



ΕΛΛΗΝΙΚΗ ΔΗΜΟΚΡΑΤΙΑ
Εθνικό και Καποδιστριακό
Πανεπιστήμιο Αθηνών



ΤΜΗΜΑ
ΠΛΗΡΟΦΟΡΙΚΗΣ &
ΤΗΛΕΠΙΚΟΙΝΩΝΙΩΝ

DATA MINING TECHNIQUES

FIRST ASSIGNMENT

Zacharopoulou Lida AM 1115201100004
Seintaridis Dimitrios AM 1115201100197

Προεπεξεργασία δεδομένων - WordClouds

Για την προεπεξεργασία δεδομένων αφαιρέσαμε αρχικά από το συνολικό λεξιλόγιο(corpus) τα stopwords(λέξεις που δεν προσφέρουν κάποια σημαντική πληροφορία στον ταξινομητή καθώς είναι πολύ συχνές).

Τέτοιες λέξεις για παράδειγμα είναι το is, I, the, κτλ. Για να το κάνουμε αυτό χρησιμοποιήσαμε τα stopwords ENGLISH_STOP_WORDS της βιβλιοθήκης sklearn.feature_extraction.text

Επιπλέον, μετά από παρατηρήσεις πάνω στα wordcloud images που παράγουμε, επιλέξαμε να προσθέσουμε κάποια additional stopwords με λέξεις που δεν πρόσφεραν κάποια αξία στο αποτέλεσμα, αλλά εμφανίζονταν με μεγάλη συχνότητα όπως said, much κτλ.

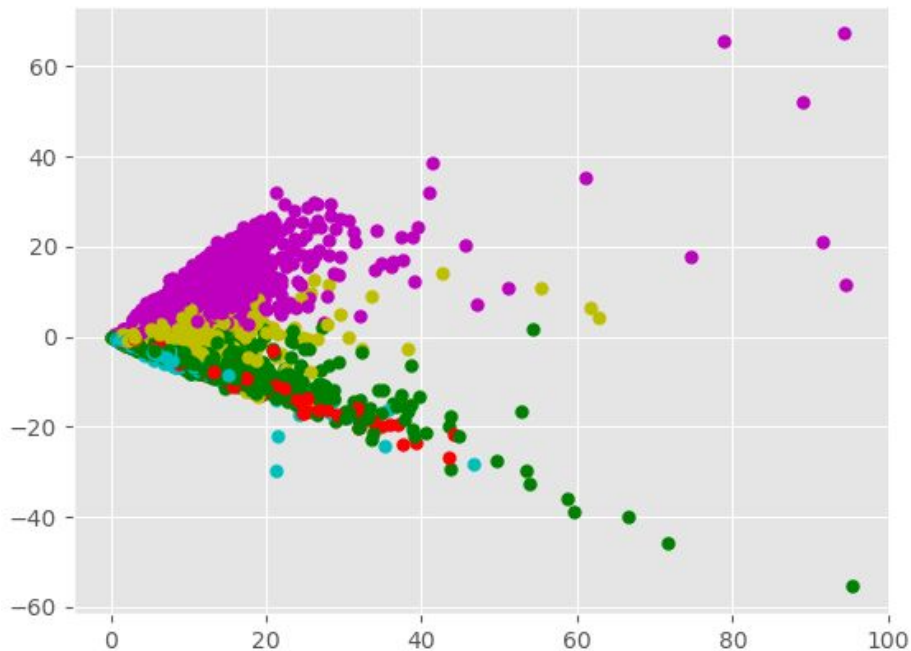
Προεπεξεργασία δεδομένων - Vectorizers

Για να μπορέσουμε να επεξεργαστούμε τα δεδομένα, τα μετατρέψαμε σε πίνακες αριθμών με την χρήση vectorizers, συγκεκριμένα τον TfidfVectorizer και τον CountVectorizer

