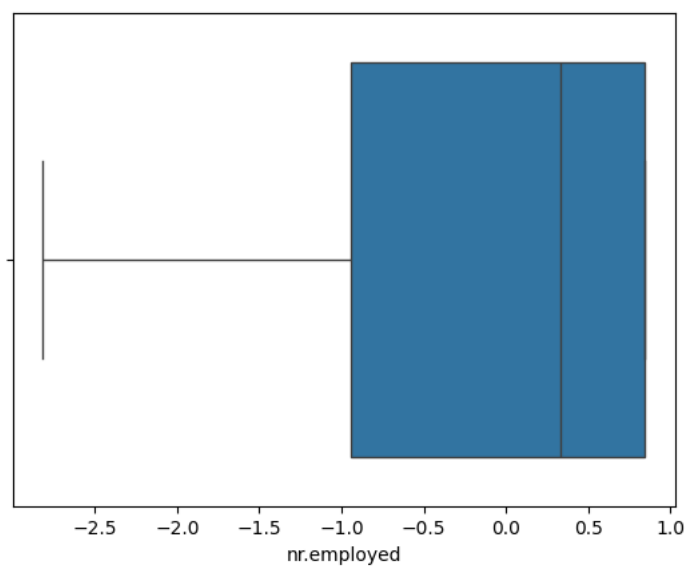
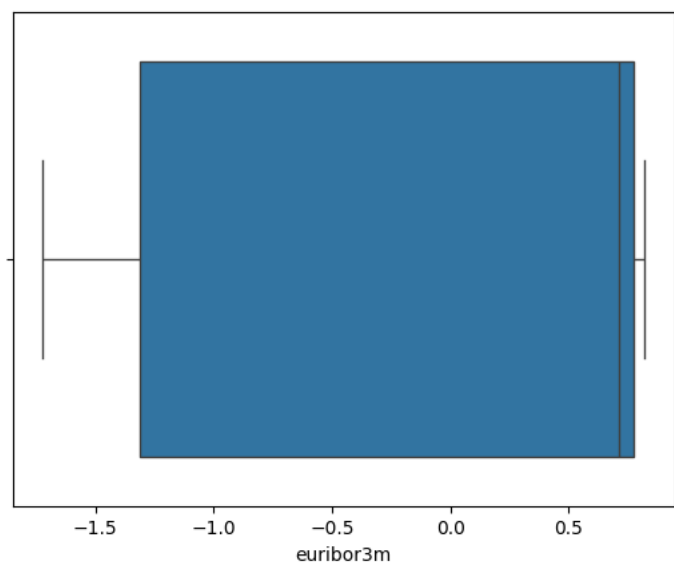
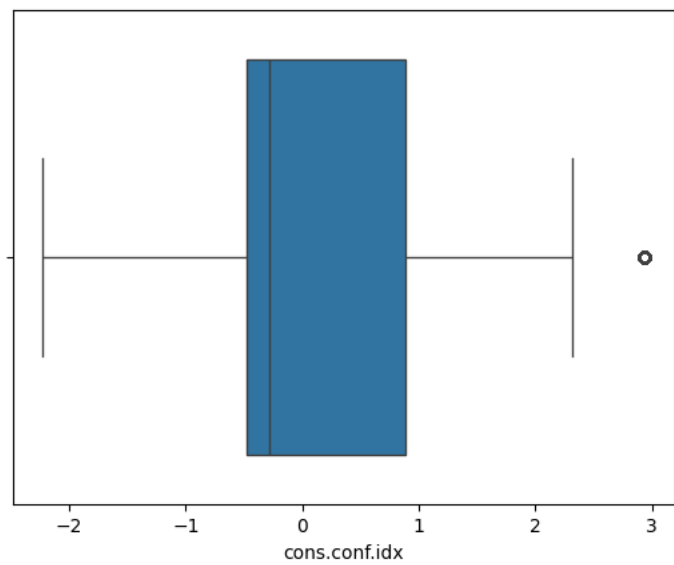



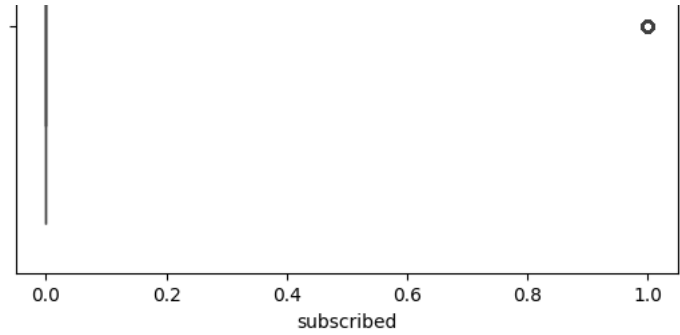












Druga varijanta dataseta

Kategoricki podaci kod kojih je jasno zastoupljenija jedna kategorija u odnosu na sve ostale transformisu se tako da imaju samo 2 moguće vrednosti

```
1 lista = ["marital", "default", "loan", "poutcome"]
2 dataset_2 = dataset.copy()
3 for ind,el in enumerate(lista):
4     dataset_2[el] = np.where(dataset_2[el] == 0, 0, 1)
5 dataset_2
```

	age	job	marital	education	default	housing	loan	contact	month	day_of_week	...	campaign	pdays	previous	p
0	4.025352	0.000000	0	-1.499673	0	0	0	0	0.000000	0	...	0.000000	0.195414	-0.349494	
1	4.043051	1.000000	0	-1.042111	1	0	0	0	0.000000	0	...	0.000000	0.195414	-0.349494	
2	3.610918	1.000000	0	-1.042111	0	1	0	0	0.000000	0	...	0.000000	0.195414	-0.349494	
3	3.688879	1.414214	0	-0.584550	0	0	0	0	0.000000	0	...	0.000000	0.195414	-0.349494	
4	4.025352	1.000000	0	-1.042111	0	0	1	0	0.000000	0	...	0.000000	0.195414	-0.349494	
...
41183	4.290459	2.236068	0	0.330573	0	1	0	1	2.236068	4	...	0.000000	0.195414	-0.349494	
41184	3.828641	1.732051	0	0.330573	0	0	0	1	2.236068	4	...	0.000000	0.195414	-0.349494	
41185	4.025352	2.236068	0	1.245696	0	1	0	1	2.236068	4	...	0.693147	0.195414	-0.349494	
41186	3.784190	2.000000	0	0.330573	0	0	0	1	2.236068	4	...	0.000000	0.195414	-0.349494	
41187	4.304065	2.236068	0	0.330573	0	1	0	1	2.236068	4	...	1.098612	0.195414	1.671136	

41188 rows × 21 columns

ALGORITMI

```
1 def eval_metrics(X,labels):
2     silhouette = silhouette_score(X, labels)
3     ch_index = calinski_harabasz_score(X,labels)
4     db_index = davies_bouldin_score(X, labels)
5     return (silhouette,db_index,ch_index)
```

```

1 def visualise_3d_data( datas, cluster_labels, n_clusters):
2     pca = PCA(n_components=3)
3     PCs = pd.DataFrame(pca.fit_transform(datas))
4     PCs.columns = ["PC1","PC2","PC3"]
5     PCs["cluster"] = cluster_labels
6     cluster = []
7     trace = []
8     for i in range(n_clusters):
9         cluster.append(PCs[PCs["cluster"]==i])
10        trace.append (go.Scatter3d(
11            x = cluster[i]["PC1"],
12            y = cluster[i]["PC2"],
13            z = cluster[i]["PC3"],
14            mode = "markers",
15            name = "Cluster "+str(i),
16            # marker = dict(color = 'rgba(255, 128, 255, 0.8)'),
17            text = None))
18
19
20 title = "Vizuelizacija klastera u 3D pomocu 3 PCA komponente"
21
22 layout = dict(title = title,
23               xaxis= dict(title= 'PC1',ticklen= 5,zeroline= False),
24               yaxis= dict(title= 'PC2',ticklen= 5,zeroline= False)
25               )
26
27 fig = dict(data = trace, layout = layout)
28
29 iplot(fig)

