

Μηχανική Μάθηση

Εργασία 2: Clustering problems

Ομάδα:

Φλάμος Κάρολος dai19200

Βασιλειάδης Παντελεήμων dai19216

15/01/2023

Ομαδοποιητής Πουλιών

Περιεχόμενα	
Εξώφυλλο.....	1
Περιεχόμενα.....	You are here!
Εισαγωγή.....	3
Μέθοδοι που εφαρμόστηκαν.....	4
Συμπεράσματα.....	7

Εισαγωγή

Οι βασικοί στόχοι της εργασίας είναι η ομαδοποίηση(clustering) ενός συνόλου εγγραφών απο το dataset fashion-mnist. Αρχικά πρέπει να βρούμε έναν τρόπο να μεταφέρουμε τα δεδομένα, στον υπολογιστή του collab. Για να δουλέψουμε σωστά τα δεδομένα και να εκπαιδεύσουμε σωστά τα μοντέλα μας, τα χωρίζουμε σε training, validation and test set. Το dataset αυτο περιέχει δεδομένα υψηλότερων διαστάσεων, και γι'αυτο θα προσπαθήσουμε να τις μειώσουμε με την χρήση του Principal component analysis(PCA). Με το PCA θα αναλυθούν οι διαστάσεις σε συνιστώσες, και σε όσα μεγέθη(στήλες) υπάρχει μεγάλη συσχέτιση, θα συγχωνευτούν. Στόχος είναι να μειώσουμε τις διαστάσεις με τρόπο τέτοιο ώστε να μην χάνεται πάνω από 15% της πληροφορίας. Παράλληλα θα πρέπει η συσταδοποίηση να γίνεται επαρκώς καλά. Για την μέτρηση της επάρκειας του PCA καλούμαστε να χρησιμοποιήσουμε κάποιες μετρικές απόδοσης των αλγοριθμων/μοντέλων συσταδοποίησης.

Συγκρίνοντας τις μετρικές πριν και μετά την μείωση διαστάσεων, θα ξέρουμε εάν έχει γίνει σωστά η διαδικασία.

Μέθοδοι που εφαρμόστηκαν

Εισάγαμε το dataset fashion-mnist στον σερβερ μας για να μπορούμε να χειριστούμε τα δεδομένα. Τα χωρίσαμε σε training, validation και test sets. Κρατήσαμε και τις πραγματικές κλασεις των εγγραφών για supervised learning.

Οι διαστάσεις του σετ είναι σε πολύ υψηλό βαθμο, για να επεξεργαστούμε δεδομένα πανω σε αυτον. Θα μειώσουμε τις διαστάσεις με Principal component Analysis. Αναλύουμε σε συνιστώσες τα μεγέθη, και έτσι μπορούμε να δούμε αν καποια απο αυτα είναι θετικά συσχετισμένα. Εάν δούμε οτι συνδιακυμαίνονται, μπορούμε να συγχωνεύσουμε τα μεγέθη σε ένα. Υλοποιήσαμε τον μετασχηματισμό με τις βιβλιοθήκες της rython κρατώντας το 85% της πληροφορίας. Παίρνουμε δειγματοληπτικά μερικές ενδείξεις απο απο τις εγγραφές του training set και όλα φαίνεται να βαίνουν καλώς. Ο μετασχηματισμός εχει γίνει, και απο την ποιότητα των εικόνων φαίνεται να έχει χαθεί λιγη ποιότητα, αλλα είναι ακομη ευδιακριτο το ποια αντικείμενα απεικονίζονται. Τον ιδιο μετασχηματισμό έχουμε κάνει και στο validation set. Απο το δείγμα που πήραμε δίνεται η διαίσθηση οτι όλα βαίνουν καλώς. Έχουμε εφαρμόσει και την αντίστροφη διαδικασία/συνάρτηση στα τροποποιημένα δεδομένα, ωστε να τα επαναφέρουμε στον αρχικό χωρο. Η επιτυχής επαναφορά μας δειχνει οτι δεν είχε χαθεί σημαντικό κομματι της πληροφορίας. Οι εκτυπώσεις των εικόνων το αποδεικνύουν και με το ματι.

Στη συνέχεια εφαρμώσαμε τον ιδιο μετασχηματισμό και στο τεστ σετ. Κάνοντας οπτικοποίηση φαίνεται πως η μετατροπή εχει γίνει και εδω, αλλα με λιγο χειρότερα αποτελέσματα απο πριν.

Μετά φτιάξαμε τα μοντέλα ομαδοποίησης (clustering). Επιλέξαμε τα: kmeans, min batch kmeans και hierarchical clustering, με υλοποιηση απο την sklearn. Εκπαιδεύσαμε το μοντέλο πάνω στο traning set.