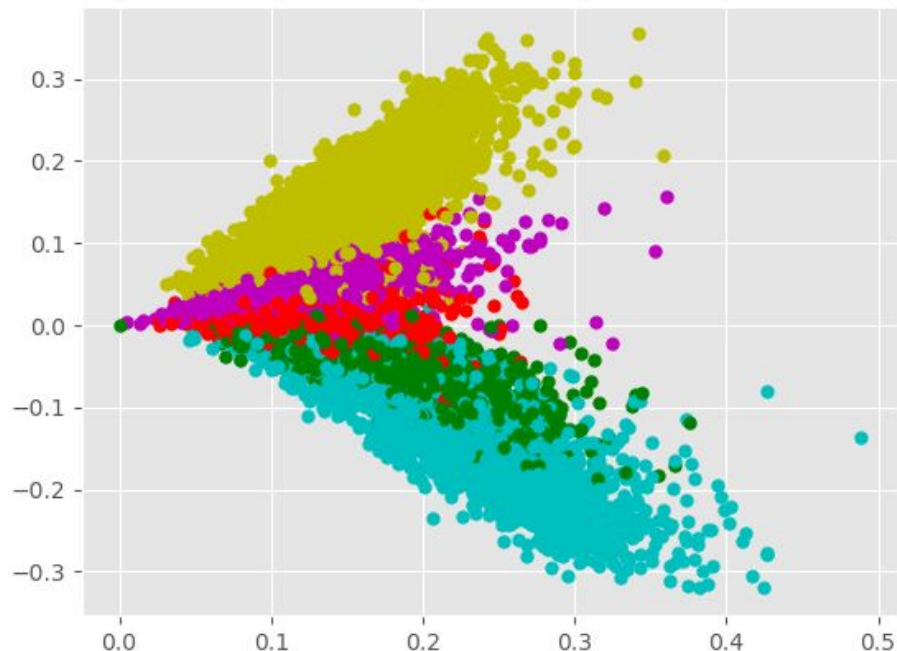
Από τις μετρήσεις μας, όπως φαίνεται αναλυτικά και στα διαγράμματα και τους πίνακες που ακολουθούν, επιλέξαμε τον TFIDF vectorizer,

