

## ΜΕΛΗ

| Όνομα   | Επώνυμο     | ΑΜ            |
|---------|-------------|---------------|
| Άννα    | Πετρίδου    | 1115201600135 |
| Μιχάλης | Ανδρουλάκης | 1115201600004 |

Το αρχείο μας αποτελείται από ένα φακέλο project3 και από 4 υποφακέλους, τους Α, Β, Γ, Δ όπου έχουν τα αντίστοιχα κομμάτια μέσα τους και τον υποφάκελο Functions\_for\_both\_algorithms που περιέχει κώδικα της πρώτης εργασίας και χρησιμοποιείται από τα Α, Β και Δ κομμάτια.

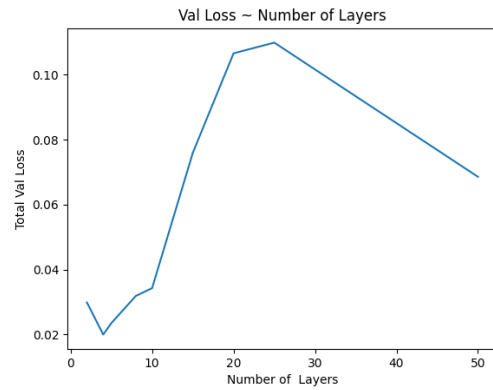
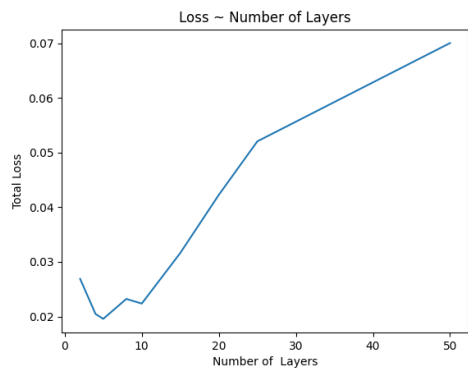
## Α' ΜΕΡΟΣ

Το Α μέρος το υλοποιήσαμε σε python. Χρησιμοποιώντας το πρότυπο του autoencoder της προηγούμενης εργασίας, φτιάξαμε έναν νέο autoencoder που παίρνει τις εξής υπερπαραμετρές: number of filters, number of hidden layers, filter size dimension d(τετραγωνή διάσταση όπως στην προηγούμενη εργασία dxd), number of epochs, batch size, learning rate και latent dimension με default=10. Προσθέσαμε Flatten layer στον encoder και FC και Reshape layer στον decoder. Έπειτα κάναμε μια σειρά πειραμάτων για να επιλέξουμε τις ιδανικές υπερπαραμετρές. Τα αποτελέσματα των πειραμάτων φαίνονται παρακάτω.

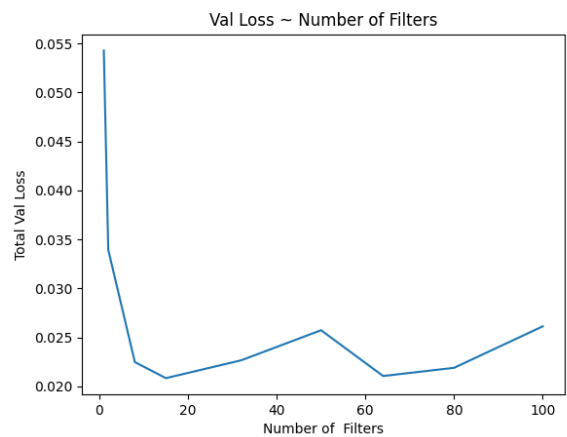
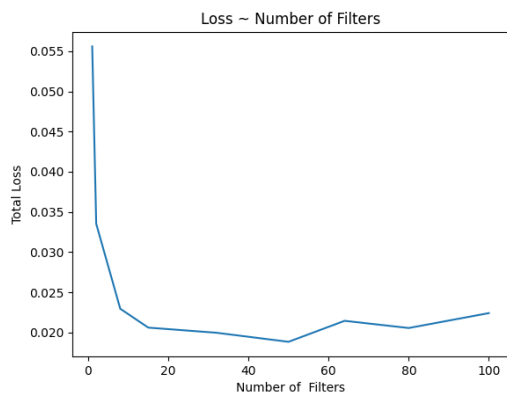
Οι μετρήσεις έγιναν με τις εξής σταθερές παραμετρές, κρατώντας τις σταθερές και αλλάζοντας μια κάθε φορά:

filters:32, layers:5, filtersize:3 ,epochs:10, batch: 100, lerning rate: 0.001, latent dimensions: 10