```

```

1 def visualise_2d_data(datas, cluster_labels, n_clusters):
2     pcas=PCA(n_components=2).fit_transform(datas)
3     df2=pd.DataFrame(pcas,columns=['PC1','PC2'])
4     sns.scatterplot(data=df2, x="PC1", y="PC2", hue=(cluster_labels+1))
5     plt.title("Vizuelizacija klastera u 2d pomocu 2 PCA komponente")
6     plt.show()

```

```

1 pcas=PCA(n_components=2).fit_transform(dataset)
2 df2=pd.DataFrame(pcas,columns=['PC1','PC2'])

```

```

1 def printFeatureImportance(df, pred, centers = None):
2
3     if centers is None:
4         cluster_means = df.groupby(pred).mean()
5     else:
6         cluster_means = centers
7
8     fig, ax = plt.subplots(figsize=(16, 6))
9     cluster_means.T.plot(kind='bar', ax=ax, cmap='plasma')
10    plt.title('Feature importance')
11    plt.show()

```

▼ Kmeans

▼ Odredjivanje broja klastera

```

1 from scipy.spatial.distance import cdist
2 inertias = []
3 mapping2 = {}
4 K = range(1, 12)
5 distortions = []
6 mapping1 = {}
7 for k in K:
8     kmeanModel = KMeans(n_clusters=k).fit(dataset)
9     kmeanModel.fit(dataset)
10    distortions.append(sum(np.min(cdist(dataset, kmeanModel.cluster_centers_,
11                                       'euclidean'), axis=1)) / dataset.shape[0])
12    mapping1[k] = sum(np.min(cdist(dataset, kmeanModel.cluster_centers_, 'euclidean'), axis=1)) / dataset.shape[0]
13    inertias.append(kmeanModel.inertia_)
14    mapping2[k] = kmeanModel.inertia_

```



```

1 print("Vrednosti distrozije po broj klastera")
2 for key, val in mapping1.items():
3     print(f'{key} : {val}')

```

Vrednosti distrozije po broj klastera

```

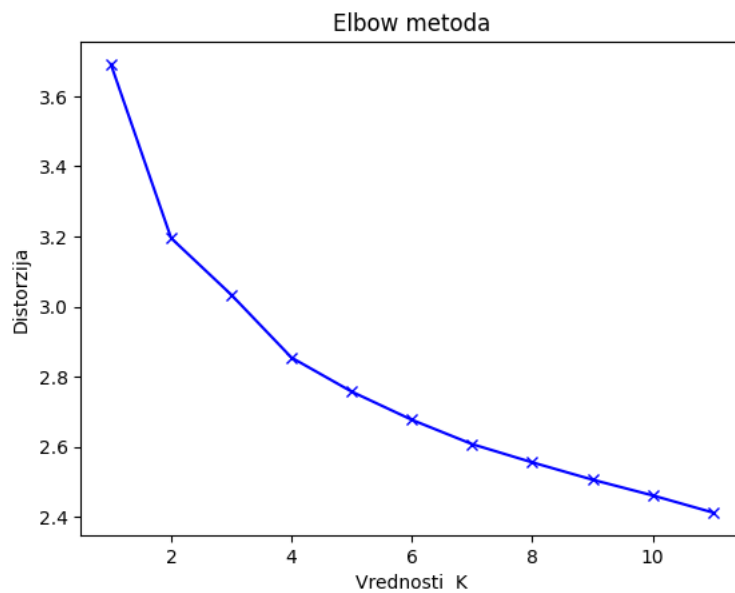
1 : 3.6913650287156656
2 : 3.195397659473886
3 : 3.0337668694656728
4 : 2.854440107706939
5 : 2.7573355324683564
6 : 2.6768714773642075
7 : 2.6076204759999184
8 : 2.5555242327853667
9 : 2.506167610990432
10 : 2.4612022069180175
11 : 2.4125094232084505

```

```

1 plt.plot(K, distortions, 'bx-')
2 plt.xlabel('Vrednosti K')
3 plt.ylabel('Distorzija')
4 plt.title('Elbow metoda')
5 plt.show()

```



```

1 print("Vrednosti inercije po broj klastera")
2 for key, val in mapping2.items():
3     print(f'{key} : {val}')

```

Vrednosti inercije po broj klastera

```

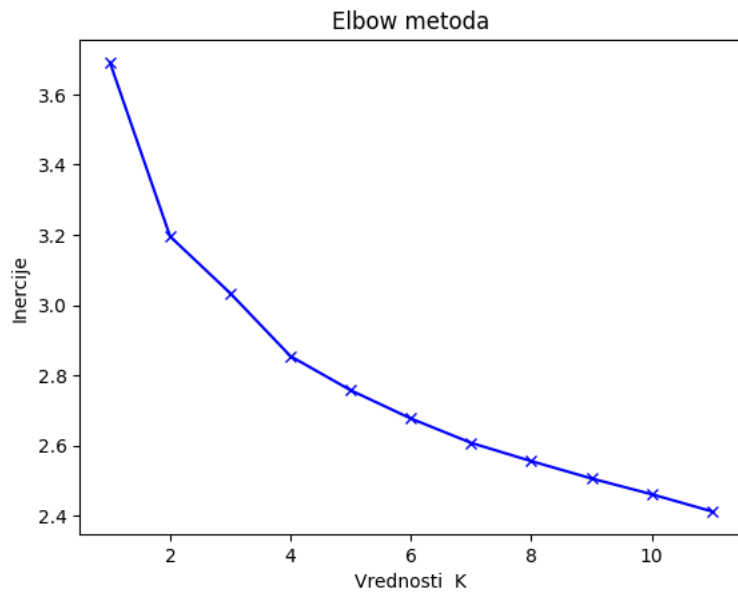
1 : 631425.3405406944
2 : 474669.1279857573
3 : 407115.17265450803
4 : 367195.09307091916
5 : 339662.1240076279
6 : 324003.9517890938
7 : 305634.40413226635
8 : 293676.15158675454
9 : 281434.2349200344
10 : 271443.97278337204
11 : 262839.81831405795

```

```

1 plt.plot(K, distortions, 'bx-')
2 plt.xlabel('Vrednosti K')
3 plt.ylabel('Inercije')
4 plt.title('Elbow metoda')
5 plt.show()

```



▼ Klasterovanje

Prvo se vrši klasterovanje neredukovanog seta podataka u 4 klastera

```
1 kmeans = KMeans(n_clusters =4)
2 cluster_labels = kmeans.fit_predict(dataset)
3 results = eval_metrics(dataset,cluster_labels)
4 print("Silhouette Score : "+str(results[0])+" Davies-Bouldin Index: "+str(results[1])+"Calinski-Harabasz Index: "+str(results[2]))
5 results_df = pd.DataFrame()
6
7
```

Silhouette Score : 0.17114653033315932 Davies-Bouldin Index: 1.65000937411862Calinski-Harabasz Index: 9862.011660070055