Counter vs TFIDF Vectorizer

nlTK kmeans, Cosine Similarity,
Counter Vectorizer

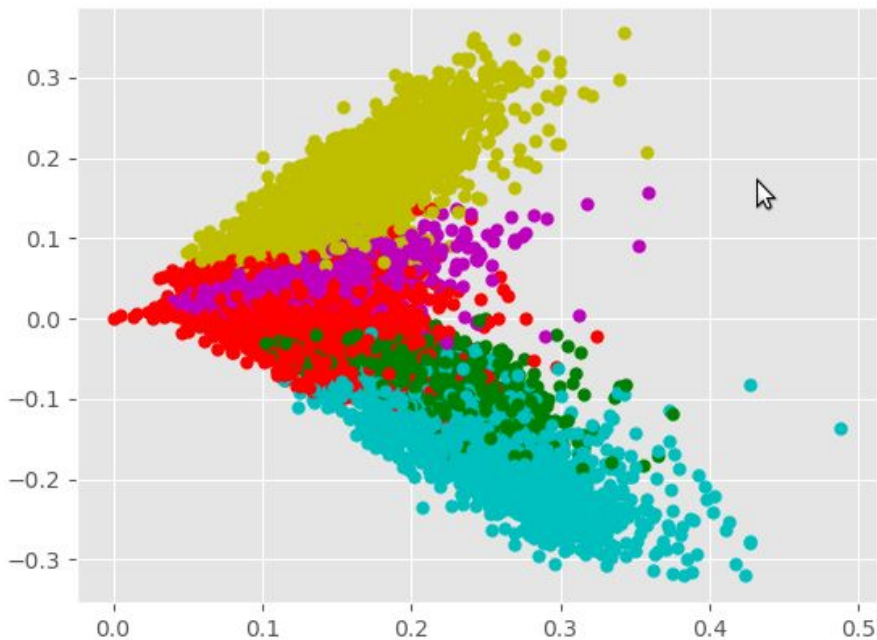


nlTK kmeans, Cosine Similarity,
TFIDF Vectorizer

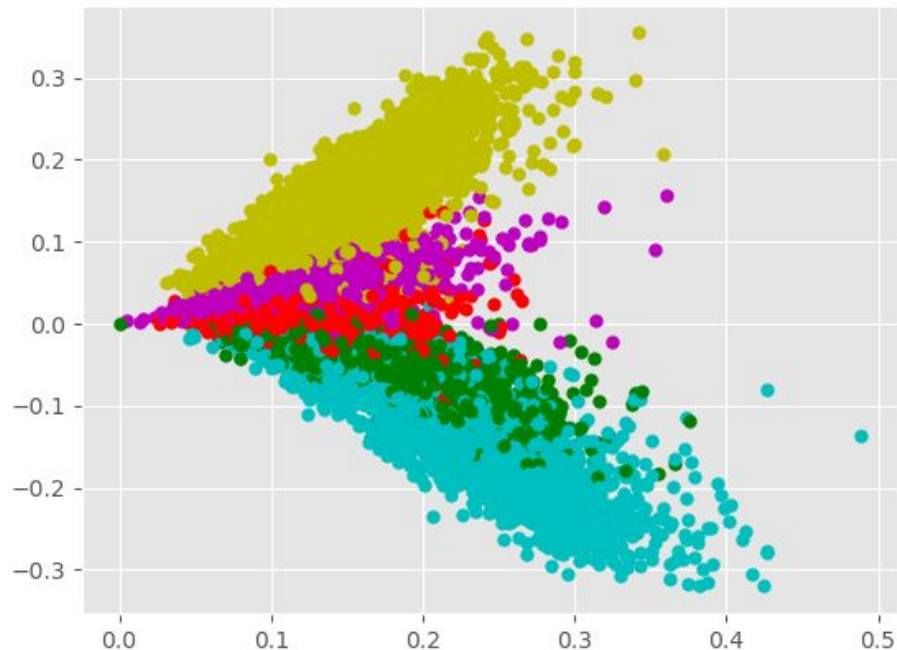


Clustering - Euclidean Distance vs Cosine Similarity

nlTK kmeans, TFIDF Vectorizer,
Euclidean Distance



nlTK kmeans, TFIDF Vectorizer,
Cosine Similarity



Clustering - Σύγκριση αποτελεσμάτων

nlTK kmeans, Cosine Similarity, Preprocessing με TfidfVectorizer και LSI

	Politics	Football	Technology	Film	Business
Cluster1	0.12	0.86	0	0.01	0.01
Cluster2	0.03	0.31	0.01	0.64	0.01
Cluster3	0.97	0.02	0	0	0.01
Cluster4	0	0	0	0	0.99
Cluster5	0.03	0.01	0.89	0.05	0.01

Μετά από μετρήσεις και δοκιμές, καταλήξαμε στο ότι η υλοποίηση με τα καλύτερα αποτελέσματα είναι αυτή, όπως φαίνεται και στον πίνακα παραπάνω (αποτελέσματα από clustering_KMeans.csv)

Clustering

nlTK kmeans, Cosine Similarity, Preprocessing με **CountVectorizer** και LSI

	Politics	Football	Technology	Film	Business
Cluster1	0.05	0.02	0.71	0.21	0.01
Cluster2	0.01	0	0.03	0.13	0.82
Cluster3	0	0	0	0	0.99
Cluster4	0.15	0.65	0.01	0.18	0.01
Cluster5	0.97	0.02	0	0	0.01

Παρατηρούμε ότι η υλοποίηση αυτή δεν είναι καθόλου αποτελεσματική καθώς στο και στο Cluster2 και στο Cluster3 συγκεντρώνονται λέξεις της κατηγορίας Business, ενώ η κατηγορία Film εμφανίζει το μεγαλύτερο ποσοστό της (μόλις 21%) στο Cluster1, το οποίο όμως κατα 71% αποτελείται από λέξεις της κατηγορίας Technology.

Clustering

nlTK kmeans, **Euclidean Distance**, Preprocessing με TfidfVectorizer και LSI

	Politics	Football	Technology	Film	Business
Cluster1	0.1	0.9	0	0	0
Cluster2	0.22	0.32	0.02	0.35	0.09
Cluster3	0.99	0.01	0	0	0
Cluster4	0	0	0	0	1
Cluster5	0.02	0	0.94	0.03	0.01

Η υλοποίηση αυτή δεν είναι η βέλτιστη καθώς η κατηγορία Film εμφανίζεται σε μεγαλύτερα ποσοστά στο Cluster2, αλλά το ποσοστό αυτό είναι πάρα πολύ μικρό (35%) και πολύ κοντά με το ποσοστό των λέξεων της κατηγορίας Football στο ίδιο cluster.

