

## ▼ 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_v
0	56	housemaid	married	basic.4y	no	no	no	telephone	may	
1	57	services	married	high.school	unknown	no	no	telephone	may	
2	37	services	married	high.school	no	yes	no	telephone	may	
3	40	admin.	married	basic.6y	no	no	no	telephone	may	
4	56	services	married	high.school	no	no	yes	telephone	may	

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"] < 0).sum()
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.api.types.is_numeric_dtype(dataset["previous"].dtype)
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, 'retiree': 5, 'married': 0, 'single': 1, 'divorced': 2, 'unknown': 3}
{'basic.4y': 0, 'high.school': 1, 'basic.6y': 2, 'basic.9y': 3, 'professional.course': 4}
{'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}
{'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
0	56	0	0	0	0	0	0	0	0	0
1	57	1	0	1	1	0	0	0	0	0
2	37	1	0	1	0	1	0	0	0	0
3	40	2	0	2	0	0	0	0	0	0
4	56	1	0	1	0	0	1	0	0	0
...	...	...	...	...	...	...	...	...	...	...
41183	73	5	0	4	0	1	0	1	5	4
41184	46	3	0	4	0	0	0	1	5	4
41185	56	5	0	6	0	1	0	1	5	4
41186	44	4	0	4	0	0	0	1	5	4
41187	74	5	0	4	0	1	0	1	5	4

41188 rows × 21 columns

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

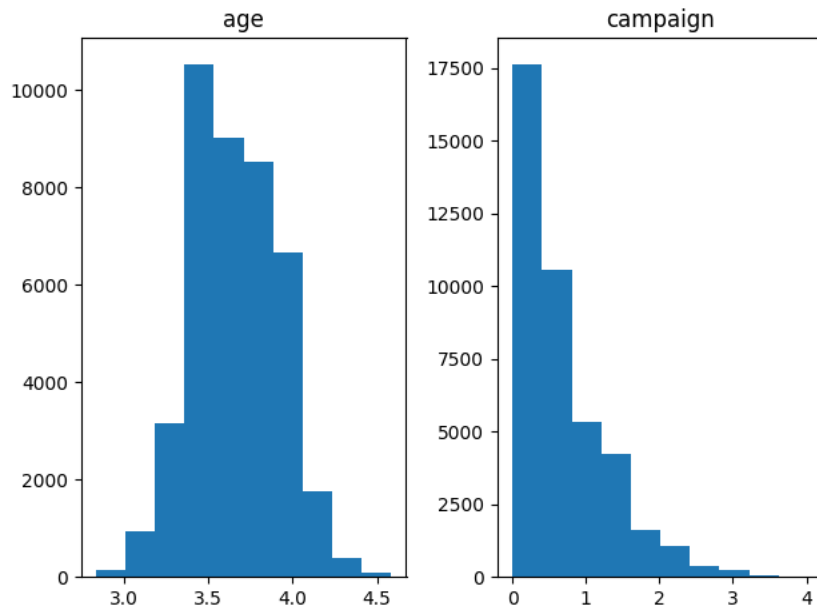
```
1 dataset = dataset_bckup.copy()
```

Iskoseni podaci normalizuju se primenom logaritma

```

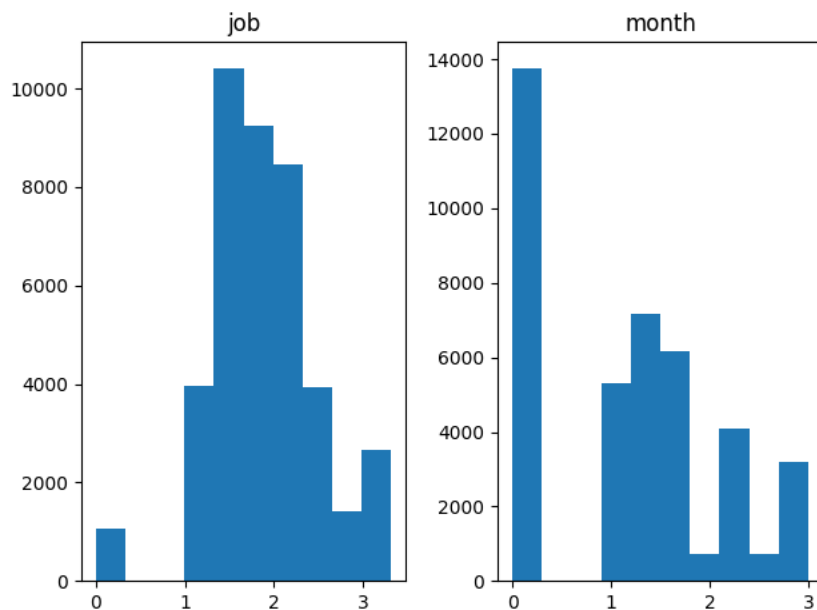
1 skewed_columns = ["age","campaign"]
2 fig, axes = plt.subplots(1,2)
3
4 for ind,col in enumerate(skewed_columns):
5     dataset[col] = np.log(dataset[col])
6
7     axes[ind].hist(dataset[col])
8     axes[ind].set_title(col)
9     axes[ind].set_xlabel('')
10    axes[ind].set_ylabel('')
11
12 plt.tight_layout()
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", "euribc
3 # normalize_columns = ["education", "duration"]
4 normalized_data = scaler.fit_transform(dataset[normalize_columns])
```

```

4 normalized_data= scaler.fit_transform(dataset[normalize_columns])
5 dataset[normalize_columns] = normalized_data
6

