
PROJEKT 2 MILESTONE3

PRZEMYSŁAW OLENDER, DOMINIK PAWLAK

EDA

- Usunięcie skrótowców (brak)
- Lematyzacja
- Usunięcie 'stopwords'
- Usunięcie słów bardzo długich i bardzo krótkich
- Stworzenie ramki statystycznej i przeskalowanie

	len	words	avg_sen	reading_ease	grade	sentences
0	1.832013	1.549162	0.749075	0.162432	-0.298802	0.775681
1	0.208099	0.189544	0.040772	0.928372	-0.808777	0.609403
2	0.738420	0.632898	0.413880	0.768816	-0.741128	1.108236
3	0.263277	0.197989	0.296945	0.614966	-0.683885	0.609403
4	-0.785101	-0.806946	3.828118	0.500498	-0.668274	-0.554540

PRZYGOTOWANIE DANYCH - TF IDF, PRZESKALOWANIE

- Stworzyliśmy również ramkę z wykorzystaniem narzędzia TF DIF

$$w_{i,j} = tf_{i,j} \times \log\left(\frac{N}{df_i}\right) \quad tf_{i,j} = \frac{n_{i,j}}{\sum_k n_{i,j}} \quad idf(w) = \log\left(\frac{N}{df_t}\right)$$

- Otrzymaliśmy następującą ramkę danych

yellow	yes	yesterday	yield	yieldeth	yoga	yoke	young	youth	zeal
0.0	0.0	0.000000	0.000000	0.0	0.0	0.000000	0.056284	0.057832	0.000000
0.0	0.0	0.000000	0.000000	0.0	0.0	0.000000	0.000000	0.000000	0.000000
0.0	0.0	0.000000	0.000000	0.0	0.0	0.000000	0.049485	0.000000	0.000000
0.0	0.0	0.000000	0.000000	0.0	0.0	0.056522	0.000000	0.000000	0.000000



ALGORYTMY KLASTROWANIA

- KMeans
- Agglomerative Clustering: linkage ,single', ,complete', ,ward'
- GMM (Gaussian Mixture Models): covariance: ,full', ,tied', ,diag'

METRYKI

- Silhouette score
- Davies bouldin score
- Rand score

$$RI = \frac{\text{Number of Agreeing Pairs}}{\text{Number of Pairs}}$$

- Adjusted mutual info score

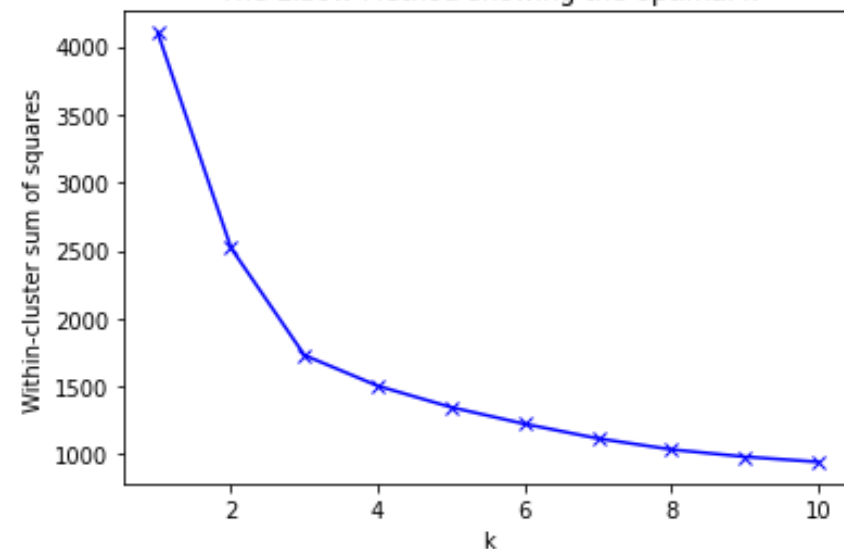
$$ARI = \frac{RI - \text{Expected RI}}{\text{Max}(RI) - \text{Expected RI}}$$

- Mutual info score

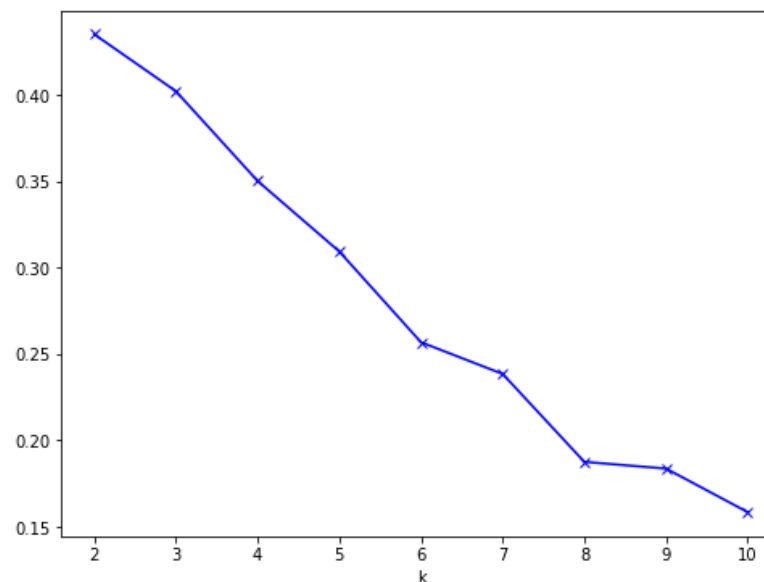
$$MI(U, V) = \sum_{i=1}^{|U|} \sum_{j=1}^{|V|} \frac{|U_i \cap V_j|}{N} \log \frac{N|U_i \cap V_j|}{|U_i||V_j|}$$