Clustering - Σχόλια / Παρατηρήσεις

- Για το clustering χρησιμοποιήσαμε την υλοποίηση του kmeans που υπάρχει στην βιβλιοθήκη nltk της Python, ώστε να μπορούμε να δώσουμε τη συνάρτηση απόστασης ως όρισμα.
- Για την cosine similarity ως συνάρτηση απόστασης χρησιμοποιήθηκε η cosine distance από το nltk.cluster ($1 - \text{cosine distance} = \text{cosine similarity}$)
- Αυξάνοντας τον αριθμό των επαναλήψεων του kmeans, παρατηρήσαμε ότι δεν πετυχαίνουμε καλύτερα αποτελέσματα, ενώ ο χρόνος εκτέλεσης αυξάνεται σημαντικά. Ενδεικτικά, οι χρόνοι εκτελέσεως ήταν:
Για $n = 500$ $t \sim 20\text{min}$,
Για $n = 100$ $t \sim 5\text{min}$
Για $n = 25$ $t \sim 1\text{min}$

Επομένως επιλέξαμε το $n=25$ για αριθμό επαναλήψεων

- Παρατηρήσαμε τέλος ότι η υλοποίηση του kmeans της βιβλιοθήκης sklearn είχε τα ίδια αποτελέσματα με αυτή της βιβλιοθήκης nltk χρησιμοποιώντας την Euclidean distance ως συνάρτηση απόστασης, αλλά ο χρόνος εκτέλεσης ήταν αρκετά μικρότερος.

Classification

Ζητούμενο:

1. Εύρεση καλύτερου ταξινομητή ύστερα από μετρήσεις
2. Κατηγοροποίηση κειμένων με την χρήση διαφορετικών μεθόδων ταξινόμησης και μετρικών.
3. Βελτίωση ταξινομητή

Περιγραφή εργασίας(1)

Για το classification δοκιμάσαμε τις παρακάτω μεθόδους:

- Support Vector Machines(SVM)
- Random Forests
- Naive Bayes
- K-Nearest Neighbor(δική μας υλοποίηση)

Για να αξιολογήσουμε και να καταγράψουμε την απόδοση κάθε μεθόδου χρησιμοποιήσαμε 10-fold Cross Validation χρησιμοποιώντας τις παρακάτω μετρικές:

- Precision / Recall / F-Measure
- Accuracy
- AUC
- ROC plot

Περιγραφή εργασίας(2)

Κατά την προ-επεξεργασία των δεδομένων χρησιμοποιήσαμε την τεχνική “Latent Semantic Indexing (LSI)” δοκιμάσαμε διαφορετικό αριθμό components κρατώντας σταθερό τον καλύτερο μας ταξινομητή (SVM).

Το γράφημα παρουσιάζεται και παρακάτω στις διαφάνειες αλλά και στο φακελο output με το όνομα LSIcomponentsAccuracy.png το οποίο παράγεται εκτελώντας μια συνάρτηση στο πρόγραμμά μας.

Προσπαθήσαμε να χρησιμοποιήσουμε αποδοτικά τον τίτλο δίνοντας του διαφορετικά βάρη χωρίς τρομερή βελτίωση. Τα διαφορετικά βάρη μπορούμε να τα περάσουμε σαν παράμετρο στις συναρτήσεις μας.

Περιγραφή εργασίας(3)

Για K-Nearest Neighbor δεν χρησιμοποιήσαμε έτοιμη υλοποίηση αλλά τον υλοποιήσαμε μόνοι μας. Χρησιμοποιήσαμε έτοιμη την ευκλίδεια απόσταση ωστόσο για να μετράμε τις αποστάσεις των σημείων. Η επιλογή του τελικού label γίνεται με Majority voting.

Για τον SVM πειραματιστήκαμε με τις παραμέτρους kernel(rbf η default του sklearn, linear), c και gamma.

Αρχεία Εξόδου στον φάκελο output

EvaluationMetric_10fold.csv περιέχει τους μέσους όρους των μετρήσεων για όλες τις μεθόδους ταξινόμησης και όλες τις μετρικές

testSet_categories.csv περιέχει τις προβλέψεις μου για το test dataset

Όλα τα αρχεία category_roc περιέχουν τα διαγράμματα roc για κάθε κατηγορία για κάθε fold.

Cross validation results

Για το cross validation χρησιμοποίησα nfold=10. Τα αποτελέσματα παρακάτω είναι το μέσο των τιμών όλων των folds. Καλύτερος ταξινομητής βγήκε ο SVM με kernel=linear και C=2

Statistic Measure	Naive Bayes	Random Forest	SVM	KNN
Accuracy	0.895	0.945	0.961	0.957
Precision	0.886	0.942	0.957	0.953
Recall	0.888	0.939	0.957	0.954
F-Measure	0.885	0.940	0.957	0.954
AUC				

Beat the Benchmark

Τα αρχικά μας αποτελέσματα ήταν αυτά με καλύτερο ταξινομητή τον KNN:

Statistic Measure	Naive Bayes	Random Forest	SVM	KNN
Accuracy	0.892	0.934	0.935	0.941
Precision	0.879	0.929	0.932	0.936
Recall	0.891	0.928	0.927	0.936
F-Measure	0.882	0.928	0.929	0.936
AUC				

Βελτίωση 1: Αλλαγή ταξινομητή

Δοκιμάζοντας τον SVM με τις παραμέτρους `kernel(rbf,linear)`, `c`, `gamma` παρατήρησα ότι έχει καλύτερα αποτελέσματα από τον KNN.