Συγκεκριμενα:

min batch Kmeans: χρησιμοποιήσαμε 3 αρχικά clusters(κ) για την εκπαίδευση του μοντέλου.

Kmeans : και πάλι το ίδιο(κ=3)

hierarchical : του επιτρέψαμε να φτάσει στο μέγιστο επίπεδο διαχωρισμού

Στην συνέχεια προσπαθήσαμε να επιλέξουμε τις κατάλληλες παραμέτρους για την εκκίνηση των αλγορίθμων. Η επιλογή αυτή έγινε με εκπαίδευση πάνω στο validation set. Πήραμε διαδοχικές τιμές εκκίνησης, και δοκιμάσαμε μέχρι να δωθεί το καλύτερο σκορ(τα σκορ αναφαίρονται αναλυτικά στο τέλος) που μπορούσαμε να βρούμε. Τελος εκπαίδευσης του Μοντέλου. Τέλος δοκιμάσαμε τα μοντέλα πάνω στο τεστ σετ. Τα αποτελέσματα που λάβαμε μετρήθηκαν ως προς την ποιότητα με τις μετρικές που φαίνονται παρακάτω:

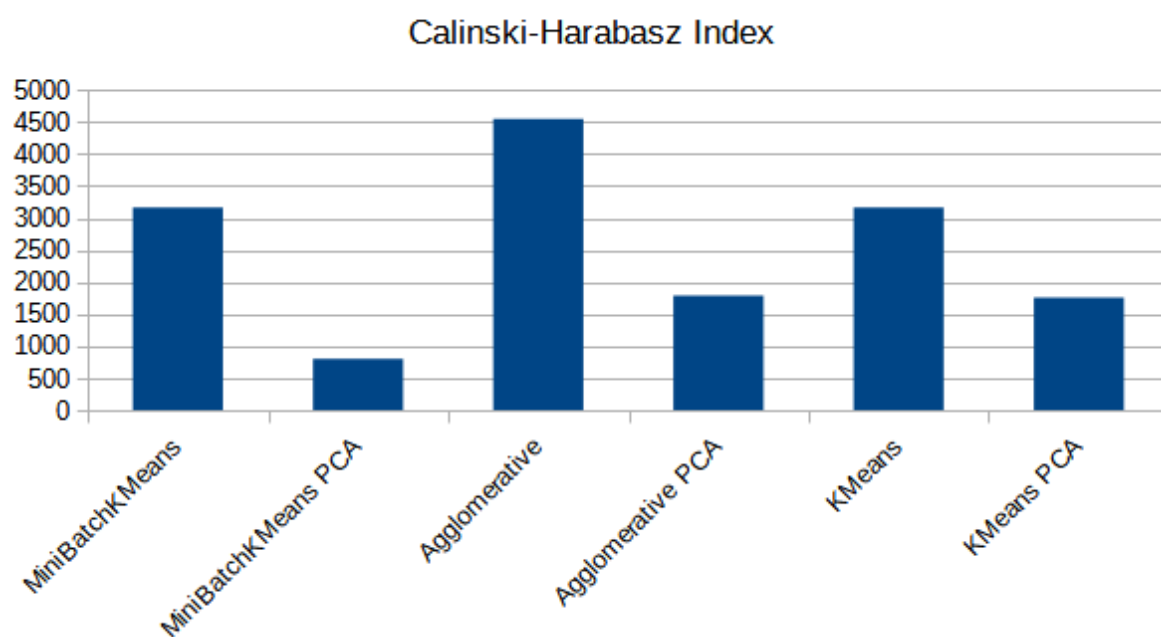
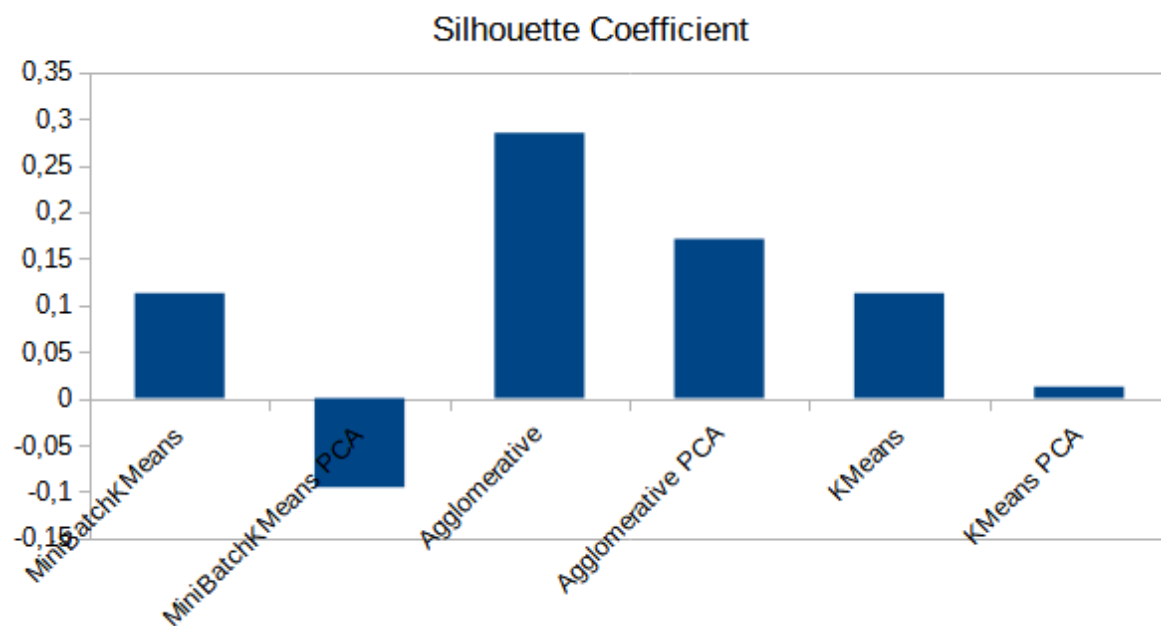
-Silhouette Coefficient: Δείχνει αν τα σημεία είναι κοντά με τα στοιχεία του cluster τους και ταυτόχρονα μακριά απο τα αλλα

