



Τοπολογική Ανάλυση Δεδομένων

Εαρινό Εξάμηνο 2026

Φροντιστηριακή Διάλεξη 4: Μηχανική Μάθηση & Τοπολογική
Ανάλυση Δεδομένων

Διδάσκοντες:

Αθανάσιος Ανδρικόπουλος

Ιωάννης Γουναρίδης

Επικουρικό Έργο: Αλέξανδρος - Μάριος Αφράτης
Πάτρα, 26 Μαρτίου 2026

Πίνακας Περιεχομένων

Πλαίσιο Επιβλεπόμενης Μάθησης

Σύνολο Δεδομένων MNIST & Προεπεξεργασία

Μετρικές Αξιολόγησης

Αλγόριθμοι Ταξινόμησης

Σύγκριση Μοντέλων

Ανάλυση Σημασίας Χαρακτηριστικών

Τροχιές Δυναμικών Συστημάτων

Αναπαραστάσεις Διαγραμμάτων Επιμονής

Πυρήνες Διαγραμμάτων Επιμονής

Ολοκληρωμένη Προσέγγιση

Σύνοψη

Κύριες Εισαγωγές

```
1 import numpy as np
2 import pandas as pd
3 import matplotlib.pyplot as plt
4 import seaborn as sns
5 from sklearn.datasets import fetch_openml
6 from sklearn.model_selection import train_test_split
7 from sklearn.preprocessing import StandardScaler
8 from sklearn.metrics import accuracy_score, classification_report, confusion_matrix
9 import time
10 import warnings
11 warnings.filterwarnings('ignore')
12
13 # Ορισμός παραμέτρου ψευδοστοχαστικότητας
14 np.random.seed(42)
15 plt.style.use('seaborn-v0_8-whitegrid')
16 plt.rcParams['figure.figsize'] = [10, 6]
17
```

Ορισμός

Δίνεται σύνολο εκπαίδευσης $\mathcal{D} = \{(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_n, y_n)\}$ όπου $\mathbf{x}_i \in \mathbb{R}^d$ είναι το **διάνυσμα χαρακτηριστικών** (feature vector) και $y_i \in \mathcal{Y}$ η **ετικέτα** (label). Στόχος: εκμάθηση συνάρτησης $f: \mathbb{R}^d \rightarrow \mathcal{Y}$.

Ορισμός

Ταξινόμηση (Classification): επιβλεπόμενη μάθηση όπου $\mathcal{Y} = \{c_1, c_2, \dots, c_k\}$ είναι πεπερασμένο διακριτό σύνολο. Κάθε c_k ονομάζεται **κλάση** (class).

Ορισμός

Το **χώρος υποθέσεων** (hypothesis space) \mathcal{H} είναι το σύνολο όλων των δυνατών συναρτήσεων που μπορεί να εκμάθει ο αλγόριθμος:

- Γραμμικοί ταξινομητές:
 $f(\mathbf{x}) = \text{sign}(\mathbf{w}^T \mathbf{x} + b)$
- Δέντρα απόφασης: τμηματικά σταθερές συναρτήσεις

Πίνακας Περιεχομένων

Πλαίσιο Επιβλεπόμενης Μάθησης

Σύνολο Δεδομένων MNIST & Προεπεξεργασία

Μετρικές Αξιολόγησης

Αλγόριθμοι Ταξινόμησης

Σύγκριση Μοντέλων

Ανάλυση Σημασίας Χαρακτηριστικών

Τροχιές Δυναμικών Συστημάτων

Αναπαραστάσεις Διαγραμμάτων Επιμονής

Πυρήνες Διαγραμμάτων Επιμονής

Ολοκληρωμένη Προσέγγιση

Σύνοψη

Σύνολο Δεδομένων MNIST (MNIST Dataset)

Περιγραφή (Description)

Το MNIST αποτελείται από **γκρι εικόνες** 28×28 **pixels** χειρόγραφων ψηφίων (0–9).

Κάθε εικόνα αναπαρίσταται ως:

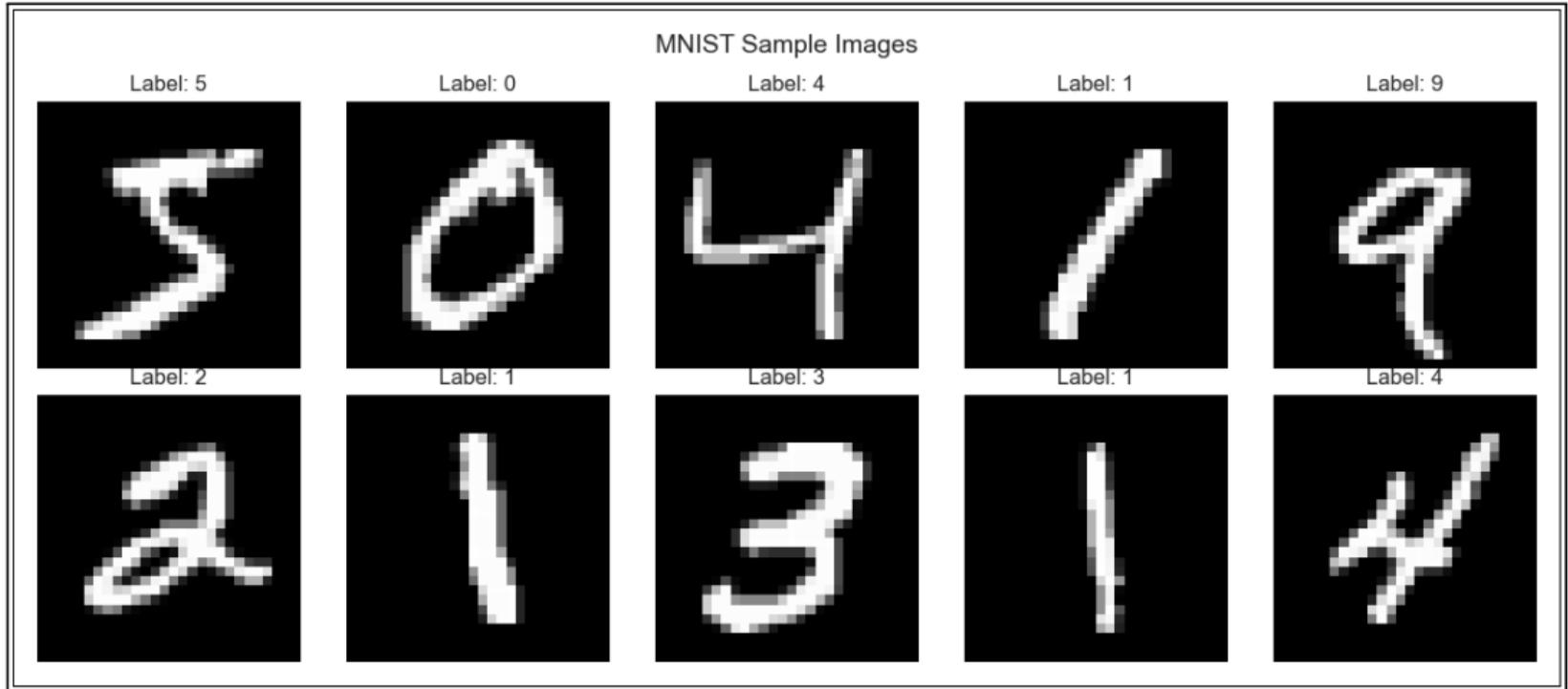
- Μητρώο $\mathbf{X} \in \mathbb{R}^{28 \times 28}$, $X_{ij} \in [0, 255]$ (ένταση pixel)
- ή ισοδύναμα ως εκτεταμένο διάνυσμα $\mathbf{x} \in \mathbb{R}^{784}$

Φόρτωση (Loading)

```
1 mnist = fetch_openml('mnist_784', version=1, as_frame=False)
2 X_full = mnist.data # shape: (70000, 784)
3 y_full = mnist.target.astype(int)
4 print(f"Dataset loaded: {X_full.shape[0]} samples, {X_full.shape[1]} features")
5
```

Σημείωση

Σύνολο δεδομένων: **70.000 δείγματα**, 784 χαρακτηριστικά (pixels), 10 κλάσεις (ψηφία 0–9).



Σχ. 4.1: Σύνολο Δεδομένων MNIST

Επιλογή Υποσυνόλου (Subset Selection)

Επιλογή 5 Οπτικά Όμοιων Ψηφίων

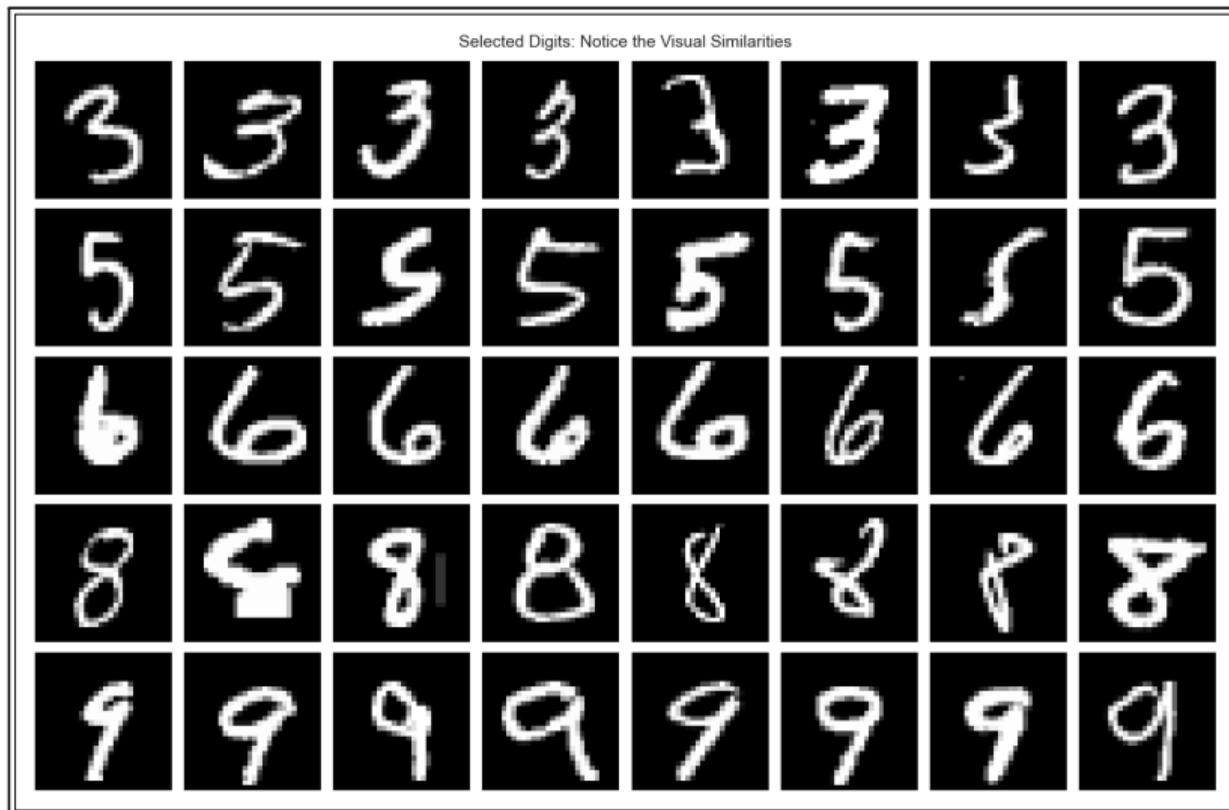
Επιλέγουμε $\mathcal{Y} = \{3, 5, 6, 8, 9\}$ – ψηφία που μοιάζουν οπτικά και δυσκολεύουν την ταξινόμηση.

```
1  selected_digits = [3, 5, 6, 8, 9]
2  samples_per_class = 1000 # 1000 δείγματα ανά κλάση
3
4  X_list, y_list = [], []
5  for digit in selected_digits:
6  mask = y_full == digit
7  indices = np.random.choice(np.where(mask)[0], samples_per_class, replace=False)
8  X_list.append(X_full[indices])
9  y_list.append(y_full[indices])
10
11 X = np.vstack(X_list) # shape: (5000, 784)
12 y = np.hstack(y_list) # shape: (5000,)
13 print(f"Selected dataset: {X.shape[0]} samples, Classes: {selected_digits}")
14
```

Κίνητρο Χρήσης TDA

Τα ψηφία 6 και 9 έχουν από μία οπή (similar topology), ενώ τα 3 και 5 δεν έχουν οπές. Αυτή η δυσκολία ταξινόμησης επιλύεται αποδοτικά μέσω TDA.

Οπτικοποίηση (Visualization)



Σχ. 4.2: Επιλογή Υποσυνόλου

Κανονικοποίηση Χαρακτηριστικών (Feature Normalization)

Ορισμός

Κανονικοποίηση Min-Max (Min-Max Normalization): για χαρακτηριστικό x με εύρος $[x_{\min}, x_{\max}]$:

$$\tilde{x} = \frac{x - x_{\min}}{x_{\max} - x_{\min}}$$

Για τιμές pixel: $x \in [0, 255] \rightarrow \tilde{x} \in [0, 1]$.

Γιατί Κανονικοποίηση;

- Πολλοί αλγόριθμοι (SVM, K-NN) χρησιμοποιούν μετρικές αποστάσεων ευαίσθητες στις κλίμακες
- Η βελτιστοποίηση βαθμίδας (gradient-based optimization) συγκλίνει ταχύτερα με κανονικοποιημένα χαρακτηριστικά
- Αποτρέπει χαρακτηριστικά με μεγάλο εύρος να κυριαρχούν

Παράδειγμα Κανονικοποίησης:

```
1 x_normalized = x / 255.0 # [0,255] -> [0,1]
```

Ορισμός

Βλάβη Γενικότητας – Σφάλμα Γενίκευσης
(Generalization Error): αναμενόμενο σφάλμα σε προηγουμένως άγνωστες εισόδους (δεδομένα):

$$R(\hat{f}) = \mathbb{E}_{(\mathbf{x}, y) \sim P} [\mathcal{L}(f(\mathbf{x}), y)]$$

Ορισμός

Διαχωρισμός εκπαίδευσης-ελέγχου:

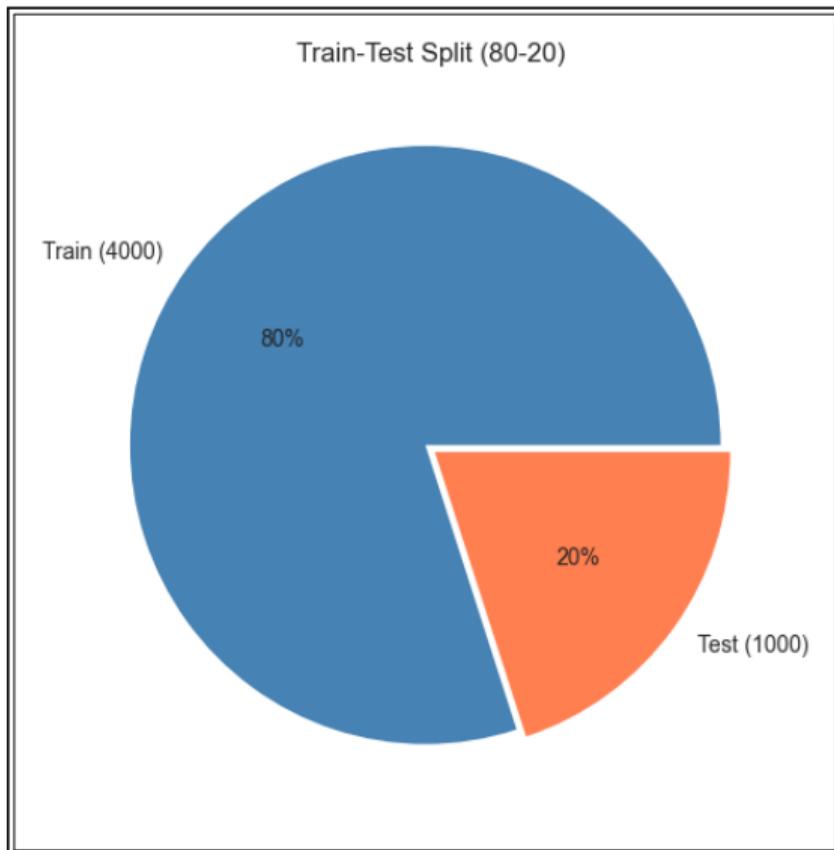
$$\mathcal{D} = \mathcal{D}_{\text{train}} \cup \mathcal{D}_{\text{test}}, \quad \mathcal{D}_{\text{train}} \cap \mathcal{D}_{\text{test}} = \emptyset.$$

Κώδικας (80% – 20%)

```
1 X_train, X_test, y_train, y_test =  
2 train_test_split(  
3 X_normalized, y,  
4 test_size=0.2,  
5 random_state=42,  
6 stratify=y)  
7 # Training set: 4000 samples  
8 # Test set:      1000 samples  
9
```

Διαισθητικά

Η εκπαίδευση θυμίζει μελέτη παλαιών θεμάτων εξεταστικής και ο έλεγχος είναι σαν τελική εξέταση με νέες ερωτήσεις.