Υστερα από αρκετές δοκιμές με τις παραμέτρους είδαμε ότι για `linear kernel` και `c=2` έχουμε τις καλύτερες τιμές γενικά, καθώς το `accuracy` μας με `cross validation 10 fold` έβγαινε σχεδόν 0.967

Παρατήρηση : Είδαμε από το clustering ότι τα δεδομένα μας είναι αρκετά γραμμικά διαχωρίσιμα επομένως θεωρούμε πως για αυτό το λόγο ο svm δίνει καλύτερα αποτελέσματα με `linear kernel`.

Βελτίωση Random Forest

Δοκιμάζοντας τον αλγόριθμο Random Forest για διαφορετικό αριθμό estimators, παρατηρήσαμε ότι όσο το n αυξανόταν, τα αποτελέσματα ήταν ολοένα και καλύτερα, για $n=128$ ήταν τα βέλτιστα, ενώ για $n>128$ η απόδοση μειωνόταν και πάλι.

Έτσι βελτιώσαμε την απόδοση του Random Forest, ωστόσο ο SVM παρέμεινε ο καλύτερος classifier.

	n=5	n=50	n=128	n=150
Accuracy	0.934	0.943	0.946	0.944
Precision	0.928	0.938	0.941	0.938
Recall	0.926	0.939	0.941	0.939
F-Measure	0.927	0.938	0.941	0.939

Δοκιμές για βελτίωση άλλων ταξινομητών

Αρχικά δοκιμάσαμε τον KNN με διαφορετικές συναρτήσεις απόστασης από την ευκλείδια (cosine similarity, linear kernel) χωρίς όμως να πετύχουμε καλύτερη απόδοση.

Τέλος προσπαθήσαμε να βελτιώσουμε την απόδοση του Naive Bayes, δοκιμάζοντας τη μέθοδο Bernoulli αντί για την Gaussian, χωρίς όμως να επιτύχουμε καλύτερα αποτελέσματα.

	GaussianNB	BernoulliNB
Accuracy	0.892	0.862
Precision	0.879	0.865
Recall	0.891	0.852
F-Measure	0.882	0.853

Βελτίωση 2: Preprocess LSI με διαφορετικό αριθμό components

Κατά την διάρκεια της εργασίας παρατηρήσαμε ότι η προεπεξεργασία των δεδομένων παίζει πολύ σημαντικό ρόλο στην απόδοση του ταξινομητή μας, καθώς αρχικά χρησιμοποιήσαμε count vectorizer και είχαμε σχετικά κακό accuracy

Στην συνέχεια χρησιμοποιήσαμε tfidf vectorizer με πολύ καλύτερα αποτελέσματα.

Αυτό που έκανε όμως μεγάλη διαφορά στην απόδοση μας είναι όταν κάναμε μετρήσεις στον ταξινομητή μας αλλάζοντας τον αριθμό διάστασης του LSI.

Παρακάτω παρουσιάζονται οι μετρήσεις μας:

Latent Semantic Indexing (LSI) με διαφορετικό αριθμό components.