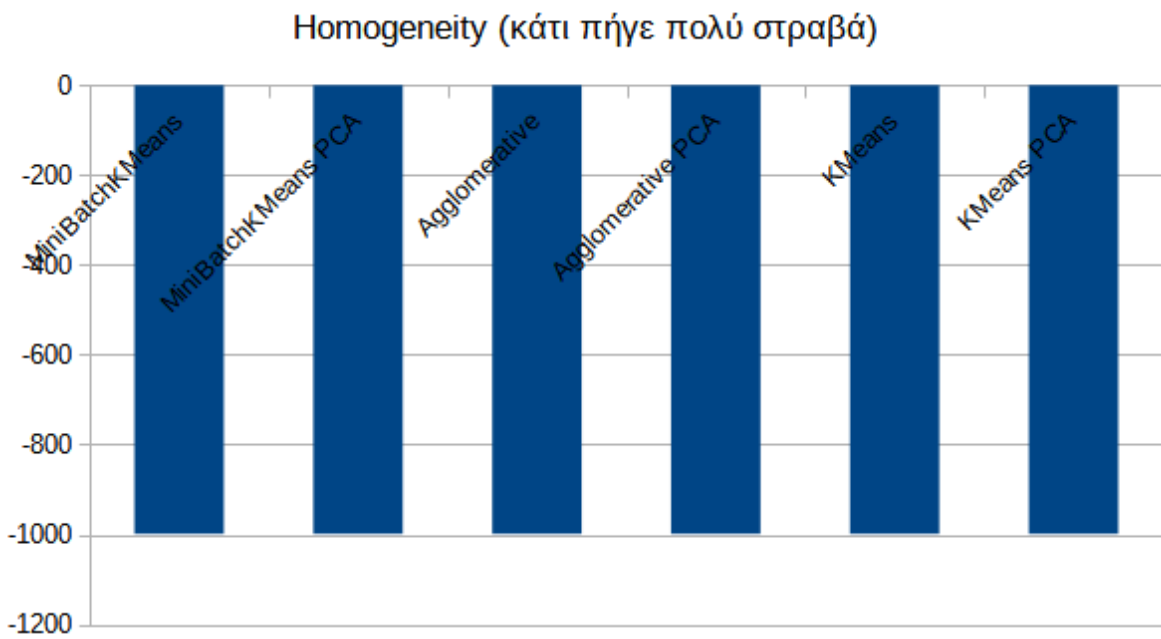
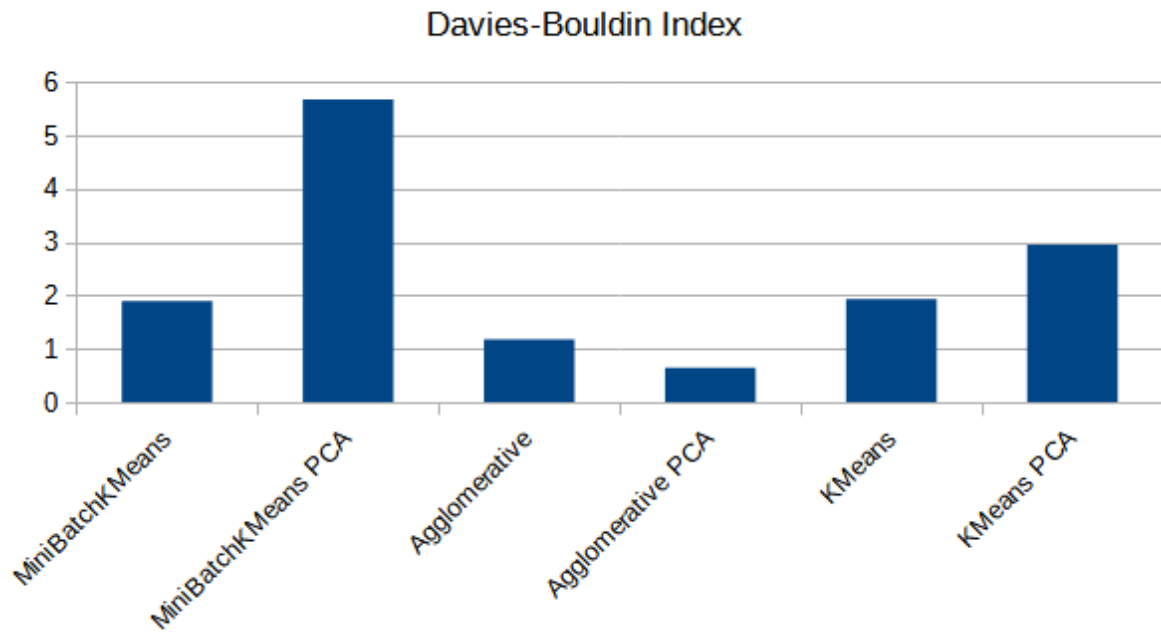
-Calinski-Harabasz Index : Δίχνει το πόσο καλά εχουν διαχωριστεί τα clusters με βάση τις ιδιότητες των συνλικών εγγραφών. Δεδομένου οτι χρησιμοποιήσαμε unsupervised learning, θεωρουμε οτι είναι πολύ βοηθητικό κριτήριο