```

```

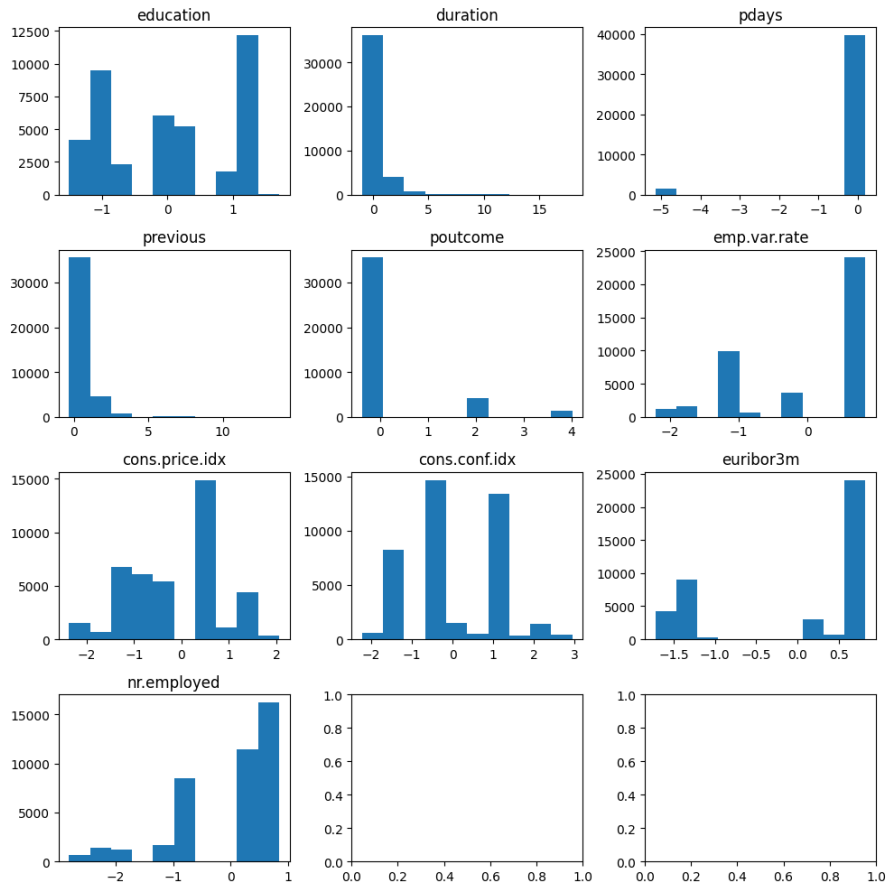
1

```

```

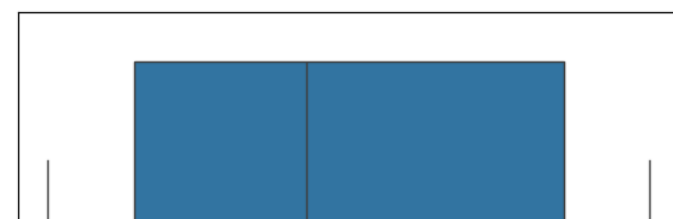
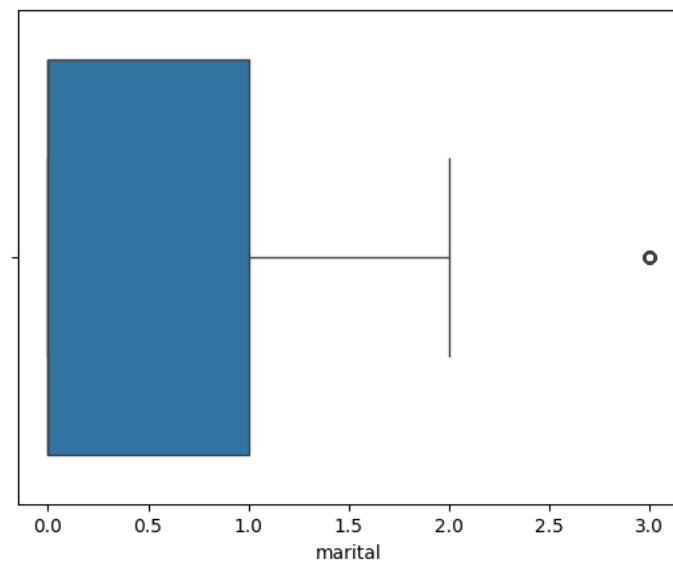
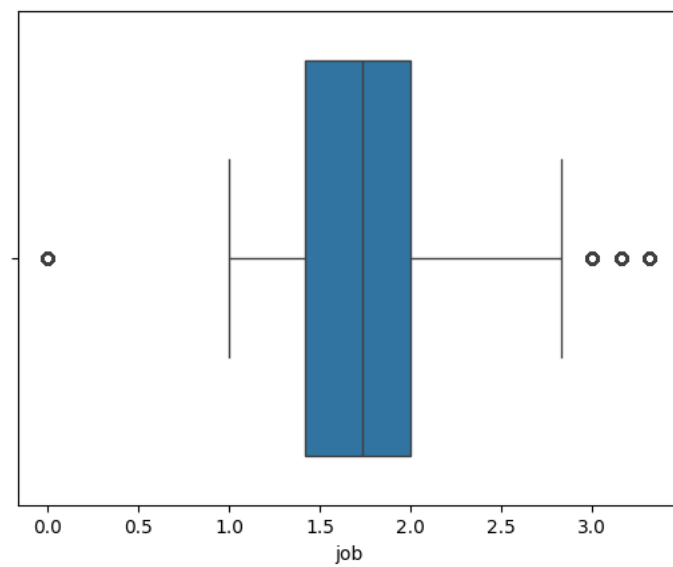
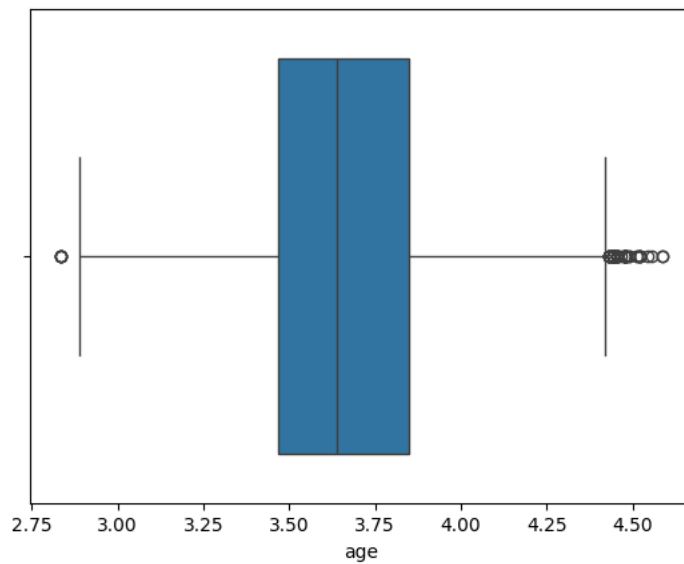
1 fig, axes = plt.subplots(4,3,figsize=(10,10))
2 axes = axes.flatten()
3 for ind,col in enumerate(normalize_columns):
4     axes[ind].hist(dataset[col])
5     axes[ind].set_title(col)
6     axes[ind].set_xlabel('')
7     axes[ind].set_ylabel('')
8
9 plt.tight_layout()
10 plt.show()
11