WYBÓR LICZBY KLASTRÓW - NA PEŁNYM ZBIORZE

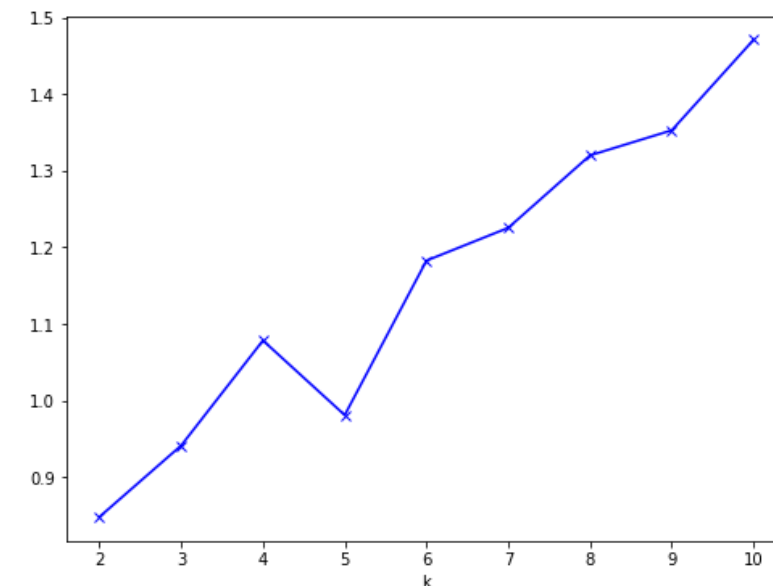
The Elbow Method showing the optimal k



Wykres łokciowy



Metryka Silhouette

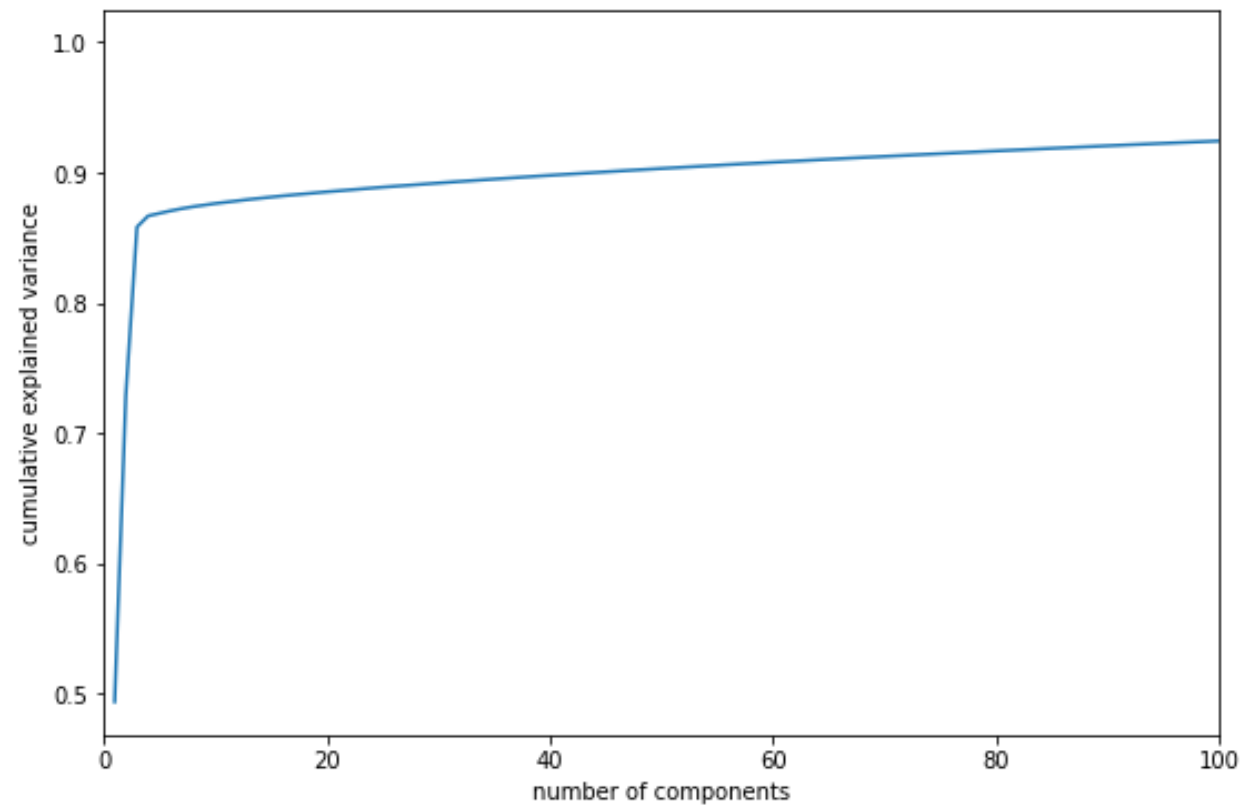


Metryka Davies Bouldin

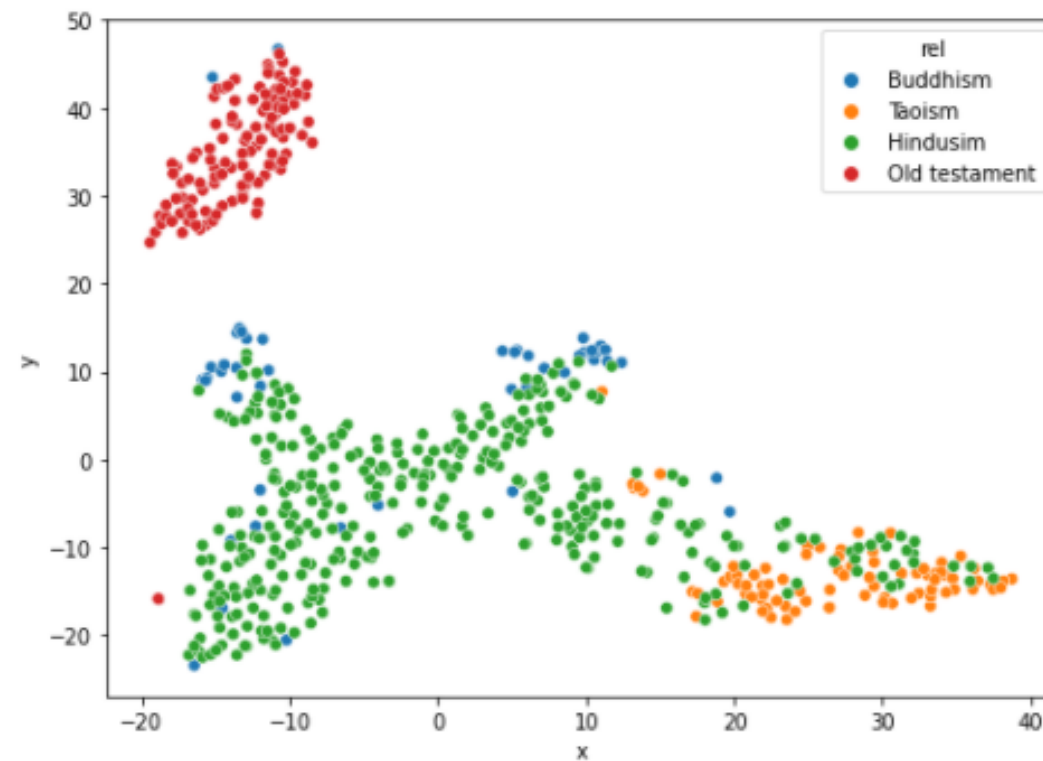
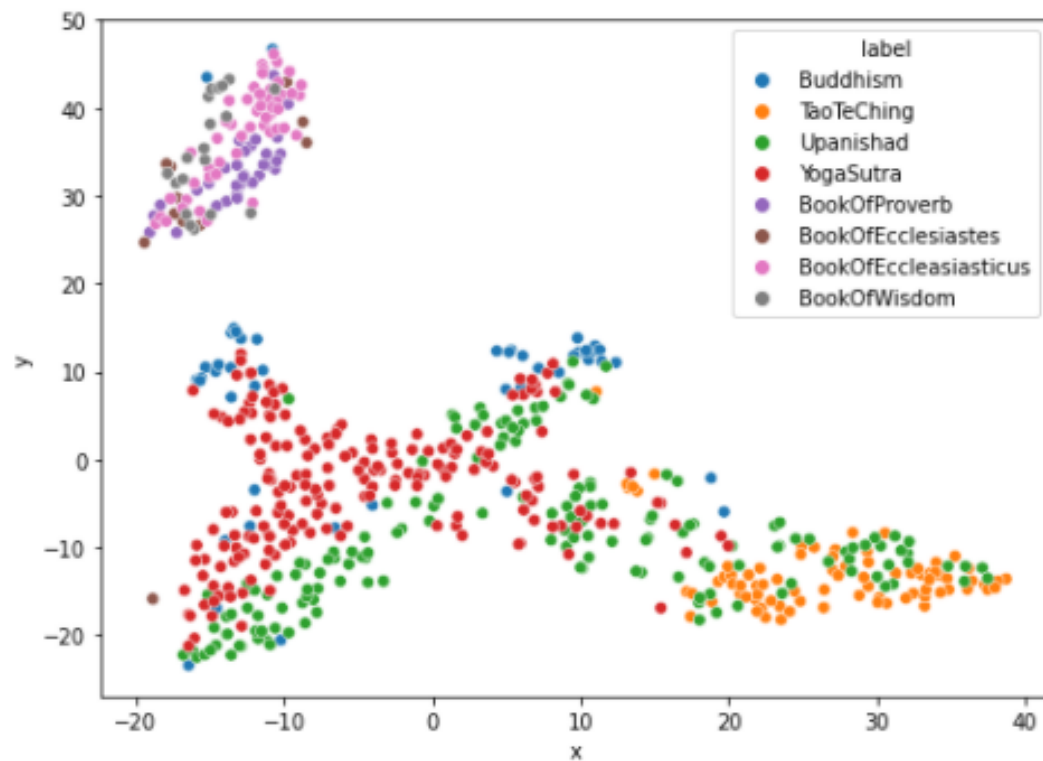
$$DB = \frac{1}{n} \sum_{i=1}^K \max_{j \neq i} \frac{\sigma_i + \sigma_j}{d(c_i, c_j)},$$

gdzie σ_i jest średnią odległością wszystkich punktów ze skupienia i do jego środka, a $d(c_i, c_j)$ jest odległością pomiędzy środkami skupień i oraz j .