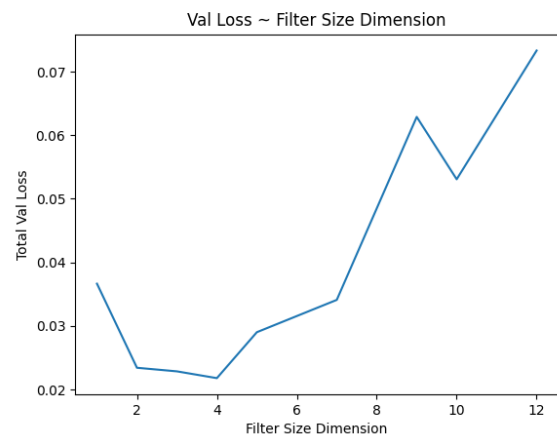
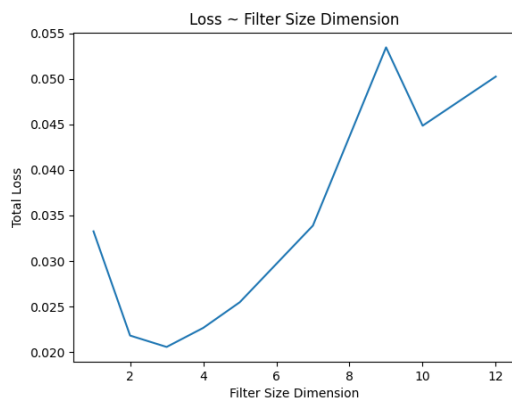
Στα



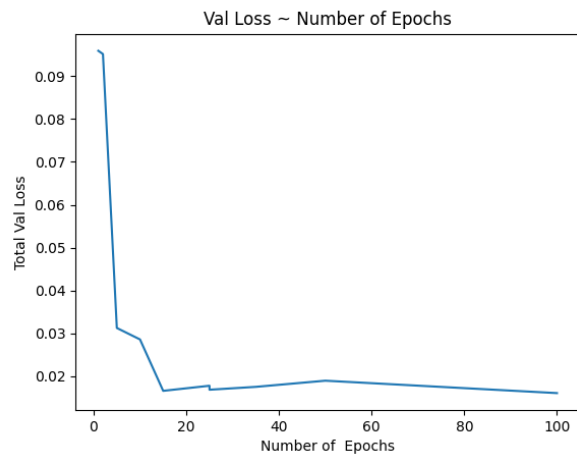
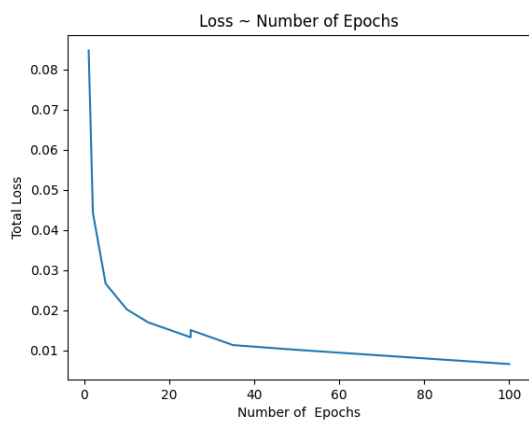
layers παρατηρούμε πως η ιδανική τιμή είναι γύρω στο 10, καθώς μετά αυξάνονται και το val loss και το loss.