```

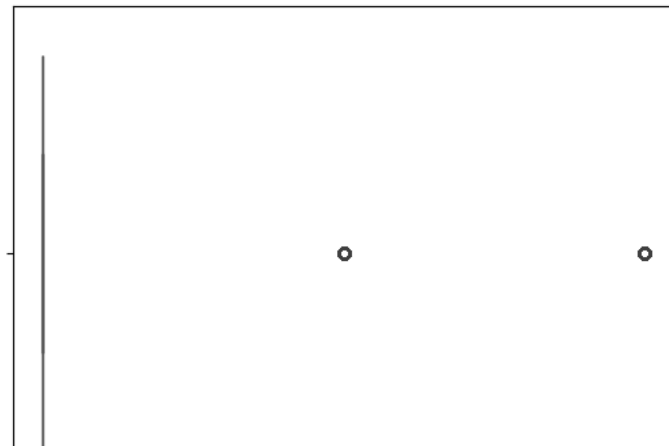
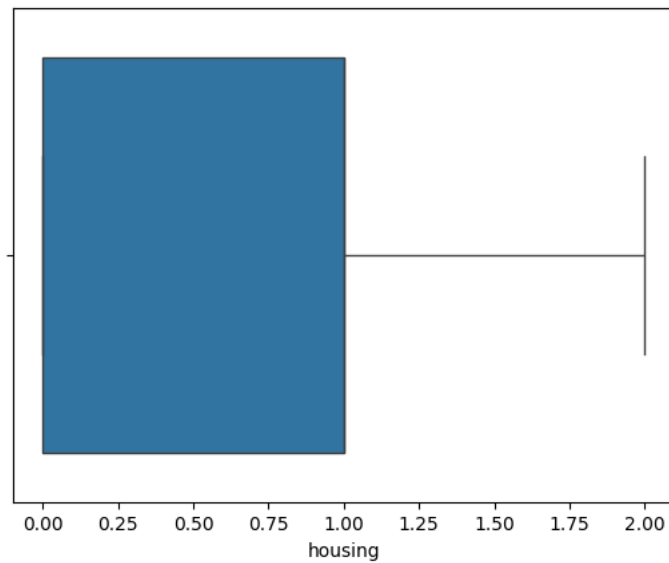
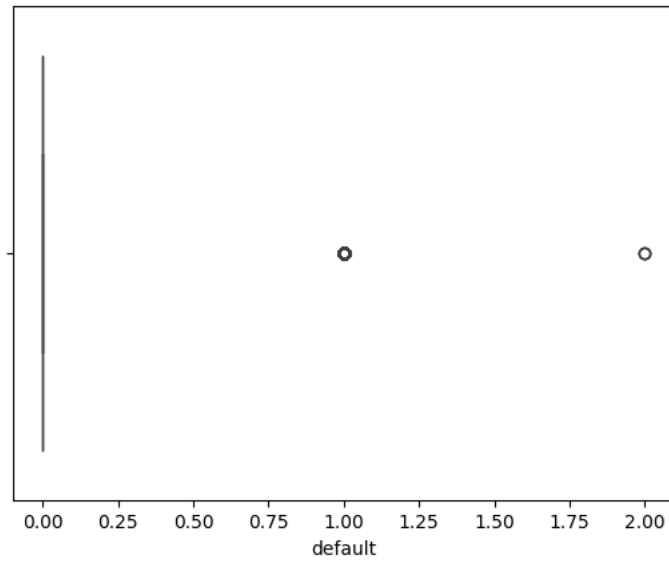
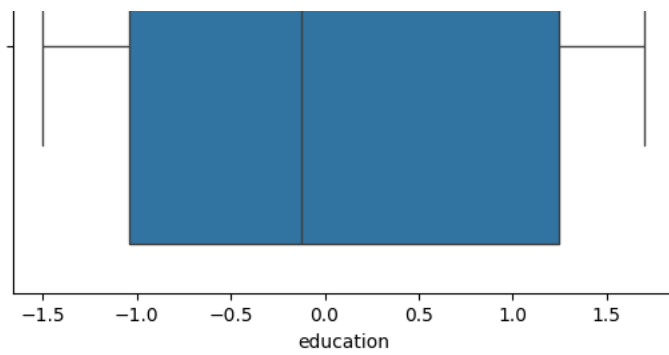


