

# Ανάπτυξη Λογισμικού για Αλγοριθμικά Προβλήματα

## Εργασία 3

Δημητρακόπουλος Κωνσταντίνος

ΑΜ: 1115201500034

Κακαβάς Αντώνιος

ΑΜ: 1115201500052

### Ανάπτυξη

Η εργασία αναπτύχθηκε στο περιβάλλον Google Collaboratory με τύπο αρχείου python notebook, όπου παραδίδεται. Επιπρόσθετα παραδίδεται ο πηγαίος κώδικας της εργασίας σε αρχεία .py σύμφωνα με τις προϋποθέσεις της εκφώνησης. Η συνεργασία κατά την ανάπτυξη πραγματοποιήθηκε στο [private repository](#).

### Εκτέλεση

Τα παραδοτέα αρχεία εκτελούνται ως εξής:

- `python3 forecast.py -d <dataset path> -n <number of time series selected> -t <"offline_all" or "online_all" or "online_self">`
- `python3 detect.py -d <dataset path> -n <number of time series selected> -t <"offline_all" or "online_all"> -mae < error value as double >`
- `python3 reduce.py -d <dataset> -q <queryset> -od <output_dataset_file> -oq <output_query_file> -t <"offline_all" or "online_all" >`

Σχετικά με την παράμετρο '-t':

- "offline\_all": Φόρτωση του offline εκπαιδευμένου ανά σύνολο χρονοσειρών μοντέλου.
- "online\_all": Εκπαίδευση ανά σύνολο χρονοσειρών κατά την εκτέλεση.
- "online\_self": Εκπαίδευση ανά χρονοσειρά κατά την εκτέλεση.

Η διαρρύθμιση των καταλόγων της εργασίας δεν πρέπει να μεταβληθεί για επιτυχή εκτέλεση της.

Κατά την εκτέλεση των αρχείων .py και όχι του αρχείου .ipynb εμφανίζονται warnings σχετικά με την μελλοντική μη υποστήριξη ορισμένων συναρτήσεων από τη πλατφόρμα tensorflow και μπορούν να αγνοηθούν, αφού δεν αφορούν την τωρινή εκτέλεση της εργασίας.

## Κατάλογος αρχείων

Η εργασία που παραδίδεται αποτελείται από τα εξής αρχεία:

- `split_dataset:`  
Δέχεται ως όρισμα το μονοπάτι του dataset της εργασίας και το διαμελίζει στα σετ δεδομένων `dataset.csv` και `queryset.csv` τοποθετώντας αυτά στο ίδιο μονοπάτι.
- `common_utils:`  
Περιέχει κοινές βοηθητικές συναρτήσεις για όλες τις υποενότητες της εργασίας.
- `forecast.py:`  
Το αντίστοιχο αρχείο που περιγράφει η εκφώνηση. Πραγματοποιείται έλεγχος των ορισμάτων, φόρτωση του σετ δεδομένων και για διαφορετικές τιμές του ορίσματος “-t” πραγματοποιείται είτε φόρτωση του offline εκπαιδευμένου ανά σύνολο χρονοσειρών μοντέλου, είτε εκπαίδευση ανά σύνολο χρονοσειρών κατά την εκτέλεση με αποθήκευση του παραγόμενου μοντέλου, είτε εκπαίδευση ανά χρονοσειρά κατά την εκτέλεση. Η ανάλυση του μοντέλου πραγματοποιείται παρακάτω. Τελικά πραγματοποιείται πρόγνωση τιμών χρονοσειρών.
- `detect.py:`  
Το αντίστοιχο αρχείο που περιγράφει η εκφώνηση. Πραγματοποιείται έλεγχος των ορισμάτων, φόρτωση του σετ δεδομένων και για διαφορετικές τιμές του ορίσματος “-t” πραγματοποιείται είτε φόρτωση του offline εκπαιδευμένου ανά σύνολο χρονοσειρών μοντέλου, είτε εκπαίδευση ανά σύνολο χρονοσειρών κατά την εκτέλεση με αποθήκευση του παραγόμενου μοντέλου. Η ανάλυση του μοντέλου πραγματοποιείται παρακάτω. Τελικά πραγματοποιείται ανίχνευση ανωμαλιών χρονοσειρών.
- `reduce.py:`  
Το αντίστοιχο αρχείο που περιγράφει η εκφώνηση. Πραγματοποιείται έλεγχος των ορισμάτων, φόρτωση των σετ δεδομένων και για διαφορετικές τιμές του ορίσματος “-t” πραγματοποιείται είτε φόρτωση του offline εκπαιδευμένου ανά σύνολο χρονοσειρών μοντέλου, είτε εκπαίδευση ανά σύνολο χρονοσειρών κατά την εκτέλεση με αποθήκευση του παραγόμενου μοντέλου και εμφάνιση των παραγόμενων χρονοσειρών (αρχική, κωδικοποιημένη, αυτοκωδικοποιημένη) κατά την εκπαίδευση σε γραφήματα. Η ανάλυση του μοντέλου πραγματοποιείται παρακάτω. Τελικά παράγονται οι συμπιεσμένες-κωδικοποιημένες χρονοσειρές και εξάγονται κατάλληλα στα αρχεία εξόδου.
- `project-algorithms.ipynb:`  
Το notebook στο οποίο αναπτύχθηκε η εργασία. Περιλαμβάνει όλα τα παραπάνω.

## Πρόγνωση Τιμών Χρονοσειρών

Όλα τα πειράματα πραγματοποιήθηκαν για ένα μεγάλο αριθμό από epochs. Τελικά επιλέχθηκε η τιμή της παραμέτρου epochs που ελαχιστοποιεί το σφάλμα πριν εμφανιστεί το φαινόμενο overfitting. Η αναγνώριση του κατάλληλου epoch προήλθε από τη γραφική παράσταση σφάλματος για τα σετ train, test κατά την εκπαίδευση (συνάρτηση plot\_training\_loss).

Τα πειράματα που πραγματοποιήθηκαν συνοψίζονται στο παρακάτω πίνακα.

| Πείραμα | Batch Size | Layers | Units/Layer | Time Steps | Train Loss | Val Loss |
|---------|------------|--------|-------------|------------|------------|----------|
| 1       | 32         | 4      | 32          | 50         | 8.3730e-04 | 0.0807   |
| 2       | 32         | 4      | 64          | 60         | 5.3747e-04 | 0.0847   |
| 3       | 32         | 4      | 128         | 50         | 5.7512e-04 | 0.0919   |
| 4       | 64         | 6      | 64          | 60         | 5.6178e-04 | 0.0865   |
| 5       | 64         | 6      | 128         | 50         | 5.5015e-04 | 0.0909   |
| 6       | 64         | 6      | 32          | 60         | 8.6226e-04 | 0.0900   |
| 7       | 64         | 8      | 64          | 50         | 5.7873e-04 | 0.0828   |
| 8       | 128        | 8      | 128         | 60         | 4.4686e-04 | 0.0857   |
| 9       | 128        | 6      | 64          | 50         | 5.7284e-04 | 0.0855   |

