

# Αναφορά για τη 1<sup>η</sup> εργασία

## Εισαγωγή

Στην εργασία αυτή συγκρινουμε τρεις κατηγορίες ταξινομητών:

- **Support Vector Machines (SVMs)**
- **k-Nearest Neighbors (kNN)**
- **Nearest Class Centroid (NCC)**

σε δυο διαφορετικά συνολα δεδομένων:

1. **CIFAR-10** (classification 10 κλάσεων)
2. **Breast cancer** (δυαδική ταξινόμηση καλοήθες / κακοήθες)

Δίνουμε ιδιαίτερη βαρύτητα στα **SVMs**, τα οποία χρησιμοποιούμε ως βασική γραμμή σύγκρισης και τα οποία σε γενικές γραμμές επιτυγχάνουν τις καλύτερες επιδόσεις, ειδικά σε συνδυασμό με κατάλληλη επιλογή χαρακτηριστικών.

Κεντρικός στόχος είναι να δείξουμε πως συμπεριφέρονται τα SVMs σε σχέση με απλούς ταξινομητές απόστασης (kNN, NCC) καθώς και πως επηρεάζει η διαστάση και ο τύπος των χαρακτηριστικών την απόδοση των μοντέλων.

## Datasets και Pretraining

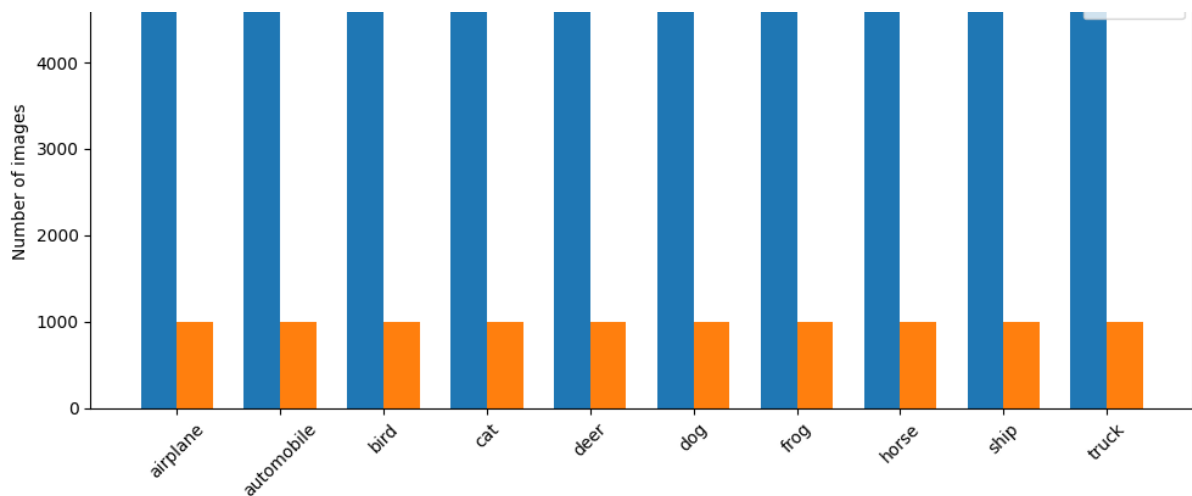
### CIFAR-10

Για την πειραματική μελέτη χρησιμοποιήθηκε το dataset CIFAR-10, το οποίο περιλαμβάνει 10 κλάσεις (airplane, automobile, bird, truck, etc) με εικόνες διαστάσεων 32x32x3 (RGB). Το σύνολο δεδομένων αποτελείται από 50 000 δείγματα για το train set και 10 000 δείγματα για το test set.

*Table 1 Τα 5 πρώτα παραδείγματα του CIFAR-10*



Table 2 Τα δεδομένα είναι ισοκατανεμημένα.