Σχ. 4.3: Διαχωρισμός Εκπαίδευσης – Ελέγχου

Μείωση Διαστατάσεων για Οπτικοποίηση (PCA)

Ορισμός

Ανάλυση Κύριων Συνιστωσών (*Principal Component Analysis – PCA*): βρίσκει ορθογώνιες κατευθύνσεις μέγιστης διακύμανσης. Για κεντραρισμένο πίνακα $\mathbf{X} \in \mathbb{R}^{n \times d}$,

$$\mathbf{X} = \mathbf{UV}^T, \quad \mathbf{Z} = \mathbf{XV}_k \in \mathbb{R}^{n \times k}$$

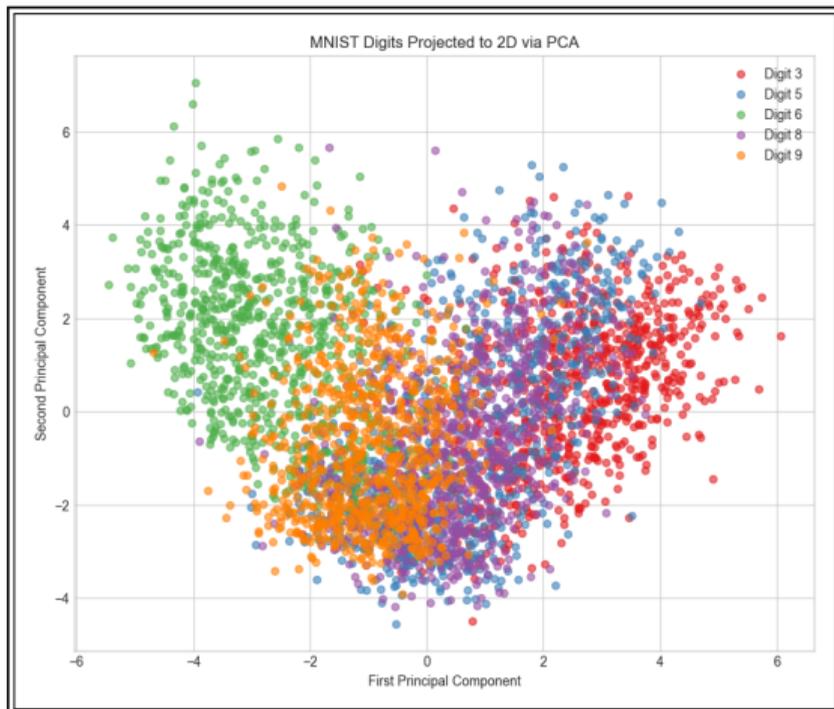
Για οπτικοποίηση: $k = 2$.

Κώδικας PCA

```
1 from sklearn.decomposition import PCA
2
3 pca = PCA(n_components=2)
4 X_train_2d = pca.fit_transform(
5 X_train)
6 print(f"Variance explained: "
7 f"{pca.explained_variance_ratio_
8 .sum()*100:.1f}%")
9
```

Ερμηνεία

Τα ψηφία 6 και 9 (μία οπή) ομαδοποιούνται μαζί, όπως και τα 3 και 5 (χωρίς οπές) – η τοπολογική δομή επηρεάζει την κατανομή!



Σχ. 4.4: Μείωση Διαστάσεων για Οπτικοποίηση (PCA)

Πίνακας Περιεχομένων

Πλαίσιο Επιβλεπόμενης Μάθησης

Σύνολο Δεδομένων MNIST & Προεπεξεργασία

Μετρικές Αξιολόγησης

Αλγόριθμοι Ταξινόμησης

Σύγκριση Μοντέλων

Ανάλυση Σημασίας Χαρακτηριστικών

Τροχιές Δυναμικών Συστημάτων

Αναπαραστάσεις Διαγραμμάτων Επιμονής

Πυρήνες Διαγραμμάτων Επιμονής

Ολοκληρωμένη Προσέγγιση

Σύνοψη

Ορισμός

Πίνακας Σύγχυσης (Confusion Matrix): $C \in \mathbb{R}^{K \times K}$ όπου C_{ij} = αριθμός δειγμάτων με αληθινή ετικέτα i που προβλέφθηκαν ως j .

- **Διαγώνιες εγγραφές** C_{ii} : σωστές προβλέψεις (True Positives κλάσης i)
- **Εκτός διαγωνίου** C_{ij} ($i \neq j$): σφάλματα ταξινόμησης

Ορισμός

Ακρίβεια (Accuracy):

$$\text{Accuracy} = \frac{\text{Αριθμός σωστών προβλέψεων}}{\text{Συνολικός αριθμός προβλέψεων}} = \frac{\sum_c TP_c}{n}$$

Προσοχή: η ακρίβεια μπορεί να είναι παραπλανητική σε μη ισορροπημένα σύνολα δεδομένων (imbalanced datasets).

Ορισμός

Ακρίβεια ανά κλάση (Precision):

$$\text{Precision}_c = \frac{TP_c}{TP_c + FP_c}$$

Από όλα τα δείγματα που προβλέφθηκαν ως c , ποιο ποσοστό όντως ανήκει σε c ;

Ορισμός

Ανάκληση (Recall):

$$\text{Recall}_c = \frac{TP_c}{TP_c + FN_c}$$

Από όλα τα δείγματα που ανήκουν σε c , ποιο ποσοστό εντοπίστηκε σωστά;

Ορισμός

F1-Score:

$$F1_c = 2 \cdot \frac{\text{Precision}_c \cdot \text{Recall}_c}{\text{Precision}_c + \text{Recall}_c}$$

Αρμονικός μέσος ακρίβειας και ανάκλησης – τιμωρεί οριακές τιμές.

Συμβολισμός

- TP_c : True Positives (σωστά ως c)
- FP_c : False Positives (λάθος ως c)
- FN_c : False Negatives (c προβλεφθέν αλλιού)

Βοηθητικές Συναρτήσεις Αξιολόγησης (Evaluation Helper Functions)

```
1  from sklearn.metrics import precision_score, recall_score, f1_score
2
3  results = {'Algorithm': [], 'Accuracy': [], 'Precision': [],
4            'Recall': [], 'F1-Score': [], 'Training Time (s)': []}
5
6  def print_metrics(y_true, y_pred, name):
7      """Εκτύπωση όλων των μετρικών αξιολόγησης για ένα μοντέλο."""
8      acc = accuracy_score(y_true, y_pred)
9      prec = precision_score(y_true, y_pred, average='weighted')
10     rec = recall_score(y_true, y_pred, average='weighted')
11     f1 = f1_score(y_true, y_pred, average='weighted')
12     print(f"\n{name} - Evaluation Metrics:")
13     print(f"Accuracy:  {acc:.4f} ({acc*100:.2f}%)")
14     print(f"Precision:  {prec:.4f}")
15     print(f"Recall:      {rec:.4f}")
16     print(f"F1-Score:   {f1:.4f}")
17     return acc, prec, rec, f1
18
```

Πίνακας Περιεχομένων

Πλαίσιο Επιβλεπόμενης Μάθησης

Σύνολο Δεδομένων MNIST & Προεπεξεργασία

Μετρικές Αξιολόγησης

Αλγόριθμοι Ταξινόμησης

Σύγκριση Μοντέλων

Ανάλυση Σημασίας Χαρακτηριστικών

Τροχιές Δυναμικών Συστημάτων

Αναπαραστάσεις Διαγραμμάτων Επιμονής

Πυρήνες Διαγραμμάτων Επιμονής

Ολοκληρωμένη Προσέγγιση

Σύνοψη

K-Πλησιέστεροι Γείτονες – Θεωρία (K-Nearest Neighbors)

Ορισμός

Για ερώτημα \mathbf{x}_q , έστω $N_k(\mathbf{x}_q)$ οι k πλησιέστεροι γείτονες στο εκπαιδευτικό σύνολο.

Η πρόβλεψη K-NN:

$$\hat{y} = \arg \max_{c \in \mathcal{Y}} \sum_{(\mathbf{x}_i, y_i) \in N_k(\mathbf{x}_q)} \begin{cases} 1 & \text{αν } y_i = c \\ 0 & \text{αν } y_i \neq c \end{cases}$$

Κοινές μετρικές αποστάσεων (*distance metrics*):

- **Ευκλείδεια:** $d(\mathbf{x}, \mathbf{x}') = \|\mathbf{x} - \mathbf{x}'\|_2$
- **Manhattan:** $d(\mathbf{x}, \mathbf{x}') = \|\mathbf{x} - \mathbf{x}'\|_1$

Ιδιότητες

- **Μη-παραμετρικός** (non-parametric): χωρίς παραδοχές για κατανομή
- **Lazy learning:** χωρίς ρητή φάση εκπαίδευσης

Διαισθητικά

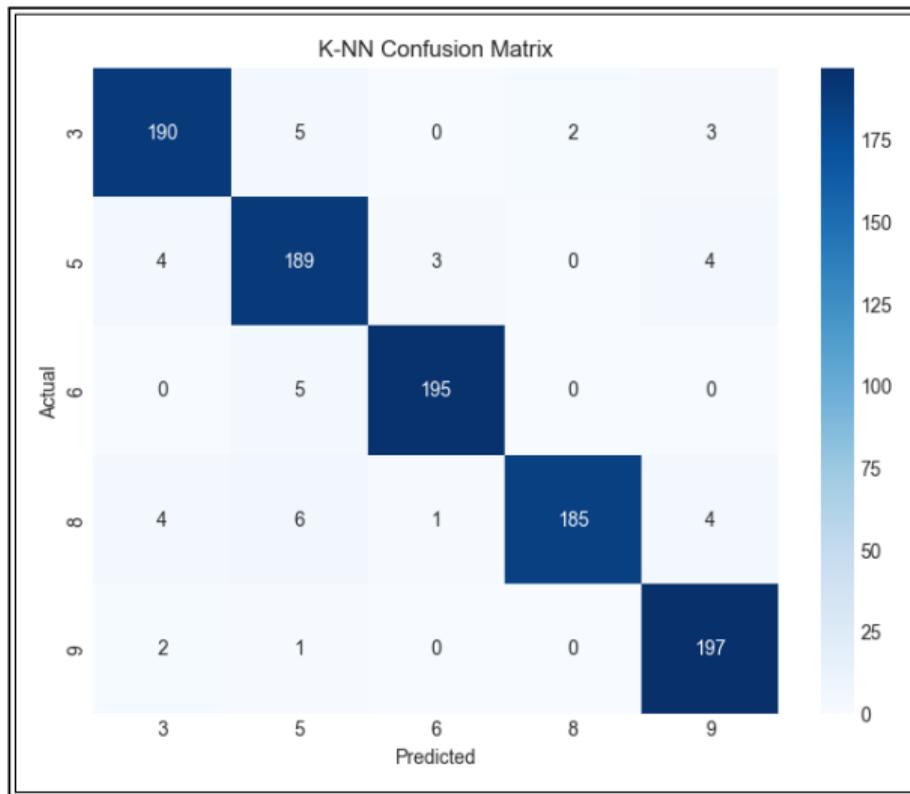
“Είσαι ο μέσος όρος των k κοντινότερων φίλων σου.” (Ταξινόμηση με πλειοψηφική ψηφοφορία γειτόνων).

K-NN: Κώδικας & Αποτελέσματα

```
1  from sklearn.neighbors import KNeighborsClassifier
2
3  # Εκπαίδευση K-NN με k=5
4  knn = KNeighborsClassifier(n_neighbors=5)
5  start = time.time()
6  knn.fit(X_train, y_train)
7  knn_time = time.time() - start
8
9  knn_pred = knn.predict(X_test)
10
11 # Εκτύπωση μετρικών
12 knn_acc, knn_prec, knn_rec, knn_f1 = print_metrics(y_test, knn_pred, 'K-NN (k=5)')
13 print(f"Training Time: {knn_time:.4f}s")
14
15 # Αποθήκευση αποτελεσμάτων για σύγκριση
16 results['Algorithm'].append('K-NN')
17 results['Accuracy'].append(knn_acc)
18 results['Precision'].append(knn_prec)
19 results['Recall'].append(knn_rec)
20 results['F1-Score'].append(knn_f1)
21 results['Training Time (s)'].append(knn_time)
22
```

Ερμηνεία

Ο K-NN επιτυγχάνει υψηλή απόδοση ($\approx 95-96\%$). Τα περισσότερα σφάλματα εμφανίζονται μεταξύ οπτικά όμοιων ζευγών: 3 \leftrightarrow 5 και 3 \leftrightarrow 8.