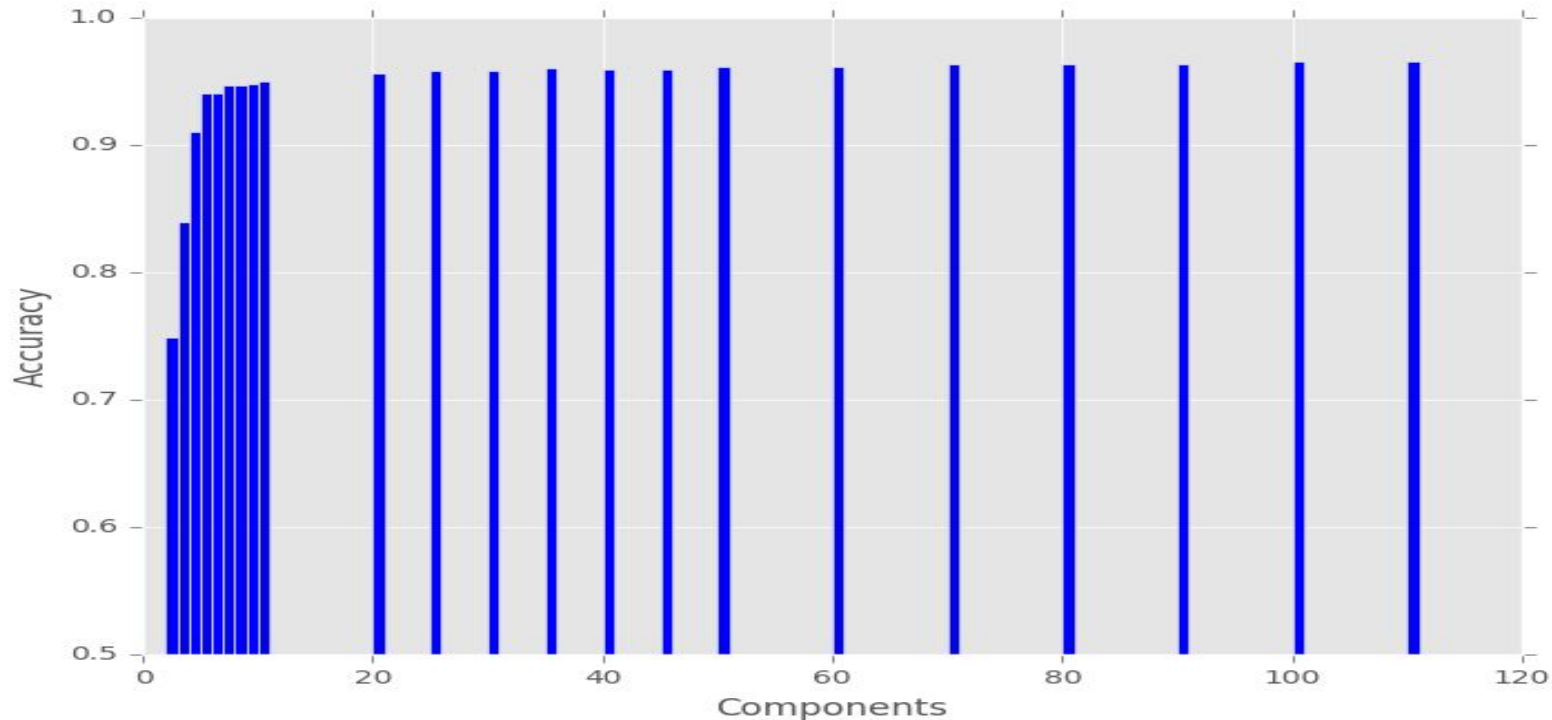
Για να παρατηρήσουμε το πώς επηρεάζονται οι μετρήσεις μας ανάλογα με τον αριθμό components που χρησιμοποιούμε, θα κρατήσουμε σταθερό τον ταξινητή μας και θα μεταβάλλουμε τον αριθμό των components.

Ταξινομητής: Support Vector Machine (SVM) με linear kernel $C=2$

Μετρική : Accuracy

Components: Από 1..120

Latent Semantic Indexing (LSI) με διαφορετικό αριθμό components.



Latent Semantic Indexing (LSI) με διαφορετικό αριθμό components.

Παρατηρούμε πως για μικρές διαστάσεις ο ταξινομητής δίνει πολύ κακά αποτελέσματα, αλλά καθώς αυξάνουμε γίνεται όλο και καλύτερος.

Στις 100 διαστάσεις βλέπουμε ότι σταθεροποιείται γύρω στο 0.965 αλλά συνεχίζει να βελτιώνετε μέχρι και τις 120 διαστάσεις στα 0.9657.