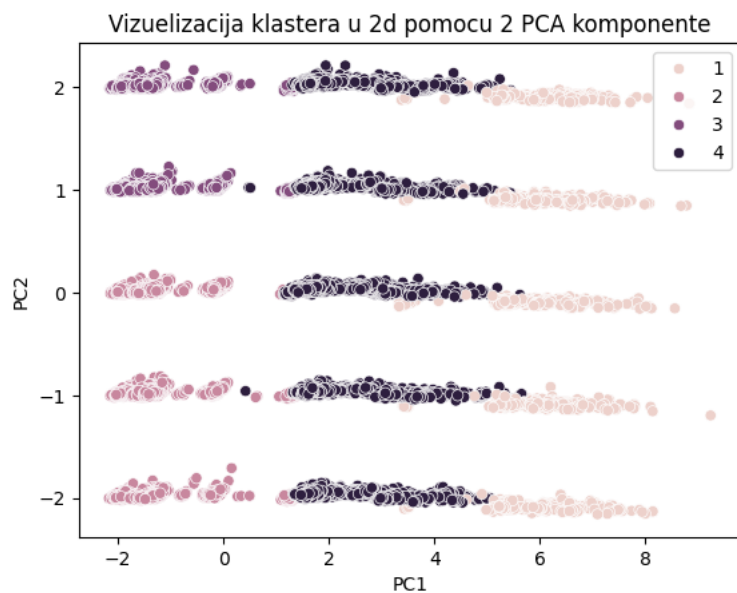
```
1 results_df.loc["KMeans","Silhouette Score"] = results[0]
```

```
1 visualise_3d_data(dataset,cluster_labels,4)
```

Vizuelizacija klastera u 3D pomocu 3 PCA komponente

- Cluster 0
- Cluster 1
- Cluster 2
- Cluster 3

```
1 visualise_2d_data(dataset, cluster_labels, 4)
```



```
1 unique_values, counts = np.unique(cluster_labels, return_counts=True)
```

```
1 clustering_dataframe = pd.DataFrame(counts)
```

```
1 clustering_dataframe.columns = ["KMeans"]
```

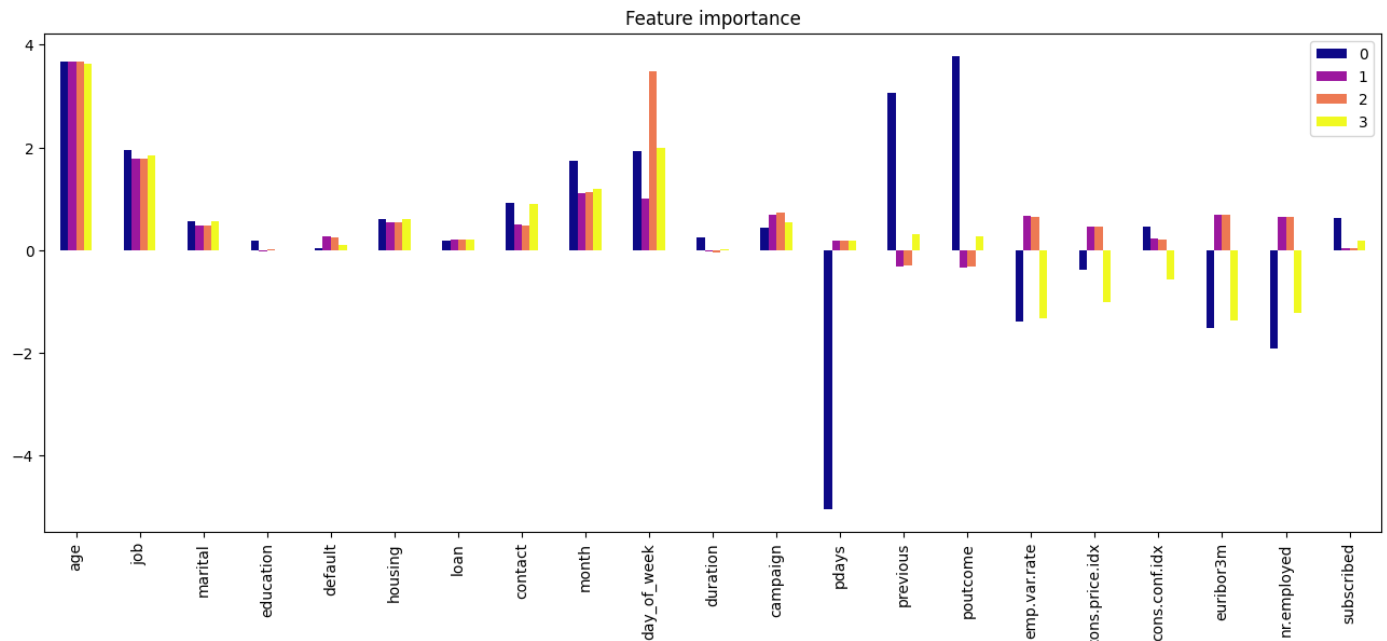
```
1 results_df
```

	Silhouette Score
KMeans	0.171147

```
1 clustering_dataframe
```

	KMeans
0	1537
1	16652
2	10845
3	12154

```
1 printFeatureImportance(dataset, cluster_labels)
```



Na osnovu uticaja svakog ficera za svaki klaster, odbacuju se kolone **education, default, loan, duration, subscribed**, i vrsi klasterovanje u 4 klastera

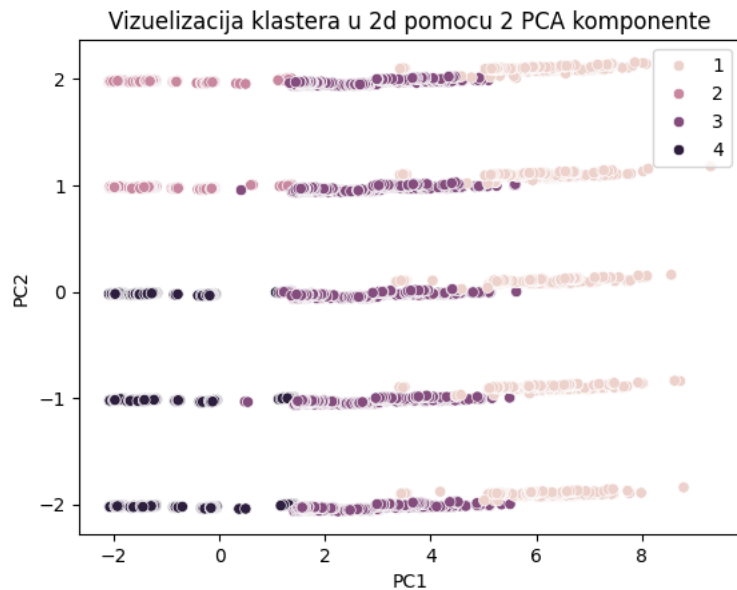
```
1 dataset_3 = dataset.copy()
2 delete_columns = ["education", "default", "loan", "duration", "subscribed"]
3 dataset_3.drop(delete_columns, axis=1, inplace=True)
4

1 kmeans = KMeans(n_clusters =4)
2 cluster_labels = kmeans.fit_predict(dataset_3)
3 results = eval_metrics(dataset_3, cluster_labels)
4 print("Silhouette Score : "+str(results[0])+" Davies-Bouldin Index: "+str(results[1])+" Calinski-Harabasz Index: "+str(results[2]))
5 results_df.loc["Najbitnije_kolone_KMeans", :] = results[0]
```

Silhouette Score : 0.21905577735050674 Davies-Bouldin Index: 1.3776859463120819 Calinski-Harabasz Index: 13567.142893603963

```
1 unique_values, counts = np.unique(cluster_labels, return_counts=True)
2 clustering_dataframe["Najbitnije_Kmeans"] = counts
```

```
1 pcas=PCA(n_components=2).fit_transform(dataset_3)
2 dataset_3_pca=pd.DataFrame(pcas, columns=['PC1', 'PC2'])
3 visualise_2d_data(dataset_3_pca, cluster_labels, 4)
```



