

Data mining - Breast Cancer

Branko Grbić - Matematički fakultet, Univerzitet u Beogradu
mi20002@alas.matf.bg.ac.rs

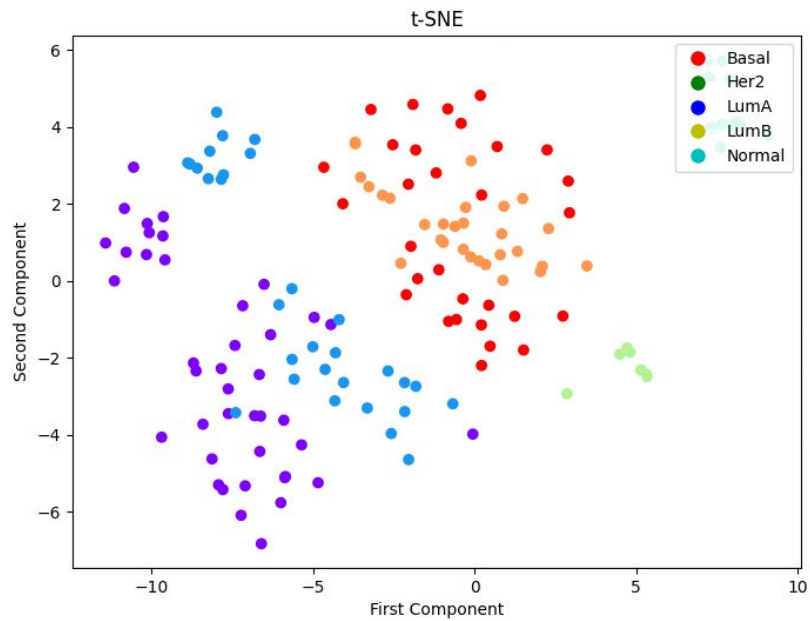
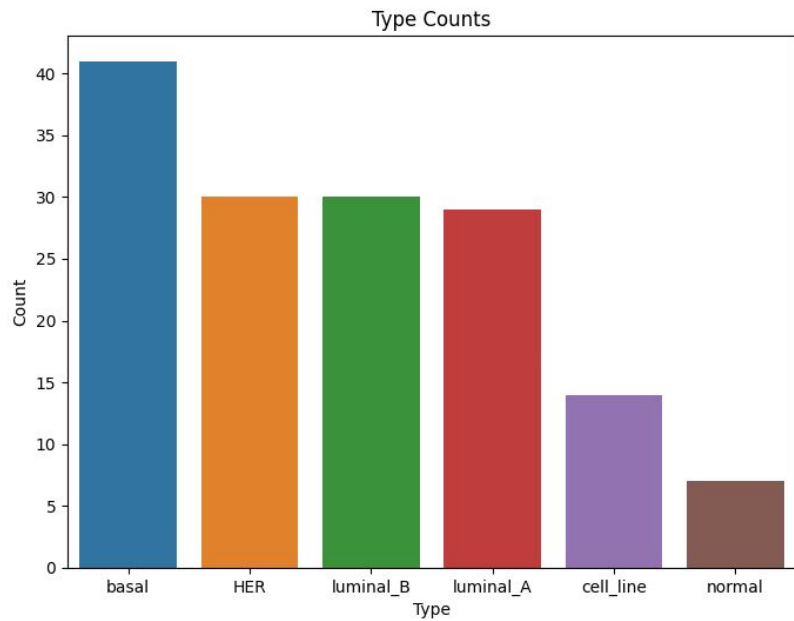
O projektu

- CuMiDa dataset - 151 entiteta, 54676 atributa
- EDA
 - Informacije o podacima
 - Prvi uvid u zavisnosti i značaj atributa
- Klasifikacija
 - XGBoost
 - SVM
 - KNN
 - Ansambl
- Klasterovanje
 - KMeans
 - Gausova mešavina
- Pravila pridruživanja
 - Apriori