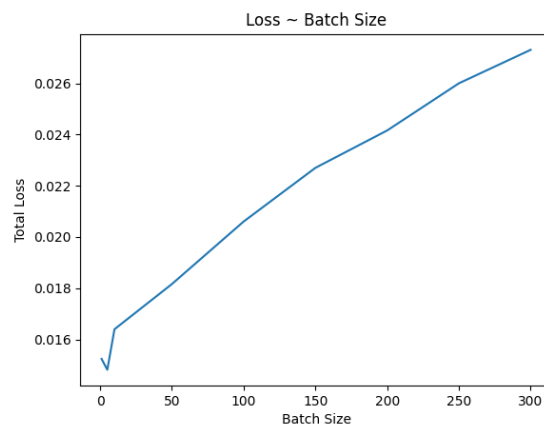
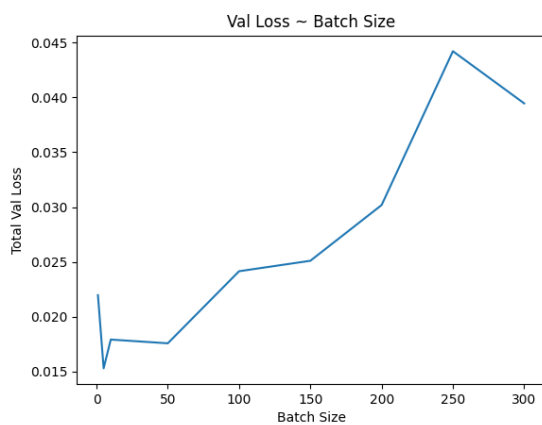
Για τα φίλτρα παρατηρούμε πως το val loss πιάνει μια κατώτατη τιμή γύρω στο 15.



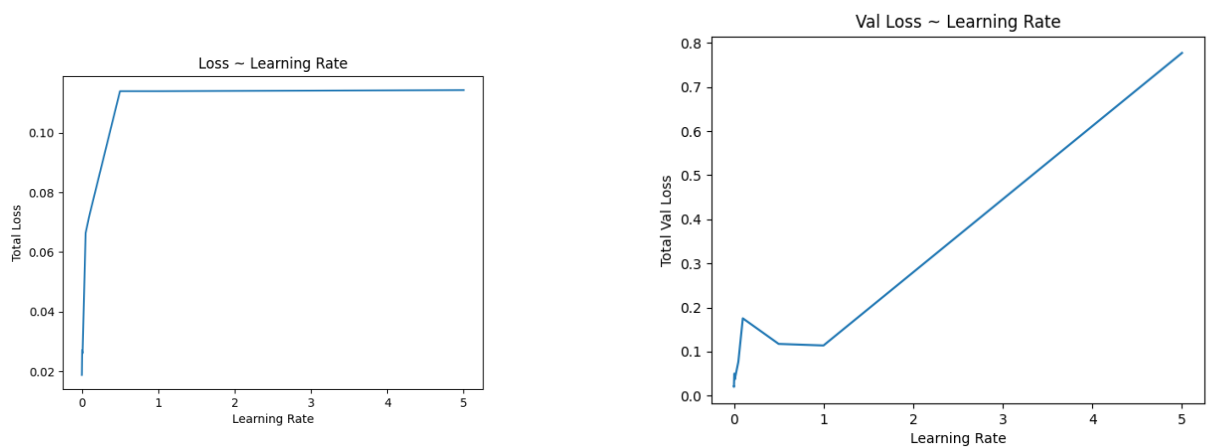
Για το μέγεθος των φίλτρων έχουμε ως βέλτιστη τιμή το 4 όπως παρατηρούμε και απ το val loss και κοντά και το loss.