Σχ. 4.5: K-NN Confusion Matrix

Συμβιβασμός Πόλωσης – Διακύμανσης (Bias-Variance Tradeoff)

- **Μικρό** k : χαμηλή πόλωση (bias), υψηλή διακύμανση (variance) – ευαίσθητο στο θόρυβο
- **Μεγάλο** k : υψηλή πόλωση, χαμηλή διακύμανση – υπερεξομάλυνση (over-smoothing)

Μέθοδος Αγκώνα (Elbow Method)

Μέθοδος Αγκώνα (Elbow Method)

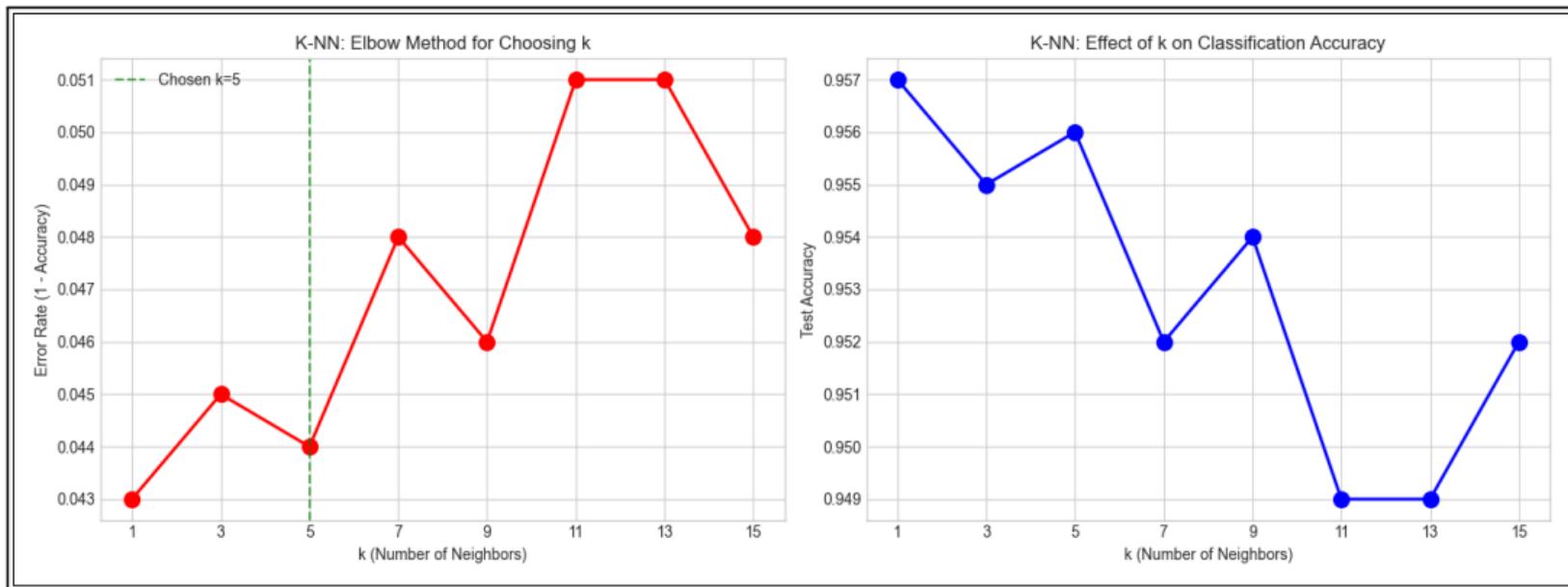
Σχεδιάζουμε το ποσοστό σφάλματος $1 - \text{accuracy}$ ως συνάρτηση του k . Αναζητούμε το σημείο "αγκώνα" όπου το σφάλμα σταματά να μειώνεται σημαντικά.

Κώδικας

```
1 k_values = [1, 3, 5, 7, 9, 11, 13, 15]
2 k_errors = []
3 for k in k_values:
4     knn_temp = KNeighborsClassifier(n_neighbors=k)
5     knn_temp.fit(X_train, y_train)
6     acc = accuracy_score(y_test, knn_temp.predict(X_test))
7     k_errors.append(1 - acc)
8     print(f"k={k:2d}: Accuracy={acc:.4f}, Error Rate={1-acc:.4f}")
9
```

Ερμηνεία

Για αυτό το σύνολο δεδομένων, $k = 1$ έως $k = 5$ δίνουν παρόμοια απόδοση ($\approx 95-96\%$) – επιλέγουμε $k = 5$ ως ισορροπία τοπικής δομής και ανθεκτικότητας στο θόρυβο.



Σχ. 4.6: K-NN: Επίδραση της Παραμέτρου k

Naive Bayes: Θεωρητικό Υπόβαθρο

Ορισμός

Θεώρημα Bayes:

$$P(y|\mathbf{x}) = \frac{P(\mathbf{x}|y)P(y)}{P(\mathbf{x})}$$

Ορισμός

Ταξινομητής Naive Bayes: υποθέτει **δεσμευμένη ανεξαρτησία** (*conditional independence*) χαρακτηριστικών δοθείσης κλάσης:

$$P(\mathbf{x}|y) = \prod_{j=1}^d P(x_j|y)$$

Πρόβλεψη: $\hat{y} = \arg \max_c P(y=c) \prod_j P(x_j|y=c)$

Ορισμός

Gaussian Naive Bayes: κάθε $P(x_j|y=c)$ ακολουθεί κατανομή Gauss:

$$P(x_j|y=c) = \frac{1}{\sqrt{2\pi\sigma_{jc}^2}} \exp\left(-\frac{(x_j - \mu_{jc})^2}{2\sigma_{jc}^2}\right)$$

Διαισθητικά

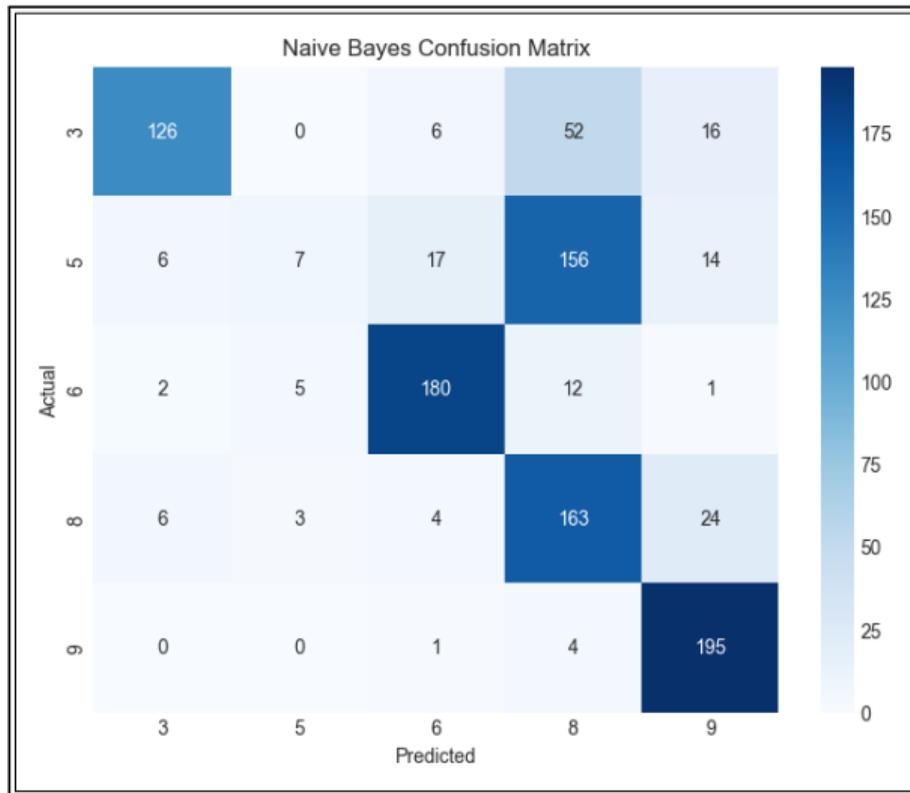
“Ποια είναι η πιθανότητα αυτό να είναι 3, δεδομένου πώς μοιάζουν τα τυπικά 3;”
Παρά την αφελή υπόθεση ανεξαρτησίας, αποδίδει καλά στην πράξη.

Naive Bayes: Κώδικας & Αποτελέσματα

```
1  from sklearn.naive_bayes import GaussianNB
2
3  # Εκπαίδευση Gaussian Naive Bayes
4  nb = GaussianNB()
5  start = time.time()
6  nb.fit(X_train, y_train)
7  nb_time = time.time() - start
8
9  nb_pred = nb.predict(X_test)
10
11 # Εκτύπωση μετρικών
12 nb_acc, nb_prec, nb_rec, nb_f1 = print_metrics(y_test, nb_pred, 'Naive Bayes')
13 print(f"Training Time: {nb_time:.4f}s (very fast!)")
14
15 # Αποθήκευση αποτελεσμάτων
16 results['Algorithm'].append('Naive Bayes')
17 results['Accuracy'].append(nb_acc)
18
```

Χαρακτηριστικά Naive Bayes

- **Πλεονέκτημα:** εξαιρετικά γρήγορη εκπαίδευση – υπολογίζει μόνο μέσους και διακυμάνσεις ανά χαρακτηριστικό και κλάση
- **Μειονέκτημα:** η υπόθεση ανεξαρτησίας χαρακτηριστικών σπάνια ισχύει για pixels (γειτονικά pixels συσχετίζονται)



Σχ. 4.7: Naive Bayes Confusion Matrix

Λογιστική Παλινδρόμηση (Logistic Regression)

Ορισμός

Λογιστική Παλινδρόμηση (δυναμική περίπτωση): μοντελοποιεί την πιθανότητα κλάσης 1 μέσω της **λογιστικής (sigmoid) συνάρτησης**:

$$P(y=1|\mathbf{x}) = \sigma(\mathbf{w}^T \mathbf{x} + b) = \frac{1}{1 + e^{-(\mathbf{w}^T \mathbf{x} + b)}}$$

Για K κλάσεις: **softmax** πολυκλαδική επέκταση.

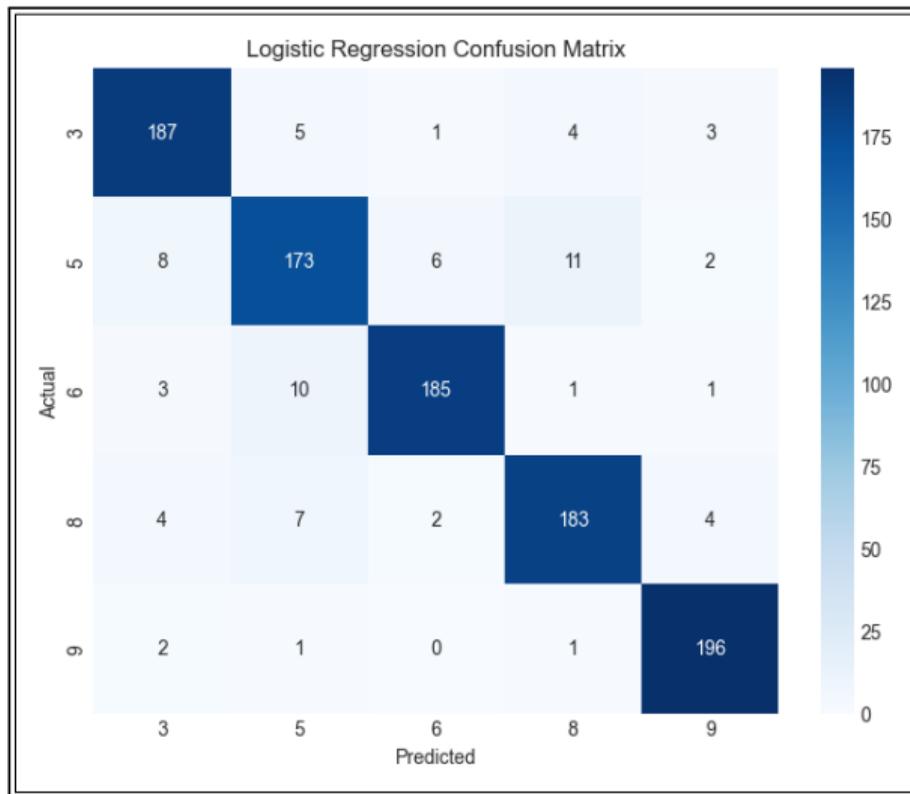
Οπτικοποίηση Συντελεστών (Coefficient Visualization)

Κάθε κλάση έχει διάνυσμα βαρών $\mathbf{w}_k \in \mathbb{R}^{784}$. Αναδιαμόρφωση σε 28×28 αποκαλύπτει ποιες περιοχές pixels συμβάλλουν σε κάθε κλάση.

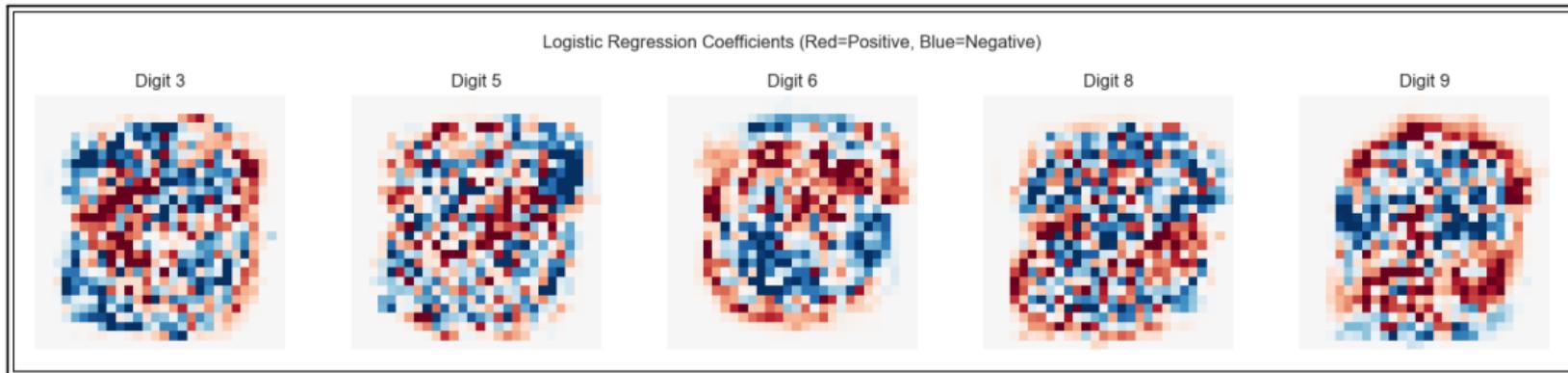
Κώδικας

```
1 from sklearn.linear_model \
2 import LogisticRegression
3
4 lr = LogisticRegression(
5     max_iter=1000,
6     random_state=42)
7 start = time.time()
8 lr.fit(X_train, y_train)
9 lr_time = time.time() - start
10 lr_pred = lr.predict(X_test)
11 lr_acc, lr_prec, lr_rec, lr_f1\
12 = print_metrics(y_test, lr_pred,
13     'Logistic Regression')
14
```

```
1 for i, digit in enumerate(selected_digits):
2     coef = lr.coef_[i].reshape(28, 28)
3     # Οπτικοποίηση με cmap='RdBu' κόκκινοθετικό(, μπλεαρνητικό)
4
```



Σχ. 4.8: Logistic Regression Confusion Matrix



**Σχ. 4.9: Συντελεστές Λογιστικής Παλινδρόμησης
(Κόκκινο = Θετικό, Μπλε = Αρνητικό)**

Δέντρα Απόφασης (Decision Trees)

Ορισμός

Δέντρο Απόφασης: αναδρομικώς διαμερίζει το χώρο χαρακτηριστικών με διαχωρισμούς παράλληλους στους άξονες. Σε κάθε εσωτερικό κόμβο:

$$(j^*, t^*) = \arg \min_{j, t} \mathcal{L}(\mathcal{D}_L) + \mathcal{L}(\mathcal{D}_R)$$

Ορισμός

Νοθεία Gini: $G(\mathcal{D}) = 1 - \sum_{k=1}^K p_k^2$

Εντροπία (Entropy): $H(\mathcal{D}) = - \sum_{k=1}^K p_k \log_2 p_k$

Ορισμός

Σημασία Χαρακτηριστικού (Feature Importance):

$$Imp(j) = \sum_{\text{κόμβοι } t \text{ που διαχωρίζουν σε } j} n_t \cdot \Delta \mathcal{L}_t$$