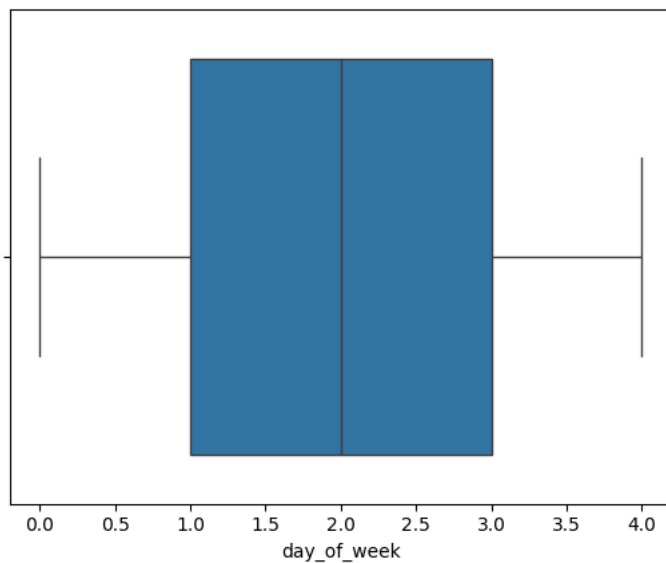
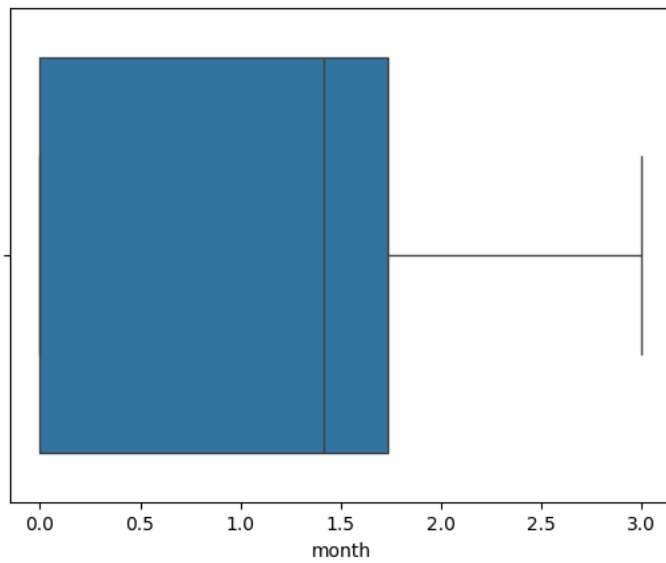
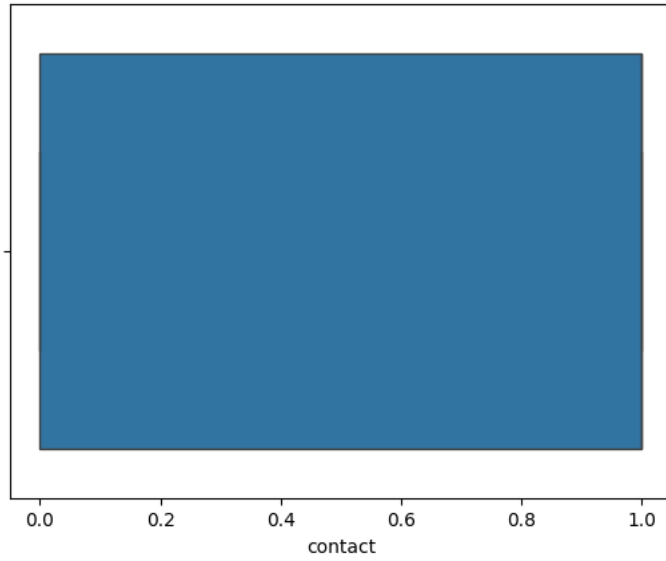
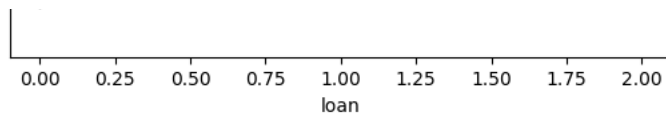
✓ Detekcija i otklanjanje outlier-a, normalizacija podataka

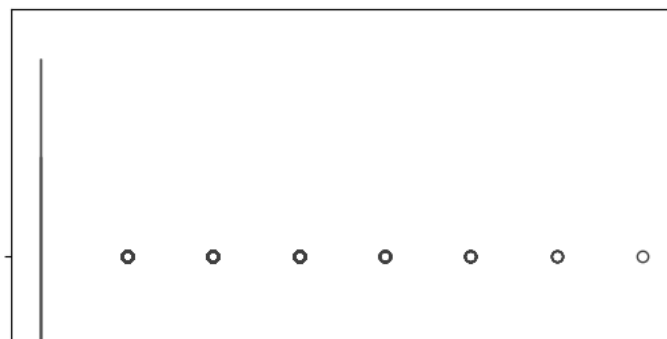
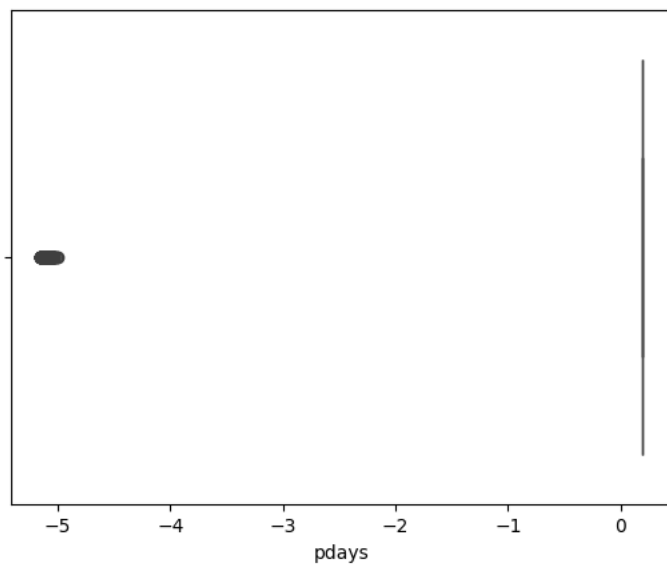
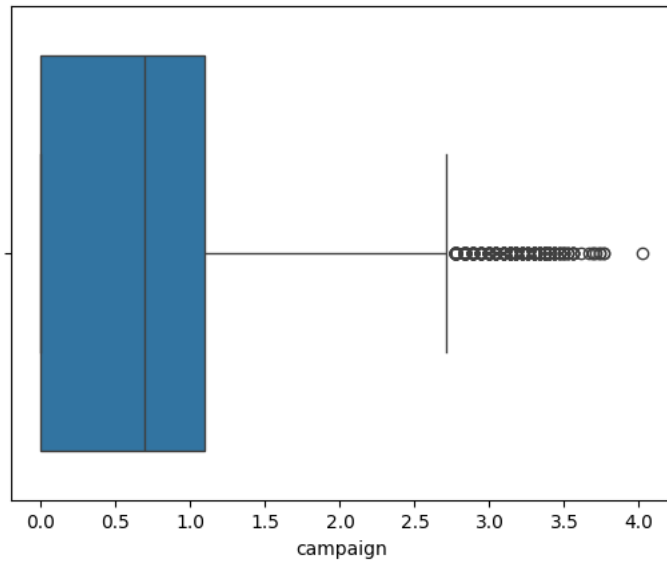
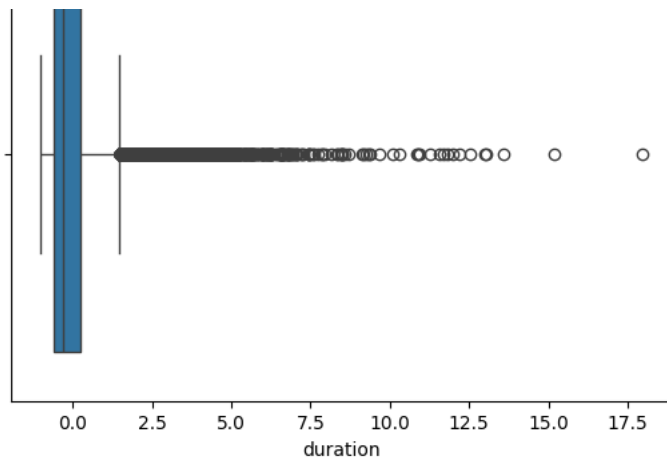
```
1 lista = ["age", "job", "marital", "education", "default", "housing", "loan", "contact", "month", "day_of_week", "duration", "campaign", "  
2 for el in lista:  
3     sns.boxplot(x = dataset[el])  
4     plt.show()
```