Izvršena je PCA redukcija za redukovani set podataka, nakon čega se vrši klasterovanje podataka u 4 grupe.

```
1 kmeans = KMeans(n_clusters =4)
2 cluster_labels = kmeans.fit_predict(dataset_3_pca)
3 results = eval_metrics(dataset_3_pca,cluster_labels)
4 print("Silhouette Score : "+str(results[0])+" Davies-Bouldin Index: "+str(results[1])+"Calinski-Harabasz Index: "+str(results[2]))
5 results_df.loc["Najbitnije_PCA_KMeans",:] = results[0]
```

Silhouette Score : 0.5198752896108609 Davies-Bouldin Index: 0.8075172260055952Calinski-Harabasz Index: 53202.08871939428

```
1 unique_values, counts = np.unique(cluster_labels, return_counts=True)
2 clustering_dataframe["Najbitnije_PCA_KMeans"] = counts
```

Na kraju je istim algoritmom klasterovan polazni set podataka nakon PCA redukcije u 4 klastera.

```
1 kmeans = KMeans(n_clusters =4)
2 cluster_labels = kmeans.fit_predict(df2)
3 results = eval_metrics(df2,cluster_labels)
4 print("Silhouette Score : "+str(results[0])+" Davies-Bouldin Index: "+str(results[1])+"Calinski-Harabasz Index: "+str(results[2]))
5 results_df.loc["PCA_KMeans",:] = results[0]
```



Silhouette Score : 0.5172239943083106 Davies-Bouldin Index: 0.8238340653763321Calinski-Harabasz Index: 52865.5266060788


```
1 unique_values, counts = np.unique(cluster_labels, return_counts=True)
2 clustering_dataframe["PCA_KMeans"] = counts
```

```
1 results_df
```

	Silhouette Score	
KMeans	0.171147	
Najbitnije_kolone_KMeans	0.219056	
Najbitnije_PCA_KMeans	0.519875	
PCA_KMeans	0.517224	

```
1 clustering_dataframe
```

	KMeans	Najbitnije_Kmeans	Najbitnije_PCA_Kmeans	PCA_Kmeans	
0	1537	1541	10821	10822	
1	16652	11102	5783	6143	
2	10845	12146	16136	16136	
3	12154	16399	8448	8087	

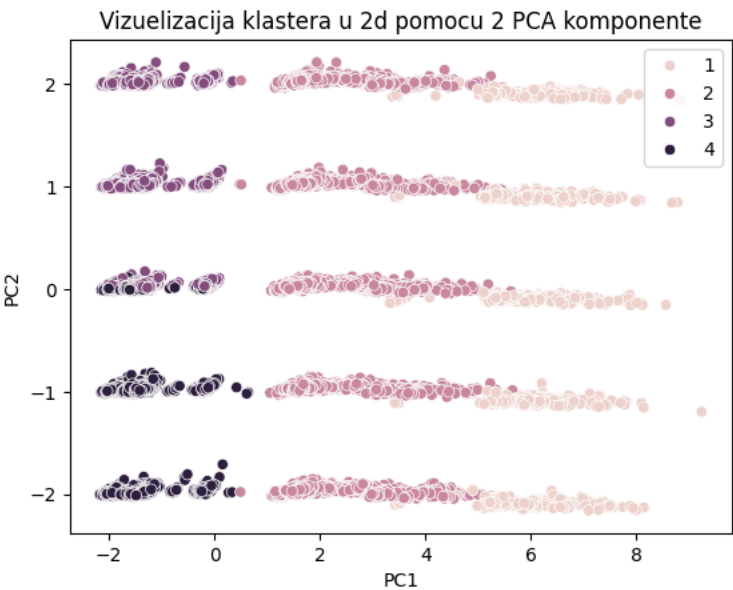
Next steps:  [View recommended plots](#)

▼ Bisecting K-means

```
1 bisectingkmeans = BisectingKMeans(n_clusters =4)
2 cluster_labels = bisectingkmeans.fit_predict(dataset)
3 results = eval_metrics(dataset,cluster_labels)
4 print("Silhouette Score : "+str(results[0])+" Davies-Bouldin Index: "+str(results[1])+"Calinski-Harabasz Index: "+str(results[2]))
5 results_df.loc["BisectingKMeans",:] = results[0]

Silhouette Score : 0.1635755840390671 Davies-Bouldin Index: 1.6770112397349854Calinski-Harabasz Index: 9796.98105156303

1 visualise_2d_data(dataset,cluster_labels,4)
2 visualise_3d_data(dataset,cluster_labels,4)
```



Vizuelizacija klastera u 3D pomocu 3 PCA komponente

- Cluster 0
- Cluster 1
- Cluster 2
- Cluster 3

```
1 unique_values, counts = np.unique(cluster_labels, return_counts=True)
2 clustering_dataframe["Bisecting Kmeans"] = counts
```

1 results_df

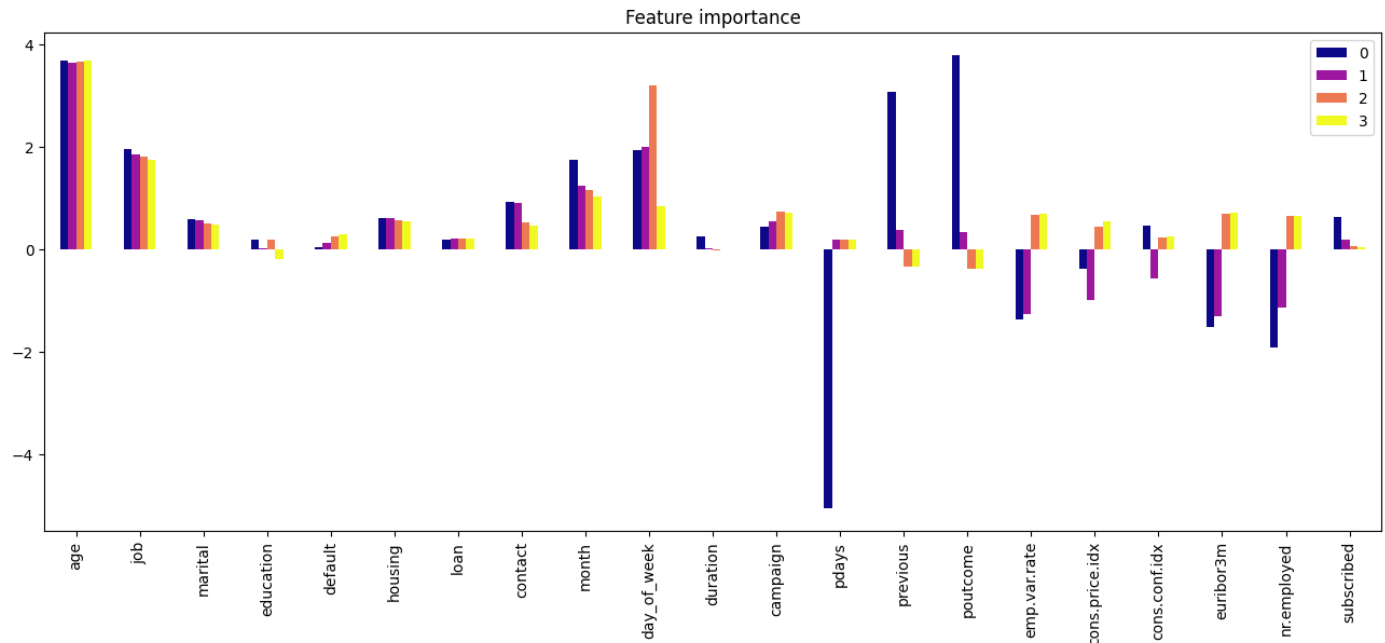
	Silhouette Score	
KMeans	0.171147	
Najbitnije_kolone_KMeans	0.219056	
Najbitnije_PCA_KMeans	0.519875	
PCA_KMeans	0.517224	
BisectingKMeans	0.163576	