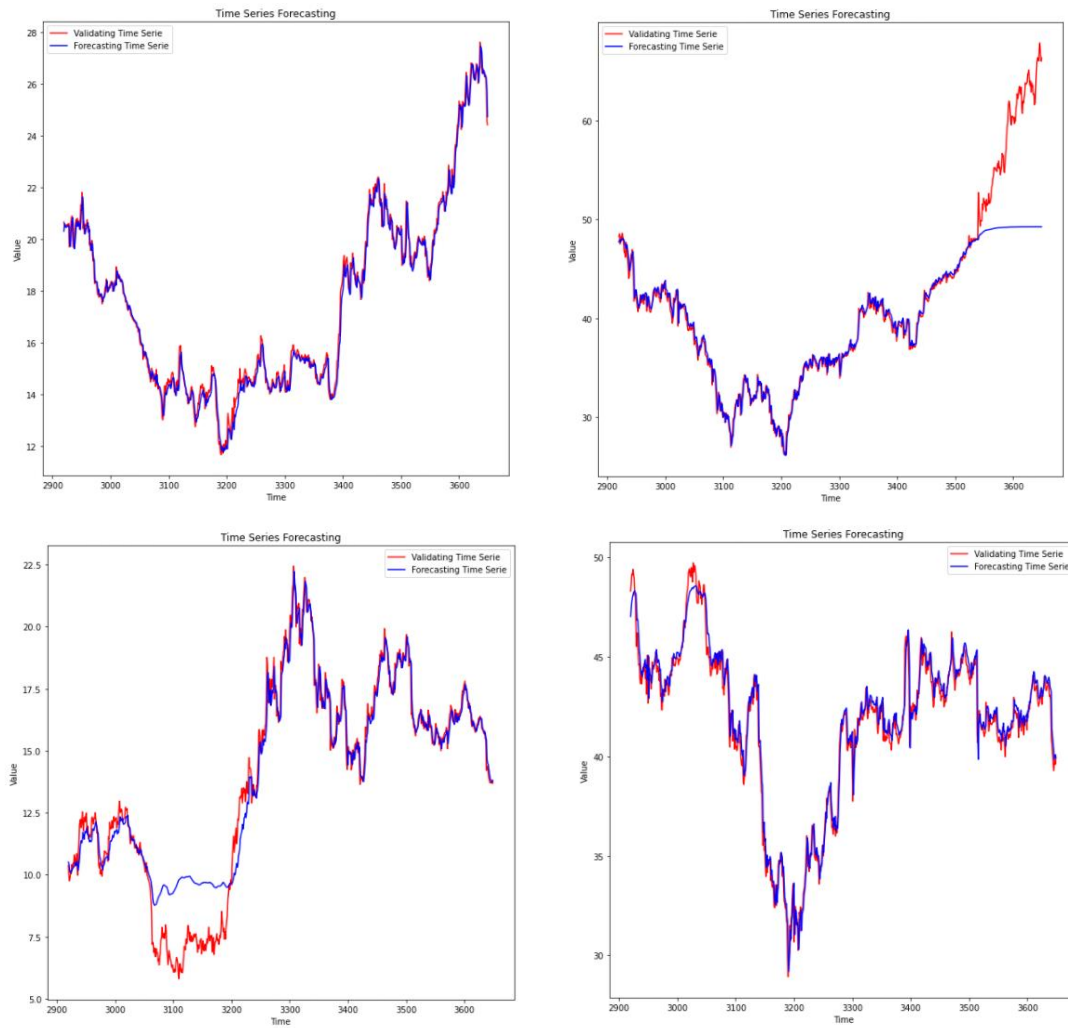
Κάθε layer αποτελείται από ένα επίπεδο LSTM και ένα επίπεδο dropout με πιθανότητα απόρριψης 0.2. Έτσι, αποφεύγεται το φαινόμενο υπερπροσαρμογής (overfitting).

Παρατηρήσεις:

- Ο χρόνος εκτέλεσης μειώνεται σημαντικά με την αύξηση του batch size, ωστόσο με χρήση πολύ μεγάλου batch size (>1024) το loss δεν ελαχιστοποιείται.
- Αυξάνοντας τον αριθμό των κρυφών επιπέδων του μοντέλου δεν παρατηρείται μείωση του σφάλματος. Τελικά το πρόβλημα δεν απαιτεί ένα πιο περίπλοκο μοντέλο από αυτά που δοκιμάστηκαν.
- Σημαντική επίδραση φαίνεται να έχει ο αριθμός των νευρώνων (units) ανά επίπεδο. Αυξάνοντας την τιμή της υπερπαραμέτρου το σφάλμα μειώνεται.
- Η υπερπαραμέτρος time-steps (look-back) δεν μεταβάλλει σημαντικά τη τιμή του σφάλματος.

Το τελικό μοντέλο περιέχει τις τιμές υπερπαραμέτρων του πειράματος 8, όπου είναι οι βέλτιστες.

## Ενδεικτικές προγνώσεις του τελικού μοντέλου



Το μοντέλο προσεγγίζει άρτια τις χρονοσειρές.

## Ανίχνευση Ανωμαλιών Χρονοσειρών

Όπως στη πρόγνωση τιμών χρονοσειρών επιλέχθηκε κατάλληλα η τιμή της υπερπαραμέτρου  $epoch$ .

Τα πειράματα που πραγματοποιήθηκαν συνοψίζονται στο παρακάτω πίνακα.

| Πείραμα | Batch Size | Layers | Units | Time Steps | Train Loss | Val Loss |
|---------|------------|--------|-------|------------|------------|----------|
| 1       | 32         | 2      | 64    | 50         | 0.2111     | 0.6069   |
| 2       | 32         | 2      | 128   | 60         | 0.2120     | 0.6148   |
| 3       | 32         | 2      | 256   | 50         | 0.2110     | 0.6074   |
| 4       | 64         | 2      | 64    | 60         | 0.2143     | 0.6132   |
| 5       | 64         | 4      | 128   | 50         | 0.2112     | 0.6121   |
| 6       | 128        | 6      | 64    | 60         | 0.1987     | 0.6241   |
| 7       | 128        | 6      | 128   | 50         | 0.1998     | 0.6218   |
| 8       | 128        | 6      | 32    | 60         | 0.1975     | 0.6243   |

Κάθε layer αποτελείται από ένα επίπεδο LSTM και ένα επίπεδο dropout με πιθανότητα απόρριψης 0.2. Έτσι, αποφεύγεται το φαινόμενο υπερπροσαρμογής (overfitting).