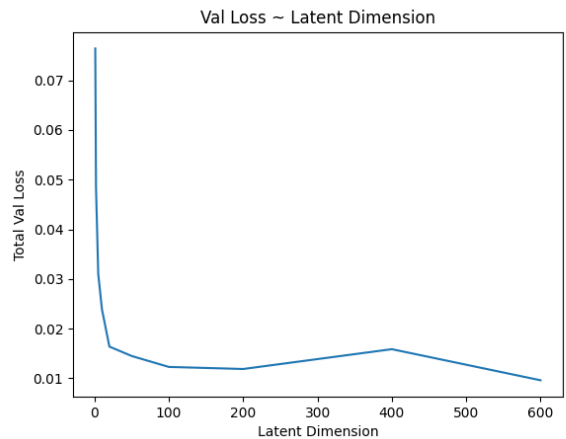
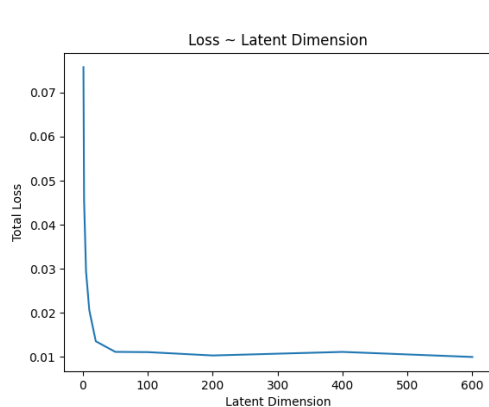
Για τις εποχές, παρατηρούμε ότι όσο περισσότερες είναι οι εποχές τόσο λιγότερο το loss και val loss. Οπότε θα επιλέγαμε μια αρκετά μεγάλη τιμή.



Για τα batches, παρατηρούμε πως το Loss πιανει και στις δυο περιπτώσεις μια κατώτατη τιμή γύρω στο 1-2.



Για το learning rate παρατηρούμε πως η βέλτιστη τιμή είναι πολύ κοντά στο 0, δηλαδή το 0,0001-0,001.



Για τον αριθμο των νεων διαστασεων, παρατηρουμε πως όσο αυξανεται, τόσο μειωνεται το loss, οποτε θα μπορούσαμε να επιλεξουμε εναν αριθμο οπως το 600 αν και η χρηση του latent δεν εξυπηρετει το σκοπο αυτο.

Επισης τα δεδομενα χωριστηκαν οπως και στη προηγουμενη εργασία σε training set και validation set. Στη συνεχεια, αφου ειχαμε το μοντελο μας, κανονικοποιησαμε τις εικονες στο (0,1) και τις μετατρεψαμε απο τις αρχικες διαστασεις στις νεες (πχ 1x10). Επειτα γραψαμε το νεο αρχειο που δημιουργησαμε που ειναι 2 bytes.

Δημιουργησαμε το reduce.py η οποια καλωντας τις καταλληλες συναρτησεις συμπιεζει καταλληλα το χωρο στις διαστασεις που

δοθηκαν (Latent) και δημιουργει το νεο αρχειο οπως εξηγησαμε παραπανω.

### **Εκτελεση προγραμματος με την εντολη:**

Για το reduce.py:

```
$python3.8 reduce.py -d t10k-images.idx3-ubyte -od t10k-images.idx3-ubyte -q train-images.idx3-ubyte -oq train-images.idx3-ubyte
```

Για τον autoencoder\_v3.py που δημιουργει το μονετλο:

```
$python3.8 autoencoder_v3.py -d t10k-images.idx3-ubyte
```

### **Τα περιεχομενα του φακελου ειναι:**

**autoencoder v3.py:** δεχεται τις υπερπαραμετρους και δημιουργει το ζητουμενο μοντελο

**autoencoding.py:** δημιουργια μοντελου

**encoder\_decoder:** υλοποιηση του encoder και decoder

**graphs.py:** για την εκτυπωση των γραφικων

**inputOutputMNIST:** για το διαβασμα και γραψιμο αρχειων MNIST

**reading\_and\_converting.py:** διαβασμα ορισματος

**reduce.py:** διαβαζει το αρχικο αρχειο κ χρησιμοποιωντας το μοντελο φτιαχνει το νεο αρχειο ανα 2 byte στο συμπιεσμενο χωρο

## Β' ΜΕΡΟΣ

Το μέρος αυτό έγινε σε c++ κωδικά. Για το Β μέρος χρησιμοποιήσαμε τα κομμάτια lsh και true από τη 1<sup>η</sup> εργασία με μικροαλλαγές και έπειτα υλοποιήσαμε τον Reduced χώρο.