1 clustering_dataframe

	KMeans	Najbitnije_Kmeans	Najbitnije_PCA_Kmeans	PCA_Kmeans	Bisecting Kmeans
0	1537	1541	10821	10822	1534
1	16652	11102	5783	6143	12700
2	10845	12146	16136	16136	13013
3	12154	16399	8448	8087	13941

Next steps: [View recommended plots](#)

```
1 printFeatureImportance(dataset, cluster_labels)
```

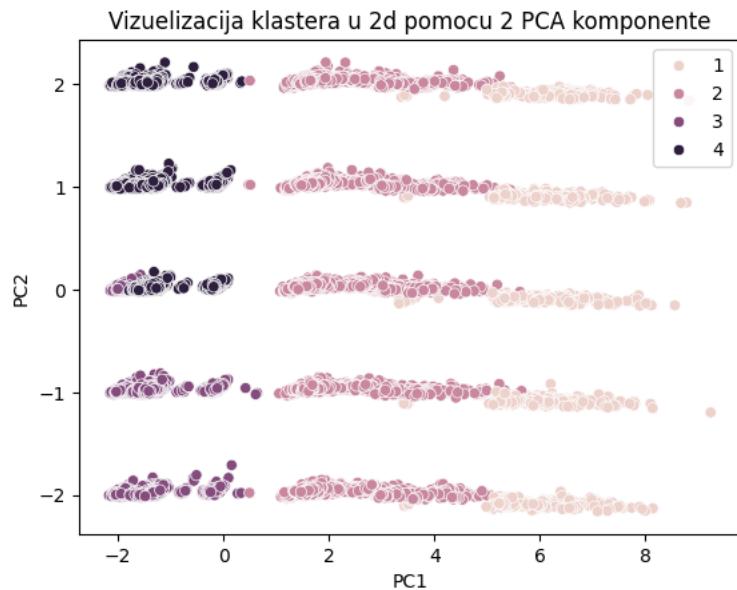


Na osnovu prikaza uticaja svakog ficera za svaki klaster izbacuju se kolone **education**, **default**, **loan** i **duration**, pa se vrši klasterovanje za dobijeni set podataka

```
1 dataset_4 = dataset.copy()
2 dataset_4.drop(["education","default","loan","duration"],axis =1,inplace = True)
3 bisectingkmeans = BisectingKMeans(n_clusters =4)
4 cluster_labels = bisectingkmeans.fit_predict(dataset_4)
5 results = eval_metrics(dataset_4,cluster_labels)
6 print("Silhouette Score : "+str(results[0])+" Davies-Bouldin Index: "+str(results[1])+"Calinski-Harabasz Index: "+str(results[2]))
7 results_df.loc["Najbitnije_BisectingKMeans",:] = results[0]
8 unique_values, counts = np.unique(cluster_labels, return_counts=True)
9 clustering_dataframe["Najbitnije_BisectingKMeans"] = counts
```

Silhouette Score : 0.20844716771499355 Davies-Bouldin Index: 1.4080440435805721Calinski-Harabasz Index: 13286.495522703544

```
1 visualise_2d_data(dataset,cluster_labels,4)
```

Izvršena je redukcija novodobijenog seta podataka na 2 glavne komponente PCA metodom.

```
1 pcas=PCA(n_components=2).fit_transform(dataset_4)
2 dataset_4_pca=pd.DataFrame(pcas,columns=['PC1','PC2'])
3 bisectingmeans = BisectingKMeans(n_clusters =4)
4 cluster_labels = bisectingmeans.fit_predict(dataset_4_pca)
5 results = eval_metrics(dataset_4_pca,cluster_labels)
6 print("Silhouette Score : "+str(results[0])+" Davies-Bouldin Index: "+str(results[1])+"Calinski-Harabasz Index: "+str(results[2]))
7 results_df.loc["Najbitnije_PCA_BisectingKMeans",:] = results[0]
8 unique_values, counts = np.unique(cluster_labels, return_counts=True)
9 clustering_dataframe["Najbitnije_PCA_BisectingKMeans"] = counts
```

Silhouette Score : 0.5158424202482829 Davies-Bouldin Index: 0.8319366775858243Calinski-Harabasz Index: 52415.54121833599

Izvršeno je klasterovanje nakon originalong redukcije primenom PCA metode na 2 komponente.

```
1 bisectingmeans = BisectingKMeans(n_clusters =4)
2 cluster_labels = bisectingmeans.fit_predict(df2)
3 results = eval_metrics(df2,cluster_labels)
4 print("Silhouette Score : "+str(results[0])+" Davies-Bouldin Index: "+str(results[1])+"Calinski-Harabasz Index: "+str(results[2]))
5 results_df.loc["PCA_BisectingKMeans",:] = results[0]
```

Silhouette Score : 0.512936777229905 Davies-Bouldin Index: 0.8340137361907445Calinski-Harabasz Index: 52175.537242845814

```
1 unique_values, counts = np.unique(cluster_labels, return_counts=True)
2 clustering_dataframe["PCA_BisectingKMeans"] = counts
```

```
1 results_df
```

	Silhouette Score	
KMeans	0.171147	
Najbitnije_kolone_KMeans	0.219056	
Najbitnije_PCA_KMeans	0.519875	
PCA_KMeans	0.517224	
BisectingKMeans	0.163576	
Najbitnije_BisectingKMeans	0.208447	
Najbitnije_PCA_BisectingKMeans	0.515842	
PCA_BisectingKMeans	0.512937	

▼ Gaussian Mixture

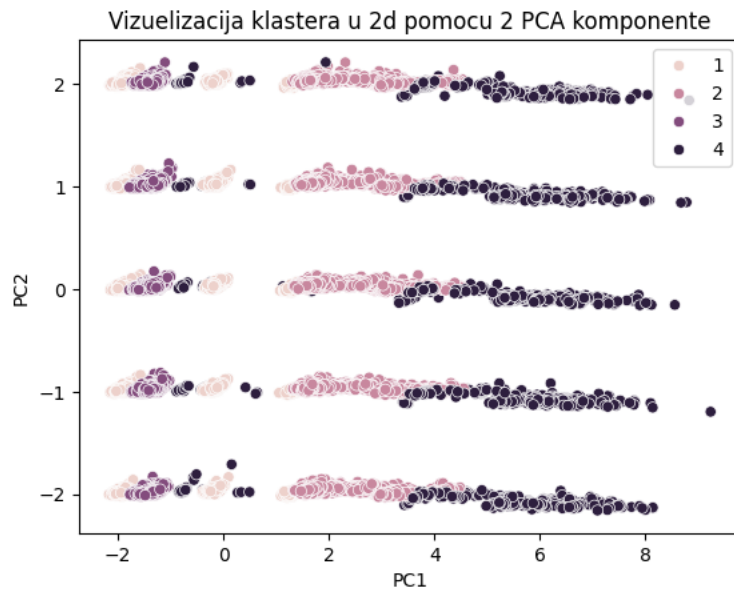
```
1 from sklearn.mixture import GaussianMixture
2 gmm = GaussianMixture(n_components = 4, random_state=42)
```

```
1 cluster_labels = gmm.fit_predict(dataset)
2 results = eval_metrics(dataset, cluster_labels)
3 print("Silhouette Score : "+str(results[0])+" Davies-Bouldin Index: "+str(results[1])+" Calinski-Harabasz Index: "+str(results[2]))
4 results_df.loc["Gaussian Mixture",:] = results[0]
```