Eksplorativna analiza podataka

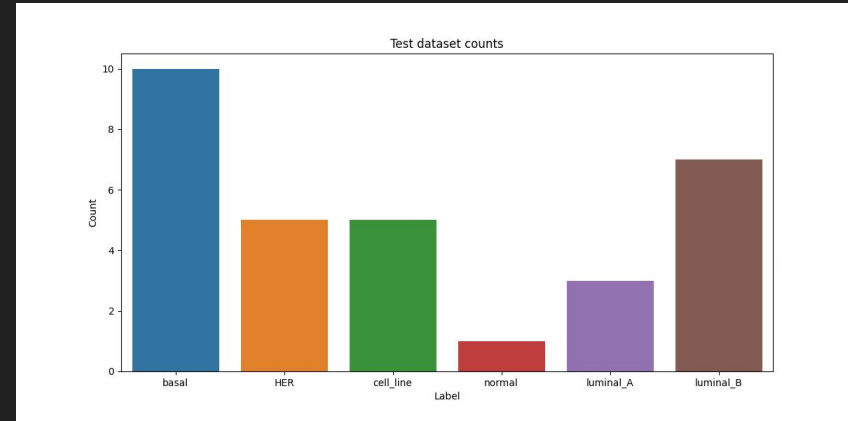
- Šest nebalansiranih klasa
- Uvid u atribute - nema naznaka zavisnosti pojedinačnih atributa od klasa
- PCA analiza
 - Objašnjiva varijansa - 90 komponenti = 85% varijanse
- t-SNE analiza
 - Demonstrira veću razliku atributa sa 2 komponente - model ima nelinearne zavisnosti
- Korelacija atributa
 - Vrlo mali broj atributa visoko korelisan
 - Ogroman broj atributa pa se ne može pokrenuti $O(n^2)$ algoritam za uklanjanje atributa

type	1007_s_at	1053_at	117_at	121_at	1255_g_at	1294_at	1316_at	1320_at	1405_i_at	...	AFFX-r2-Ec-bioD-3_at	AFFX-r2-Ec-bioD-5_at	AFFX-r2-P1-cre-3_at	AFFX-r2-P1-cre-5_at	AFFX-ThrX-3_at
basal	9.850040	8.097927	6.424728	7.353027	3.029122	6.880079	4.963740	4.408328	8.870780	...	12.229711	11.852955	13.658701	13.477698	6.265781
basal	9.861357	8.212222	7.062593	7.685578	3.149468	7.542283	5.129607	4.584418	7.767646	...	12.178531	11.809408	13.750086	13.470146	6.771853
basal	10.103478	8.936137	5.735970	7.687822	3.125931	6.562369	4.813449	4.425195	9.417956	...	12.125108	11.725766	13.621732	13.295080	6.346952
basal	9.756875	7.357148	6.479183	6.986624	3.181638	7.802344	5.490982	4.567956	9.022345	...	12.111235	11.719215	13.743108	13.508861	6.610284
basal	9.408330	7.746404	6.693980	7.333426	3.169923	7.610457	5.372469	4.424426	9.400056	...	12.173642	11.861296	13.797774	13.542206	6.414354



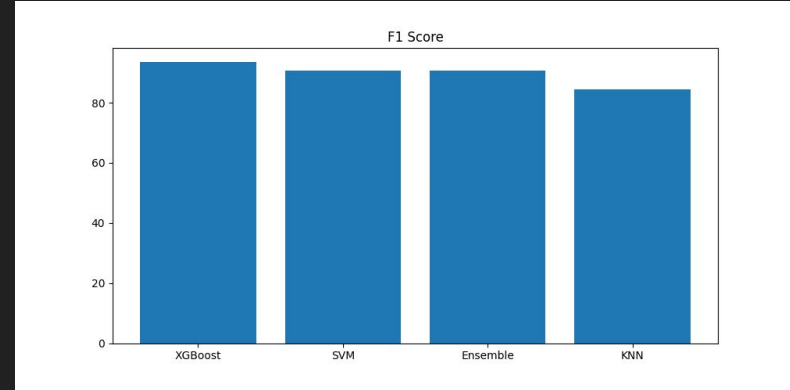
Klasifikacija

- Podela skupa - 80% / 20%
- Nadsemplovanje?
- Više varijanti skupa podataka
 - Običan dataset (Normalized dataset)
 - Dataset sa prvih 90 PCA komponenti (PCA Dataset)
 - Nadsemplovan dataset (Normalized SMOTE dataset)
 - Nadsemplovan dataset sa prvih 90 PCA komponenti (PCA+SMOTE dataset)
- Parametrizacija - Grid Search
- Rezultat - **F1 skor** (*otežinjen*), tačnost
- Analiza rezultata - Matrica konfuzije