-Davies-Bouldin Index: άλλη μία μέτρηση βασισμένη σε εσωτερικά χαρακτηριστικά.

Ιδανική για Unsupervised learning.

Παρατίθενται οι εικόνες των γραφημάτων για τις αποδόσεις των μετρικών.





Συμπεράσματα-αποτιμήσεις

Ως προς το Silhouette Coefficient, τα αποτελέσματα είναι μακράν καλύτερα πριν την μετατροπή PCA. Το ιεραρχικό φαίνεται να αντιστέκεται περισσότερο στην πτώση αποδόσεων μετά το PCA μετατροπή. Γενικά τα silhouette scores είναι πολύ κακά, με εξαίρεση τον agglomerative.

Το Calinski-Harabasz Index επιβεβαιώνει αυτά τα αποτελέσματα. Είναι φανερό ότι τα αποτελέσματα είναι πεσμένα μετά το PCA και αυτό σημαίνει πεσμένη απόδοση.

Στο Davies-Bouldin τα νούμερα είναι πιο ανεβασμένα μετά το PCA, ωστόσο για την συγκεκριμένη μετρική αυτό δεν συνεπάγεται καλύτερες αποδόσεις.

Ο δείκτης που χρησιμοποιήσαμε μόνοι μας ήταν της ομοιογένειας, αλλά δεν πήγε καλά, οπότε τον αποκλείουμε από την εξαγωγή συμπερασμάτων ούτε μιλάμε για αυτόν.

Ο καλύτερος αλγόριθμος είναι ξεκαθαρά του μοντέλου agglomerative clustering. Δεν παίρνουμε τις ίδιες μετρικές. Θα λεγάμε ότι οι αποδόσεις είναι πεσμένες.

Ο καλύτερος συνδιασμός είναι agglomerative με ανεπεξέργαστα δεδομένα.

Οι μετρικές του είναι:

```
.. Silhouette Coefficient score is 0.28
.. Calinski-Harabasz Index score is 4552.63
.. Davies-Bouldin Index score is 1.18
```

Τα αποτελέσματα θεωρούνται αποδεκτά, αφού δηλώνουν ότι τα cluster είναι διαχωρισμένα αρκετά. Δεν είναι καλά διαχωρισμένα αλλά είναι αποδεκτά. Μπορούμε να είμαστε αρκετά σίγουροι ότι το κάθε σημείο ανήκει (στατιστικά) στο cluster του. Επίσης δεν επικαλύπτονται τα clusters. Θα μπορούσαμε ενδεχομένως να κρατήσουμε μεγαλύτερο ποσοστό πληροφορίας κατά τον PCA. Έτσι ίσως να ήταν αποδεκτά τα αποτελέσματα και στους άλλους αλγορίθμους.