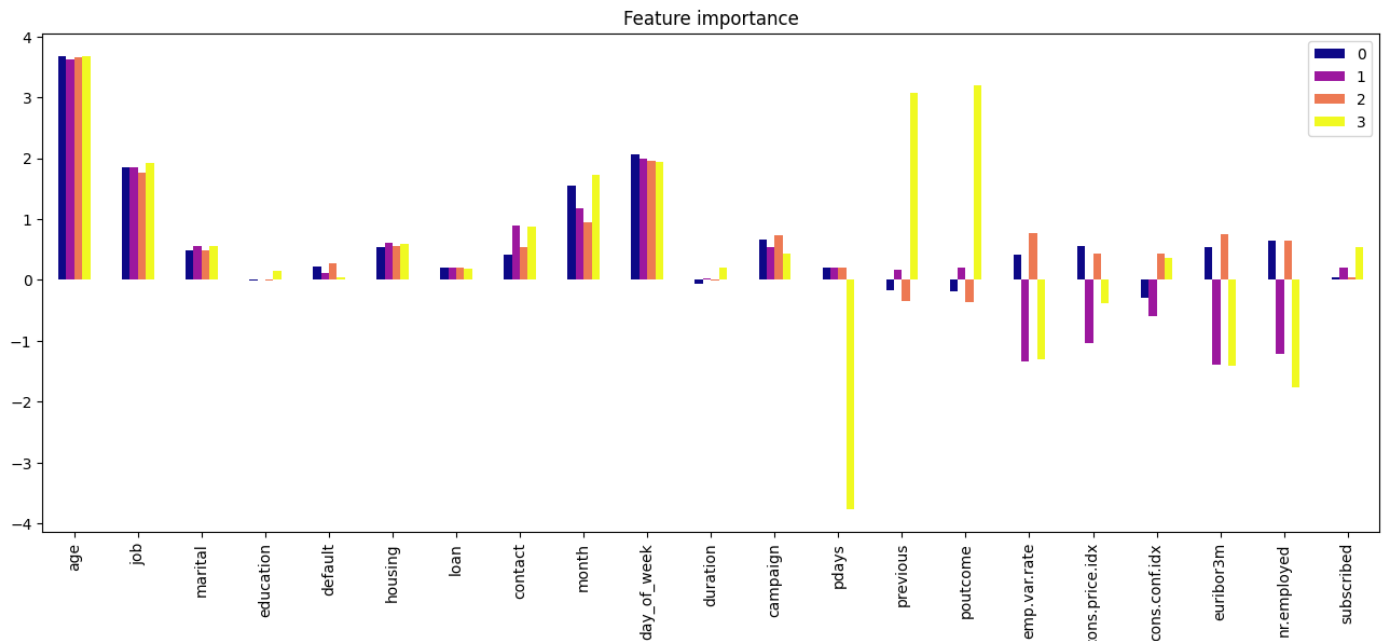
Silhouette Score : 0.09975477726149581 Davies-Bouldin Index: 3.3015872624880185 Calinski-Harabasz Index: 7460.016926791017

```
1 unique_values, counts = np.unique(cluster_labels, return_counts=True)
2 clustering_dataframe["Gaussian Mixture"] = counts
```

```
1 visualise_2d_data(dataset, cluster_labels, 4)
```



```
1 printFeatureImportance(dataset, cluster_labels)
```



Na osnovu prikazanog uticaja ficera na klastere izbacuju se **education, default, loan, duration**. Vrsi se klasterovanje novodobijenog dataseta.

```
1 dataset_5 = dataset.copy()
2 dataset_5.drop(["education","default","loan","duration"],axis =1, inplace = True)
3 cluster_labels = gmm.fit_predict(dataset_5)
4 results = eval_metrics(dataset_5,cluster_labels)
5 print("Silhouette Score : "+str(results[0])+" Davies-Bouldin Index: "+str(results[1])+"Calinski-Harabasz Index: "+str(results[2]))
6 results_df.loc["Najbitnije_Gaussian_Mixture",:] = results[0]
```

Silhouette Score : 0.12713395035214892 Davies-Bouldin Index: 2.329832381913431Calinski-Harabasz Index: 10306.446779961354

```
1 unique_values, counts = np.unique(cluster_labels, return_counts=True)
2 clustering_dataframe["Najbitnije_Gaussian_Mixture"] = counts
```

Nakon toga se vrsi redukcija dimenzionalnosti PCA metodom na 2 komponente i klasterizacija

```
1 pcas=PCA(n_components=2).fit_transform(dataset_5)
2 dataset_5_pca=pd.DataFrame(pcas,columns=['PC1','PC2'])
3 cluster_labels = gmm.fit_predict(dataset_5_pca)
4 results = eval_metrics(dataset_5_pca,cluster_labels)
5 print("Silhouette Score : "+str(results[0])+" Davies-Bouldin Index: "+str(results[1])+"Calinski-Harabasz Index: "+str(results[2]))
6 results_df.loc["Najbitnije_PCA_Gaussian_Mixture",:] = results[0]
7 unique_values, counts = np.unique(cluster_labels, return_counts=True)
8 clustering_dataframe["Najbitnije_PCA_Gaussian_Mixture"] = counts
```

Silhouette Score : 0.46328132173778425 Davies-Bouldin Index: 0.6972159159809895Calinski-Harabasz Index: 38306.47305693685



Na kraju se izvrsava klasterizacija ovom metodom za polazni dataset koji je redukovan na dve glavne komponente PCA metodom

```
1 cluster_labels = gmm.fit_predict(df2)
2 results = eval_metrics(df2,cluster_labels)
3 print("Silhouette Score : "+str(results[0])+" Davies-Bouldin Index: "+str(results[1])+"Calinski-Harabasz Index: "+str(results[2]))
4 results_df.loc["PCA_Gaussian_Mixture",:] = results[0]
```

Silhouette Score : 0.46546688598306035 Davies-Bouldin Index: 0.8838661416398568Calinski-Harabasz Index: 38707.137751568815

```
1 unique_values, counts = np.unique(cluster_labels, return_counts=True)
2 clustering_dataframe["PCA_Gaussian Mixture"] = counts
```

```
1 results_df
```

	Silhouette Score	
KMeans	0.171147	
Najbitnije_kolone_KMeans	0.219056	
Najbitnije_PCA_KMeans	0.519875	
PCA_KMeans	0.517224	
BisectingKMeans	0.163576	
Najbitnije_BisectingKMeans	0.208447	
Najbitnije_PCA_BisectingKMeans	0.515842	
PCA_BisectingKMeans	0.512937	
Gaussian Mixture	0.099755	
Najbitnije_Gaussian_Mixture	0.127134	
Najbitnije_PCA_Gaussian_Mixture	0.463281	
PCA_Gaussian Mixture	0.465467	

```
1 clustering_dataframe
```