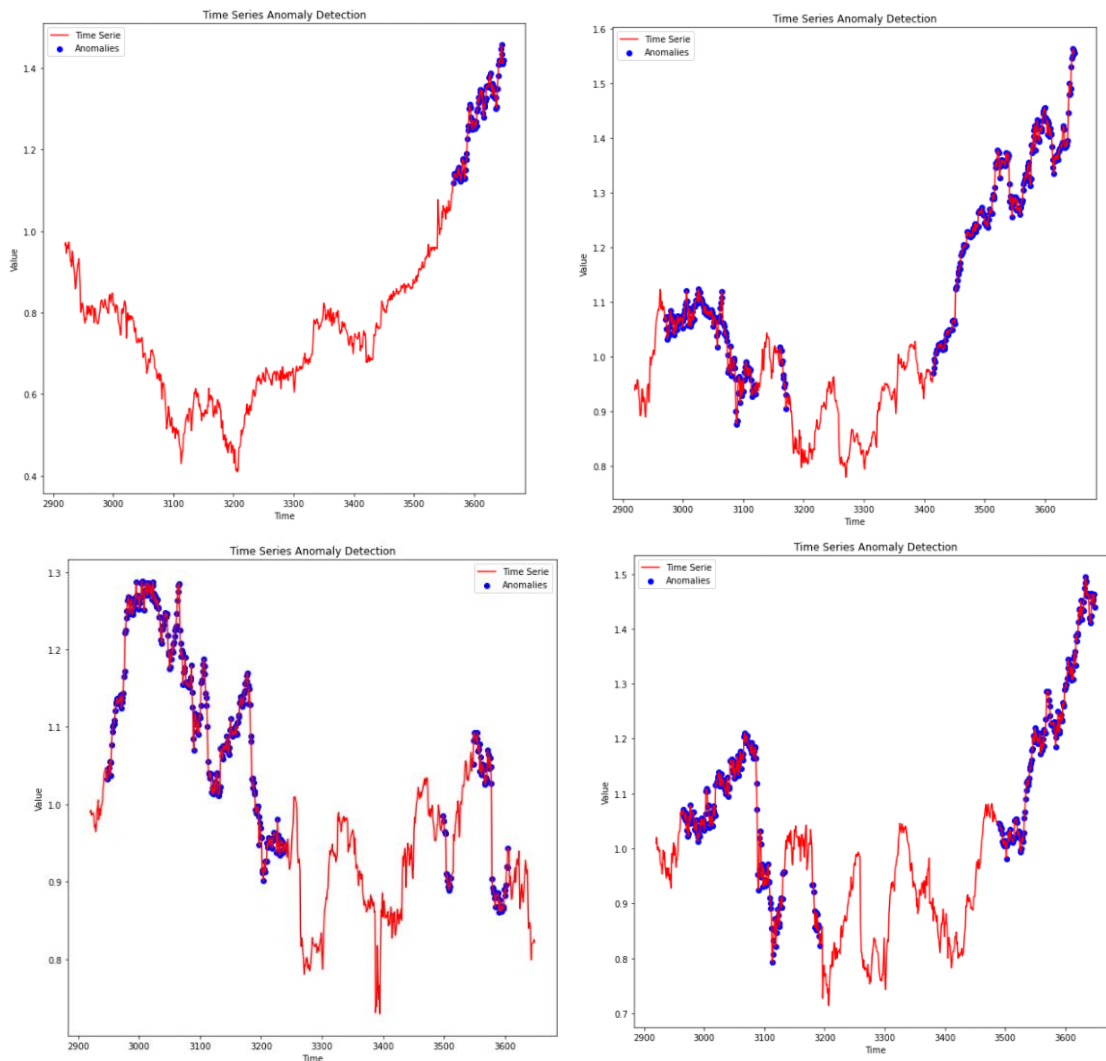
Παρατηρήσεις:

- Ο χρόνος εκτέλεσης μειώνεται σημαντικά με την αύξηση του batch size, ωστόσο με χρήση πολύ μεγάλου batch size (>1024) το loss δεν ελαχιστοποιείται.
- Αυξάνοντας τον αριθμό των κρυφών επιπέδων του μοντέλου παρατηρείται σημαντική μείωση του σφάλματος. Το πρόβλημα απαιτεί ένα αρκετά περίπλοκο μοντέλο.
- Μικρή επίδραση φαίνεται να έχει ο αριθμός των νευρώνων (units) ανά επίπεδο. Αυξάνοντας την τιμή της υπερπαραμέτρου το σφάλμα δεν μεταβάλλεται ισχυρά.
- Η υπερπαραμέτρος time-steps (look-back) δεν μεταβάλλει σημαντικά τη τιμή του σφάλματος.

Το τελικό μοντέλο περιέχει τις τιμές υπερπαραμέτρων του πειράματος 8, όπου είναι οι βέλτιστες.

## Ενδεικτικές ανιχνεύσεις του τελικού μοντέλου

MAE=0.65



Οι ανωμαλίες επισημαίνονται σε σημεία με απότομες μεταβάσεις ή απρόβλεπτη πορεία τιμών σε σύγκριση με τις υπόλοιπες χρονοσειρές, όπως αναμενόταν.

## Συνελικτική Αυτοκωδικοποίηση Χρονοσειρών

Όπως στη πρόγνωση τιμών και στην ανίχνευση ανωμαλιών χρονοσειρών επιλέχθηκε κατάλληλα η τιμή της υπερπαραμέτρου  $epoch$ .

Τα πειράματα που πραγματοποιήθηκαν συνοψίζονται στο παρακάτω πίνακα.

| Πείραμα | Batch Size | Layers | Filter Size | Window Length | Latent Dimension | Train Loss | Val Loss |
|---------|------------|--------|-------------|---------------|------------------|------------|----------|
| 1       | 32         | 2      | 16          | 10            | 3                | 0,6791     | 0,7576   |
| 2       | 32         | 2      | 16          | 50            | 13               | 0,5563     | 0,2912   |
| 3       | 32         | 2      | 32          | 50            | 13               | 0,5303     | -0,1815  |
| 4       | 32         | 3      | 16          | 50            | 7                | 0,5317     | -0,5217  |
| 5       | 32         | 3      | 16          | 10            | 2                | 0,5436     | -0,9232  |
| 6       | 64         | 3      | 32          | 50            | 7                | 0,5304     | -0,8348  |
| 7       | 64         | 3      | 64          | 10            | 2                | 0,5295     | -0,2726  |
| 8       | 32         | 4      | 16          | 50            | 4                | 0,5307     | -0,1027  |
| 9       | 32         | 4      | 32          | 50            | 4                | 0,6665     | -0,8822  |