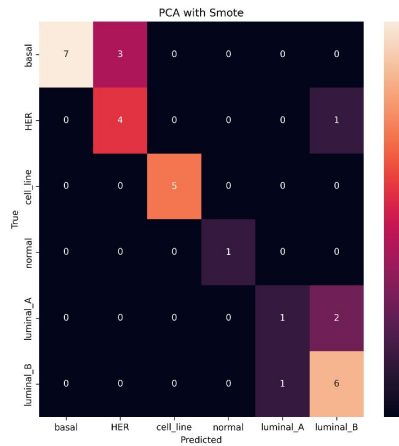
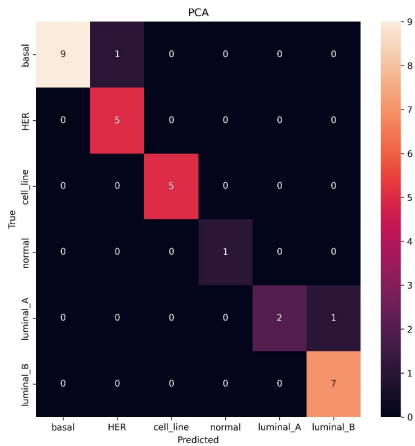
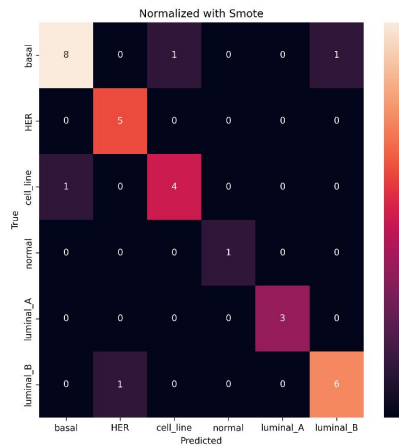
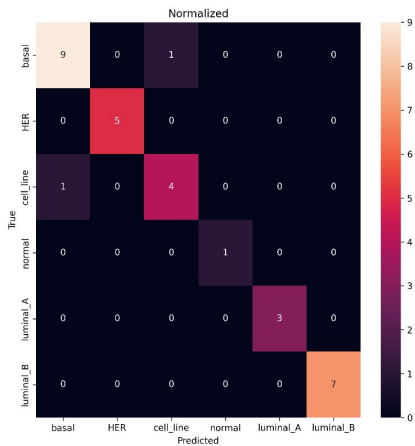


Rezultati

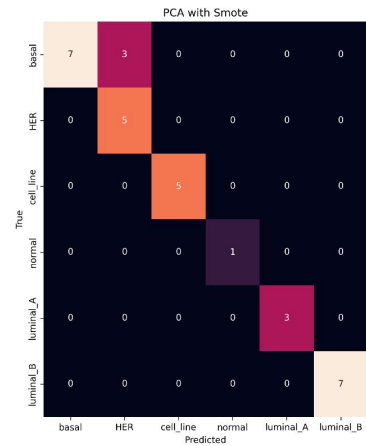
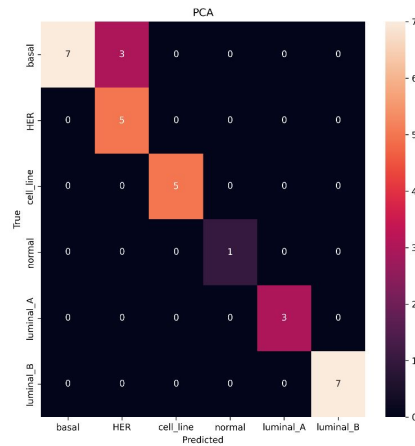
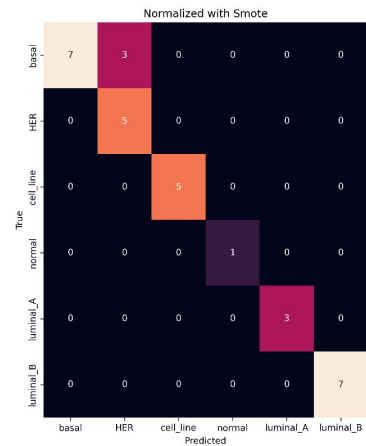
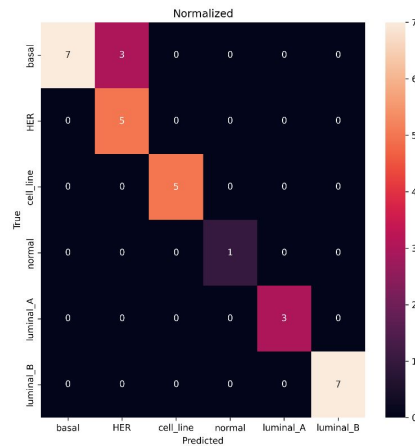
- **XGBoost (Normalized)**
 - Parametri: lr=0.1, max depth=6, estimators=100
 - F1 skor - 93.55%
- **SVM (Normalized)**
 - Parametri: C=0.1, gamma=1, kernel=linear
 - F1 skor - 90.58%
- **KNN (Normalized with SMOTE)**
 - Parametri: n_neighbors=2, p=1, weights=uniform
 - F1 skor: 84.40%
- **Ansaml:**
 - F1 skor: 90.58%