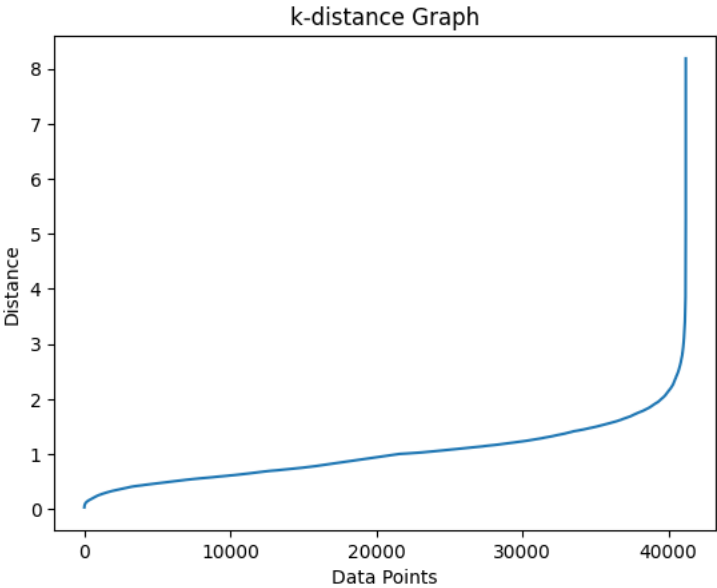
	KMeans	Najbitnije_Kmeans	Najbitnije_PCA_Kmeans	PCA_Kmeans	Bisecting Kmeans	Najbitnije_BisectingKMeans	Najbitnije_PCA_BisectingKMeans
0	1537	1541	10821	10822	1534	1530	6246
1	16652	11102	5783	6143	12700	12704	7986
2	10845	12146	16136	16136	13013	13326	10820
3	12154	16399	8448	8087	13941	13628	16136

Next steps:

 [View recommended plots](#)

▽ DBSCAN

```
1 from sklearn.neighbors import NearestNeighbors
2 neigh = NearestNeighbors(n_neighbors=4)
3 distances, indices = neigh.fit(dataset).kneighbors(dataset)
4 distances = np.sort(distances[:, -1])
5 plt.plot(distances)
6 plt.xlabel('Data Points')
7 plt.ylabel('Distance')
8 plt.title('k-distance Graph')
9 plt.show()
```



```
1 dbscan = DBSCAN(eps = 2.75, n_jobs=-1)
2 cluster_labels = dbscan.fit_predict(dataset)
3 max(cluster_labels)

3
```

```
1 results = eval_metrics(dataset,cluster_labels)
2 print("Silhouette Score : "+str(results[0])+" Davies-Bouldin Index: "+str(results[1])+"Calinski-Harabasz Index: "+str(results[2]))
3 results_df.loc["DBSCAN",:] = results[0]

Silhouette Score : 0.2389458490253637 Davies-Bouldin Index: 1.968267519982127Calinski-Harabasz Index: 3303.7085306796844
```

```
1 unique_values, counts = np.unique(cluster_labels, return_counts=True)
2 unique_values

array([-1,  0,  1,  2,  3])
```

```
1 clustering_dataframe["DBSCAN"] = counts[1:5]
```

```
1 results_df
```

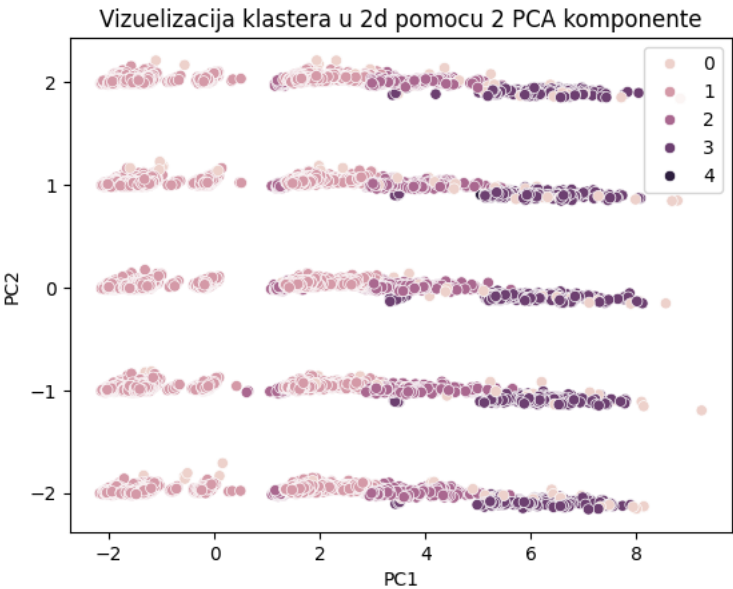
	Silhouette Score	
KMeans	0.171147	
Najbitnije_kolone_KMeans	0.219056	
Najbitnije_PCA_KMeans	0.519875	
PCA_KMeans	0.517224	
BisectingKMeans	0.163576	
Najbitnije_BisectingKMeans	0.208447	
Najbitnije_PCA_BisectingKMeans	0.515842	
PCA_BisectingKMeans	0.512937	
Gaussian Mixture	0.099755	
Najbitnije_Gaussian_Mixture	0.127134	
Najbitnije_PCA_Gaussian_Mixture	0.463281	
PCA_Gaussian Mixture	0.465467	
DBSCAN	0.238946	

```
1 clustering_dataframe
```

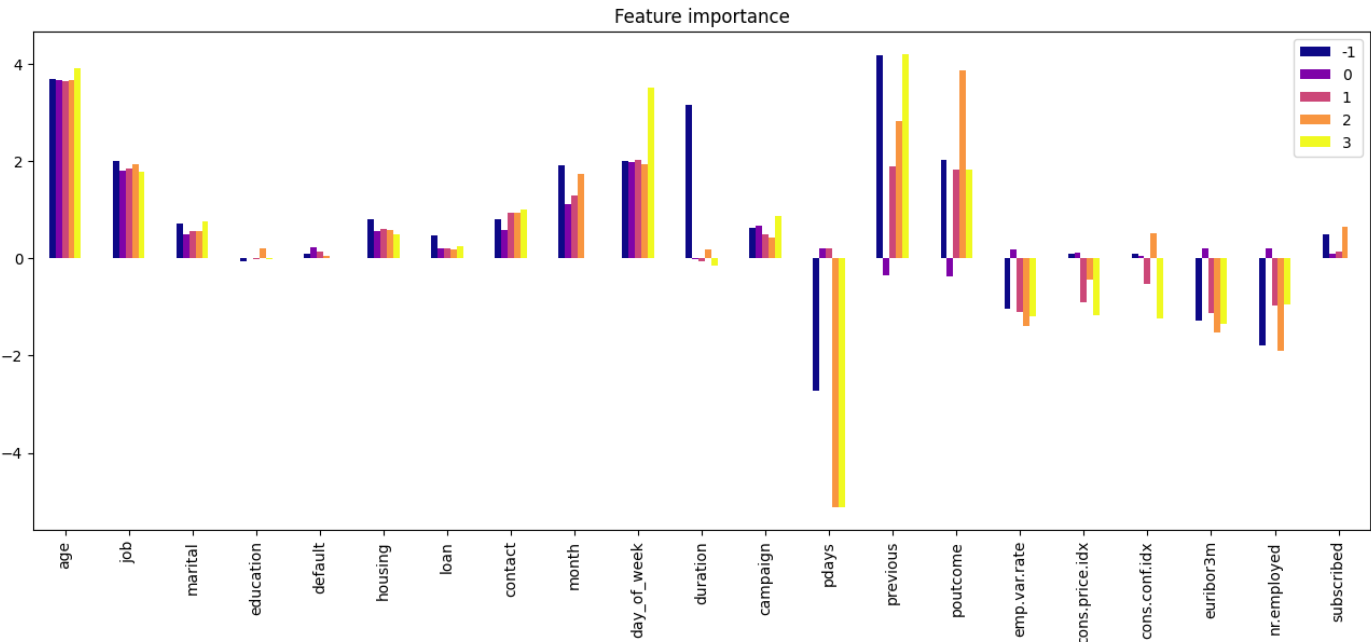
	KMeans	Najbitnije_Kmeans	Najbitnije_PCA_Kmeans	PCA_Kmeans	Bisecting Kmeans	Najbitnije_BisectingKMeans	Najbitnije_PCA_BisectingKMeans
0	1537	1541	10821	10822	1534	1530	6246
1	16652	11102	5783	6143	12700	12704	7986
2	10845	12146	16136	16136	13013	13326	10820
3	12154	16399	8448	8087	13941	13628	16136