Υπερπροσαρμογή (Overfitting)

Τα δέντρα μεγάλου βάθους απομνημονεύουν σε προβληματικό βαθμό τα δεδομένα εκπαίδευσης:

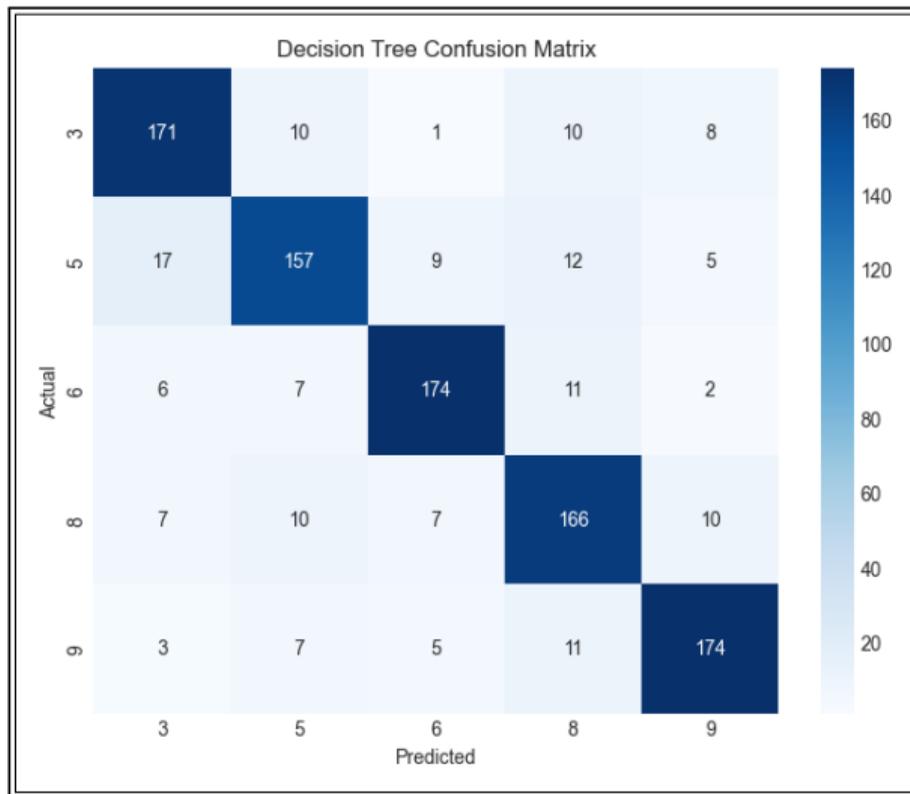
$$R_{train}(f) \ll R_{test}(f)$$

Δέντρα Απόφασης: Κώδικας & Υπερπροσαρμογή

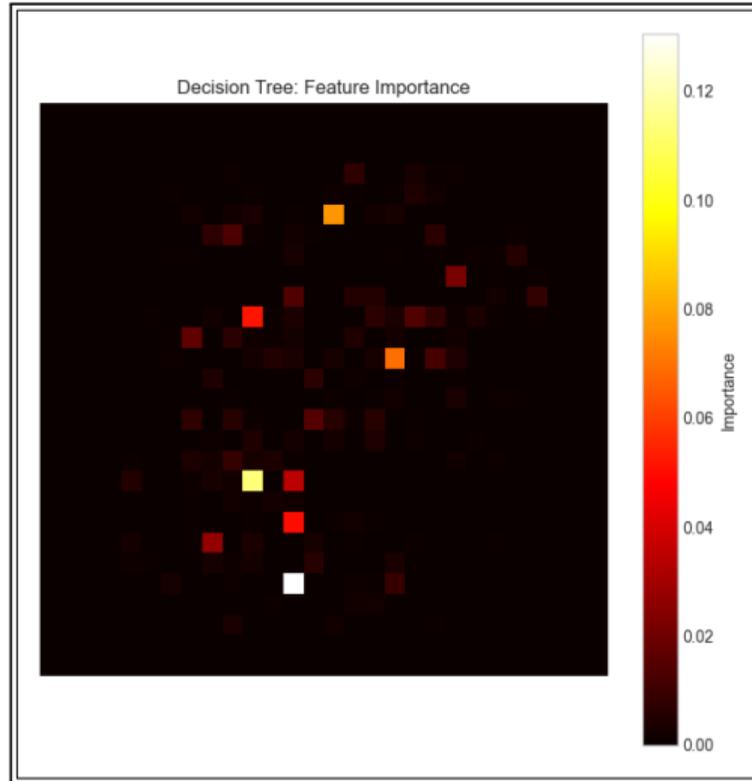
```
1 from sklearn.tree import DecisionTreeClassifier
2
3 # Εκπαίδευση Δέντρου Απόφασης (max_depth=10)
4 dt = DecisionTreeClassifier(max_depth=10, random_state=42)
5 dt.fit(X_train, y_train)
6 dt_pred = dt.predict(X_test)
7 dt_acc, dt_prec, dt_rec, dt_f1 = print_metrics(y_test, dt_pred, 'Decision Tree')
8
9 # Απόδειξη υπερπροσαρμογής για διάφορα βάθη
10 depths = [3, 5, 10, 15, None]
11 for depth in depths:
12     dt_temp = DecisionTreeClassifier(max_depth=depth, random_state=42)
13     dt_temp.fit(X_train, y_train)
14     train_acc = accuracy_score(y_train, dt_temp.predict(X_train))
15     test_acc = accuracy_score(y_test, dt_temp.predict(X_test))
16     depth_str = str(depth) if depth else 'None (unrestricted)'
17     print(f"Depth={depth_str}: Train={train_acc:.4f}, Test={test_acc:.4f}")
18
```

Παρατήρηση: Υπερπροσαρμογή

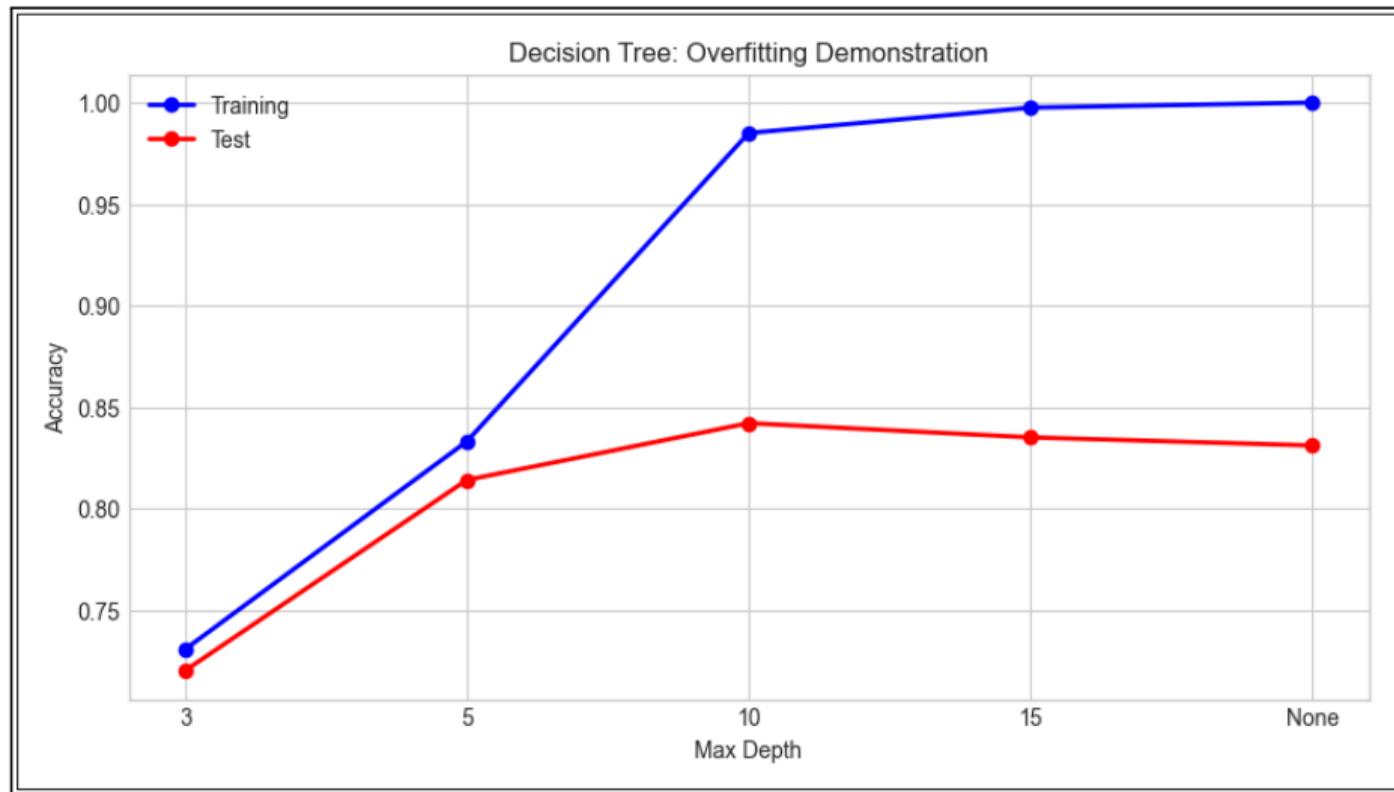
Καθώς το βάθος αυξάνεται, η ακρίβεια εκπαίδευσης προσεγγίζει 100% ενώ η ακρίβεια ελέγχου κορυφώνεται γύρω στο βάθος 10–15 και στη συνέχεια μειώνεται. Αυτό είναι το χαρακτηριστικό σύμπτωμα υπερπροσαρμογής.



Σχ. 4.10: Decision Tree Confusion Matrix



Σχ. 4.11: Decision Tree: Feature Importance



Σχ. 4.12: Decision Tree: Overfitting

Τυχαίο Δάσος (Random Forest)

Ορισμός

Μέθοδος Ensemble: συνδυάζει M βασικούς εκμαθητές για μείωση διακύμανσης:

$$\hat{y} = \arg \max_{c \in \mathcal{Y}} \sum_{m=1}^M \mathbb{1}[f_m(\mathbf{x}) = c]$$

Τυχαίο Δάσος (Random Forest): ensemble δέντρων απόφασης, κάθε δέντρο εκπαιδεύεται σε:

1. **Bootstrap δείγμα:** τυχαίο δείγμα με επανάθεση
2. **Τυχαίο υποσύνολο χαρακτηριστικών:** σε κάθε διαχωρισμό, μόνο \sqrt{d} τυχαία χαρακτηριστικά

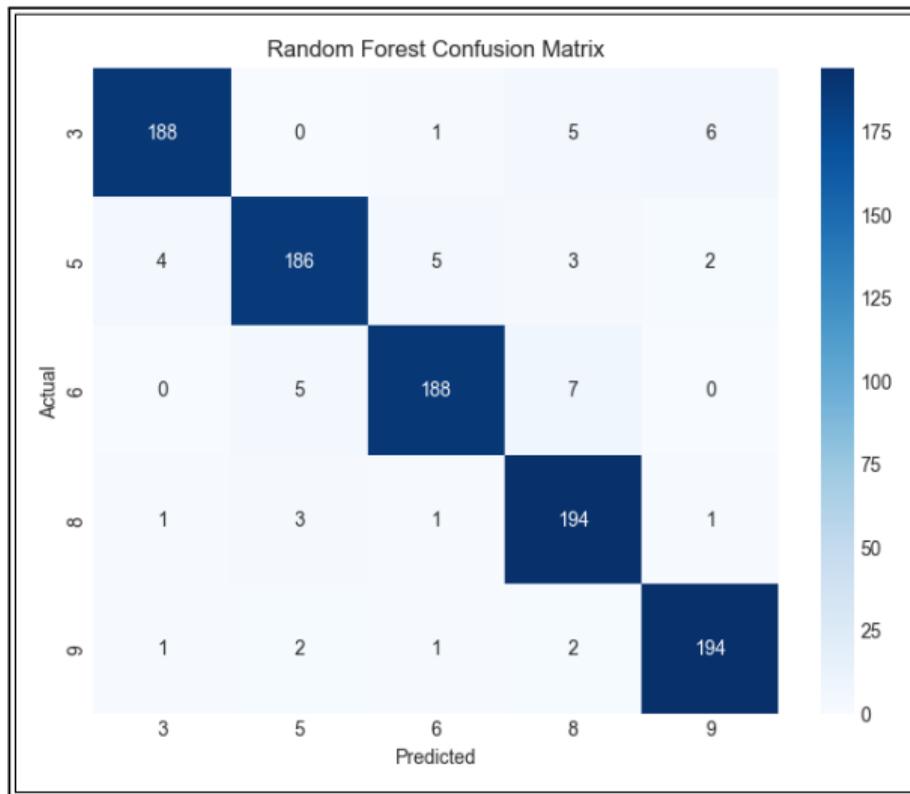
Θεώρημα

Μείωση Διακύμανσης: για M ανεξάρτητους εκτιμητές με διακύμανση σ^2 :

$$\text{Var}\left(\frac{1}{M} \sum_{m=1}^M f_m\right) = \frac{\sigma^2}{M}$$

Διαισθητικά

“Σοφία του πλήθους”: πολλά ελαφρώς διαφορετικά δέντρα ψηφίζουν, εξισορροπώντας τα μεμονωμένα σφάλματα.