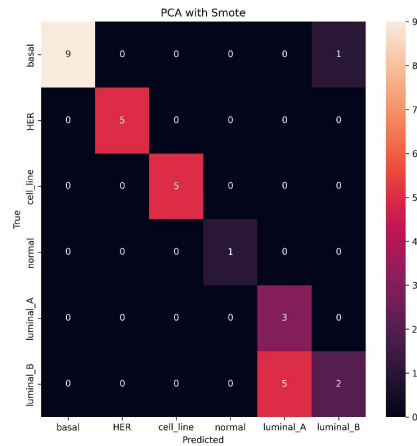
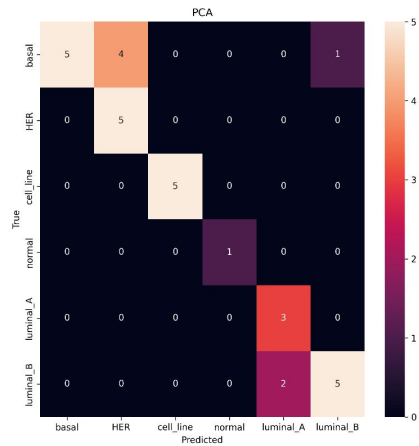
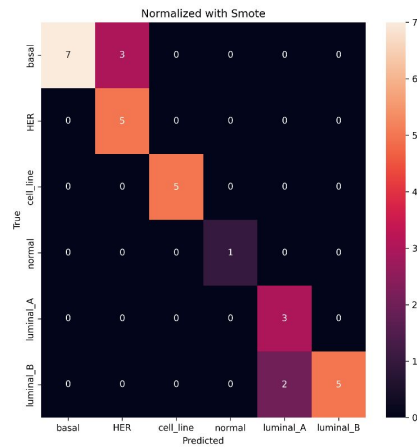
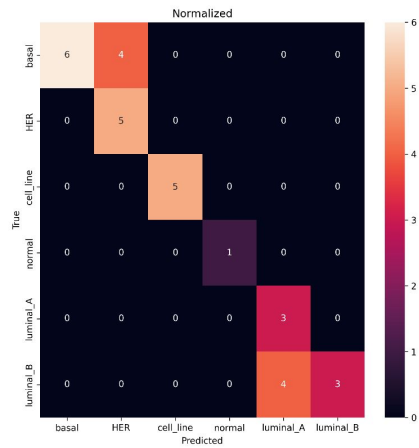
XGBoost