Ως βήματα προεπεξεργασίας, αρχικά έγινε φόρτωση των δεδομένων. Στην συνέχεια, οι εικόνες μετατράπηκαν σε διανύσματα με την εντολή

```
X_train_flat = X_train.reshape(len(X_train), -1)
```

και αντιστοίχα για το test set, έτσι ώστε κάθε εικόνα να αναπαριστάται ως ένα ενιαίο διάνυσμα χαρακτηριστικών. Οπότε από τις διαστάσεις (32, 32, 3) πάμε στις flattened 3072. Επιπλέον, οι τιμές των pixels κλιμακώθηκαν στην περιοχή [0, 1] με διαίρεση δια 255, ώστε να βελτιωθεί η αριθμητική σταθερότητα και η απόδοση των ταξινομητών.

Τα δεδομένα χρησιμοποιήθηκαν σε μορφή raw pixels για την εκπαίδευση ταξινομητών SVM, kNN και Nearest Class Centroid (NCC). Επιπλέον, εξηχθήσαν HOG χαρακτηριστικά πάνω στις ίδιες εικόνες και ξαναεκπαιδεύτηκαν πάνω σε αυτά, επιτυγχάνοντας αισθητά καλύτερες επιδόσεις σε σχέση με τα αντιστοίχα μοντέλα που εκπαιδεύτηκαν απευθείας στα raw pixels.

Table 3 Παραδειγμα HOG μετασχηματισμου.



Το **HOG** (*Histogram of Oriented Gradients*) υπολογίζει για κάθε μικρή περιοχή της εικόνας (cell) gradients και φτιάχνει ιστογράμματα των κατευθύνσεών τους. Τα histograms αυτά ομαδοποιούνται και κανονικοποιούνται σε blocks, έτσι κωδικοποιούνται οι ακμές και τα σχήματα της εικόνας, και έτσι δεν δίνεται βαρύτητα στις ακριβείς τιμές των pixels. Με αυτό τον τρόπο, τα HOG features είναι πιο συμπαγή, πιο ανθεκτικά σε αλλαγές φωτισμού και θορύβο και αναδεικνύουν τη δομή του

αντικειμενου, κατι που επιτρεπει στα μοντελα να βρουν ενα καθαροτερο και πιο γραμμικο οριο αποφασης σε σχεση με την εκπαιδευση πανω σε raw pixels.

## Breast Cancer

Για το δευτερο μερος της εργασιας χρησιμοποιηθηκε το dataset Breast Cancer, το οποιο αντιστοιχει σε προβλημα δυαδικης ταξινομησης μεταξυ καλοηθων και κακοηθων ογκων (κλασεις 2 για benign και 4 για malignant). Καθε δειγμα απο τα 683, που ειναι στο συνολο, περιγραφεται απο 9 χαρακτηριστικα, τα οποια ποσοτικοποιουν ιδιοτητες των κυτταρων σε κλιμακα απο 1 εως 10: Clump Thickness, Uniformity of Cell Size, Uniformity of Cell Shape, Marginal Adhesion, Single Epithelial Cell Size, Bare Nuclei, Bland Chromatin, Normal Nucleoli και Mitoses. Στα χαρακτηριστικα αυτα χρησιμοποιουμε την κλιμακωμενη (scaled) εκδοχη τους στο διαστημα  $[-1,1]$ , και για τα class labels απο 2,4 μετατραπηκαν στις τιμες 0,1, γεγονος που διευκολυνει την εκπαιδευση των μοντελων και την συγκρισιμοτητα μεταξυ των διαφορων διαστασεων του χωρου χαρακτηριστικων.

Για τον διαχωρισμο των δεδομενων σε συνολα εκπαιδευσης και ελεγχου χρησιμοποιηθηκε η συναρτηση `train_test_split` με ορισμα `test_size=0.3`, ωστε το 30% των δειγματων να αποτελει το test set και το υπολοιπο 70% το train set.

```
Train Shape: (478, 9)
```

```
Test Shape: (205, 9)
```

Επιπλεον, χρησιμοποιησαμε `random_state=0` για αναπαραγωγιμοτητα και χρησιμοποιηθηκε η παραμετρος `stratify=y`, ωστε να διασφαλιστει οτι η αναλογια των δυο κλασεων διατηρειται περιπου ιδια και στα συνολα train και test. Με τον τροπο αυτο αποτρεπεται η δημιουργια μη αντιπροσωπευτικων υποσυνολων και βελτιωνεται η αξιοπιστια της αξιολογησης των ταξινομητων.