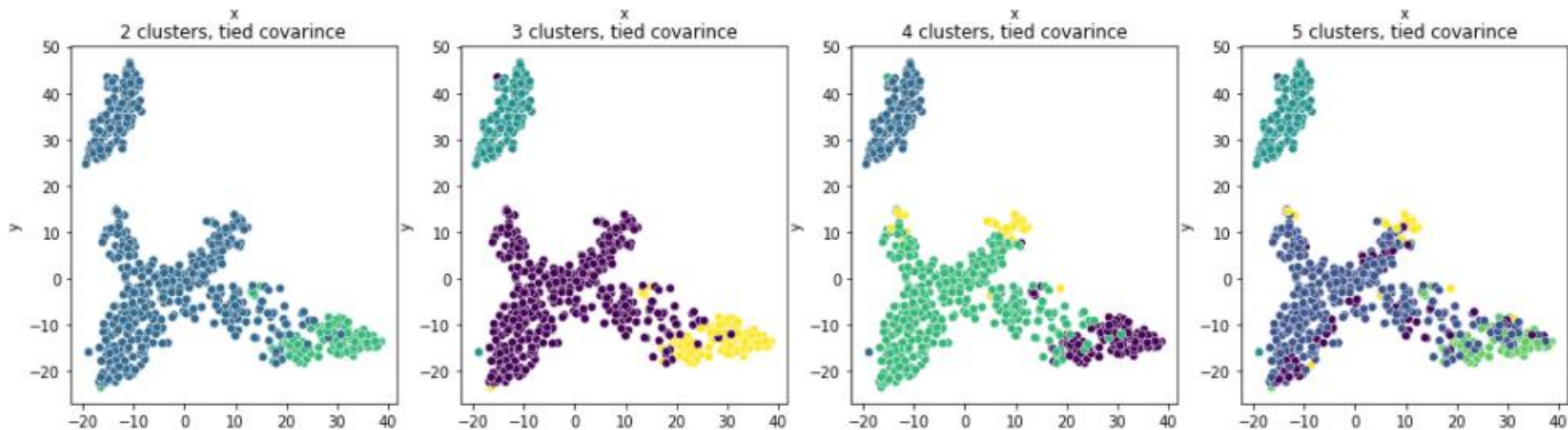
KOMPONENTY PCA



PCA DLA 45 KOMPONENTÓW

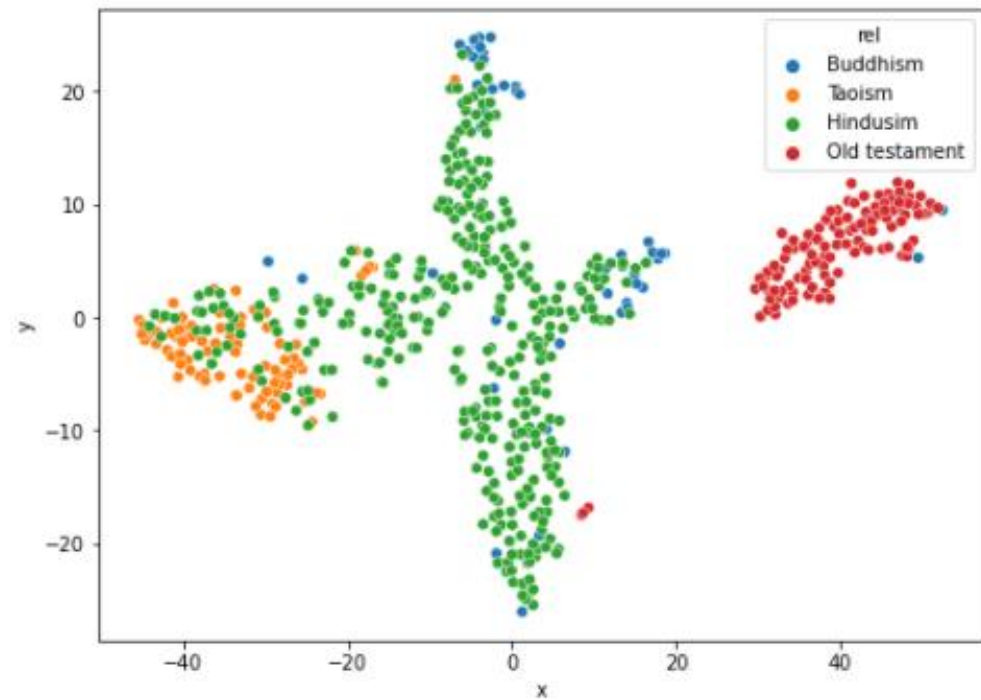
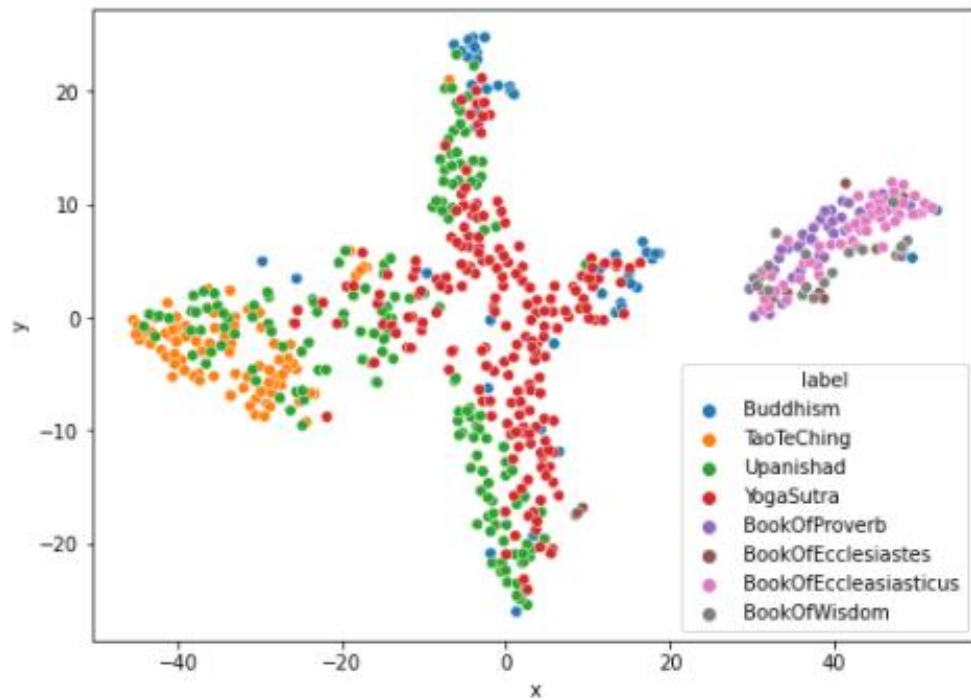