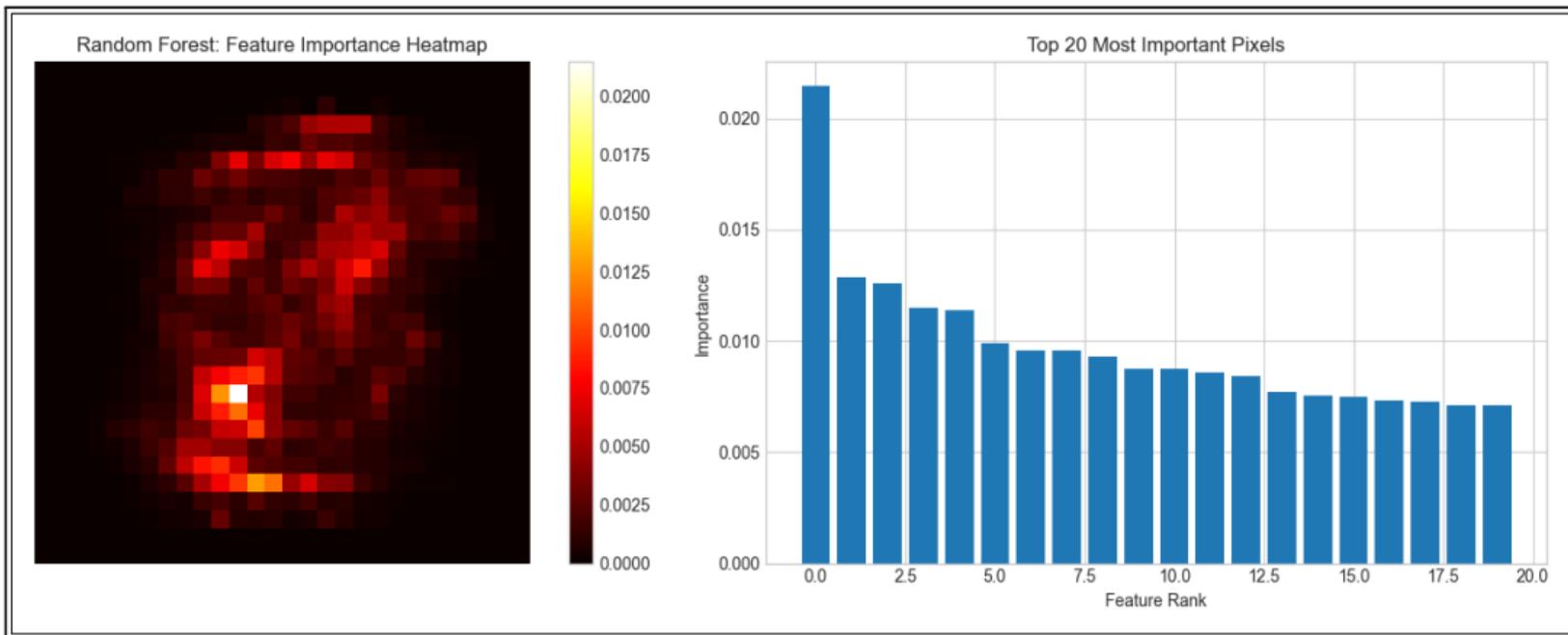
Σχ. 4.13: Random Forest Confusion Matrix

Τυχαίο Δάσος: Κώδικας & Σημασία Χαρακτηριστικών

```
1  from sklearn.ensemble import RandomForestClassifier
2
3  # Εκπαίδευση Τυχαίου Δάσους (100 δέντρα)
4  rf = RandomForestClassifier(n_estimators=100, random_state=42, n_jobs=-1)
5  start = time.time()
6  rf.fit(X_train, y_train)
7  rf_time = time.time() - start
8
9  rf_pred = rf.predict(X_test)
10 rf_acc, rf_prec, rf_rec, rf_f1 = print_metrics(y_test, rf_pred, 'Random Forest')
11
12 # Σύγκριση μονό δέντρο vs δάσος
13 print(f"Single Decision Tree:      {dt_acc:.4f}")
14 print(f"Random Forest (100 trees): {rf_acc:.4f}")
15 print(f"Improvement:                {(rf_acc - dt_acc)*100:.2f}%")
16
17 # Ανάκτηση σημασίας χαρακτηριστικών
18 rf_importance = rf.feature_importances_ # shape: (784,)
19 # Αναδιαμόρφωση σε (28,28) για οπτικοποίηση
20
```

Σημασία Χαρακτηριστικών

Η χαρτογράφηση σημασίας του Τυχαίου Δάσους είναι πιο ομαλή σε σχέση με μονό δέντρο, με κέντρο βάρους στην κεντρική περιοχή της εικόνας όπου γράφεται το ψηφίο. Ένα μικρό σύνολο pixels (~ 20–30) φέρει το μεγαλύτερο μέρος της προβλεπτικής ισχύος.



Σχ. 4.14: Σημασία Χαρακτηριστικών & Κορυφαία 20 Σημαντικά Pixels

Ορισμός

SVM Σκληρού Περιθωρίου (Hard-Margin SVM): εύρεση υπερεπιπέδου που μεγιστοποιεί το περιθώριο $\frac{2}{\|\mathbf{w}\|}$:

$$\min_{\mathbf{w}, b} \frac{1}{2} \|\mathbf{w}\|^2 \quad \text{υπό } y_i(\mathbf{w}^T \mathbf{x}_i + b) \geq 1$$

Ορισμός

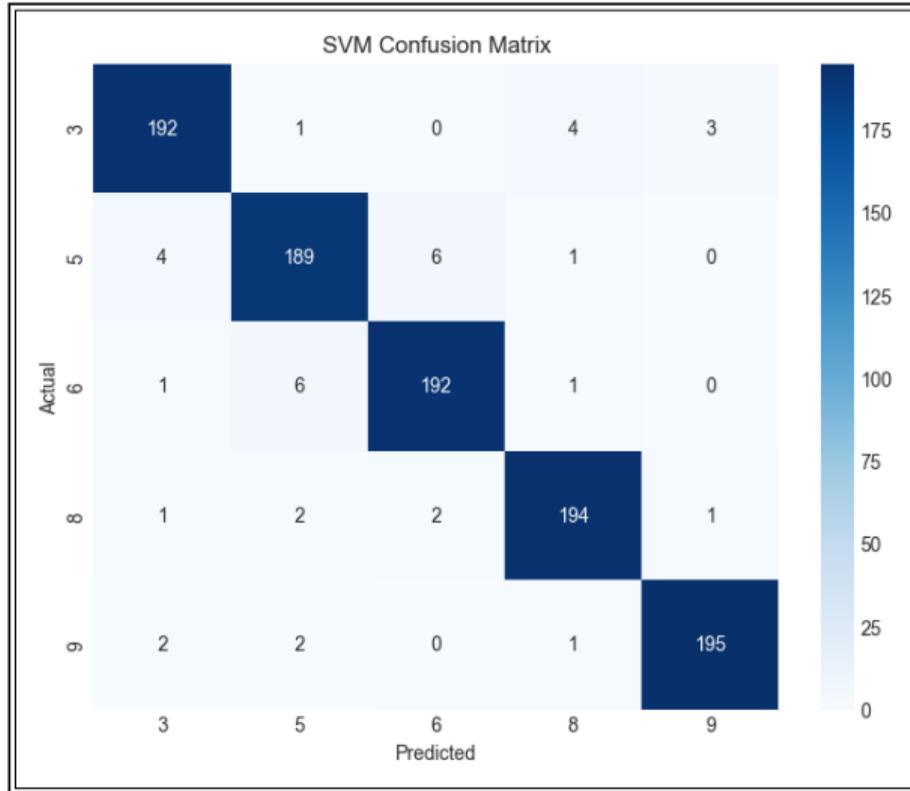
SVM Μαλακού Περιθωρίου (Soft-Margin SVM) με μεταβλητές χαλάρωσης $\xi_i \geq 0$:

$$\min_{\mathbf{w}, b, \xi} \frac{1}{2} \|\mathbf{w}\|^2 + C \sum_i \xi_i$$

Ορισμός

Τέχνασμα Πυρήνα (Kernel Trick): αντιστοίχιση σε υψηλότερες διαστάσεις μέσω συνάρτησης πυρήνα $k(\mathbf{x}, \mathbf{x}') = \langle \phi(\mathbf{x}), \phi(\mathbf{x}') \rangle$.

- **Γραμμικός:** $k = \mathbf{x}^T \mathbf{x}'$
- **RBF:** $k = \exp(-\gamma \|\mathbf{x} - \mathbf{x}'\|^2)$
- **Πολυωνυμικός:** $k = (\mathbf{x}^T \mathbf{x}' + c)^d$



Σχ. 4.15: SVM Confusion Matrix

SVM: Κώδικας & Αποτελέσματα

```
1  from sklearn.svm import SVC
2
3  # Εκπαίδευση SVM με πυρήνα RBF
4  svm = SVC(kernel='rbf', random_state=42)
5  start = time.time()
6  svm.fit(X_train, y_train)
7  svm_time = time.time() - start
8
9  svm_pred = svm.predict(X_test)
10 svm_acc, svm_prec, svm_rec, svm_f1 = print_metrics(y_test, svm_pred, 'SVM')
11 print(f"Training Time: {svm_time:.4f}s")
12
13 # Αποθήκευση αποτελεσμάτων
14 results['Algorithm'].append('SVM')
15 results['Accuracy'].append(svm_acc)
16 results['Precision'].append(svm_prec)
17 results['Recall'].append(svm_rec)
18 results['F1-Score'].append(svm_f1)
19 results['Training Time (s)'].append(svm_time)
20
```

Ερμηνεία

Το SVM με πυρήνα RBF επιτυγχάνει κορυφαία απόδοση σε αυτό το σύνολο δεδομένων. Η ανά-κλάση ανάλυση δείχνει F1-scores 94–98%, με το ψηφίο 8 (δύο οπές – διακριτή τοπολογία) να επιτυγχάνει τις υψηλότερες βαθμολογίες.

Πίνακας Περιεχομένων

Πλαίσιο Επιβλεπόμενης Μάθησης

Σύνολο Δεδομένων MNIST & Προεπεξεργασία

Μετρικές Αξιολόγησης

Αλγόριθμοι Ταξινόμησης

Σύγκριση Μοντέλων

Ανάλυση Σημασίας Χαρακτηριστικών

Τροχιές Δυναμικών Συστημάτων

Αναπαραστάσεις Διαγραμμάτων Επιμονής

Πυρήνες Διαγραμμάτων Επιμονής

Ολοκληρωμένη Προσέγγιση

Σύνοψη

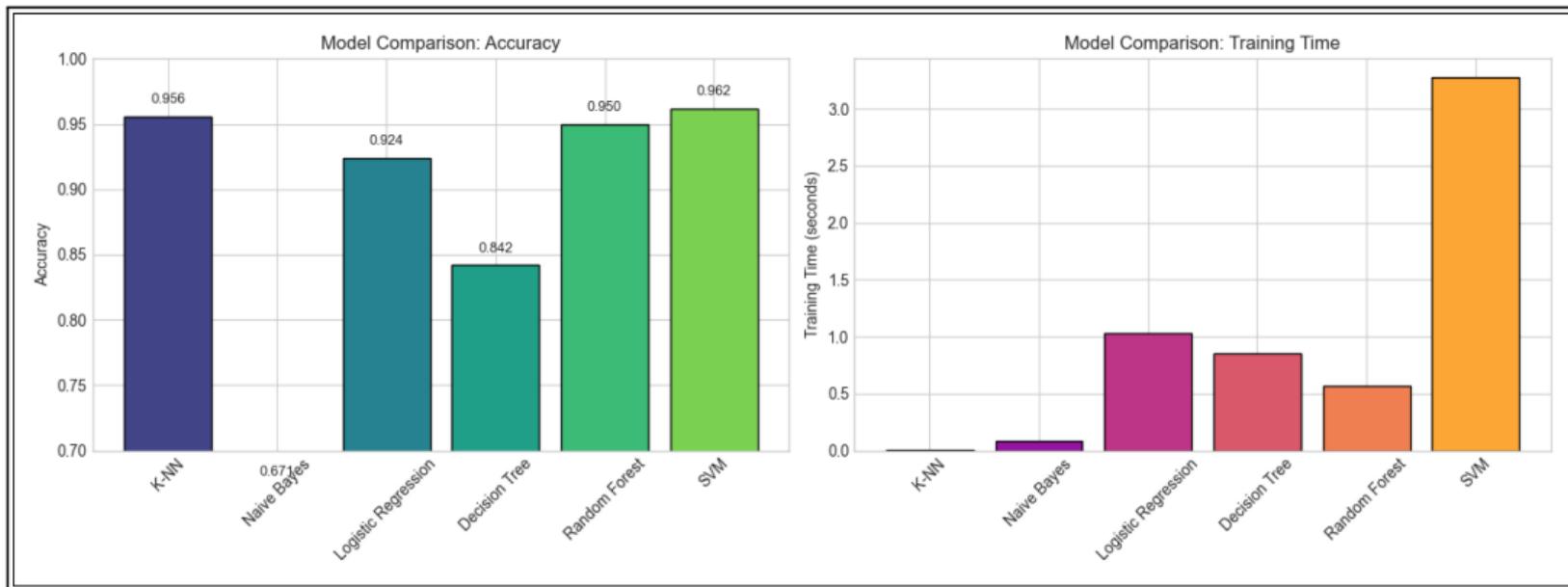
Σύγκριση Μοντέλων (Model Comparison)

Πίνακας Σύνοψης Απόδοσης

```
1 comparison_df = pd.DataFrame(results)
2 comparison_df = comparison_df.sort_values('Accuracy', ascending=False)
3 print(comparison_df.to_string(index=False))
4
```

Αλγόριθμος	Accuracy	Precision	Recall	F1-Score
SVM (RBF)	≈ 0.97	≈ 0.97	≈ 0.97	≈ 0.97
Random Forest	≈ 0.96	≈ 0.96	≈ 0.96	≈ 0.96
K-NN ($k = 5$)	≈ 0.96	≈ 0.96	≈ 0.96	≈ 0.96
Logistic Regression	≈ 0.91	≈ 0.91	≈ 0.91	≈ 0.91
Decision Tree	≈ 0.88	≈ 0.88	≈ 0.88	≈ 0.88
Naive Bayes	≈ 0.78	≈ 0.79	≈ 0.78	≈ 0.78

Οπτικοποίηση (Visualization)



Σχ. 4.16: Σύγκριση Μοντέλων

Οδηγός Επιλογής Αλγορίθμου (Algorithm Selection Guidelines)

Κριτήρια Επιλογής

Κριτήριο	Καλύτερη Επιλογή
Υψηλότερη ακρίβεια	SVM, Random Forest
Ταχύτερη εκπαίδευση	Naive Bayes
Ερμηνευσιμότητα (interpretability)	Decision Tree, Logistic Regression
Σημασία χαρακτηριστικών	Random Forest, Decision Tree
Μη-γραμμικά όρια απόφασης	SVM (RBF), Random Forest
Μεγάλα δεδομένα	Logistic Regression, Naive Bayes

Κεντρικό Συμπέρασμα

Δεν υπάρχει καθολικά “καλύτερος” αλγόριθμος (No Free Lunch Theorem). Η επιλογή εξαρτάται από τη φύση των δεδομένων, τους υπολογιστικούς πόρους και τις εκάστοτε απαιτήσεις.