Οι καλύτερες τιμές που πέτυχα ήταν με διαστάσεις 300,400 με απόδοση 0.9676 και 0.9677 αντίστοιχα που συμφωνούν με την wikipedia

https://en.wikipedia.org/wiki/Latent_semantic_analysis#Latent_semantic_indexing

Στην παράγραφο Challenges to LSI

Βελτίωση 3: Αλλαγή βάρους των λέξεων του τίτλου

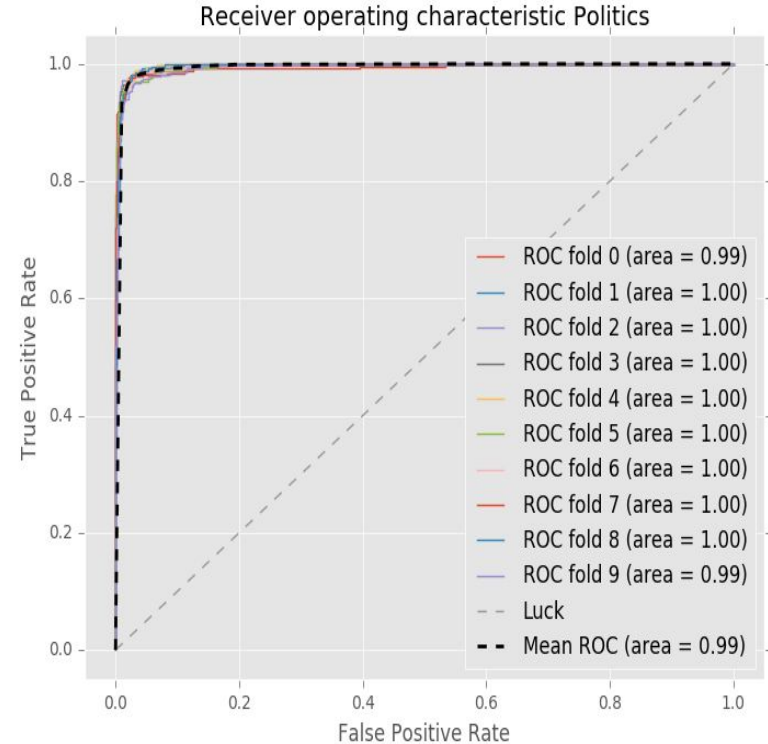
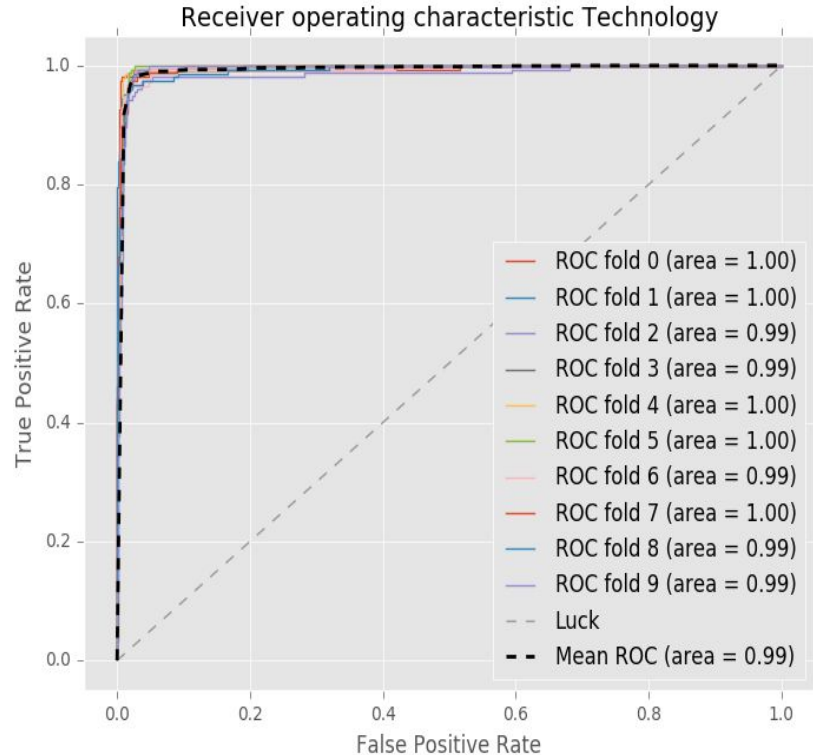
Δοκίμασα διαφορετικά βάρη για τις λέξεις του τίτλου πολλαπλασιάζοντας κάθε φορά με διαφορετικό αριθμό το vector των λέξεων του τίτλου και παρατήρησα τα εξής:

Για μεγάλα βάρη(5,100,1000) έπαιρνα χειρότερα αποτελέσματα και αυτό είναι λογικό καθώς σημαίνει ότι πλέον αγνοούσαμε τις λέξεις του κειμένου.