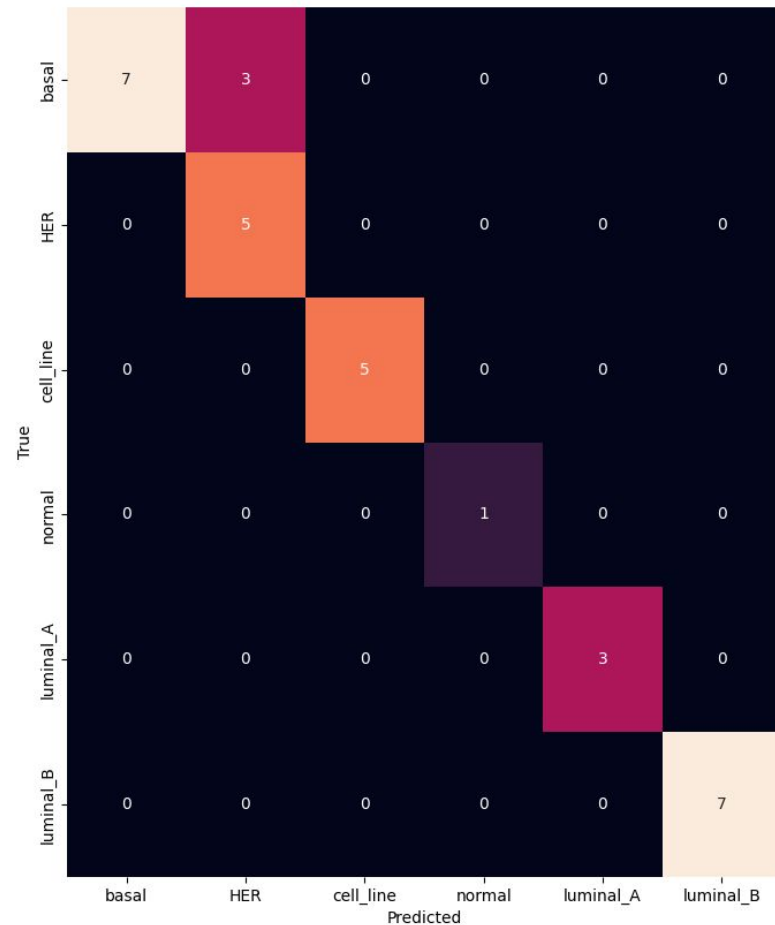
SVM



KNN

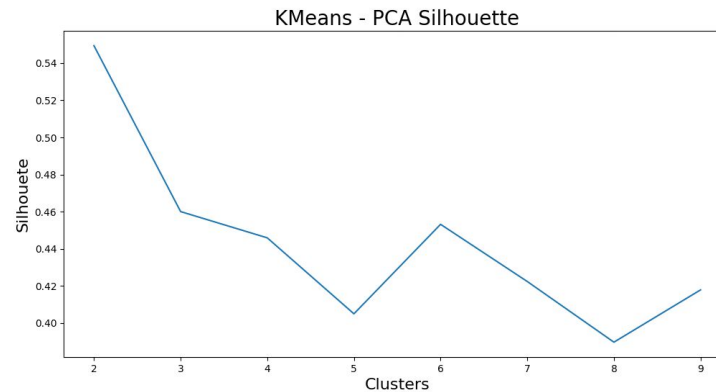
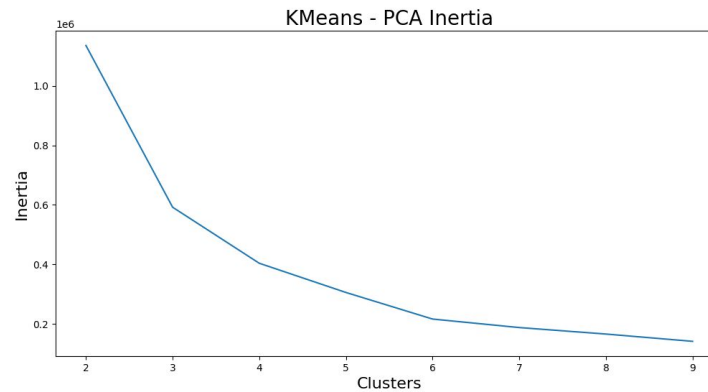
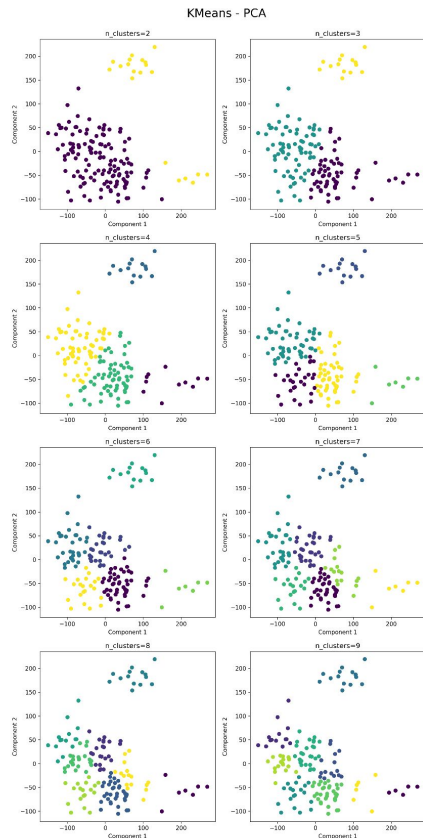