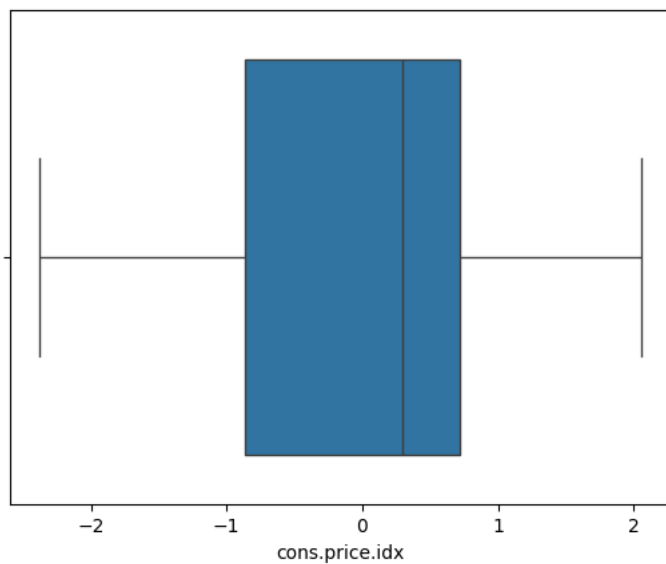
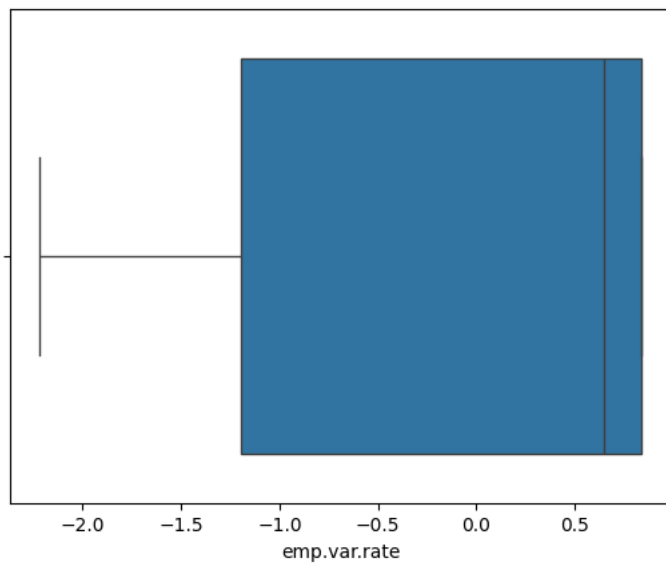
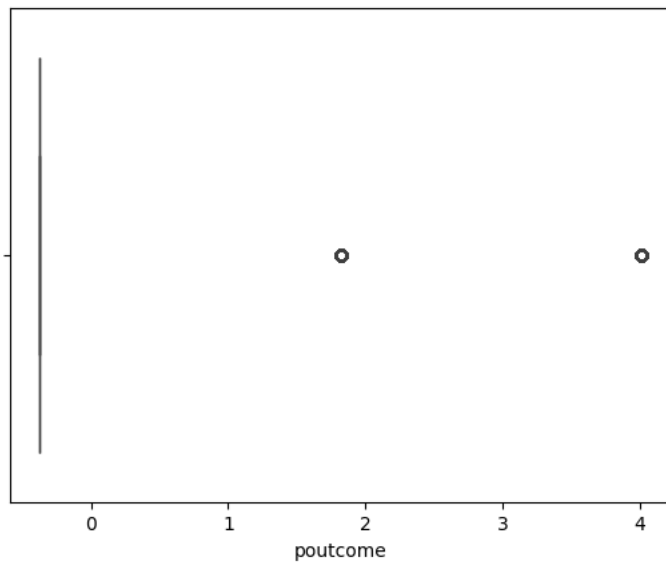
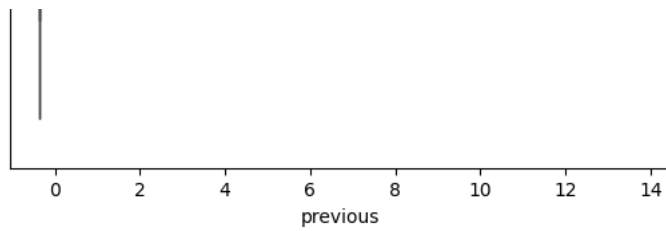