Πίνακας Περιεχομένων

Πλαίσιο Επιβλεπόμενης Μάθησης

Σύνολο Δεδομένων MNIST & Προεπεξεργασία

Μετρικές Αξιολόγησης

Αλγόριθμοι Ταξινόμησης

Σύγκριση Μοντέλων

Ανάλυση Σημασίας Χαρακτηριστικών

Τροχιές Δυναμικών Συστημάτων

Αναπαραστάσεις Διαγραμμάτων Επιμονής

Πυρήνες Διαγραμμάτων Επιμονής

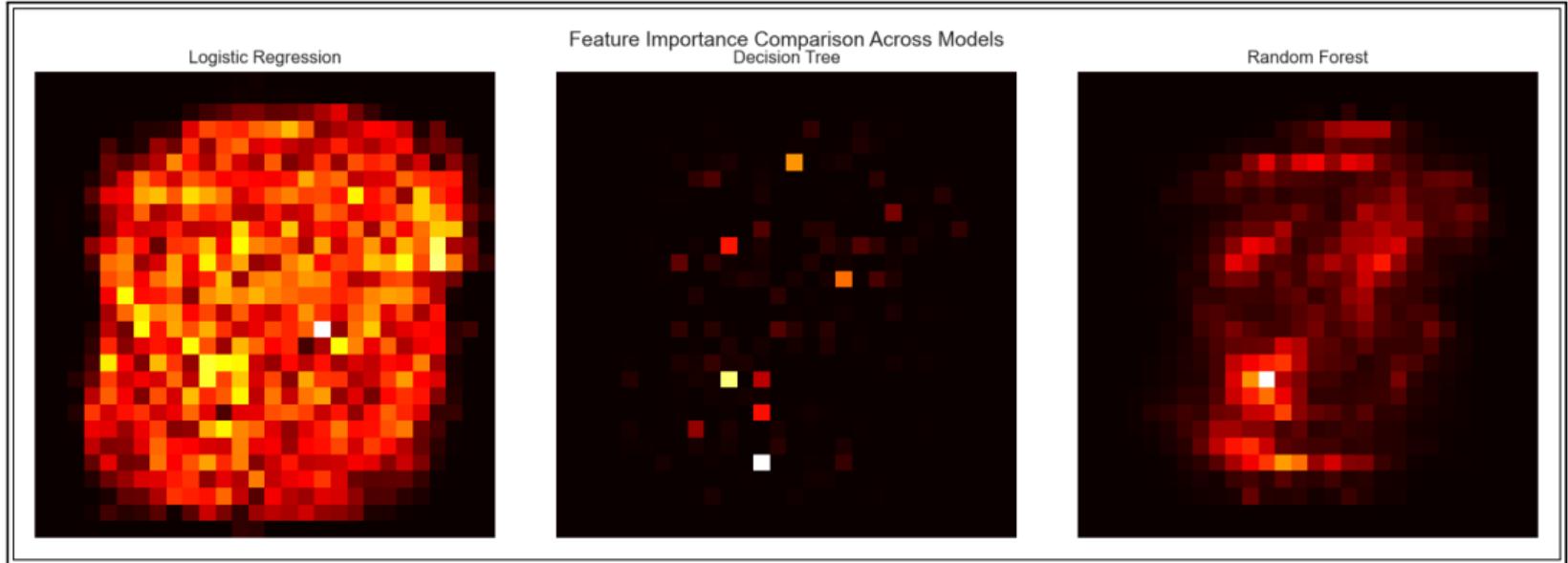
Ολοκληρωμένη Προσέγγιση

Σύνοψη

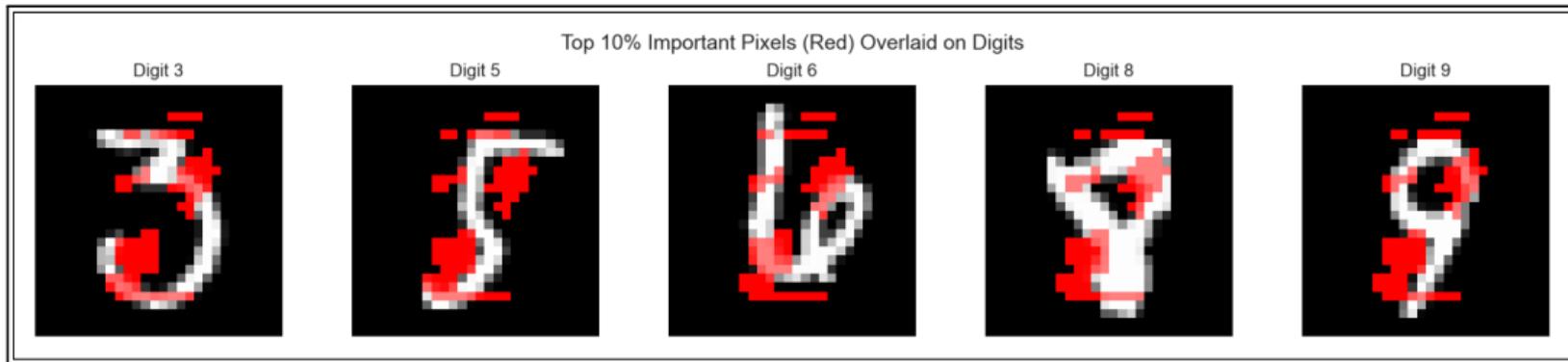
Σύγκριση Σημασίας Χαρακτηριστικών (Feature Importance Comparison)

Ανάλυση για τρία μοντέλα

```
1 # Logistic Regression: μέσος όρος απόλυτης τιμής συντελεστών ανά pixel
2 lr_imp = np.abs(lr.coef_).mean(axis=0) # shape: (784,)
3
4 # Decision Tree: ενσωματωμένη σημασία
5 dt_importance = dt.feature_importances_ # shape: (784,)
6
7 # Random Forest: μέσος όρος σημασίας ανά δέντρο
8 rf_importance = rf.feature_importances_ # shape: (784,)
9
10 # Επικάλυψη σημαντικότερων pixels σε δείγματα
11 importance_mask = rf_importance.reshape(28, 28) > np.percentile(rf_importance, 90)
12 # Τα κόκκινα pixels: κορυφαίο 10% σε σημασία
13
```



Σχ. 4.17: Σύγκριση Τριών Μοντέλων



Σχ. 4.18: Κορυφαία 10% Σημαντικών Pixels (Κόκκινο) Επικαλυμμένο στα Ψηφία

Σύνοψη Βασικών Ευρημάτων

- **Κεντρική εστίαση:** όλα τα μοντέλα συμφωνούν ότι η κεντρική περιοχή (όπου γράφεται το ψηφίο) φέρει τη μεγαλύτερη πληροφορία. Γωνίες και άκρες (μαύρο φόντο) συμβάλλουν ελάχιστα.
- **Διαφορές μοντέλων:** η Logistic Regression κατανέμει ομαλά, το Decision Tree είναι αραιό, το Random Forest δίνει ομαλή χαρτογράφηση.
- **Διακριτές περιοχές:** σημαντικά pixels αντιστοιχούν σε περιοχές όπου τα ψηφία διαφέρουν – περιοχές βρόχων (6, 8, 9) και καμπύλες (3, 5).

Βασικά Συμπεράσματα

1. **Επιβλεπόμενη μάθηση** απαιτεί: διαχωρισμό εκπαίδευσης/ελέγχου, ομαλοποίηση, αξιολόγηση με σωστές μετρικές.
2. **Χαρακτηριστικά δεδομένων καθορίζουν τον αλγόριθμο**: τα pixels MNIST έχουν συσχετίσεις – το Naive Bayes υποφέρει, τα SVM/RF ευδοκιμούν.
3. **Υπερπροσαρμογή** ελέγχεται με περιορισμό πολυπλοκότητας (βάθος δέντρου, regularization).
4. **Ανάλυση σημασίας χαρακτηριστικών** αποκαλύπτει ποια pixels “βλέπουν” τα μοντέλα.

Σύνδεση με TDA

Η ανάλυση σημασίας ρixel υπονοεί ότι τα μοντέλα ανακαλύπτουν εμμέσως **τοπολογικά χαρακτηριστικά** (βρόχοι, οπές). Η TDA μπορεί να **εξάγει ρητά** αυτά τα χαρακτηριστικά και να τα χρησιμοποιήσει ως είσοδο στους παραπάνω ταξινομητές.

Κύριες Εισαγωγές

```
1 import numpy as np
2 import gudhi as gd
3 import gudhi.representations
4 import matplotlib.pyplot as plt
5
```

Πίνακας Περιεχομένων

Πλαίσιο Επιβλεπόμενης Μάθησης

Σύνολο Δεδομένων MNIST & Προεπεξεργασία

Μετρικές Αξιολόγησης

Αλγόριθμοι Ταξινόμησης

Σύγκριση Μοντέλων

Ανάλυση Σημασίας Χαρακτηριστικών

Τροχιές Δυναμικών Συστημάτων

Αναπαραστάσεις Διαγραμμάτων Επιμονής

Πυρήνες Διαγραμμάτων Επιμονής

Ολοκληρωμένη Προσέγγιση

Σύνοψη

Τροχιές Δυναμικών Συστημάτων (Orbits of Dynamical Systems)

Περιγραφή Συστήματος (System Description)

Χρησιμοποιούμε δυναμικό σύστημα που εξαρτάται από παράμετρο $r > 0$:

$$x_{n+1} = x_n + r y_n (1 - y_n) \pmod{1}$$

$$y_{n+1} = y_n + r x_{n+1} (1 - x_{n+1}) \pmod{1}$$

Για διαφορετικές τιμές r , η τροχιά εμφανίζει **διαφορετική τοπολογική δομή**.

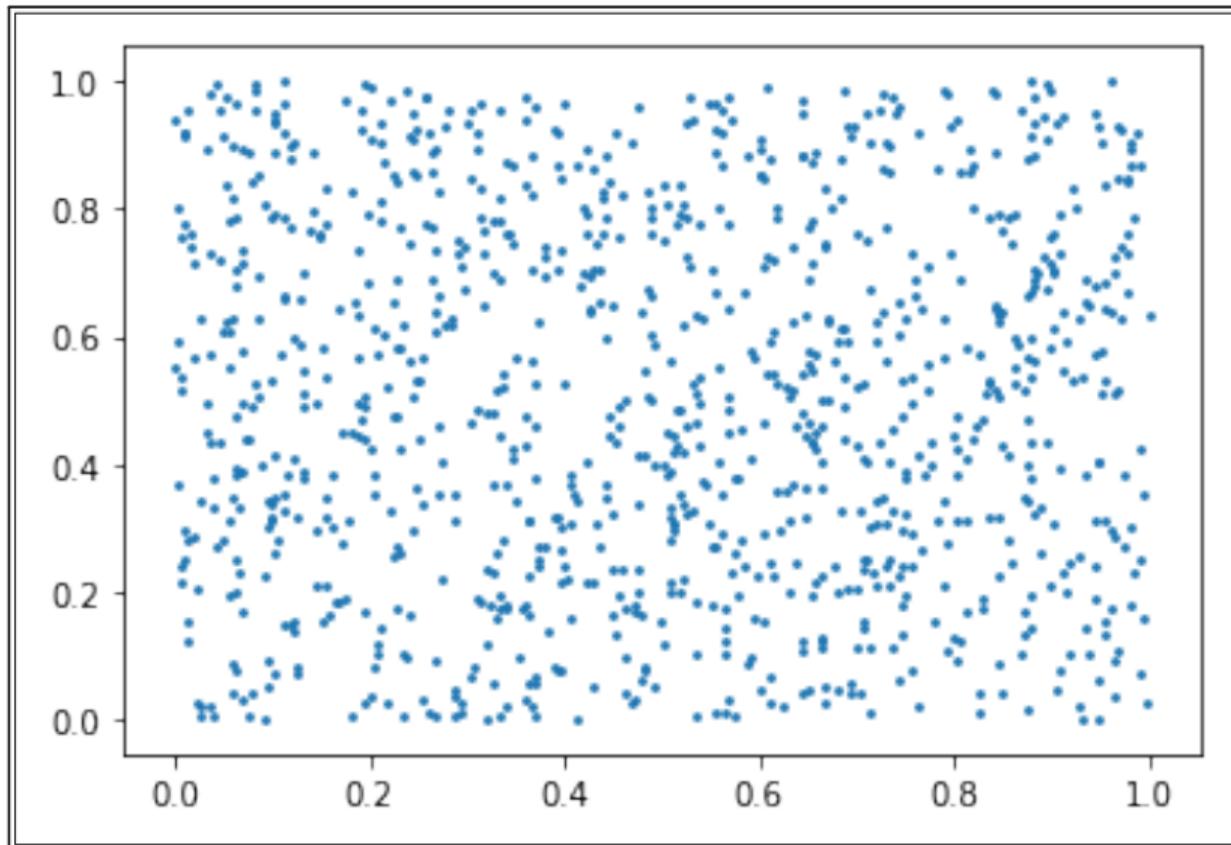
Παραγωγή Νέφους Σημείων

```
1 num_pts = 1000
2 r = 3.5
3
4 X = np.empty([num_pts, 2])
5 x, y = np.random.uniform(), np.random.uniform()
6 for i in range(num_pts):
7     X[i,:] = [x, y]
8     x = (X[i,0] + r*X[i,1]*(1-X[i,1])) % 1.
9     y = (X[i,1] + r*x*(1-x)) % 1.
10
```

Βασική Παρατήρηση

Για $r = 3.5$: το νέφος φαίνεται τυχαίο.

Για $r = 4.1$: εμφανίζεται **οπή στο κέντρο** – σημαντικό τοπολογικό χαρακτηριστικό. Η τιμή r καθορίζει την τοπολογία!



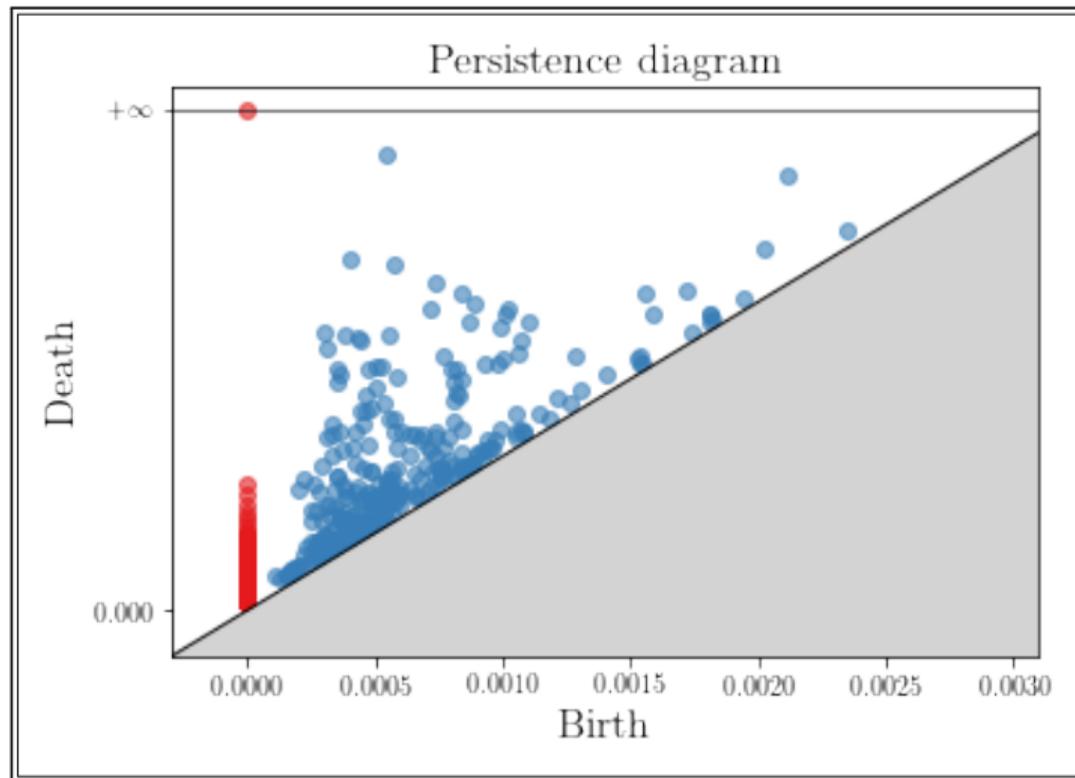
Σχ. 4.19: Παραγωγή Νέφους Σημείων

Παραγωγή Διαγραμμάτων Επιμονής (Computing Persistence Diagrams)

Alpha Φιλτράρισμα με GUDHI

```
1 # Κατασκευή Alpha συμπλέγματος και υπολογισμός επιμονής
2 acX = gd.AlphaComplex(points=X).create_simplex_tree()
3 dgmX = acX.persistence()
4
5 # Οπτικοποίηση διαγράμματος επιμονής
6 gd.plot_persistence_diagram(dgmX)
7
8 # Δεύτερο νέφος σημείων (r=4.1 -- έχει τοπολογία!)
9 r = 4.1
10 Y = np.empty([num_pts, 2])
11 x, y = np.random.uniform(), np.random.uniform()
12 for i in range(num_pts):
13     Y[i,:] = [x, y]
14     x = (Y[i,0] + r*Y[i,1]*(1-Y[i,1])) % 1.
15     y = (Y[i,1] + r*x*(1-x)) % 1.
16
17 acY = gd.AlphaComplex(points=Y).create_simplex_tree()
18 dgmY = acY.persistence()
```

Σύγκριση: Για $r = 3.5$: τυχαίο νέφος, χωρίς έντονα τοπολογικά χαρακτηριστικά στο διάγραμμα επιμονής. Για $r = 4.1$: εμφανής οπή στη διάσταση 1, υποδηλώνει **δακτύλιο** στη δομή της τροχιάς.