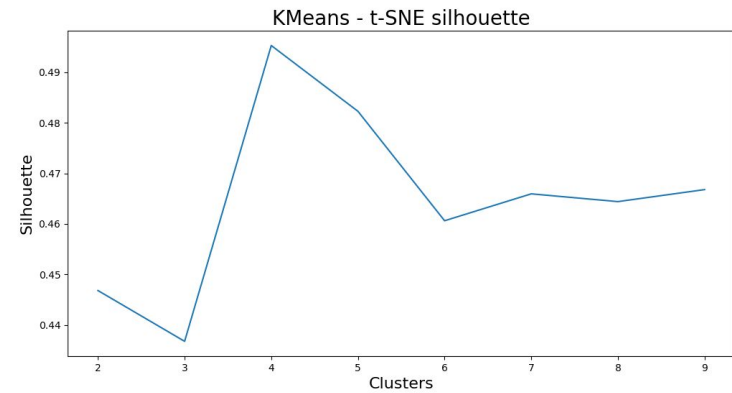
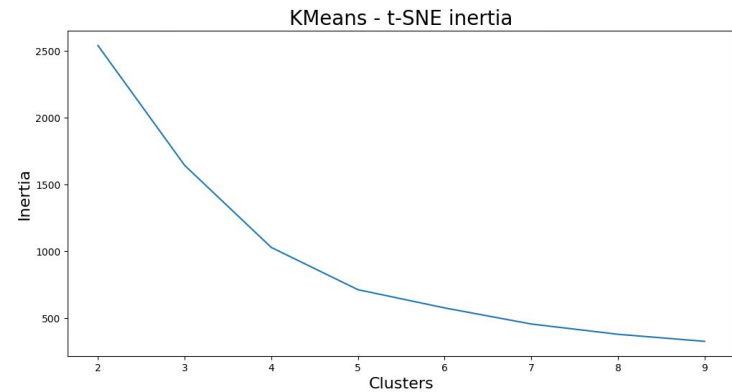
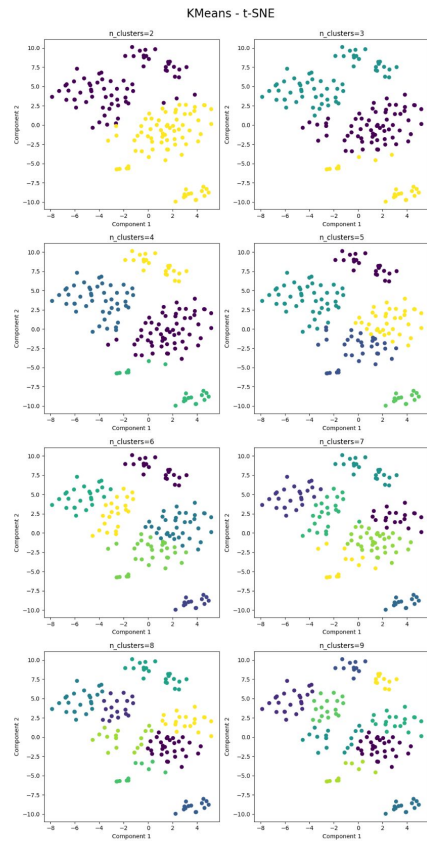
Ensemble Confusion Matrix