PCA DLA 45 KOMPONENTÓW

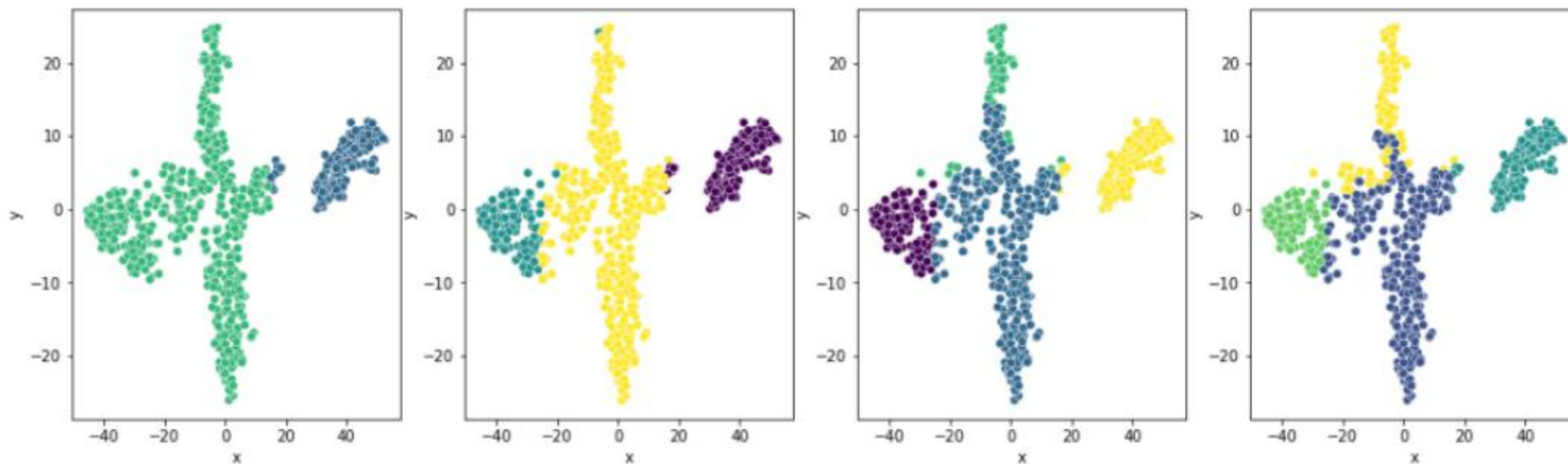


clusters	covariance	silhouette_score	davies_bouldin_score	rand_score	adjusted_mutual_info_score	mutual_info_score
2	tied	0.319456	1.031972	0.617830	0.391053	0.306398
3	tied	0.434745	0.820805	0.860551	0.741807	0.745419
4	tied	0.412935	0.948276	0.895347	0.767429	0.838767
5	tied	0.182694	2.393333	0.838698	0.731193	0.868887

SPARSE PCA DLA 45 KOMPONENTÓW

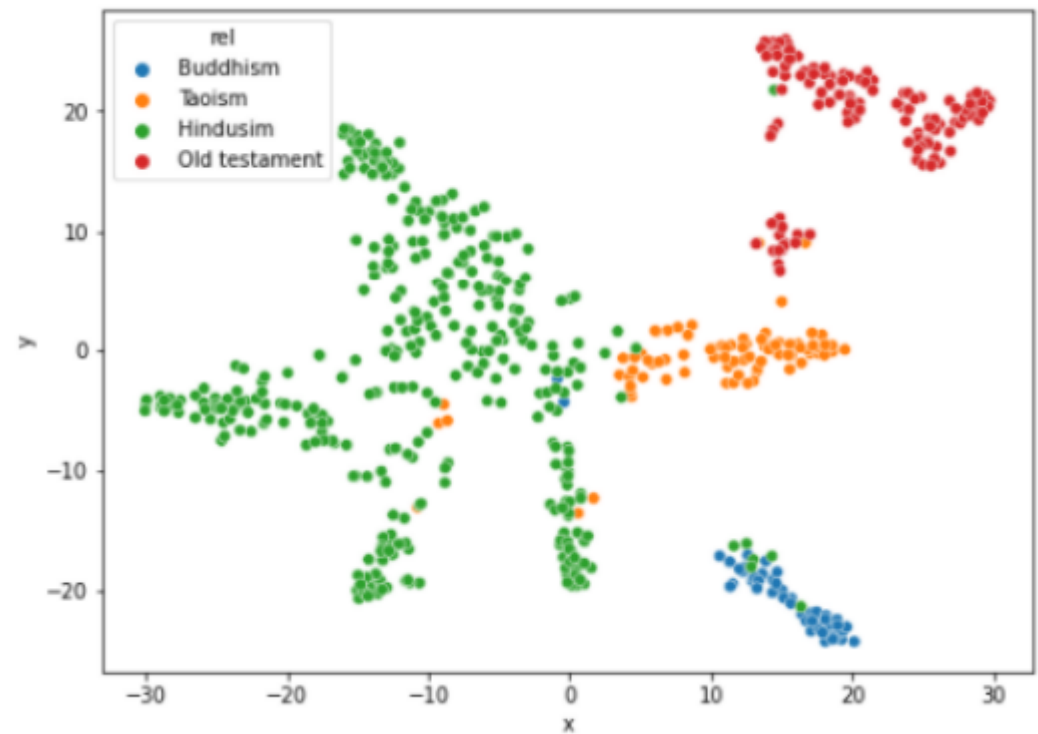
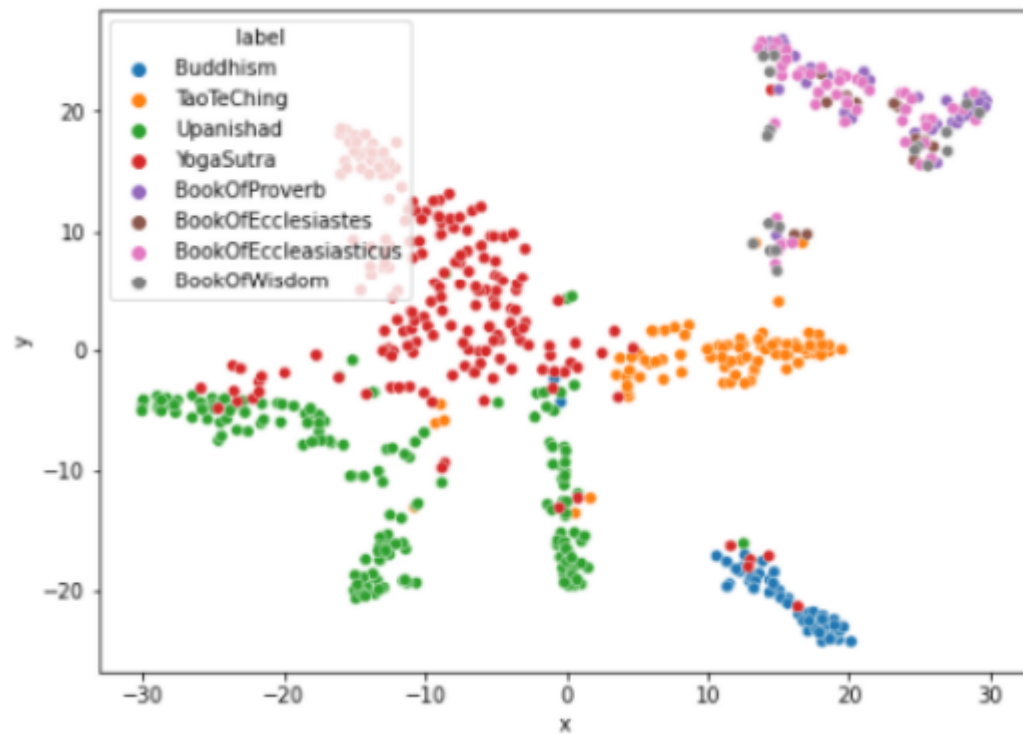