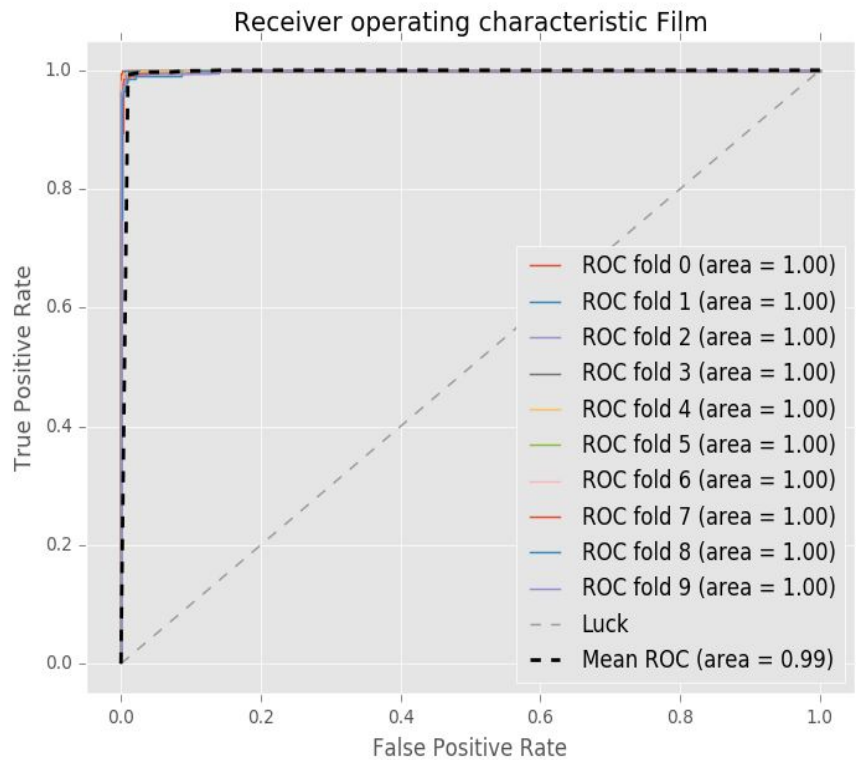
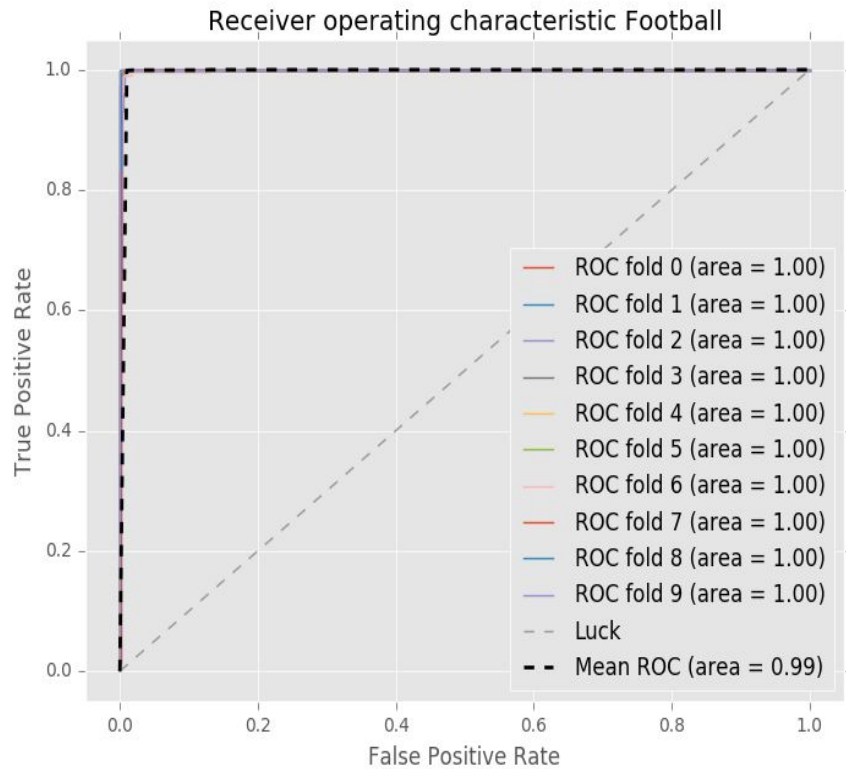
Για μικρά βάρη(2,3) έπαιρνα σχεδόν ίδιες τιμές

Για βάρος 1.1 πήραμε μια πολύ μικρή βελτίωση.

Receiver operating characteristic



Receiver operating characteristic



Receiver operating characteristic