Klasterovanje

- KMeans
 - Broj klastera - [2, 9]
 - Inertia
 - Skor siluete
- Gausova mešavina
 - Broj klastera - [2, 9]
 - Kovarijantni tip?
 - BIC skor
 - Skor siluete
- Korišćeni su PCA i t-SNE
- Rezultati nezadovoljavajući po metrikama kvaliteta





Gausova mešavina

- Kovarijantni tipovi

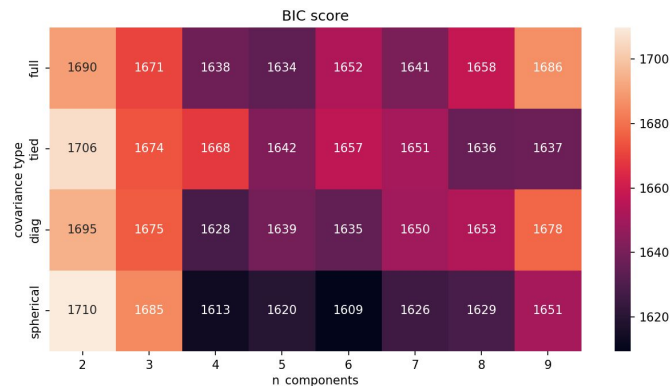
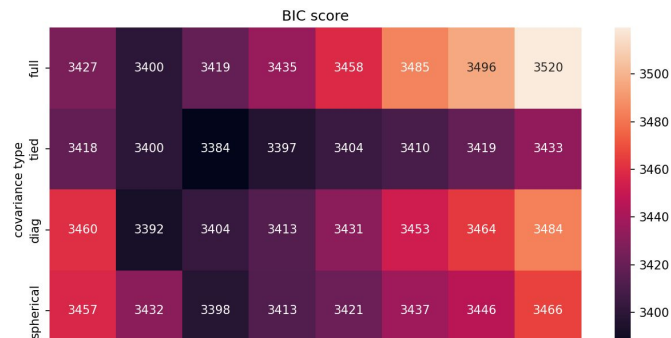
- Full
- Tied
- Diag
- Spherical

PCA

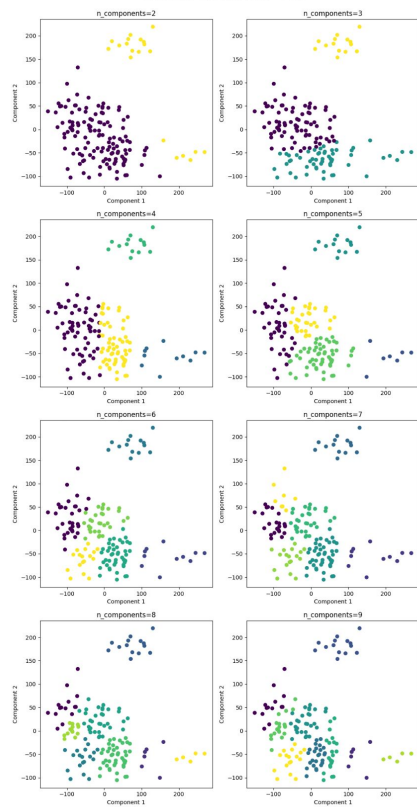
- Kako odrediti najboljeg?