Next steps: [View recommended plots](#)

```
1 visualise_2d_data(dataset,cluster_labels,4)
```



```
1 printFeatureImportance(dataset, cluster_labels)
```



```

1 dataset_6 = dataset.copy()
2 dataset_6.drop(["education","default"],axis =1, inplace = True)
3 dbscan = DBSCAN(eps = 2.75, n_jobs=-1)
4 cluster_labels = dbscan.fit_predict(dataset_6)
5 results = eval_metrics(dataset_6,cluster_labels)
6 print("Silhouette Score : "+str(results[0])+" Davies-Bouldin Index: "+str(results[1])+"Calinski-Harabasz Index: "+str(results[2]))
7 results_df.loc["Najbitnije_DBSCAN",:] = results[0]

```

Silhouette Score : 0.26144058353883676 Davies-Bouldin Index: 1.719416723737258Calinski-Harabasz Index: 3683.1650027445903

```

1 unique_values, counts = np.unique(cluster_labels, return_counts=True)
2 unique_values

array([-1,  0,  1,  2,  3])

```

```
1 clustering_dataframe["DBSCAN"] = counts[1:5]
```

```
1 results_df
```

	Silhouette Score	
KMeans	0.171147	
Najbitnije_kolone_KMeans	0.219056	
Najbitnije_PCA_KMeans	0.519875	
PCA_KMeans	0.517224	
BisectingKMeans	0.163576	
Najbitnije_BisectingKMeans	0.208447	
Najbitnije_PCA_BisectingKMeans	0.515842	
PCA_BisectingKMeans	0.512937	
Gaussian Mixture	0.099755	
Najbitnije_Gaussian_Mixture	0.127134	
Najbitnije_PCA_Gaussian_Mixture	0.463281	
PCA_Gaussian Mixture	0.465467	
DBSCAN	0.238946	
Najbitnije_DBSCAN	0.261441	

Na osnovu dobijenih tabela zaključuje se da primena KMeans algoritma za klasterovanje u 4 klastera, daje najbolje rezultate, sa malom razlikom u odnosu na Bisecting KMeans i Gaussian Mixture, sve nad podacima redukovane dimenzionalnosti. KMeans i Gaussian Mixture bice testirani i na drugoj varijanti dataseta.

▼ Druga varijanta

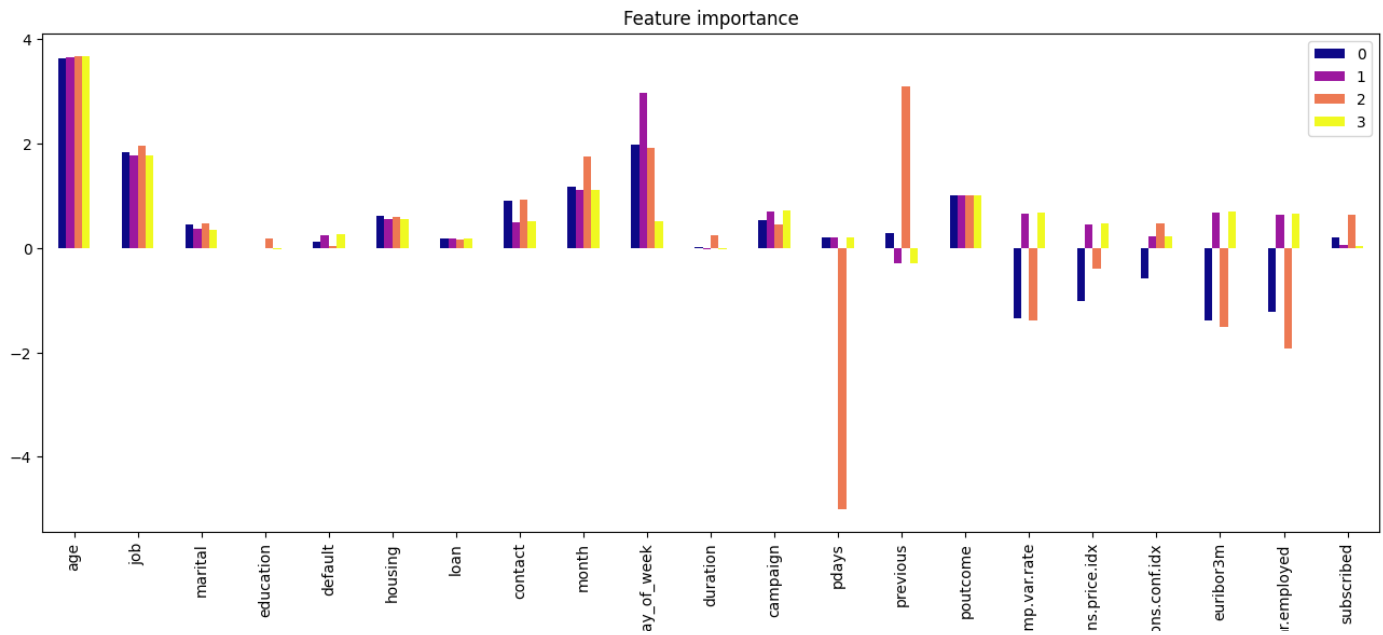
```

1 kmeans = KMeans(n_clusters =4)
2 cluster_labels = kmeans.fit_predict(dataset_2)
3 results = eval_metrics(dataset_2,cluster_labels)
4 print("Silhouette Score : "+str(results[0])+" Davies-Bouldin Index: "+str(results[1])+"Calinski-Harabasz Index: "+str(results[2]))
5 results_df_2 = pd.DataFrame()
6 results_df_2.loc["KMeans","Silhouette Score"] = results[0]

```

Silhouette Score : 0.17896779164819634 Davies-Bouldin Index: 1.6419446085231888Calinski-Harabasz Index: 9668.543718253422

```
1 printFeatureImportance(dataset_2, cluster_labels)
```



1 results_df_2

Silhouette Score	
KMeans	0.178968

```

1 dataset_7 = dataset_2.copy()
2 dataset_7.drop(["education", "default", "loan", "duration"], axis=1, inplace=True)
3 cluster_labels = kmeans.fit_predict(dataset_7)
4 results = eval_metrics(dataset_7, cluster_labels)
5 print("Silhouette Score : "+str(results[0])+" Davies-Bouldin Index: "+str(results[1])+" Calinski-Harabasz Index: "+str(results[2]))
6
7 results_df_2.loc["Najbitniji_KMeans", "Silhouette Score"] = results[0]

```

Silhouette Score : 0.22921389181397597 Davies-Bouldin Index: 1.3724886463730646 Calinski-Harabasz Index: 13398.876434245802

Nakon toga odredjuju se dva glavne komponente sa PCA metodom i vrši klasterovanje

```

1 pcas=PCA(n_components=2).fit_transform(dataset_7)
2 dataset_7_pca = pd.DataFrame(pcas, columns=['PC1', 'PC2'])
3 cluster_labels = kmeans.fit_predict(dataset_7_pca)

```