Από το Α έπειτα από μετρήσεις πήραμε τα 2 καλύτερα αρχεία εικόνων, ένα train και ένα test, τα οποία τα έχουμε επισυναψεί στον Β φάκελο και είναι τα train-images.idx3-ubyte και test-images.idx3-ubyte, τα οποία περιέχουν εικόνες μεγέθους 10 διαστάσεων όπου η πληροφορία για το κάθε pixel προκύπτει από το διάστημα 2 bytes από το αρχείο και συνένωση τους. Βρήκαμε εξαντλητικά το καλύτερο γείτονα όλων των queries στον Reduced χώρο, και έπειτα υπολογίσαμε τη απόσταση του γείτονα που βρήκαμε με το query στον αρχικό χώρο.

Για τα ενδεικτικά αποτελέσματα μας διατηρήσαμε όλο το input file (60.000 εικόνες) και όλο το query file (10.000 εικόνες) και βρήκαμε:

Για k=4 και L=5:

tReduced: 327.008 secs

tLSH: 1439.16 secs

tTrue: 5849.43 secs

Approximation Factor LSH: 1.12371

Approximation Factor Reduced: 1.23968

Για k=4 και L=6:

tReduced: 328.551 secs

tLSH: 1812.94 secs

tTrue: 5447.67 secs

Approximation Factor LSH: 1.11478

Approximation Factor Reduced: 1.23968

Για k=5 και L=5:

tReduced: 315.035 secs

tLSH: 1352.87 secs

tTrue: 7623.44 secs



Approximation Factor LSH: 1.14763  
Approximation Factor Reduced: 1.23968

Για k=5 και L=6:

tReduced: 317.238 secs  
tLSH: 1860.103 secs  
tTrue: 5335.53 secs  
Approximation Factor LSH: 1.12804  
Approximation Factor Reduced: 1.23968

**Τα περιεχόμενα του φακέλα είναι:**

**Hash.cpp - Hash.h** : όπου βρίσκεται η δομή του Lsh

**Lsh.cpp - Lsh.h** : όπου βρίσκονται οι συναρτήσεις για το hashing του train file καθώς και τον υπολογισμό των Lsh αλλά και των πραγματικών γειτονών των queries

**Nearest\_Neighbor.cpp - Nearest\_Neighbor.h** : όπου βρίσκεται η συνάρτηση για τον υπολογισμό του πραγματικού κοντινότερου γείτονα ενός query

**Reading.cpp - Reading.h** : όπου βρίσκονται οι συναρτήσεις , για το διαβάσμα των ορισμάτων της γραμμής εντολών , καθώς και το διαβάσμα αρχικού και νέο χώρου αρχείων

**Reduced\_Space.cpp - Reduced\_Space.h** : όπου βρίσκεται η συνάρτηση για την εύρεση του πραγματικά πλησιέστερου γείτονα στον νέο χώρο αλλά και μια συνάρτηση για τον υπολογισμό των average σωστών εικόνων που δίνει ο Reduced χώρος για τις συγκρίσεις με το Γ

**Search.cpp** : όπου εδώ υπάρχει η main και καλούνται όλες οι υπολοιπές συναρτήσεις

**t10k-images.idx3-ubyte** : όπου είναι το test αρχείο στις 10 διαστάσεις που προέκυψε από το A

**train-images.idx3-ubyte** : όπου είναι το train αρχείο στις 10 διαστάσεις που προέκυψε από το A

**makefile** : για να κανουμε compile το κωδικα μας

**Εκτελεση του προγραμματος με την εντολη:** `./search -d ../train-images.idx3-ubyte -q ../t10k-images.idx3-ubyte -i train-images.idx3-ubyte -s t10k-images.idx3-ubyte -o output.txt -k 5 -L 5`

### **Γ' ΜΕΡΟΣ**

Το μερος αυτο εγινε σε python κωδικα.Υλοποιησαμε τη μετρικη Manhattan που απλα μετατρεψαμε τη συναρτηση της c++ σε python αλλα και την EMD χρησιμοποιοντας το linprog της scipy βιβλιοθηκης.Για την EMD,επειτα απο κανονικοποιηση που καναμε στα βαρη των εικονων τις αποστασεις μεταξυ του καθε cluster input και query εικονας και επειτα φτιαχνουμε την linprog δηλαδη την objective function αλλα και τις εξισωσεις γεμιζοντας τα καταλληλα vectors.