- Niži BIC - bolji kvalitet
- Visok skor siluete - bolji kvalitet
- PCA vs t-SNE

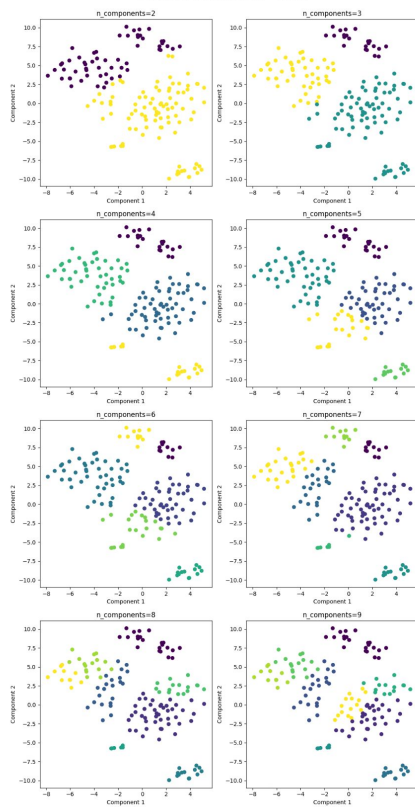
t-SNE



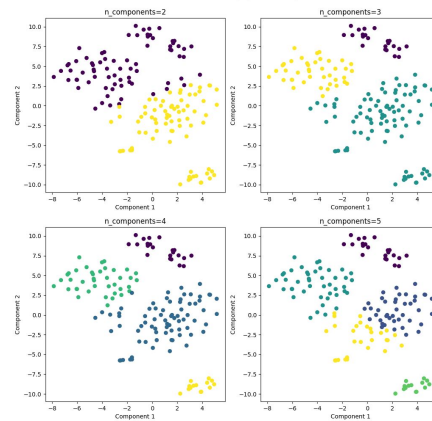
Gaussian Mixture (full) - PCA



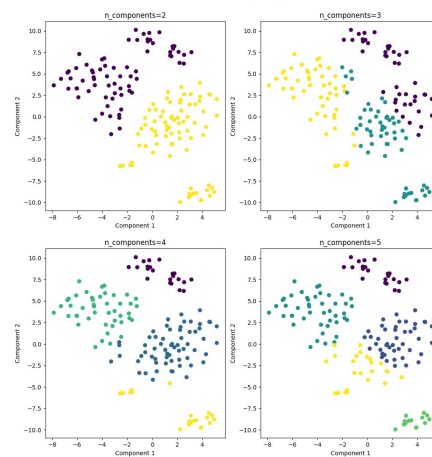
Gaussian Mixture (full) - t-SNE



Gaussian Mixture (spherical) - t-SNE



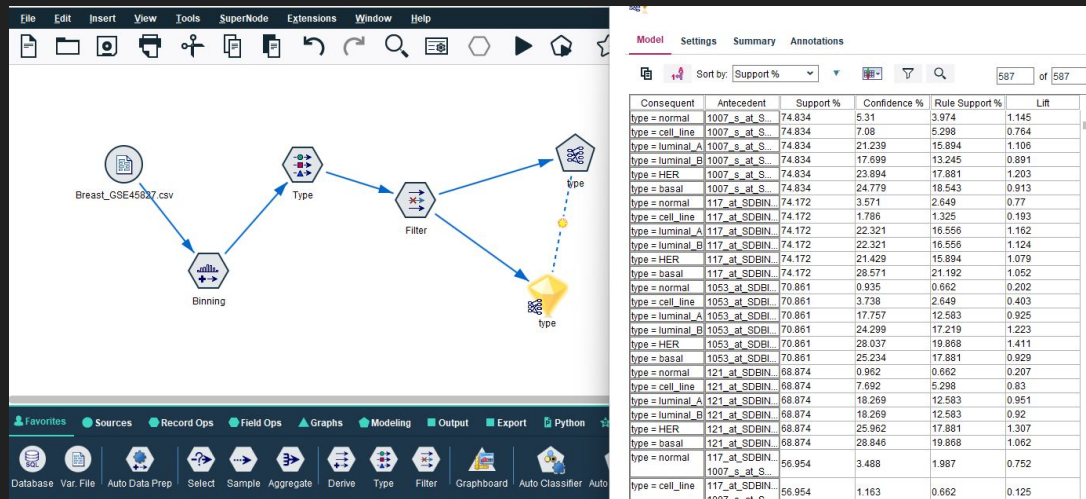
Gaussian Mixture (tied) - t-SNE



Pravila pridruživanja

- Apriori algoritam
- Zahtevaju se kategoričke klase
 - Odraditi pridruživanje na par atributa
 - Razdvojiti kontinualne atribute u kategoričke
 - Binning - po +2 STD

- **Rezultat**
 - Mali lift - Velika podrška
 - Veliki lift - Mala podrška
 - Ne nalazi dobra pravila



Hvala na pažnji

Branko Grbić 2/2020
mi20002@alas.matf.bg.ac.rs
[GitHub link](#)