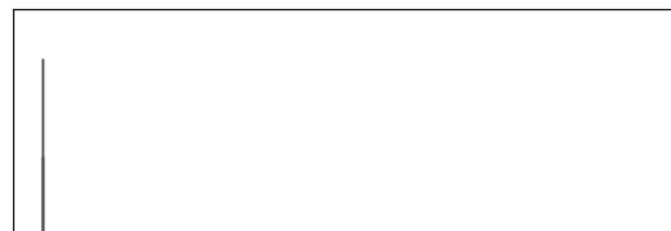
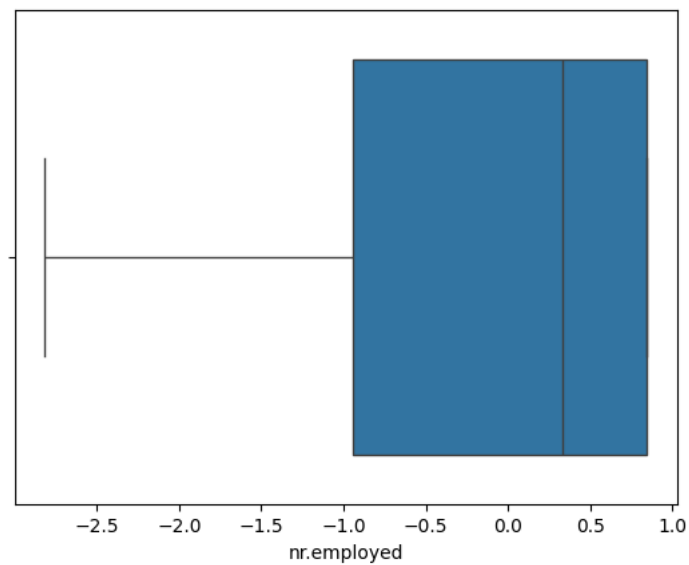
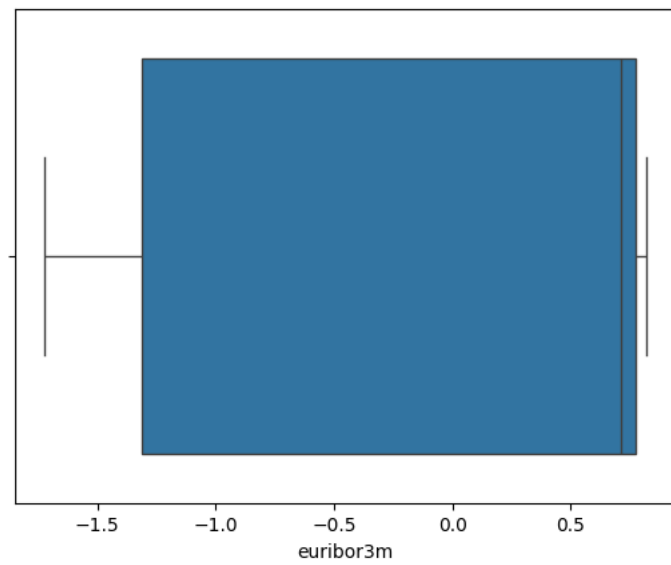
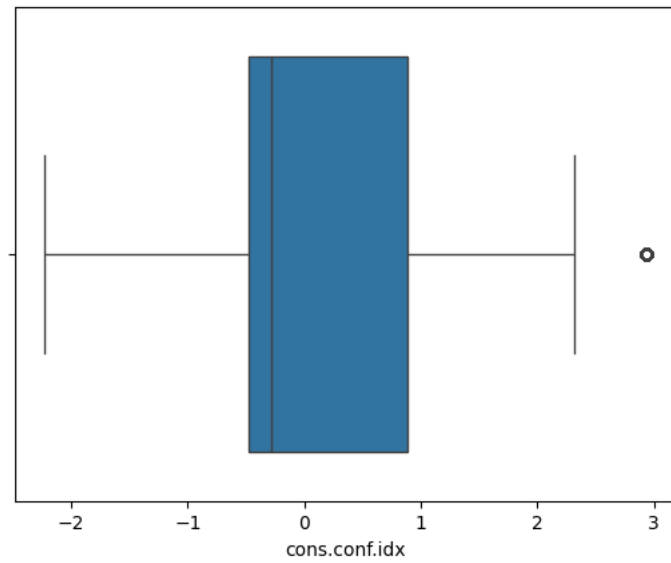


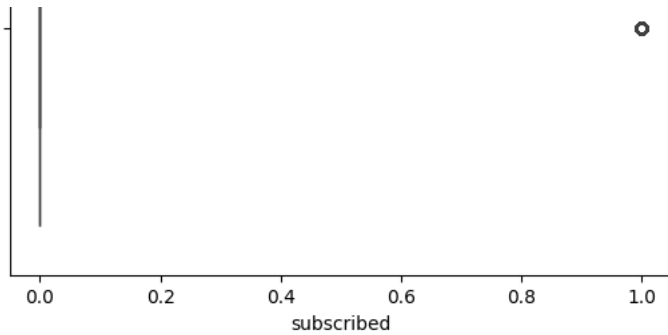












## ▼ Druga varijanta dataseta

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
0	4.025352	0.000000	0	-1.499673	0	0	0	0	0.000000
1	4.043051	1.000000	0	-1.042111	1	0	0	0	0.000000
2	3.610918	1.000000	0	-1.042111	0	1	0	0	0.000000
3	3.688879	1.414214	0	-0.584550	0	0	0	0	0.000000
4	4.025352	1.000000	0	-1.042111	0	0	1	0	0.000000
...	...	...	...	...	...	...	...	...	...
41183	4.290459	2.236068	0	0.330573	0	1	0	1	2.236068
41184	3.828641	1.732051	0	0.330573	0	0	0	1	2.236068
41185	4.025352	2.236068	0	1.245696	0	1	0	1	2.236068
41186	3.784190	2.000000	0	0.330573	0	0	0	1	2.236068
41187	4.304065	2.236068	0	0.330573	0	1	0	1	2.236068

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)
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'])

```

## ▼ Kmeans

```

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

```

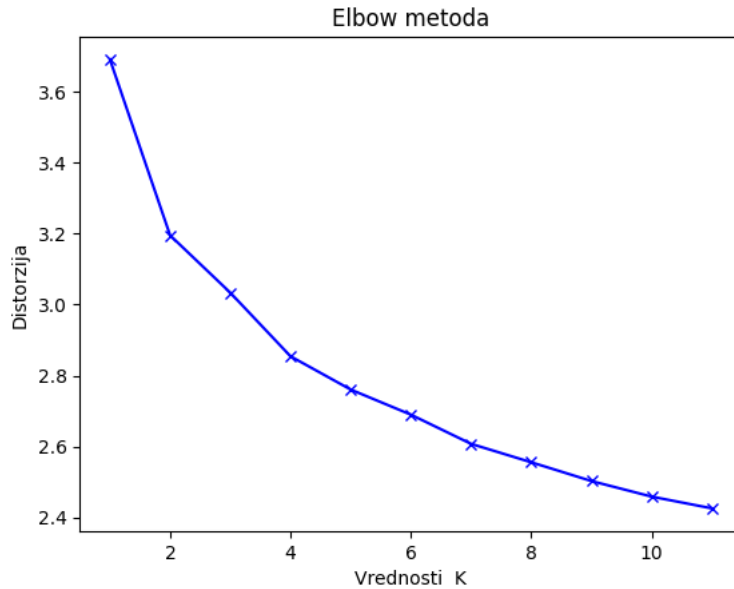
1 : 3.6913650287156656
2 : 3.195397659473886
3 : 3.0337668694656728
4 : 2.854440107706939
5 : 2.7606542867947423
6 : 2.6898759944248862
7 : 2.6076150932171145
8 : 2.5556592662526554
9 : 2.5027815746300863
10 : 2.4589177197818985
11 : 2.425843796089128