Τα αποτελεσματα υστερα απο πειραματα που πηραμε ειναι :

#### **Για clusters 2x2**

Για 10 queries και 100 input:

Time\_EMD : 8.61 secs και Average Correct EMD: 0.279999

Time\_Manhattan : 7.55 secs και Average Correct Manhattan: 0.56

Για 10 queries και 1000 input:

Time\_EMD : 78.76 secs και Average Correct EMD: 0.4

Time\_Manhattan : 38.87 secs και Average Correct Manhattan: 0.73

Για 10 queries και 10000 input:

Time\_EMD : 392.59 seconds και Average Correct EMD: 0.3299

Time\_Manhattan : 341.52 και Average Correct Manhattan : 0.89

Για 100 queries και 100 input:

Time\_EMD: 97.20 και Average Correct EMD: 0.23

Time\_Manhattan : 74.40 και Average Correct Manhattan: 0.48

Για 100 queries και 1000 input:

Time\_EMD : 413.96 secs και Average Correct EMD:0.314

Time\_Manhattan : 349.12 secs και Average Correct Manhattan: 0.75

Για 100 queries και 10000 input:

Time\_EMD : 4073.42 secs και Average Correct EMD:0.354

Time\_Manhattan :3407.39 secs και Average Correct Manhattan:0.907

#### **Για clusters 4x4:**

Για 10 queries και 100 input:

Time\_EMD : 48.60 secs και Average Correct EMD: 0.41

Time\_Manhattan : 3.55 secs και Average Correct Manhattan: 0.56

Για 10 queries και 1000 input:

Time\_EMD : 825.48 secs και Average Correct EMD: 0.59

Time\_Manhattan : 36.27 secs και Average Correct Manhattan: 0.73

Για 10 queries και 10000 input:

Time\_EMD : 5754.47 secs και Average Correct EMD: 0.68

Time\_Manhattan : 340.90 secs και Average Correct Manhattan: 0.89

Για 100 queries και 100 input:

Time\_EMD : 782.27 secs και Average Correct EMD: 0.37

Time\_Manhattan : 36.96 secs και Average Correct Manhattan: 0.48

Για 100 queries και 1000 input:

Time\_EMD : 6694.10 secs και Average Correct EMD: 0.602

Time\_Manhattan : 353.80 secs και Average Correct Manhattan: 0.75

#### **Για clusters 7x7:**

Για queries 10 και input 100:

Time\_EMD : 5329.80 secs και Average Correct EMD: 0.49

Time\_Manhattan : 3.65 secs και Average Correct Manhattan 0.56

Για queries 100 και input 100:

Time\_EMD : 44309.18 secs και Average Correct EMD: 0.47

Time\_Manhattan : 35.65 secs και Average Correct Manhattan : 0.48

Για τα clusters 7x7 και 4x4 δε μπορεσαμε να κανουμε μεγαλυτερες μετρησεις καθως μας παιρνανε πανω απο μια μερα.

Εκτος απο αυτες τις μετρησεις εχουμε και μετρησεις συγκρισης μεταξυ του B κομματιου του Reduced χωρου.Για τις συγκρισεις χρησιμοποιησαμε 1 πλησιεστερο γειτονα και στα 2 κομματια και μετρησαμε τους χρονους αλλα και τα αποτελεσματα. Για να κανουμε τις συγκρισεις αυτες φτιαξαμε μια συναρτηση στο B που διαβαζει τα labels και για τα πρωτα 10 queries υπολογιζει το ποσοστο των σωστων labels που βρηκε η reduced μετρικη για N=1. Αναλογως για το Γ παλι κανοντας μετρησεις για N=1 βρηκαμε το ποσοστο των 10 πρωτων queries.Χρησιμοποιηθηκαν οι πρωτες 10000 εικονες του dataset.