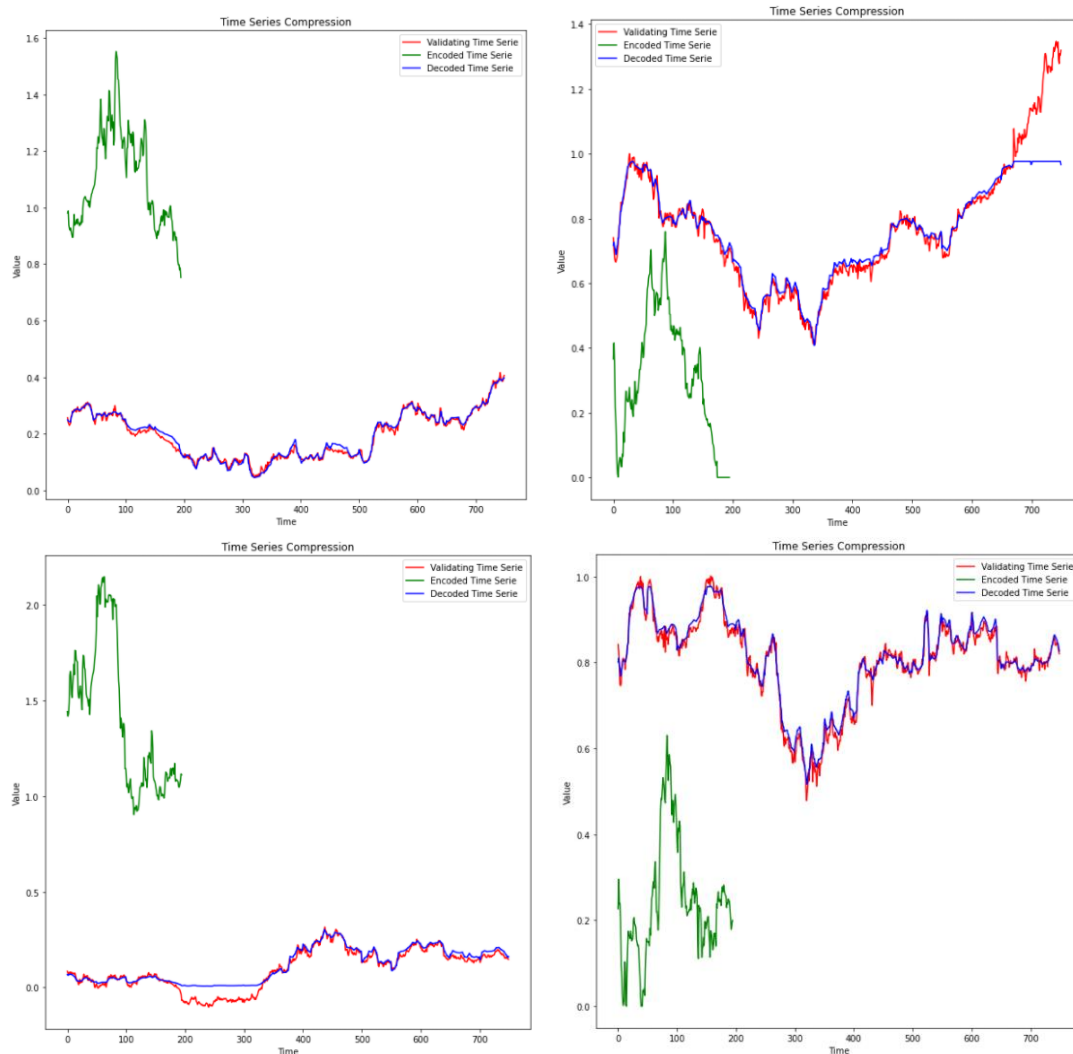
Κάθε layer αποτελείται από ένα επίπεδο conv1d και ένα επίπεδο αύξησης/μείωσης διάστασης. Κατά την αναφορά layers=3, υποδεικνύεται η ύπαρξη 3 επιπέδων downsampling (conv1d και maxpooling1d) και τριών επιπέδων upsampling (conv1d και upsampling).

Παρατηρήσεις:

- Ο χρόνος εκτέλεσης μειώνεται σημαντικά με την αύξηση του batch size, ωστόσο με χρήση πολύ μεγάλου batch size (>1024) το loss δεν ελαχιστοποιείται.
- Σημαντική αύξηση ή μείωση του αριθμού των επιπέδων αποκλείει την ελαχιστοποίηση του σφάλματος.
- Μικρή επίδραση φαίνεται να έχει ο αριθμός των νευρώνων (units) ανά επίπεδο. Αυξάνοντας την τιμή της υπερπαραμέτρου το σφάλμα δεν μεταβάλλεται ισχυρά.
- Η υπερπαραμέτρος time-steps (look-back) δεν μεταβάλλει σημαντικά τη τιμή του σφάλματος.

Το τελικό μοντέλο περιέχει τις τιμές υπερπαραμέτρων του πειράματος 4, όπου είναι οι βέλτιστες.

## Ενδεικτικές συμπίεσεις του τελικού μοντέλου



Η προσέγγιση της αρχικής χρονοσειράς είναι ικανοποιητική. Η encoded χρονοσειρά εμφανίζει κοινά μοτίβα με τις υπόλοιπες δύο, όπως αναμενόταν.

## Σύγκριση σε Αναζήτηση και Συσταδοποίηση

Οι παραγόμενες συμπιεσμένες χρονοσειρές από το μοντέλο συνελικτικής αυτοκωδικοποίησης συγκρίνονται με τις αρχικές χρονοσειρές του σετ δεδομένων στα συστήματα αναζήτησης πλησιέστερων γειτόνων (search) και συσταδοποίησης (cluster) διανυσμάτων και χρονοσειρών της εργασίας 2. Τα αποτελέσματα ανά εκτέλεση δύναται να διαφέρουν επειδή οι αλγόριθμοι search και cluster είναι πιθανοκρατικοί.



## Αναζήτηση Διανυσμάτων με LSH ( k=8 , L=5 )

Τα αποτελέσματα:

```
Query: xel
Algorithm: LSH_Vector
Approximate Nearest neighbor: cop
True Nearest neighbor: t
distanceApproximate: 1358.71
distanceTrue: 211.418
Query: xl
Algorithm: LSH_Vector
Approximate Nearest neighbor: phm
True Nearest neighbor: sti
distanceApproximate: 1436.84
distanceTrue: 561.222
Query: xlnx
Algorithm: LSH_Vector
Approximate Nearest neighbor: payx
True Nearest neighbor: snps
distanceApproximate: 357.046
distanceTrue: 307.226
Query: xom
Algorithm: LSH_Vector
Approximate Nearest neighbor: pep
True Nearest neighbor: slb
distanceApproximate: 792.258
distanceTrue: 586.457
Query: xray
Algorithm: LSH_Vector
Approximate Nearest neighbor: pcg
True Nearest neighbor: pcg
distanceApproximate: 274.728
distanceTrue: 274.728
Query: xrx
Algorithm: LSH_Vector
Approximate Nearest neighbor: ge
True Nearest neighbor: kss
distanceApproximate: 1114.69
distanceTrue: 627.702
Query: yum
Algorithm: LSH_Vector
Approximate Nearest neighbor: tmk
True Nearest neighbor: pls
distanceApproximate: 415.387
distanceTrue: 403.546
Query: zbh
Algorithm: LSH_Vector
Approximate Nearest neighbor: omc
True Nearest neighbor: ups
distanceApproximate: 1807.93
distanceTrue: 669.383
Query: zion
Algorithm: LSH_Vector
Approximate Nearest neighbor: sti
True Nearest neighbor: sti
distanceApproximate: 445.729
distanceTrue: 445.729
tApproximateAverage: 4822 microseconds
tTrueAverage: 34357 microseconds
MAF: 6.42665
```

Vector LSH initial

```
Approximate Nearest neighbor: fitb
True Nearest neighbor: fitb
distanceApproximate: 3.9533
distanceTrue: 3.9533
Query: xlnx
Algorithm: LSH_Vector
Approximate Nearest neighbor: mmm
True Nearest neighbor: adi
distanceApproximate: 2.40091
distanceTrue: 1.84224
Query: xom
Algorithm: LSH_Vector
Approximate Nearest neighbor: mon
True Nearest neighbor: cvx
distanceApproximate: 4.3812
distanceTrue: 4.04647
Query: xray
Algorithm: LSH_Vector
Approximate Nearest neighbor: t
True Nearest neighbor: t
distanceApproximate: 3.38407
distanceTrue: 3.38407
Query: xrx
Algorithm: LSH_Vector
Approximate Nearest neighbor: hig
True Nearest neighbor: rrd
distanceApproximate: 4.86928
distanceTrue: 3.44404
Query: yum
Algorithm: LSH_Vector
Approximate Nearest neighbor: orcl
True Nearest neighbor: orcl
distanceApproximate: 2.48487
distanceTrue: 2.48487
Query: zbh
Algorithm: LSH_Vector
Approximate Nearest neighbor: mdt
True Nearest neighbor: mdt
distanceApproximate: 3.0961
distanceTrue: 3.0961
Query: zion
Algorithm: LSH_Vector
Approximate Nearest neighbor: rf
True Nearest neighbor: key
distanceApproximate: 2.96483
distanceTrue: 2.30908
tApproximateAverage: 13272 microseconds
tTrueAverage: 4715 microseconds
MAF: 1.41383
```