```

```

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 : 340396.85283842846
6 : 321577.40636241285
7 : 305634.37154894805
8 : 293727.4616294227
9 : 283660.95170666
10 : 270095.10949484084
11 : 264163.24255030265

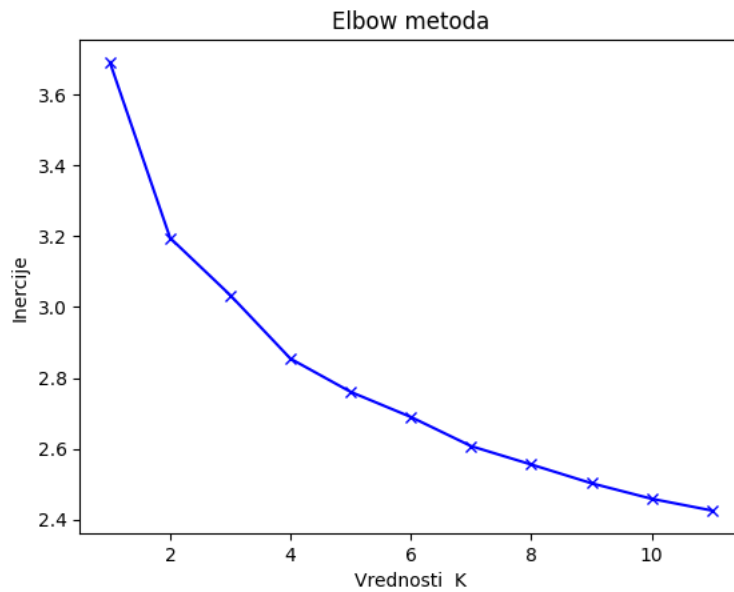
```

```

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

```





```

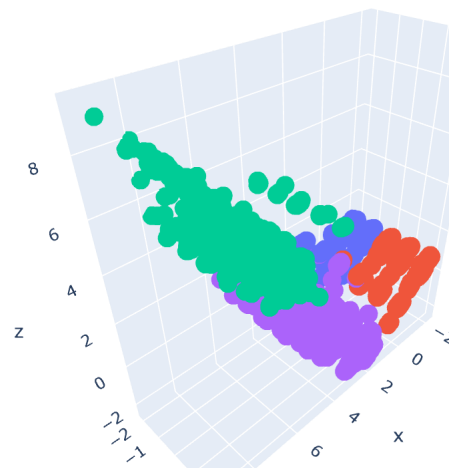
1 kmeans = KMeans(n_clusters =4)
2 cluster_labels = kmeans.fit_predict(dataset)
3 results = eval_metrics(dataset,cluster_labels)
4 print("Silhouette Score : "+str(results[0])+" Davies-Bouldin Index: "+str(results[1])+"Calinski-Harabasz Index: "+str(results[2]))
5 results_df = pd.DataFrame()
6 results_df.loc[1,"KMeans"]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

```

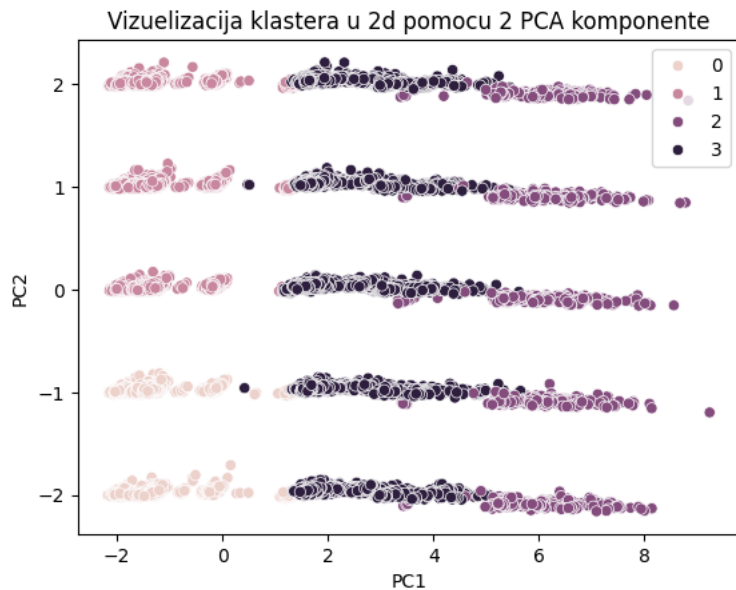
Silhouette Score : 0.16992386604518878 Davies-Bouldin Index: 1.6498489644815606Calinski-Harabasz Index: 9878.544076019489

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

Vizuelizacija klastera u 3D pomocu 3 PCA komponente



```
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
```

	KMeans	
1	0.169924	

```
1 clustering_dataframe
```

	KMeans	
0	11101	
1	16398	
2	1537	
3	12152	

```
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[1,"PCA_KMeans"] = results[0]
```

Silhouette Score : 0.516524642626481 Davies-Bouldin Index: 0.8108890221200817Calinski-Harabasz Index: 52734.52377635431



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

```
2 clustering_dataframe["PCA_Kmeans"] = counts
```

```
1 results_df
```

	KMeans	PCA_KMeans	
1	0.169924	0.516525	

```
1 clustering_dataframe
```

	KMeans	PCA_Kmeans	
0	11101	16136	
1	16398	5783	
2	1537	8448	
3	12152	10821	

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[1,"BisectingKMeans"] = results[0]