SPARSE PCA DLA 45 KOMPONENTÓW

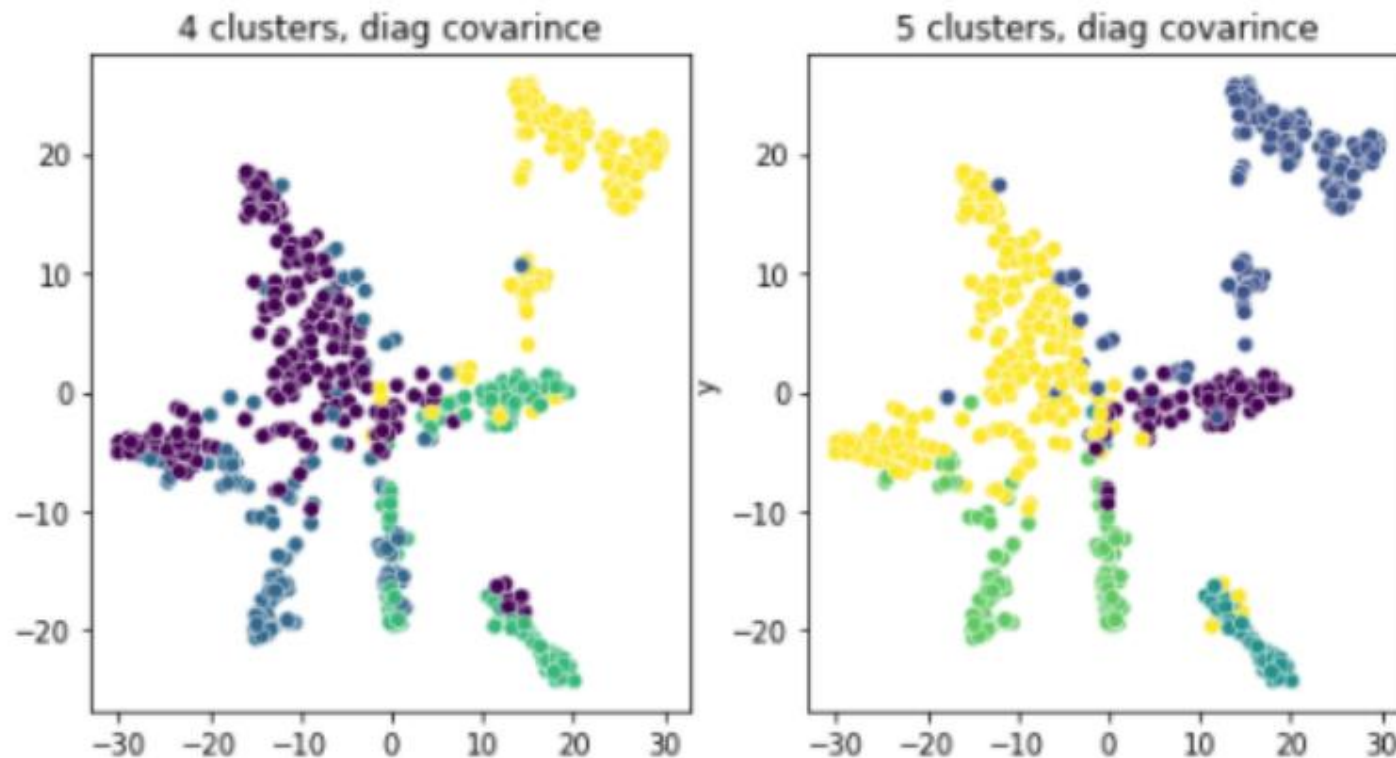


clusters	silhouette_score	davies_bouldin_score	rand_score	adjusted_mutual_info_score	mutual_info_score
2	0.507076	0.708318	0.724555	0.525175	0.422229
3	0.506940	0.710972	0.807856	0.581368	0.591890
4	0.451944	0.816602	0.802066	0.556585	0.638565
5	0.406508	0.725202	0.768271	0.527087	0.630927

NMF DLA 8 KOMPONENTÓW

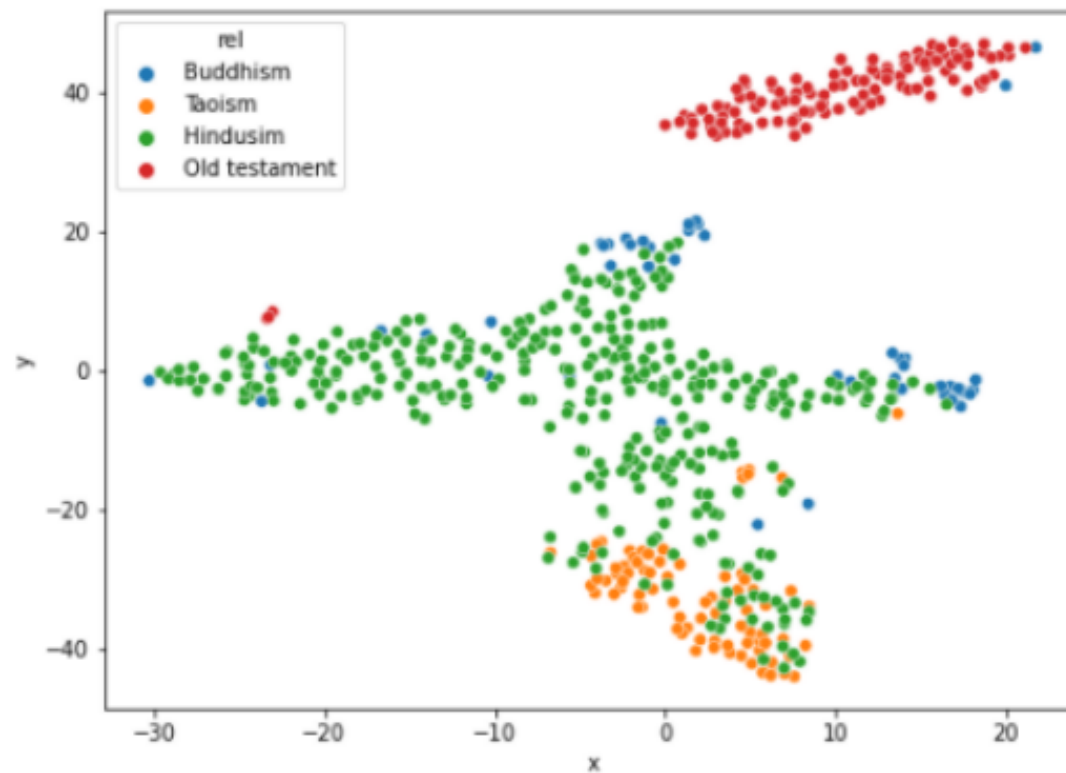
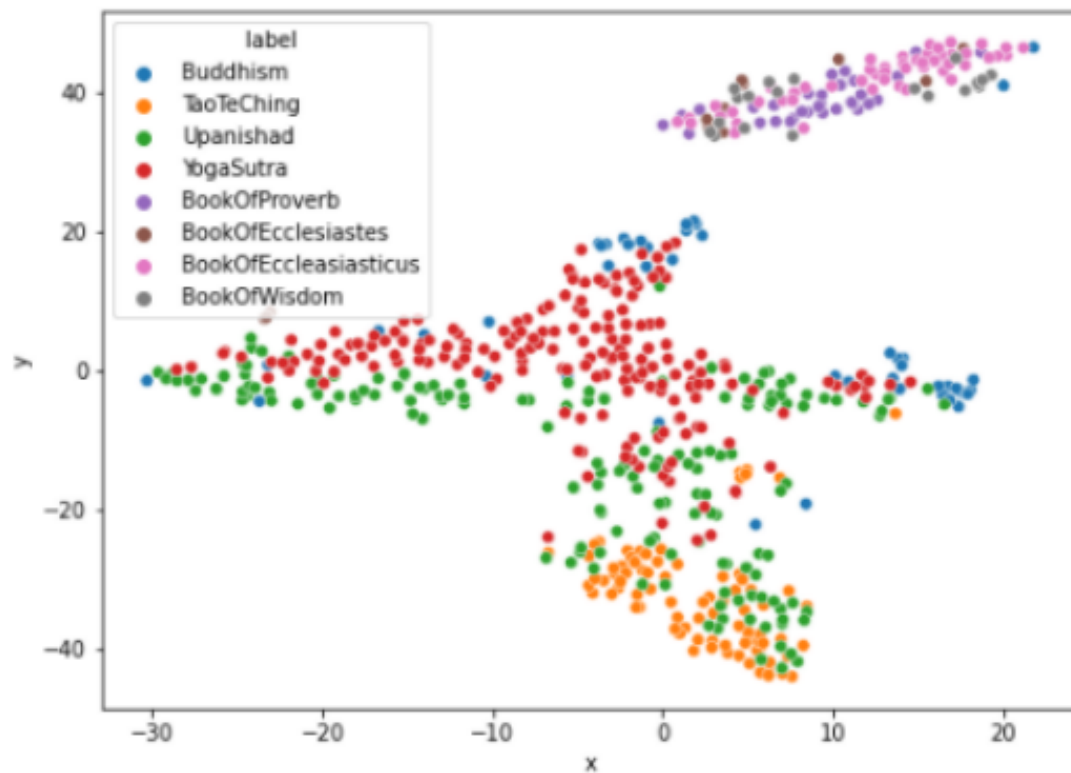