Vector LSH Compressed

Οι πραγματικές αποστάσεις των γειτόνων στη συμπιεσμένη αναπαράσταση είναι μικρότερες, από της αντίστοιχης αρχικής. Αυτό είναι αναμενόμενο, αφού πραγματοποιείται αναπαραγωγή της χρονοσειράς σε μικρότερη διάσταση και οι αναλογίες πραγματικής και σχετικής απόστασης παραμένουν όμοιες στις δύο αναπαραστάσεις. Συγκρίνοντας τις δύο αναπαραστάσεις ξεχωρίζεται ο κοινός γείτονας “t” υποδεικνύοντας μας πως η συμπιεσμένη αναπαράσταση αποτελεί μια καλή προσέγγιση της αρχικής.

## Αναζήτηση Διανυσμάτων με Hypercube ( probes=10 , k=10 )

Τα αποτελέσματα:

```
Query: xel
Algorithm: Hypercube
Approximate Nearest neighbor: wor
True Nearest neighbor: t
distanceApproximate: 306.496
distanceTrue: 211.418
Query: xl
Algorithm: Hypercube
Query: xlnx
Algorithm: Hypercube
Approximate Nearest neighbor: msft
True Nearest neighbor: snps
distanceApproximate: 335.431
distanceTrue: 307.226
Query: xom
Algorithm: Hypercube
Approximate Nearest neighbor: pg
True Nearest neighbor: slb
distanceApproximate: 626.793
distanceTrue: 586.457
Query: xray
Algorithm: Hypercube
Approximate Nearest neighbor: wfc
True Nearest neighbor: pcg
distanceApproximate: 481.973
distanceTrue: 274.728
Query: xrx
Algorithm: Hypercube
Approximate Nearest neighbor: kss
True Nearest neighbor: kss
distanceApproximate: 627.702
distanceTrue: 627.702
Query: yum
Algorithm: Hypercube
Approximate Nearest neighbor: txn
True Nearest neighbor: gis
distanceApproximate: 558.55
distanceTrue: 403.546
Query: zbh
Algorithm: Hypercube
Approximate Nearest neighbor: ph
True Nearest neighbor: ups
distanceApproximate: 1350.49
distanceTrue: 669.383
Query: zion
Algorithm: Hypercube
Approximate Nearest neighbor: ncr
True Nearest neighbor: sti
distanceApproximate: 1245.41
distanceTrue: 445.729
tApproximateAverage: 3342 microseconds
tTrueAverage: 37103 microseconds
MAF: 2.79409
```

Vector Hypercube

```
Query: xel
Algorithm: Hypercube
Approximate Nearest neighbor: adp
True Nearest neighbor: ed
distanceApproximate: 1.9617
distanceTrue: 1.62213
Query: xl
Algorithm: Hypercube
Approximate Nearest neighbor: mxim
True Nearest neighbor: fitb
distanceApproximate: 9.05043
distanceTrue: 3.9533
Query: xlnx
Algorithm: Hypercube
Approximate Nearest neighbor: mco
True Nearest neighbor: adi
distanceApproximate: 4.11979
distanceTrue: 1.84224
Query: xom
Algorithm: Hypercube
Approximate Nearest neighbor: flr
True Nearest neighbor: cvx
distanceApproximate: 7.82497
distanceTrue: 4.04647
Query: xray
Algorithm: Hypercube
Approximate Nearest neighbor: duk
True Nearest neighbor: t
distanceApproximate: 3.4706
distanceTrue: 3.38407
Query: xrx
Algorithm: Hypercube
Approximate Nearest neighbor: ati
True Nearest neighbor: rrd
distanceApproximate: 9.86862
distanceTrue: 3.44404
Query: yum
Algorithm: Hypercube
Approximate Nearest neighbor: mcd
True Nearest neighbor: orcl
distanceApproximate: 2.51684
distanceTrue: 2.48487
Query: zbh
Algorithm: Hypercube
Approximate Nearest neighbor: ca
True Nearest neighbor: mdt
distanceApproximate: 3.38861
distanceTrue: 3.0961
Query: zion
Algorithm: Hypercube
Approximate Nearest neighbor: key
True Nearest neighbor: key
distanceApproximate: 2.30908
distanceTrue: 2.30908
tApproximateAverage: 237 microseconds
tTrueAverage: 4746 microseconds
MAF: 2.86542
```

Vector Hypercube Compressed

Ξανά συγκρίνοντας τις δύο αναπαραστάσεις ξεχωρίζεται ο κοινός γείτονας “t” υποδεικνύοντας μας πως η συμπιεσμένη αναπαράσταση αποτελεί μια καλή προσέγγιση της αρχικής.

## Αναζήτηση Χρονοσειρών με Discrete Frechet ( $k=8$ , $L=5$ )

Τα αποτελέσματα:

```
Query: xel
Algorithm: LSH_Frechet_Discrete
Approximate Nearest neighbor: fitb
True Nearest neighbor: cms
distanceApproximate: 38.892
distanceTrue: 7.2957
Query: xl
Algorithm: LSH_Frechet_Discrete
Approximate Nearest neighbor: mkc
True Nearest neighbor: sti
distanceApproximate: 72.1053
distanceTrue: 34.366
Query: xlnx
Algorithm: LSH_Frechet_Discrete
Query: xom
Algorithm: LSH_Frechet_Discrete
Approximate Nearest neighbor: nue
True Nearest neighbor: pg
distanceApproximate: 52.6524
distanceTrue: 24.085
Query: xray
Algorithm: LSH_Frechet_Discrete
Approximate Nearest neighbor: etr
True Nearest neighbor: pcg
distanceApproximate: 51.8454
distanceTrue: 15.8576
Query: xrx
Algorithm: LSH_Frechet_Discrete
Approximate Nearest neighbor: jhg
True Nearest neighbor: kss
distanceApproximate: 29.9598
distanceTrue: 22.7972
Query: yum
Algorithm: LSH_Frechet_Discrete
Approximate Nearest neighbor: gps
True Nearest neighbor: tmk
distanceApproximate: 53.9371
distanceTrue: 16.6603
Query: zbh
Algorithm: LSH_Frechet_Discrete
Query: zion
Algorithm: LSH_Frechet_Discrete
Approximate Nearest neighbor: unh
True Nearest neighbor: sti
distanceApproximate: 167.14
distanceTrue: 16.953
tApproximateAverage: 1311267 microseconds
tTrueAverage: 453160309 microseconds
MAF: 9.85902
```