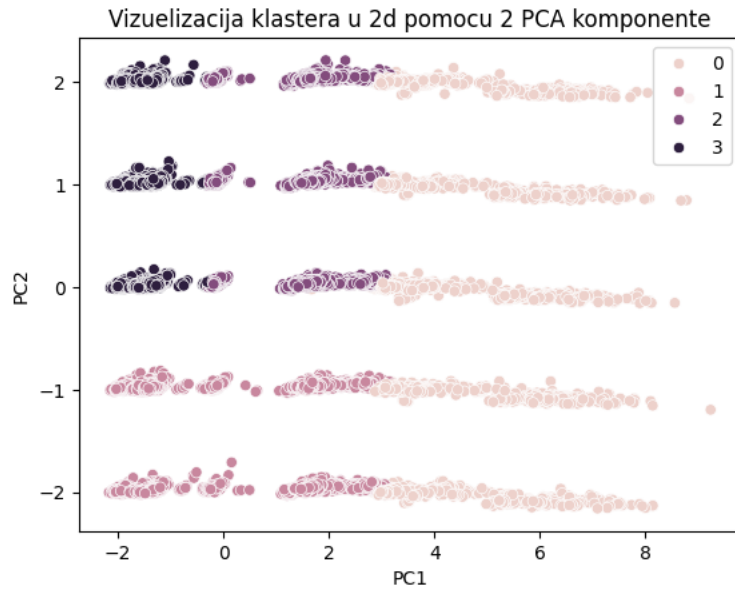
```

Silhouette Score : 0.13733267375248376 Davies-Bouldin Index: 1.8845460378968608Calinski-Harabasz Index: 7462.77928516293

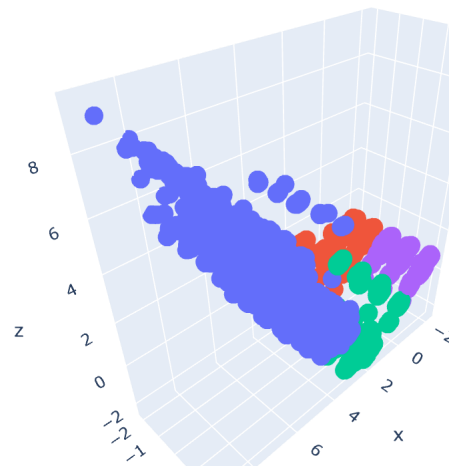
```

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

```



Vizuelizacija klastera u 3D pomocu 3 PCA komponente



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

```
1 results_df
```

	KMeans	PCA_KMeans	BisectingKMeans
1	0.169924	0.516525	0.137333

```
1 clustering_dataframe
```

	KMeans	PCA_Kmeans	Bisecting Kmeans
0	11101	16136	4927
1	16398	5783	14641
2	1537	8448	7285
3	12152	10821	14335

Next steps: [View recommended plots](#)

```

1 bisectingkmeans = BisectingKMeans(n_clusters =4)
2 cluster_labels = bisectingkmeans.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[1,"PCA_BisectingKMeans"] = results[0]

```

Silhouette Score : 0.5105757848760489 Davies-Bouldin Index: 0.7468223050428169Calinski-Harabasz Index: 50709.23961838911

```

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

```

## ✓ 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[1,"Gaussian Mixture"] = results[0]

```

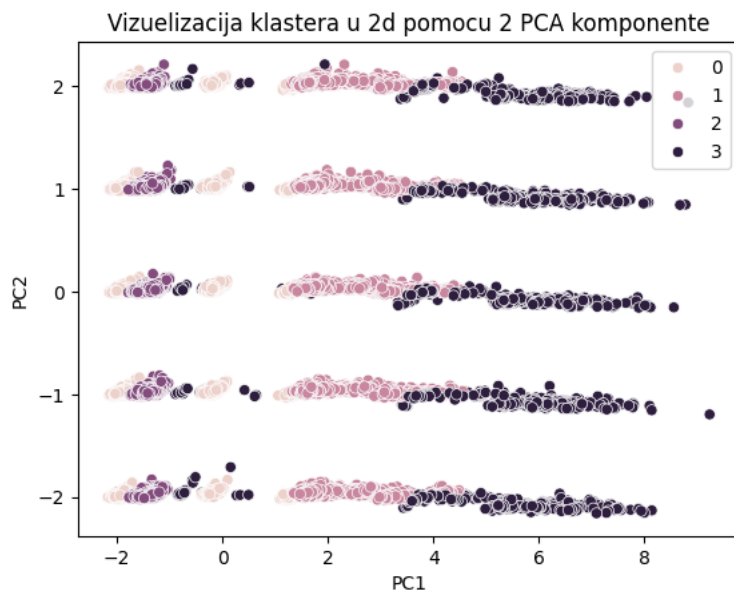
Silhouette Score : 0.09975477726149581 Davies-Bouldin Index: 3.3015872624880185Calinski-Harabasz Index: 7460.016926791017

```

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 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[1,"PCA_Gaussian Mixture"] = results[0]

```

Silhouette Score : 0.4654669162073473 Davies-Bouldin Index: 0.8838660347437242Calinski-Harabasz Index: 38707.1393547403

```



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

```

```
1 results_df
```

	KMeans	PCA_KMeans	BisectingKMeans	PCA_BisectingKMeans	Gaussian Mixture	PCA_Gaussian Mixture
Silhouette Score	0.5105757848760489	0.5105757848760489	0.09975477726149581	0.09975477726149581	0.4654669162073473	0.4654669162073473
Davies-Bouldin Index	0.7468223050428169	0.7468223050428169	3.3015872624880185	3.3015872624880185	0.8838660347437242	0.8838660347437242
Calinski-Harabasz Index	50709.23961838911	50709.23961838911	7460.016926791017	7460.016926791017	38707.1393547403	38707.1393547403

```
1 clustering_dataframe
```

	KMeans	PCA_Kmeans	Bisecting Kmeans	PCA_BisectingKMeans	Gaussian Mixture	PCA_Gaussian Mixture	
0	11101	16136	4927	12223	7931	9756	
1	16398	5783	14641	2008	11603	7315	
2	1537	8448	7285	16136	19623	14241	
3	12152	10821	14335	10821	2031	9876	

Next steps:  [View recommended plots](#)

## ✓ Hierarchical Clustering

```
1 from sklearn.cluster import AgglomerativeClustering
2
```

```
1 from scipy.cluster.hierarchy import dendrogram,linkage
```

```
1 # linkage_data = linkage(df2,method = "ward",metric = "euclidean")
```

```
1 # dendrogram(linkage_data)
2 # plt.show()
```

```
1 # hierarchical_cluster = AgglomerativeClustering(n_clusters = 4, affinity = "euclidean",linkage = "ward")
2 # cluster_labels = hierarchical_cluster.fit_predict(df2)
```

```
1 # results = eval_metrics(df2,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[1,"Hierarchical Clustering(ward)"] = results[0]
```

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