NMF DLA 8 KOMPONENTÓW

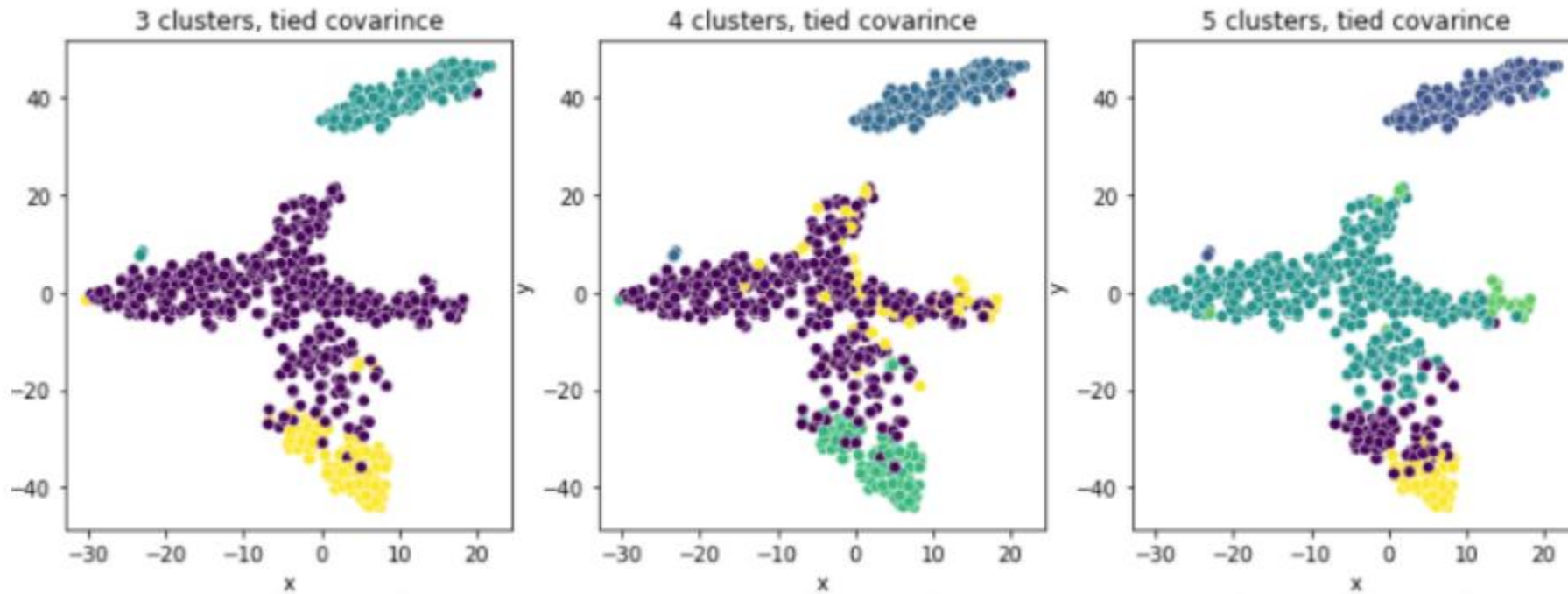


clusters	covariance	silhouette_score	davies_bouldin_score	rand_score	adjusted_mutual_info_score	mutual_info_score
4	diag	0.176353	1.989555	0.744358	0.527929	0.648687
5	diag	0.210749	1.409956	0.775385	0.616757	0.796627