Time-Serie Discrete Frechet

```
Query: xel
Algorithm: LSH_Frechet_Discrete
Approximate Nearest neighbor: cms
True Nearest neighbor: cms
distanceApproximate: 0.199831
distanceTrue: 0.199831
Query: xl
Algorithm: LSH_Frechet_Discrete
Approximate Nearest neighbor: rf
True Nearest neighbor: rf
distanceApproximate: 0.600322
distanceTrue: 0.600322
Query: xlnx
Algorithm: LSH_Frechet_Discrete
Approximate Nearest neighbor: a
True Nearest neighbor: a
distanceApproximate: 0.295518
distanceTrue: 0.295518
Query: xom
Algorithm: LSH_Frechet_Discrete
Approximate Nearest neighbor: cvx
True Nearest neighbor: cvx
distanceApproximate: 0.427988
distanceTrue: 0.427988
Query: xray
Algorithm: LSH_Frechet_Discrete
Approximate Nearest neighbor: apd
True Nearest neighbor: apd
distanceApproximate: 0.406797
distanceTrue: 0.406797
Query: xrx
Algorithm: LSH_Frechet_Discrete
Approximate Nearest neighbor: rrd
True Nearest neighbor: rrd
distanceApproximate: 0.495628
distanceTrue: 0.495628
Query: yum
Algorithm: LSH_Frechet_Discrete
Approximate Nearest neighbor: wat
True Nearest neighbor: wat
distanceApproximate: 0.327785
distanceTrue: 0.327785
Query: zbh
Algorithm: LSH_Frechet_Discrete
Approximate Nearest neighbor: ca
True Nearest neighbor: ca
distanceApproximate: 0.369852
distanceTrue: 0.369852
Query: zion
Algorithm: LSH_Frechet_Discrete
Approximate Nearest neighbor: key
True Nearest neighbor: key
distanceApproximate: 0.354276
distanceTrue: 0.354276
tApproximateAverage: 37597657 microseconds
tTrueAverage: 8634456 microseconds
MAF: 1
```

Time-Serie Discrete Frechet Compressed

Ξανά συγκρίνοντας τις δύο αναπαραστάσεις ξεχωρίζεται ο κοινός γείτονας “cms” υποδεικνύοντας μας πως η συμπιεσμένη αναπαράσταση αποτελεί μια καλή προσέγγιση της αρχικής.

## Αναζήτηση Χρονοσειρών με Continuous Frechet ( k=8 , L=5 )

Τα δειγματοληπτικά αποτελέσματα:

```
Query: xel
Algorithm: LSH_Frechet_Continuous
Approximate Nearest neighbor: aep
True Nearest neighbor: adp
distanceApproximate: 0.373937
distanceTrue: 0.365148
Query: xl
Algorithm: LSH_Frechet_Continuous
Approximate Nearest neighbor: afl
True Nearest neighbor: aee
distanceApproximate: 1.00081
distanceTrue: 0.765189
Query: xlnx
Algorithm: LSH_Frechet_Continuous
Query: xom
Algorithm: LSH_Frechet_Continuous
Approximate Nearest neighbor: aes
True Nearest neighbor: adm
distanceApproximate: 0.915197
distanceTrue: 0.675205
Query: xray
Algorithm: LSH_Frechet_Continuous
Approximate Nearest neighbor: aes
True Nearest neighbor: a
distanceApproximate: 0.913134
distanceTrue: 0.593819
Query: xrx
Algorithm: LSH_Frechet_Continuous
Query: yum
Algorithm: LSH_Frechet_Continuous
Query: zbh
Algorithm: LSH_Frechet_Continuous
Query: zion
Algorithm: LSH_Frechet_Continuous
Approximate Nearest neighbor: abc
True Nearest neighbor: aa
distanceApproximate: 1.40266
distanceTrue: 0.651523
tApproximateAverage: 11457641 microseconds
tTrueAverage: 203219369 microseconds
MAF: 2.1529
```

Time-Serie Continuous Frechet

```
Query: xel
Algorithm: LSH_Frechet_Continuous
Query: xl
Algorithm: LSH_Frechet_Continuous
Approximate Nearest neighbor: afl
True Nearest neighbor: aee
distanceApproximate: 1.00081
distanceTrue: 0.765189
Query: xlnx
Algorithm: LSH_Frechet_Continuous
Query: xom
Algorithm: LSH_Frechet_Continuous
Approximate Nearest neighbor: aes
True Nearest neighbor: adm
distanceApproximate: 0.915197
distanceTrue: 0.675205
Query: xray
Algorithm: LSH_Frechet_Continuous
Approximate Nearest neighbor: aes
True Nearest neighbor: a
distanceApproximate: 0.913134
distanceTrue: 0.593819
Query: xrx
Algorithm: LSH_Frechet_Continuous
Query: yum
Algorithm: LSH_Frechet_Continuous
Approximate Nearest neighbor: adsk
True Nearest neighbor: aapl
distanceApproximate: 0.498318
distanceTrue: 0.426215
Query: zbh
Algorithm: LSH_Frechet_Continuous
Query: zion
Algorithm: LSH_Frechet_Continuous
Approximate Nearest neighbor: abc
True Nearest neighbor: aa
distanceApproximate: 1.40266
distanceTrue: 0.651523
tApproximateAverage: 7989246 microseconds
tTrueAverage: 204221841 microseconds
MAF: 2.1529
```

Time-Serie Continuous Frechet

Πολλά ερωτήματα δεν επέστρεψαν κοντινότερο γείτονα εξ αιτίας του μικρού αριθμού δείγματος. Παρόλα αυτά αναγνωρίζονται πολλές κοινές χρονοσειρές ως κοντινότεροι γείτονες. Συγκεκριμένα οι afl, aee, aes, adm, a, aa, abc. Το μοντέλο φαίνεται αποδοτικό.

## Συσταδοποίηση Διανυσμάτων με LSH

Τα αποτελέσματα:

```
Algorithm: Range Search LSH
CLUSTER-1 {size: 195, centroid: 1.24663 1.23366 1.20245 1.17767 1.17414}
CLUSTER-2 {size: 97, centroid: 1.65038 1.63671 1.61237 1.59603 1.59075}
CLUSTER-3 {size: 23, centroid: 0.780348 0.772161 0.730144 0.701004 0.695065}
CLUSTER-4 {size: 13, centroid: 0.489088 0.483121 0.44748 0.417972 0.42334}
CLUSTER-5 {size: 22, centroid: 1.20719 1.19369 1.16687 1.14768 1.14214}
clustering time: 30
Silhouette: [0.0271573, 0.240173, 0.163257, 0.537754, 0.426017, 0.108912]
CLUSTER-1 { 1.24663 1.23366 1.20245 1.17767 1.17414 1.16538 1.1819 1.18808}
CLUSTER-2 { 1.65038 1.63671 1.61237 1.59603 1.59075 1.57903 1.58861 1.53223}
CLUSTER-3 { 0.780348 0.772161 0.730144 0.701004 0.695065 0.672031 0.69982}
CLUSTER-4 { 0.489088 0.483121 0.44748 0.417972 0.42328 0.410627 0.43334}
CLUSTER-5 { 1.20719 1.19369 1.16687 1.14768 1.14214 1.12705 1.15049 1.14214}
```