Σχ. 4.20: Παραγωγή Διαγράμματος Επιμονής

Πίνακας Περιεχομένων

Πλαίσιο Επιβλεπόμενης Μάθησης

Σύνολο Δεδομένων MNIST & Προεπεξεργασία

Μετρικές Αξιολόγησης

Αλγόριθμοι Ταξινόμησης

Σύγκριση Μοντέλων

Ανάλυση Σημασίας Χαρακτηριστικών

Τροχιές Δυναμικών Συστημάτων

Αναπαραστάσεις Διαγραμμάτων Επιμονής

Πυρήνες Διαγραμμάτων Επιμονής

Ολοκληρωμένη Προσέγγιση

Σύνοψη

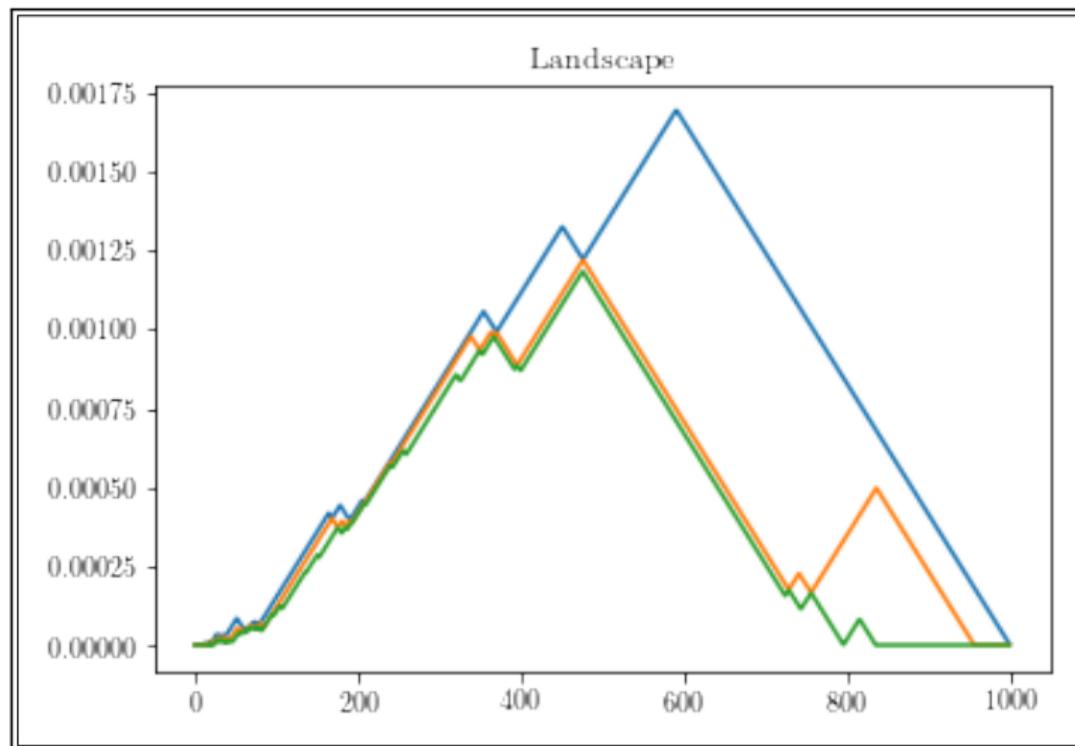
Ορισμός

Το **τοπίο επιμονής** (*persistence landscape*) λαμβάνεται:

1. Περιστροφή του διαγράμματος επιμονής κατά $-\pi/4$ (η διαγώνιος γίνεται άξονας x)
2. Τοποθέτηση **σκηνικών συναρτήσεων** (*tent functions*) σε κάθε σημείο
3. Το k -ιοστό τοπίο (*landscape*) είναι η k -ιοστή μεγαλύτερη τιμή μεταξύ όλων των σκηνικών συναρτήσεων
4. Μετατροπή σε διάνυσμα με δειγματοληψία σε ομοιόμορφα σημεία του άξονα x

Κώδικας

```
1 LS = gd.representations.Landscape(resolution=1000)
2 L = LS.fit_transform([acX.persistence_intervals_in_dimension(1)])
3 # L[0] έχει μήκος 3000 (3 τοπία x 1000 σημεία)
4 # Οπτικοποίηση:
5 plt.plot(L[0][:1000]) # ο1 τοπίο
6 plt.plot(L[0][1000:2000]) # ο2 τοπίο
7 plt.plot(L[0][2000:3000]) # ο3 τοπίο
8 plt.title("Landscape")
```



Σχ. 4.21: Παραγωγή Τοπίου Επιμονής

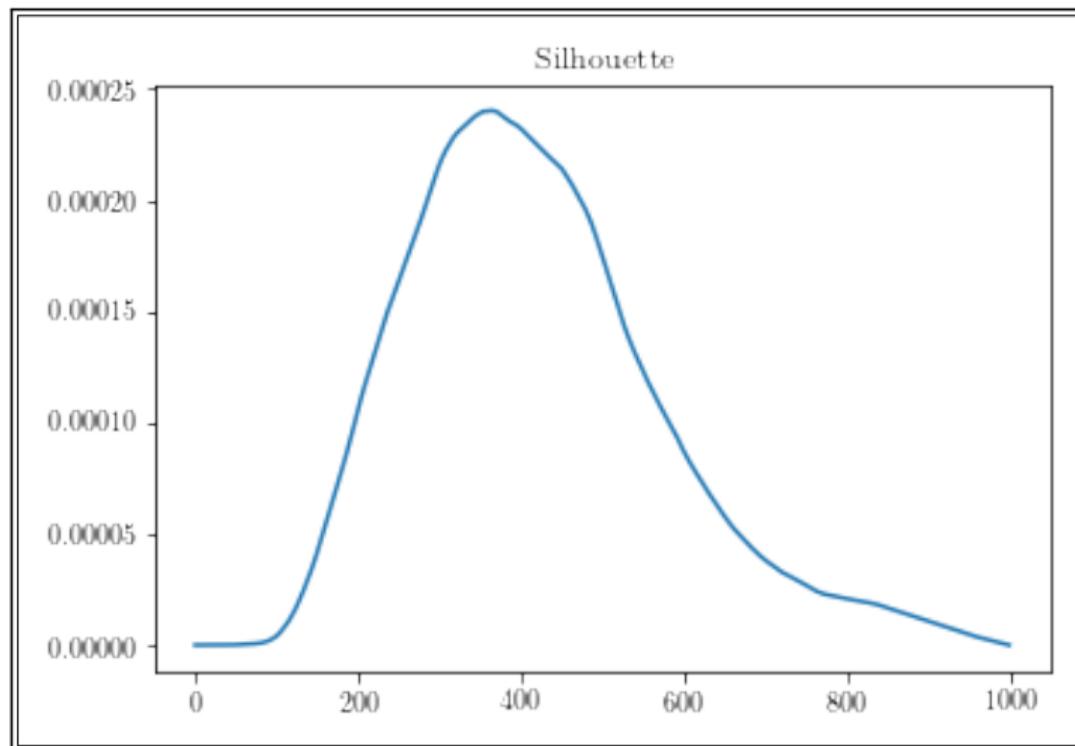
Σιλουέτα Επιμονής (Persistence Silhouette)

Ορισμός

Η **σιλουέτα επιμονής** (*persistence silhouette*) αποτελεί παραλλαγή του τοπίου: λαμβάνει **σταθμισμένο μέσο** όρο των σκηνικών συναρτήσεων αντί του k -ιοστού μέγιστου. Κάθε σκηνική συνάρτηση σταθμίζεται ανάλογα με την απόσταση του αντίστοιχου σημείου από τη διαγώνιο.

Κώδικας

```
1 # Στάθμιση: τετράγωνο της επιμονής (death - birth)^1
2 SH = gd.representations.Silhouette(
3     resolution=1000,
4     weight=lambda x: np.power(x[1]-x[0], 1)
5 )
6 sh = SH.fit_transform([acX.persistence_intervals_in_dimension(1)])
7 plt.plot(sh[0])
8 plt.title("Silhouette")
9
```



Σχ. 4.22: Παραγωγή Σιλουέτας Επιμονής

Εικόνα Επιμονής (Persistence Image)

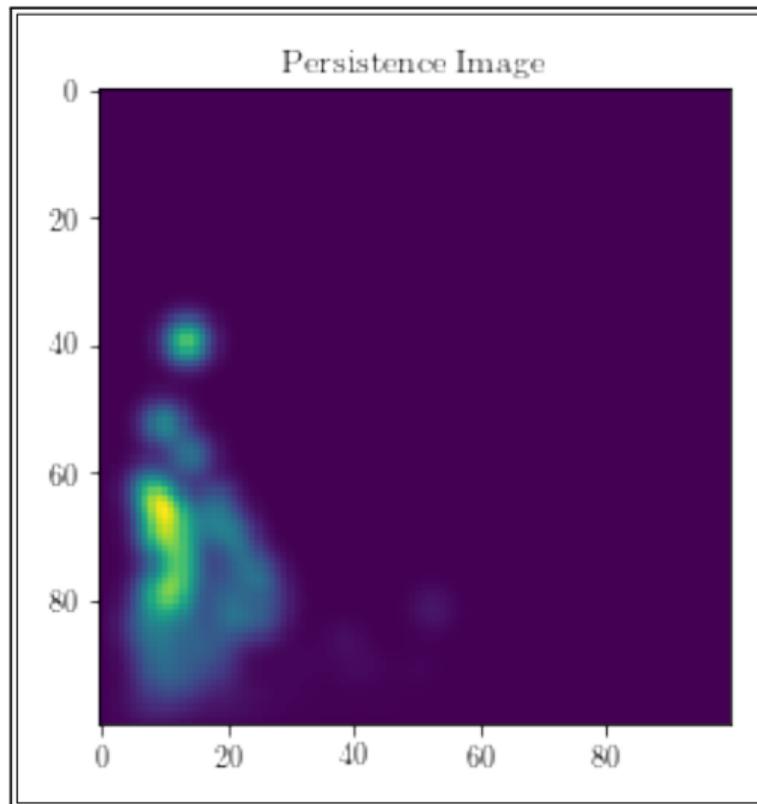
Ορισμός

Η **εικόνα επιμονής** (*persistence image*) λαμβάνεται:

1. Περιστροφή κατά $-\pi/4$
2. Κέντρωση Γκαουσιανών συναρτήσεων σε κάθε σημείο (σταθμισμένες από παραμετρική συνάρτηση – συνήθως τετράγωνο απόστασης από τη διαγώνιο)
3. Άθροιση όλων των Γκαουσιανών – δίνει 2D συνάρτηση
4. *Pixelization* σε εικόνα (*rasterization*)

Κώδικας

```
1 PI = gd.representations.PersistenceImage(  
2 bandwidth=1e-4,  
3 weight=lambda x: x[1]**2, # τετράγωνο death  
4 im_range=[0, .004, 0, .004],  
5 resolution=[100, 100] # 100x100 pixels  
6 )  
7 pi = PI.fit_transform([acX.persistence_intervals_in_dimension(1)])  
8 # pi[0] shape: (10000,) -- αναδιαμόρφωση σε (100,100)  
9
```



Σχ. 4.23: Παραγωγή Εικόνας Επιμονής

Διαθέσιμες Μέθοδοι στο `gudhi.representations`

- **Καμπύλη Betti** (Betti Curve): μετρά τον αριθμό ζωντανών χαρακτηριστικών σε κάθε τιμή του φίλτρου
- **Μιγαδικό Πολυώνυμο** (Complex Polynomial): χαρτογράφηση σε χώρο πολυωνύμων
- **Τοπολογικό Διάνυσμα** (Topological Vector): συμπίεση πληροφορίας σε διάνυσμα σταθερού μήκους

Κοινό Πλαίσιο `scikit-learn`

Όλες οι μέθοδοι ακολουθούν το API του `scikit-learn`: `fit()`, `transform()`, `fit_transform()`. Αυτό επιτρέπει απρόσκοπτη ενσωμάτωση με μηχανική μάθηση.

Πίνακας Περιεχομένων

Πλαίσιο Επιβλεπόμενης Μάθησης

Σύνολο Δεδομένων MNIST & Προεπεξεργασία

Μετρικές Αξιολόγησης

Αλγόριθμοι Ταξινόμησης

Σύγκριση Μοντέλων

Ανάλυση Σημασίας Χαρακτηριστικών

Τροχιές Δυναμικών Συστημάτων

Αναπαραστάσεις Διαγραμμάτων Επιμονής

Πυρήνες Διαγραμμάτων Επιμονής

Ολοκληρωμένη Προσέγγιση

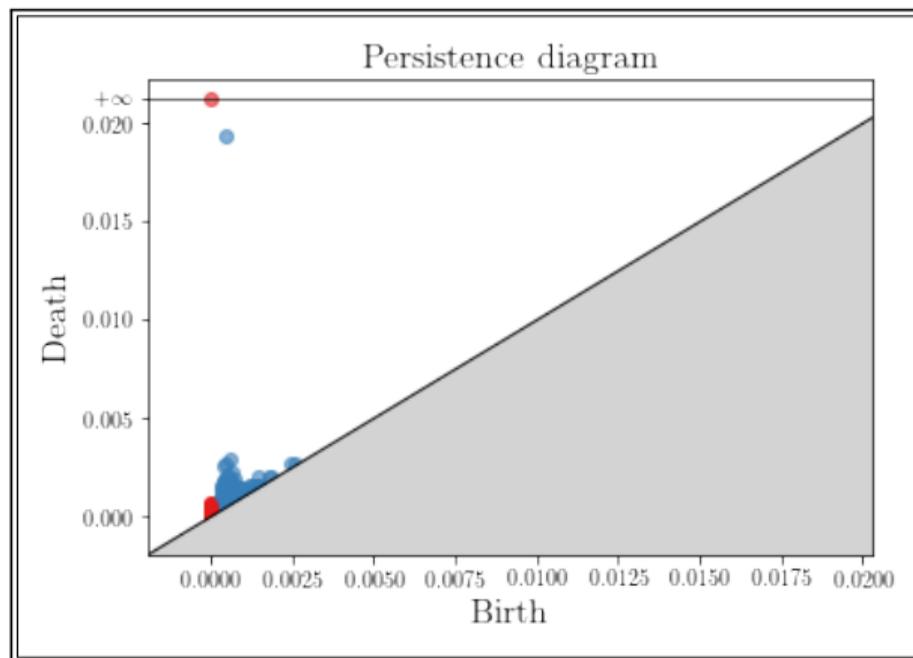
Σύνοψη

Ορισμός

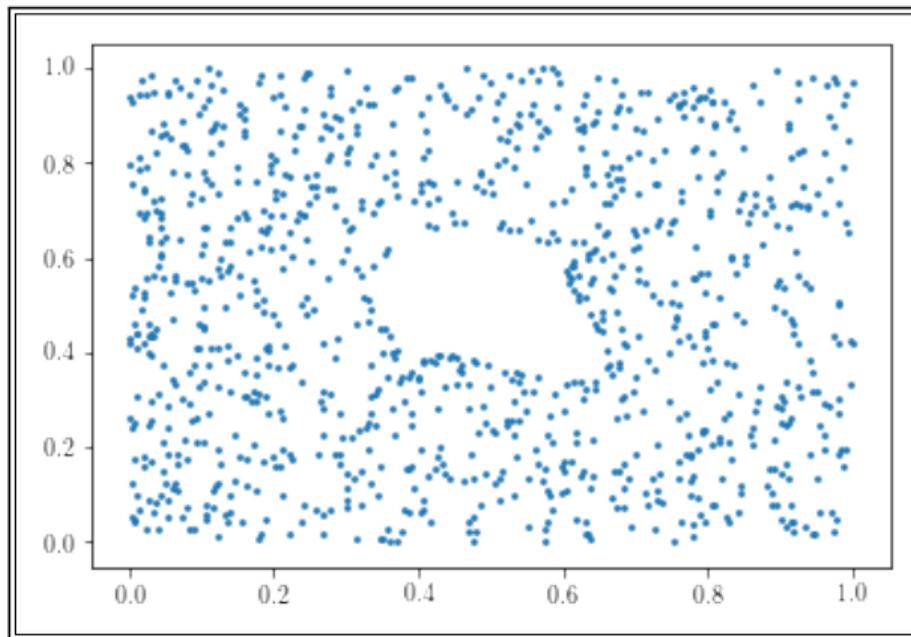
Ένας **πυρήνας** (kernel) k δέχεται ζεύγος διαγραμμάτων D, D' και επιστρέφει πραγματικό αριθμό ώστε:

$$k(D, D') = \langle \Phi(D), \Phi(D') \rangle_{\mathcal{H}}$$

για κάποιον **υπονοούμενο χώρο Hilbert** (implicit Hilbert space) \mathcal{H} και συνεχή συνάρτηση $\Phi : \mathcal{D} \rightarrow \mathcal{H}$.



Σχ. 4.24: Διάγραμμα Επιμονής



Σχ. 4.25: Επίδειξη Ομολογίας Δεδομένων

Διαθέσιμοι Πυρήνες στο GUDHI

- **Sliced Wasserstein Kernel (SW)**: βασισμένος σε προβολές Wasserstein
- **Persistence Weighted Gaussian Kernel (PWG)**: σταθμισμένες Γκαουσιανές σε χώρο διαγραμμάτων
- **Persistence Scale Space Kernel (PSS)**: stable multi-scale kernel
- **Persistence Fisher Kernel (PF)**: βασισμένος σε απόκλιση Fisher

Πλεονέκτημα

Πολλοί αλγόριθμοι (SVM, kernel PCA) απαιτούν μόνο ζευγαρωτά εσωτερικά γινόμενα. Με τους πυρήνες, μπορούμε να τους εφαρμόσουμε **απευθείας σε διαγράμματα επιμονής!**

Υπολογισμός Πυρήνων & Αποστάσεων (Computing Kernels & Distances)

```
1 # Persistence Weighted Gaussian Kernel
2 PWG = gd.representations.PersistenceWeightedGaussianKernel(
3 bandwidth=0.01, kernel_approx=None,
4 weight=lambda x: np.arctan(np.power(x[1], 1)))
5 PWG.fit([acX.persistence_intervals_in_dimension(1)])
6 pwg = PWG.transform([acY.persistence_intervals_in_dimension(1)])
7 print("PWG kernel is " + str(pwg[0][0]))
8
9 # Persistence Scale Space Kernel
10 PSS = gd.representations.PersistenceScaleSpaceKernel(bandwidth=1.)
11 PSS.fit([acX.persistence_intervals_in_dimension(1)])
12 pss = PSS.transform([acY.persistence_intervals_in_dimension(1)])
13
14 # Persistence Fisher Kernel
15 PF = gd.representations.PersistenceFisherKernel(
16 bandwidth_fisher=.001, bandwidth=.001, kernel_approx=None)
17 PF.fit([acX.persistence_intervals_in_dimension(1)])
18
19 # Sliced Wasserstein Kernel
20 SW = gd.representations.SlicedWassersteinKernel(bandwidth=.1, num_directions=100)
21 SW.fit([acX.persistence_intervals_in_dimension(1)])
22
```

Αποστάσεις Bottleneck & Wasserstein

```
1 # Απόσταση Bottleneck μέσω( gudhi.representations)
2 BD = gd.representations.BottleneckDistance(epsilon=.001)
3 BD.fit([acX.persistence_intervals_in_dimension(1)])
4 bd = BD.transform([acY.persistence_intervals_in_dimension(1)])
5 print("Bottleneck distance is " + str(bd[0][0]))
6
7 # Απόσταση Wasserstein
8 WD = gd.representations.WassersteinDistance(internal_p=2, order=2)
9 WD.fit([acX.persistence_intervals_in_dimension(1)])
10 wd = WD.transform([acY.persistence_intervals_in_dimension(1)])
11 print("Wasserstein distance is " + str(wd[0][0]))
12
```

Συνοπτικός Πίνακας Μεθόδων

Μέθοδος	Τύπος	Παραμέτρος
Persistence Landscape (LS)	Vectorization	<code>resolution</code>
Silhouette (SH)	Vectorization	<code>resolution, weight</code>
Persistence Image (PI)	Vectorization	<code>bandwidth, resolution</code>
Sliced Wasserstein (SW)	Kernel	<code>bandwidth, num_directions</code>
PWG	Kernel	<code>bandwidth, weight</code>
Bottleneck (BD)	Distance	<code>epsilon</code>
Wasserstein (WD)	Distance	<code>internal_p, order</code>

Πίνακας Περιεχομένων

Πλαίσιο Επιβλεπόμενης Μάθησης

Σύνολο Δεδομένων MNIST & Προεπεξεργασία

Μετρικές Αξιολόγησης

Αλγόριθμοι Ταξινόμησης

Σύγκριση Μοντέλων

Ανάλυση Σημασίας Χαρακτηριστικών

Τροχιές Δυναμικών Συστημάτων

Αναπαραστάσεις Διαγραμμάτων Επιμονής

Πυρήνες Διαγραμμάτων Επιμονής

Ολοκληρωμένη Προσέγγιση

Σύνοψη

5 Κλάσεις – Διαφορετικές Τιμές r

```
1 num_diag_per_class = 10 # Διαγράμματα ανά κλάση
2 dgms, labs = [], []
3
4 for idx, radius in enumerate([2.5, 3.5, 4., 4.1, 4.3]):
5     for _ in range(num_diag_per_class):
6         labs.append(idx)
7         X = np.empty([num_pts, 2])
8         x, y = np.random.uniform(), np.random.uniform()
9         for i in range(num_pts):
10            X[i,:] = [x, y]
11            x = (X[i,0] + radius * X[i,1] * (1-X[i,1])) % 1.
12            y = (X[i,1] + radius * x * (1-x)) % 1.
13            ac = gd.AlphaComplex(points=X).create_simplex_tree(max_alpha_square=1e12)
14            dgm = ac.persistence()
15            dgms.append(ac.persistence_intervals_in_dimension(1))
16
```

Δομή Συνόλου Δεδομένων

5 κλάσεις \times 10 διαγράμματα = 50 διαγράμματα επιμονής συνολικά. Κάθε κλάση αντιστοιχεί σε διαφορετική τιμή r (διαφορετική τοπολογική δομή τροχιάς).

Διαχωρισμός Εκπαίδευσης – Ελέγχου (Train-Test Split)

```
1 test_size          = 0.2
2 perm              = np.random.permutation(len(labs))
3 limit            = int(test_size * len(labs))
4 test_sub, train_sub = perm[:limit], perm[limit:]
5
6 train_labs       = np.array(labs)[train_sub]
7 test_labs        = np.array(labs)[test_sub]
8 train_dgms       = [dgms[i] for i in train_sub]
9 test_dgms        = [dgms[i] for i in test_sub]
10
```

Σύνθεση Συνόλων

- **Εκπαιδευτικό:** 80% – 40 διαγράμματα
- **Ελέγχου:** 20% – 10 διαγράμματα

Στάδια – Βήματα

Η προσέγγιση δομείται ως εξής:

1. **DiagramSelector**: εξαγωγή σημείων με πεπερασμένες συντεταγμένες (non-essential points)
2. **DiagramScaler**: κλιμάκωση ή μη στο μοναδιαίο τετράγωνο
3. **TDA (Αναπαράσταση)**: vectorization ή kernel method με GUDHI
4. **Estimator (Ταξινομητής)**: αλγόριθμος ταξινόμησης scikit-learn

Βασική Δομή Προσέγγισης

```
1 from sklearn.preprocessing import MinMaxScaler
2 from sklearn.pipeline import Pipeline
3 from sklearn.svm import SVC
4
5 pipe = Pipeline([
6     ("Separator", gd.representations.DiagramSelector(limit=np.inf, point_type="finite")),
7     ("Scaler", gd.representations.DiagramScaler(scalers=[([0,1], MinMaxScaler())])),
8     ("TDA", gd.representations.PersistenceImage()),
9     ("Estimator", SVC())
10 ])
```

Πλέγμα Παραμέτρων (Parameter Grid)

```
11 param = [  
12 # Sliced Wasserstein kernel SVM  
13 {"Scaler__use": [False],  
14  "TDA":          [gd.representations.SlicedWassersteinKernel()],  
15  "TDA__bandwidth": [0.1, 1.0], "TDA__num_directions": [20],  
16  "Estimator":    [SVC(kernel="precomputed", gamma="auto")]},  
17  
18 # Persistence Weighted Gaussian kernel SVM  
19 {"Scaler__use": [False],  
20  "TDA":          [gd.representations.PersistenceWeightedGaussianKernel()],  
21  "TDA__bandwidth": [0.1, 0.01],  
22  "TDA__weight":  [lambda x: np.arctan(x[1]-x[0])],  
23  "Estimator":    [SVC(kernel="precomputed", gamma="auto")]},  
24  
25 # Persistence Image + SVM  
26 {"Scaler__use": [True],  
27  "TDA":          [gd.representations.PersistenceImage()],  
28  "TDA__resolution": [[5,5], [6,6]], "TDA__bandwidth": [0.01, 0.1, 1.0, 10.0],  
29  "Estimator":    [SVC()]},  
30  
31 # Persistence Landscape + Random Forest  
32 {"Scaler__use": [True],  
33  "TDA":          [gd.representations.Landscape()], "TDA__resolution": [100],  
34  "Estimator":    [RandomForestClassifier()]},  
35  
36 # Bottleneck Distance + K-NN  
37 {"Scaler__use": [False],  
38  "TDA":          [gd.representations.BottleneckDistance()], "TDA__epsilon": [0.1],  
39  "Estimator":    [KNeighborsClassifier(metric="precomputed")]},  
40 ]
```

Διασταυρούμενη Επικύρωση & Αποτελέσματα (Cross-Validation & Results)

GridSearchCV – 3-fold Cross-Validation

```
1 from sklearn.model_selection import GridSearchCV
2
3 model = GridSearchCV(pipe, param, cv=3)
4
5 # Εκπαίδευση -- μπορεί να αργεί ειδικά( k-NN + Bottleneck)
6 model = model.fit(train_dgms, train_labs)
7
8 # Βέλτιστες παράμετροι
9 print(model.best_params_)
10
11 # Αξιολόγηση
12 print("Train accuracy = " + str(model.score(train_dgms, train_labs)))
13 print("Test accuracy = " + str(model.score(test_dgms, test_labs)))
14
```

Αποτέλεσμα

Random Forest + Persistence Landscape ανέδειξε την καλύτερη απόδοση για αυτό το μικρό σύνολο δεδομένων. Ακρίβεια ελέγχου $\approx 70\%$ – αξιόλογο για μόλις 50 διαγράμματα εκπαίδευσης! Η απόδοση βελτιώνεται σημαντικά με μεγαλύτερα σύνολα δεδομένων.

Γιατί Random Forest + Landscape;

- Το τοπίο επιμονής παρέχει **σταθερή, πλούσια αναπαράσταση** της τοπολογικής δομής σε πολλαπλές κλίμακες
- Το Τυχαίο Δάσος εκμεταλλεύεται αποτελεσματικά **πολυδιάστατα διανύσματα** (1000 σημεία ανά τοπίο)
- Ισχυρός συνδυασμός για **μικρά σύνολα δεδομένων** με πολυδιάστατες αναπαραστάσεις

Συμβουλές Πρακτικής Εφαρμογής

- Χρήση GridSearchCV για αυτόματη βελτιστοποίηση υπερπαραμέτρων
- Αρχικά δοκιμάστε **Persistence Image + SVM** (καλή ισορροπία ταχύτητας/απόδοσης)
- Για εφαρμογές πραγματικού χρόνου: **Silhouette** ή **Betti Curve** (ταχύτεροι)
- Αποφυγή k-NN + Bottleneck σε μεγάλα σύνολα (υπολογιστικά ακριβό)

Πίνακας Περιεχομένων

Πλαίσιο Επιβλεπόμενης Μάθησης

Σύνολο Δεδομένων MNIST & Προεπεξεργασία

Μετρικές Αξιολόγησης

Αλγόριθμοι Ταξινόμησης

Σύγκριση Μοντέλων

Ανάλυση Σημασίας Χαρακτηριστικών

Τροχιές Δυναμικών Συστημάτων

Αναπαραστάσεις Διαγραμμάτων Επιμονής

Πυρήνες Διαγραμμάτων Επιμονής

Ολοκληρωμένη Προσέγγιση

Σύνοψη

Σύνοψη Εννοιών (Summary of Concepts)

- **Supervised Learning:** MNIST, preprocessing, train-test split, PCA
- **Μετρικές:** accuracy, precision, recall, F1-score, confusion matrix
- **Αλγόριθμοι:** K-NN, Naive Bayes, Logistic Regression, Decision Trees, Random Forest, SVM
- **Σημασία χαρακτηριστικών:** κεντρική εστίαση σε pixels – υπονοούμενη τοπολογική πληροφορία
- **Τροχιές δυναμικών συστημάτων** ως γεννήτρια διαγραμμάτων επιμονής
- **Αναπαραστάσεις:** Landscape, Silhouette, Persistence Image
- **Πυρήνες:** SW, PWG, PSS, PF – για kernel-SVM
- **Ολοκληρωμένη Προσέγγιση Μέσω Scikit-Learn + GridSearchCV** για αυτόματη επιλογή μεθόδου

Κεντρικό Μήνυμα

Η TDA παρέχει **ρητές τοπολογικές αναπαραστάσεις** δεδομένων που μπορούν να ενσωματωθούν **πλήρως** σε παραδοσιακούς αλγορίθμους μηχανικής μάθησης.

Σας ευχαριστώ θερμά για την προσοχή σας!

Ερωτήσεις;



ΠΑΝΕΠΙΣΤΗΜΙΟ
ΠΑΤΡΩΝ
UNIVERSITY OF PATRAS