Time\_Reduced=317.65 seconds και average\_reduced=0.9 δηλαδη 9/10 queries ο πλησιεστερος γειτονας του reduced χωρου εχει το ιδιο label με το query

Για clusters=4x4 time\_EMD=5554.47 seconds και average\_EMD=0.6

Για clusters=2x2 time\_EMD=392.57 seconds και average\_EMD=0.3299

**Τα περιεχομενα του φακελου ειναι:**

**EMD.py** : Που περιεχει την συναρτηση EMD\_Distance για τον υπολογισμο της EMD αποστασης μεταξυ 2 εικονων

**Manhattan.py** : Που περιεχει τη συναρτηση Manhattan\_Distance για τον υπολογισμο της Manhattan αποστασης 2 εικονων

**Quicksort.py** : Που περιεχει τις συναρτησεις quickSort και partition για για το sort και την ευρεση των 10 πλησιεστερων γειτονων

**reading\_and\_converting.py** : Που περιεχει τις συναρτησεις readingArguments για διαβασμα της γραμμης εντολων και file\_normalization ωστε να κανονικοποιησουμε τα αρχεια

**Writing.py** : Που περιεχει τη συναρτηση για το γραψιμο του αποτελεσματος σε ενα αρχειο

**search.py** : Που περιεχει τη main

**Εκτελεση προγραμματος με την εντολη:** python3 search.py -d ../train-images.idx3-ubyte -q ../t10k-images.idx3-ubyte -l1 ../train-labels.idx1-ubyte -l2 ../t10k-labels.idx1-ubyte -o output.txt

## Δ' ΜΕΡΟΣ

Για το Δ μερος καναμε αρχικα το Σ1, δηλαδη το clustering για τις εικονες στο νεο χωρο, διαβαζοντας το αρχειο που δημιουργησαμε στην Α ερωτηση, καταλληλα ανα 2 bytes. Αφου ολοκληρωσαμε το clustering, καναμε την αντιστοιχια των clusters με τα σημεια που προεκυψαν, στον αρχικο μας χωρο, υπολογισαμε τα νεα centroids και ετσι βρηκαμε και εκτυπωσαμε το silhouette. Επειτα καναμε το Σ2, το οποιο ειναι το clustering που ειχαμε και στη 1η εργασια, στον αρχικο χωρο και βρηκαμε και παλι το silhouette.

Τελος, παιρνοντας το αρχειο απ την κατηγοριοποιηση στη 2η εργασια, καναμε και παλι την αντιστοιχια και βρηκαμε τα αντιχτοιχα σημεια των clusters στον αρχικο χωρο. Πηραμε τις συντεταγμενες τους, υπολογισαμε τα νεα centroids και καναμε τελος silhouette.

Ενδεικτικα αποτελεσματα μετα απο μια εκτελεση στα 10000 εχουμε:

## NEW SPACE

CLUSTER-0 {size: 1339, centroid:...

CLUSTER-1 {size: 642, centroid:...

CLUSTER-2 {size: 1073, centroid:...

CLUSTER-3 {size: 1060, centroid:..

CLUSTER-4 {size: 1008, centroid:...

CLUSTER-5 {size: 1408, centroid:...

CLUSTER-6 {size: 2324, centroid:...

CLUSTER-7 {size: 393, centroid:...

CLUSTER-8 {size: 221, centroid:...

CLUSTER-9 {size: 532, centroid:...

clustering\_time: 1 sec

Silhouette:

[0.00469718,0.00214029,0.019,.....,0,04157820]

## ORIGINAL SPACE

CLUSTER-0 {size: 7551, centroid:...

CLUSTER-1 {size: 143, centroid:...

CLUSTER-2 {size: 572, centroid:...

CLUSTER-3 {size: 11, centroid:..

CLUSTER-4 {size: 91, centroid:...

CLUSTER-5 {size: 231, centroid:...

CLUSTER-6 {size: 4, centroid:...

CLUSTER-7 {size: 501, centroid:...

CLUSTER-8 {size: 286, centroid:...

CLUSTER-9 {size: 610, centroid:...

clustering\_time: 26 sec

Silhouette:

[0.129548,0.143949,...,0.144479]

**Εκτέλεση προγράμματος με την εντολή:**

```
$ ./cluster -i t10k-images.idx3-ubyte -d t10k-images.idx3-ubyte -c  
config.txt -o output.txt -n class.txt
```

**Τα περιεχόμενα του φακέλου είναι ίδια με της προηγούμενης  
εργασίας. Ο,τι προστεθηκε, ήταν στα ήδη υπάρχοντα αρχεία.**

**Συνδεμος Github:<https://github.com/Omada53/project3>**