Vector LSH Clustering

```
Algorithm: Range Search LSH
CLUSTER-1 {size: 208, centroid: 1.57852 1.65455 1.63364 1.63788 1.63431 1.63788}
CLUSTER-2 {size: 42, centroid: 1.14606 1.18834 1.16051 1.16944 1.17064 1.16944}
CLUSTER-3 {size: 44, centroid: 1.42836 1.49425 1.46378 1.46465 1.46588 1.46465}
CLUSTER-4 {size: 37, centroid: 1.07346 1.09215 1.05638 1.08097 1.07392 1.07392}
CLUSTER-5 {size: 19, centroid: 1.15025 1.20201 1.18134 1.18466 1.17788 1.17788}
clustering time: 5
Silhouette: [0.105965, 0.0808833, 0.249008, 0.252602, 0.260356, 0.124341]
CLUSTER-1 { 1.57852 1.65455 1.63364 1.63788 1.63431 1.62765 1.53283 1.54258}
CLUSTER-2 { 1.14606 1.18834 1.16051 1.16944 1.17064 1.16656 1.09753 1.10028}
CLUSTER-3 { 1.42836 1.49425 1.46378 1.46465 1.46588 1.46246 1.37165 1.37698}
CLUSTER-4 { 1.07346 1.09215 1.05638 1.08097 1.07392 1.07596 1.01319 1.02847}
CLUSTER-5 { 1.15025 1.20201 1.18134 1.18466 1.17788 1.17149 1.09971 1.10189}
```

Vector LSH Clustering Compressed

Τα κεντροειδή και τα περιεχόμενα των συστάδων δεν αναγράφονται γιατί αποτελούν μεγάλο όγκο δεδομένων. Τα αποτελέσματα είναι ικανοποιητικά. Η μετρική silhouette επιβεβαιώνει τη σωστή επιλογή των τεσσάρων συστάδων. Υπάρχουν αρκετά όμοια διανύσματα εντός των ίδιων συστάδων στις δύο αναπαραστάσεις.

## Συσταδοποίηση Διανυσμάτων με Hypercube

Τα αποτελέσματα:

```
Algorithm: Range Search Hypercube
CLUSTER-1 {size: 64, centroid: 1.08475 1.11384 1.08236 1.1035 1.09597 1.08962}
CLUSTER-2 {size: 112, centroid: 1.56475 1.64111 1.62011 1.63209 1.62709 1.61606}
CLUSTER-3 {size: 56, centroid: 1.51252 1.58228 1.56104 1.56371 1.56306 1.56223}
CLUSTER-4 {size: 60, centroid: 1.22917 1.28487 1.25187 1.25247 1.24805 1.23082}
CLUSTER-5 {size: 50, centroid: 1.62102 1.69544 1.67197 1.67818 1.67477 1.67263}
clustering time: 1
Silhouette: [0.191946, 0.241572, 0.0706752, 0.198269, 0.0801534, 0.170981]
CLUSTER-1 { 1.08475 1.11384 1.08236 1.1035 1.09597 1.08962 1.02157 1.03874 1.10384}
CLUSTER-2 { 1.56475 1.64111 1.62011 1.63209 1.62709 1.61805 1.51804 1.52703 1.10384}
CLUSTER-3 { 1.51252 1.58228 1.56104 1.56371 1.56306 1.56223 1.47057 1.47747 1.10384}
CLUSTER-4 { 1.22917 1.28487 1.25187 1.25247 1.24805 1.23002 1.15941 1.16639 1.10384}
CLUSTER-5 { 1.62102 1.69544 1.67197 1.67818 1.67477 1.67263 1.56837 1.57734 1.10384}
```

Vector Hypercube Clustering

```
Algorithm: Range Search Hypercube
CLUSTER-1 {size: 57, centroid: 0.468225 0.487603 0.420437 0.405219 0.433314}
CLUSTER-2 {size: 163, centroid: 1.64243 1.63055 1.58338 1.59814 1.62033 1.62033}
CLUSTER-3 {size: 37, centroid: 1.49365 1.44402 1.41148 1.38917 1.38843 1.38843}
CLUSTER-4 {size: 72, centroid: 0.885255 0.884329 0.850824 0.853727 0.855939}
CLUSTER-5 {size: 21, centroid: 1.00842 1.0028 1.00042 0.967309 0.954006 0.954006}
clustering time: 0
Silhouette: [0.0315847, 0.200756, 0.227509, 0.44541, 0.307664, 0.232777]
CLUSTER-1 { 0.468225 0.487603 0.420437 0.405219 0.433314 0.485783 0.507896}
CLUSTER-2 { 1.64243 1.63055 1.58338 1.59814 1.62033 1.61304 1.60250 1.5971}
CLUSTER-3 { 1.49365 1.44402 1.41148 1.38917 1.38843 1.36443 1.39935 1.3946}
CLUSTER-4 { 0.885255 0.884329 0.850824 0.853727 0.855939 0.87259 0.878248}
CLUSTER-5 { 1.00842 1.0028 1.00042 0.967309 0.954006 0.928023 0.929351 0.929351}
```

Vector Hypercube Clustering Compressed

Όπως στη συσταδοποίηση διανυσμάτων με LSH, τα αποτελέσματα είναι ικανοποιητικά. Η μετρική silhouette επιβεβαιώνει τη σωστή επιλογή των τεσσάρων συστάδων. Υπάρχουν αρκετά όμοια διανύσματα εντός των ίδιων συστάδων στις δύο αναπαραστάσεις.

## Συσταδοποίηση Διανυσμάτων με Lloyds

### Τα αποτελέσματα:

```

Algorithm: Lloyds
CLUSTER-1 (size: 103, centroid: 1.58704 1.66431 1.64509 1.64918 1.64447 1.64129 1.54456 1.55616)
CLUSTER-2 (size: 26, centroid: 1.39688 1.46356 1.44487 1.44787 1.44959 1.36625 1.37676 1.38983)
CLUSTER-3 (size: 74, centroid: 1.37805 1.42564 1.39756 1.40274 1.39999 1.31212 1.30595 1.31212)
CLUSTER-4 (size: 43, centroid: 0.889698 0.883268 0.841861 0.869441 0.855672 0.831213 0.786768 0.786768)
CLUSTER-5 (size: 104, centroid: 1.55614 1.63554 1.61391 1.61583 1.61783 1.61822 1.52076 1.52076)
clustering time: 0
Silhouette: [0.278927, 0.303299, 0.0947371, 0.274256, 0.178512, 0.211383]
CLUSTER-1 (1.58704 1.66431 1.64509 1.64918 1.64447 1.64129 1.54456 1.55616)
CLUSTER-2 (1.39688 1.46356 1.44487 1.44787 1.44959 1.36625 1.37676 1.38983)
CLUSTER-3 (1.37805 1.42564 1.39756 1.40274 1.39999 1.31212 1.30595 1.31212)
CLUSTER-4 (0.889698 0.883268 0.841861 0.869441 0.855672 0.831213 0.786768 0.786768)
CLUSTER-5 (1.55614 1.63554 1.61391 1.61583 1.61783 1.61822 1.52076 1.52076)

```

## Vector Lloyds Clustering

```
Algorithm: Lloyd's
CLUSTER-1 (size: 72, centroid: 1.3899 1.36992 1.33857 1.315 1.31196 1.31196)
CLUSTER-2 (size: 119, centroid: 1.19408 1.18218 1.15253 1.1258 1.12076 1.12076)
CLUSTER-3 (size: 94, centroid: 1.66543 1.65503 1.62939 1.61307 1.60825 1.60825)
CLUSTER-4 (size: 24, centroid: 1.21117 1.19313 1.17805 1.16189 1.15625 1.15625)
CLUSTER-5 (size: 41, centroid: 0.640432 0.636703 0.590106 0.563442 0.563442 0.563442)
clustering_time = 0
Silhouette: [0.183495 0.205231 0.225303 0.341203 0.253859 0.221171]
CLUSTER-1 (1.3899 1.36992 1.33857 1.315 1.31196 1.3021 1.31661 1.3196)
CLUSTER-2 (1.19408 1.18218 1.15253 1.1258 1.12076 1.11133 1.12834 1.13)
CLUSTER-3 (1.66543 1.65503 1.62939 1.61307 1.60825 1.59587 1.61006 1.60825)
CLUSTER-4 (1.21117 1.19313 1.17805 1.16189 1.15625 1.14783 1.15213 1.15)
CLUSTER-5 (0.640432 0.636703 0.590106 0.563442 0.563442 0.563864 0.547475 0.57)
```

## Vector Lloyds Clustering Compressed

Τα μεγέθη εντός των συστάδων διαφέρουν ελαφρώς. Παρόλα αυτά υπάρχουν αρκετά κοινά διανύσματα εντός αυτών.