TRUNCATED SVD DLA 50 KOMPONENTÓW



TRUNCATED SVD DLA 50 KOMPONENTÓW

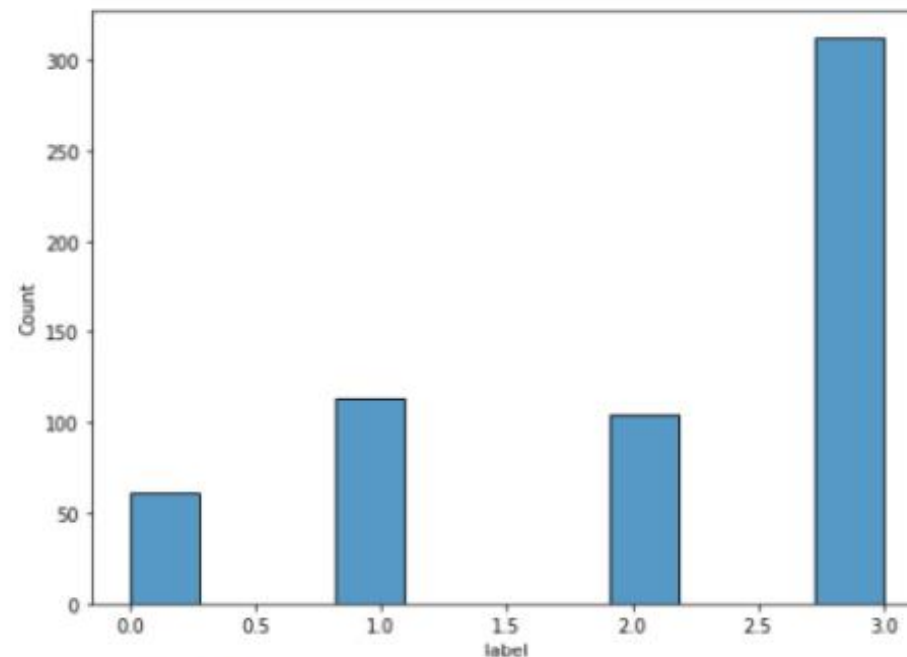
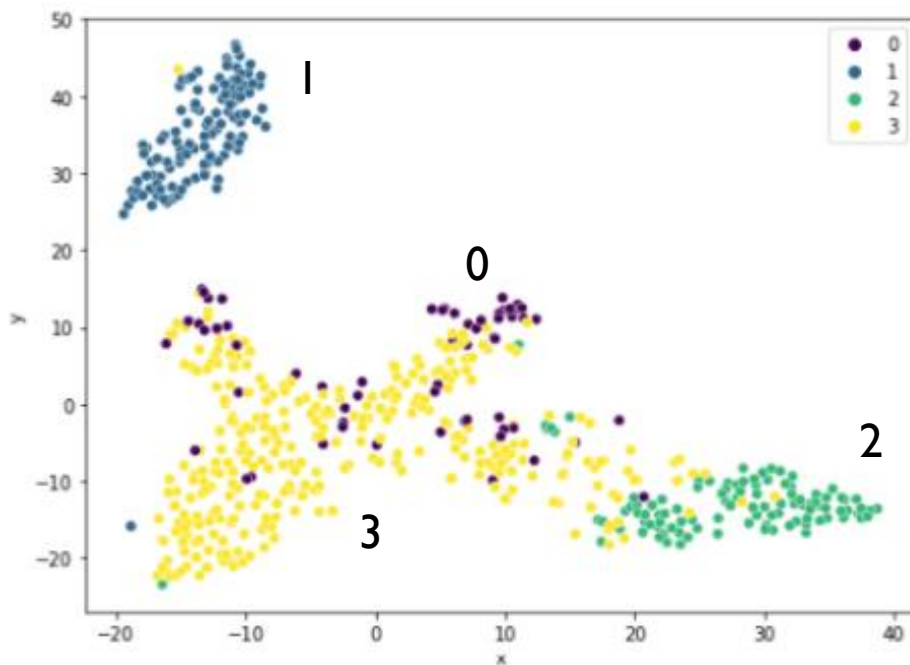


clusters	covariance	silhouette_score	davies_bouldin_score	rand_score	adjusted_mutual_info_score	mutual_info_score
3	tied	0.432942	0.825195	0.860551	0.741807	0.745419
4	tied	0.281062	1.504503	0.840902	0.680186	0.768905
5	tied	0.349039	1.082496	0.833559	0.626290	0.746307



NAJLEPSZE KLASTROWANIA

GMM DLA 4 KLASTRÓW Z COVARIANCJĄ 'TIED', NA ZBIORZE ZREDUKOWANYM PRZY POMOCY PCA Z 45 KOMPONENTAMI.



silhouette_score	0.412935
davies_bouldin_score	0.948276
rand_score	0.895347
adjusted_mutual_info_score	0.767429
mutual_info_score	0.838767
calinski_harabasz_score	364.309993

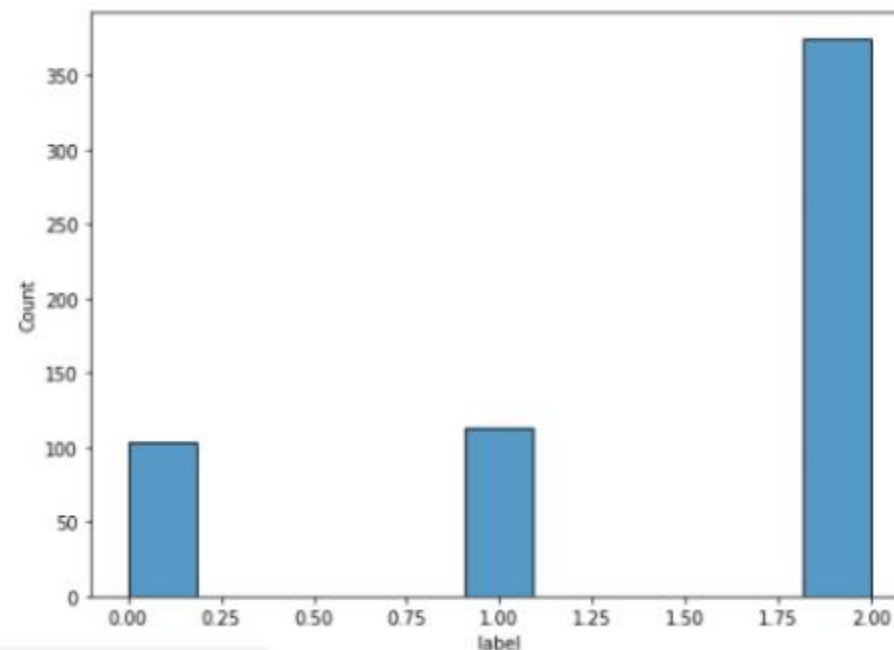
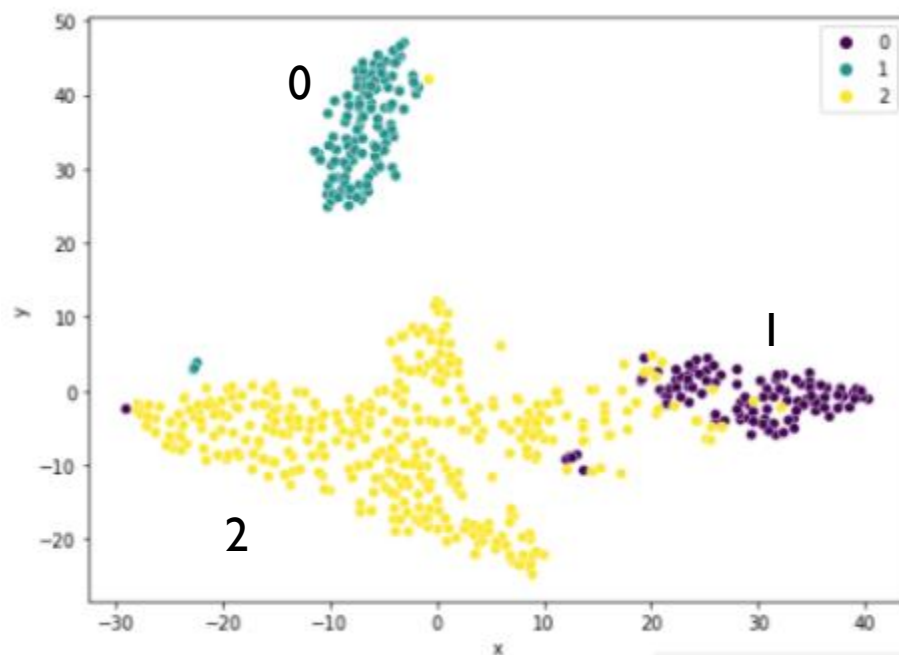
A word cloud of Sanskrit and English terms related to yoga and spirituality. The words are arranged in a circular pattern, with some words being larger and more prominent than others. The colors of the words vary, including shades of green, yellow, orange, and red.