## Συσταδοποίηση Χρονοσειρών με LSH Discrete Frechet

### Τα δειγματοληπτικά αποτελέσματα:

```
Algorithm: LSH Frechet_Discrete
CLUSTER-1 {size: 12, centroid:( 1 1.44054 ) ( 2 1.42211 ) ( 3 1.37516 )}
CLUSTER-2 {size: 2, centroid:( 1 1.3743 ) ( 2 1.36902 ) ( 3 1.36431 )}
CLUSTER-3 {size: 1, centroid:( 1 0.756473 ) ( 2 0.723891 ) ( 3 0.699379 )}
CLUSTER-4 {size: 1, centroid:( 1 0.981648 ) ( 2 0.98652 ) ( 3 0.997958 )}
CLUSTER-5 {size: 4, centroid:( 1 1.64981 ) ( 2 1.64336 ) ( 3 1.61365 )}
clustering_time: 683
Silhouette: [0.491238,0.20537,1,1,0.491479,0.513576]
```

Time-Series LSH Discrete Frechet Clustering

```

Algorithm: LSH Frechet Discrete
CLUSTER-1 {size: 13, centroid:( 1 1.36119 ) ( 2 1
CLUSTER-2 {size: 4, centroid:( 1 1.64912 ) ( 2 1.
CLUSTER-3 {size: 1, centroid:( 1 1.76472 ) ( 2 1.
CLUSTER-4 {size: 1, centroid:( 1 0.756473 ) ( 2 0
CLUSTER-5 {size: 1, centroid:( 1 1.76175 ) ( 2 1.
clustering_time: 2103
Silhouette: [0.216118,-1,1,1,1,0.0904769]
CLUSTER-1 {( 1 1.36119 ) ( 2 1.34036 ) ( 3 1.2941
CLUSTER-2 {( 1 1.64912 ) ( 2 1.64334 ) ( 3 1.6067
CLUSTER-3 {( 1 1.76472 ) ( 2 1.76463 ) ( 3 1.7653
CLUSTER-4 {( 1 0.756473 ) ( 2 0.723891 ) ( 3 0.69
CLUSTER-5 {( 1 1.76175 ) ( 2 1.75234 ) ( 3 1.7475

```

Time-Serie LSH Discrete Frechet Clustering Compressed

Τα μεγέθη των συστάδων συγκλίνουν στους ίδιους αριθμούς. Υπάρχουν αρκετές κοινές χρονοσειρές εντός των συστάδων των δύο αναπαραστάσεων.

## Συσταδοποίηση Χρονοσειρών με Lloyds Discrete Frechet

Τα δειγματοληπτικά αποτελέσματα:

```
Algorithm: Lloyds
CLUSTER-1 {size: 4, centroid:( 1 1.6917 ) ( 2 1.68812 ) ( 3 1.67302 ) ( 4 1.6537 ) ( 5 1.6302 )}
CLUSTER-2 {size: 7, centroid:( 1 1.1586 ) ( 2 1.14377 ) ( 3 1.10754 ) ( 4 1.1012 ) ( 5 1.0851 ) ( 6 1.06937 ) ( 7 1.05851 )}
CLUSTER-3 {size: 5, centroid:( 1 1.65414 ) ( 2 1.64389 ) ( 3 1.58773 ) ( 4 1.55851 ) ( 5 1.54671 )}
CLUSTER-4 {size: 2, centroid:( 1 1.35846 ) ( 2 1.35065 ) ( 3 1.35969 ) ( 4 1.34671 ) ( 5 1.34671 )}
CLUSTER-5 {size: 2, centroid:( 1 1.00558 ) ( 2 0.972113 ) ( 3 0.894748 ) ( 4 0.862415 ) ( 5 0.862415 )}
clustering_time: 55
Silhouette: [0.394106,0.156633,0.324456,0.429248,0.282979,0.28598]
```

Time-Serie Lloyds Discrete Frechet Clustering

```
Algorithm: Lloyds
CLUSTER-1 {size: 5, centroid:( 1 1.46966 ) ( 2 1.4547 ) ( 3 1.42179 ) ( 4 1.39658 ) ( 5 1.3658 )}
CLUSTER-2 {size: 3, centroid:( 1 1.07983 ) ( 2 1.05615 ) ( 3 0.981994 ) ( 4 0.975112 ) ( 5 0.95112 )}
CLUSTER-3 {size: 1, centroid:( 1 0.756473 ) ( 2 0.723891 ) ( 3 0.699379 ) ( 4 0.656542 ) ( 5 0.656542 )}
CLUSTER-4 {size: 9, centroid:( 1 1.70961 ) ( 2 1.71178 ) ( 3 1.6985 ) ( 4 1.68552 ) ( 5 1.68552 ) ( 6 1.68552 ) ( 7 1.68552 ) ( 8 1.68552 ) ( 9 1.68552 )}
CLUSTER-5 {size: 2, centroid:( 1 1.35846 ) ( 2 1.35065 ) ( 3 1.35969 ) ( 4 1.34671 ) ( 5 1.34671 )}
clustering_time: 55
Silhouette: [0.083162,0.401113,1,0.519625,0.482735,0.413062]
```

Time-Serie Lloyds Discrete Frechet Clustering Compressed

Τα μεγέθη των συστάδων διαφέρουν σχετικά με τις προηγούμενες μετρικές. Ωστόσο κάποιες χρονοσειρές είναι κοινές μεταξύ των συστάδων.

## Συμπεράσματα

Οι πειραματικές εκτελέσεις των δύο αναπαραστάσεων επιφέρουν άρτια αποτελέσματα κατά τη συσταδοποίηση. Στη πλειοψηφία των περιπτώσεων υπάρχουν αρκετά κοινά μέλη μεταξύ των συστάδων των δύο αναπαραστάσεων. Η συγκεκριμένη συμπεριφορά δεν αναγνωρίζεται εύκολα κατά την αναζήτηση του πλησιέστερου γείτονα. Η αναγνώριση θα ήταν ευκολότερη στην εμφάνιση των  $N$  πλησιέστερων γειτόνων. Η πιθανότητα όμοιου κοντινότερου γείτονα δεν κυριαρχεί, συνεπώς κατά τους  $N$  κοντινότερους γείτονες θα υπάρχουν αρκετές κοινές αναφορές.