Words included in the word cloud:

- Sanskrit: *brahman*, *soul*, *nachiketa*, *man*, *thing*, *consciousness*, *high*, *psychical*, *death*, *mean*, *object*, *force*, *boon*, *god*, *form*, *sense*, *nature*, *power*, *divine*, *realm*, *pure*, *truth*, *immortal*, *master*, *wise*, *knowledge*, *worship*, *great*, *rest*, *physical*, *meditation*, *thee*, *enter*, *past*, *dwell*, *thou*, *long*, *father*, *meaning*, *obedience*, *exist*, *stage*, *remain*, *vision*, *work*, *live*, *teach*, *visible*, *teacher*, *mortal*, *shall*, *perception*, *true*, *yama*, *perceive*, *atman*, *fine*, *future*, *mind*, *image*, *like*, *sorrow*, *personal*, *heart*, *bring*, *light*, *bear*, *effort*, *inner*, *desire*, *realize*, *stress*, *vision*, *work*, *live*, *teach*, *visible*, *teacher*, *mortal*, *shall*, *perception*, *true*, *yama*, *perceive*, *atman*, *fine*, *future*, *mind*, *image*, *like*, *sorrow*, *personal*, *heart*, *bring*, *light*, *bear*, *effort*, *inner*, *desire*, *realize*, *stress*, *vision*, *work*, *live*, *teach*, *visible*, *teacher*, *mortal*, *shall*, *perception*, *true*, *yama*, *perceive*, *atman*, *fine*, *future*, *mind*, *image*, *like*, *sorrow*, *personal*, *heart*, *bring*, *light*, *bear*, *effort*, *inner*, *desire*, *realize*, *stress*, *vision*, *work*, *live*, *teach*, *visible*, *teacher*, *mortal*, *shall*, *perception*, *true*, *yama*, *perceive*, *atman*, *fine*, *future*, *mind*, *image*, *like*, *sorrow*, *personal*, *heart*, *bring*, *light*, *bear*, *effort*, *inner*, *desire*, *realize*, *stress*, *vision*, *work*, *live*, *teach*, *visible*, *teacher*, *mortal*, *shall*, *perception*, *true*, *yama*, *perceive*, *atman*, *fine*, *future*, *mind*, *image*, *like*, *sorrow*, *personal*, *heart*, *bring*, *light*, *bear*, *effort*, *inner*, *desire*, *realize*, *stress*, *vision*, *work*, *live*, *teach*, *visible*, *teacher*, *mortal*, *shall*, *perception*, *true*, *yama*, *perceive*, *atman*, *fine*, *future*, *mind*, *image*, *like*, *sorrow*, *personal*, *heart*, *bring*, *light*, *bear*, *effort*, *inner*, *desire*, *realize*, *stress*, *vision*, *work*, *live*, *teach*, *visible*, *teacher*, *mortal*, *shall*, *perception*, *true*, *yama*, *perceive*, *atman*, *fine*, *future*, *mind*, *image*, *like*, *sorrow*, *personal*, *heart*, *bring*, *light*, *bear*, *effort*, *inner*, *desire*, *realize*, *stress*, *vision*, *work*, *live*, *teach*, *visible*, *teacher*, *mortal*, *shall*, *perception*, *true*, *yama*, *perceive*, *atman*, *fine*, *future*, *mind*, *image*, *like*, *sorrow*, *personal*, *heart*, *bring*, *light*, *bear*, *effort*, *inner*, *desire*, *realize*, *stress*, *vision*, *work*, *live*, *teach*, *visible*, *teacher*, *mortal*, *shall*, *perception*, *true*, *yama*, *perceive*, *atman*, *fine*, *future*, *mind*, *image*, *like*, *sorrow*, *personal*, *heart*, *bring*, *light*, *bear*, *effort*, *inner*, *desire*, *realize*, *stress*, *vision*, *work*, *live*, *teach*, *visible*, *teacher*, *mortal*, *shall*, *perception*, *true*, *yama*, *perceive*, *atman*, *fine*, *future*, *mind*, *image*, *like*, *sorrow*, *personal*, *heart*, *bring*, *light*, *bear*, *effort*, *inner*, *desire*, *realize*, *stress*, *vision*, *work*, *live*, *teach*, *visible*, *teacher*, *mortal*, *shall*, *perception*, *true*, *yama*, *perceive*, *atman*, *fine*, *future*, *mind*, *image*, *like*, *sorrow*, *personal*, *heart*, *bring*, *light*, *bear*, *effort*, *inner*, *desire*, *realize*, *stress*, *vision*, *work*, *live*, *teach*, *visible*, *teacher*, *mortal*, *shall*, *perception*, *true*, *yama*, *perceive*, *atman*, *fine*, *future*, *mind*, *image*, *like*, *sorrow*, *personal*, *heart*, *bring*, *light*, *bear*, *effort*, *inner*, *desire*, *realize*, *stress*, *vision*, *work*, *live*, *teach*, *visible*, *teacher*, *mortal*, *shall*, *perception*, *true*, *yama*, *perceive*, *atman*, *fine*, *future*, *mind*, *image*, *like*, *sorrow*, *personal*, *heart*, *bring*, *light*, *bear*, *effort*, *inner*, *desire*, *realize*, *stress*, *vision*, *work*, *live*, *teach*, *visible*, *teacher*, *mortal*, *shall*, *perception*, *true*, *yama*, *perceive*, *atman*, *fine*, *future*, *mind*, *image*, *like*, *sorrow*, *personal*, *heart*, *bring*, *light*, *bear*, *effort*, *inner*, *desire*, *realize*, *stress*, *vision*, *work*, *live*, *teach*, *visible*, *teacher*, *mortal*, *shall*, *perception*, *true*, *yama*, *perceive*, *atman*, *fine*, *future*, *mind*, *image*, *like*, *sorrow*, *personal*, *heart*, *bring*, *light*, *bear*, *effort*, *inner*, *desire*, *realize*, *stress*, *vision*, *work*, *live*, *teach*, *visible*, *teacher*, *mortal*, *shall*, *perception*, *true*, *yama*, *perceive*, *atman*, *fine*, *future*, *mind*, *image*, *like*, *sorrow*, *personal*, *heart*, *bring*, *light*, *bear*, *effort*, *inner*, *desire*, *realize*, *stress*, *vision*, *work*, *live*, *teach*, *visible*, *teacher*, *mortal*, *shall*, *perception*, *true*, *yama*, *perceive*, *atman*, *fine*, *future*, *mind*, *image*, *like*, *sorrow*, *personal*, *heart*, *bring*, *light*, *bear*, *effort*, *inner*, *desire*, *realize*, *stress*, *vision*, *work*, *live*, *teach*, *visible*, *teacher*, *mortal*, *shall*, *perception*, *true*, *yama*, *perceive*,

[illegible][illegible]

GMM DLA 3 KLASTRÓW Z COVARIANCJĄ 'TIED', NA ZBIORZE ZREDUKOWANYM PRZY POMOCY TRUNCATEDSVD Z 50 KOMPONENTAMI



silhouette_score	0.432961
davies_bouldin_score	0.825250
rand_score	0.860551
adjusted_mutual_info_score	0.741807
mutual_info_score	0.745419
calinski_harabasz_score	442.245899

[illegible]

A word cloud of Sanskrit terms related to yoga and spirituality. The words are arranged in a dense, overlapping manner, with some words being significantly larger than others. The colors of the words vary, including shades of blue, green, yellow, and red. The words include: death, stage, mind, tell, mortal, pure, knowledge, psychological, fine, wise, live, visible, birth, exist, dwell, shall, force, future, lead, past, brahman, high, vision, mental, personal, psychic, nature, time, power, worship, thee, immortal, spirit, stress, long, self, come, great, form, work, divine, view, world, thou, yama, god, perceive, master, pain, rest, atman, teacher, effort, right, eternal, physical, meaning, meditation, light, remain, like, sorrow, mean, desire, control, boon, man, thing, realm, truth, enter, soul, heart, gain, able, nachiketa, intellect, fabrication, realize, perception, think, true, attain, bear.