## Models

### SVMs

Τα Support Vector Machines χρησιμοποιηθηκαν ως κυριο μοντελο αναφορας για την συγκριση των υπολοιπων μεθοδων. Η βασικη ιδεα των SVMs ειναι η αναζητηση ενος υπερεπιτεδου που διαχωριζει τις κλασεις με μεγαστο περιθωριο, μεσα απο την λυση ενος προβληματος βελτιστοποιησης που ισορροπει μεταξυ μεγιστοποιησης του περιθωριου (ορος regularization) και μειωσης του σφαλματος ταξινομησης, μεσω της παραμετρου C που ελεγχει το trade off μεταξυ των δυο στοχων. Στην εργασια δοκιμαστηκαν τοσο γραμμικα οσο και μη γραμμικα SVMs, χρησιμοποιωντας τις υλοποιησεις LinearSVC και SVC της βιβλιοθηκης scikit learn.

### CIFAR 10

Στο CIFAR 10, λογω του υψηλου υπολογιστικου κοστους σε ολοκληρο το dataset, ο κυριος πειραματισμος με τις υπερπαραμετρους εγινε σε ενα stratified subset (σταθερος αριθμος δειγματων ανα κλαση), οπου δοκιμαστηκαν γραμμικα SVMs (LinearSVC), SVM με πολυωνυμικο kernel βαθμου 3, καθως και RBF SVM με διαφορετικους συνδυασμους των C και gamma το οποιο ελεγχει ποσο "κοντα" πρεπει να ειναι δυο δειγματα για να επηρεαζουν πολυ το ενα το αλλο. Αν το gamma ειναι μεγαλο καθε δειγμα εχει μικρη εμβελεια και το SVM μπορει να φτιαξει πιο πολυπλοκα συνορα με κινδυνο, ομως, να παθει overfitting. Ενω αν ειναι μικρο το gamma, καθε δειγμα επηρεαζει μεγαλυτερη περιοχη, τα συνορα γινονται πιο λεια και "απλα".

Τα αποτελεσματα ειναι τα εξης:

Table 4 CIFAR 10 – SVM σε raw pixels (stratified subset, 10 000 train / 2 000 test)

Μοντελο	Χρονος train	Test accuracy	Macro F1
LinearSVC_C1	27.46 sec	0.3385	0.3191
SVC_poly_deg3_C1	3 min 35 sec	0.4400	0.4412
SVC_linear_C1	6 min 2 sec	0.2540	0.2167
SVC_rbf_C1_gamma1e-3	6 min 19 sec	0.2430	0.1965
SVC_rbf_C10_gamma1e-3	6 min 18 sec	0.2430	0.1965
SVC_rbf_C100_gamma1e-3	377.75	0.2430	0.1965

Στην συνεχεια τρεξαμε ξανα στο full dataset για το γραμμικο μοντελο λογω του οτι είναι το πιο γρηγορο αλλα και για το πολυωνυμικο το οποιο ειχε τα καλυτερα αποτελεσματα:

Table 5 CIFAR 10 – SVM σε raw pixels (full dataset, 50 000 train / 10 000 test)

Μοντελο	Χρονος train	Test accuracy	Macro F1
LinearSVC_C1	3 min 5 sec	0.3854	0.3707
SVC_poly_deg3_C1	1 hr 31 min	0.5244	0.5254

Στη συνεχεια, υπολογιστηκαν **HOG** χαρακτηριστικα για τις εικονες και εκπαιδευτηκαν εκ νεου SVMs πανω σε αυτο το χωρο χαρακτηριστικων. Στο pipeline χρησιμοποιειται αρχικα το StandardScaler, το οποιο αφαιρει τη μεση τιμη και διαιρει με την τυπικη αποκλιση καθε χαρακτηριστικου, ωστε ολα τα features να βρισκονται σε συγκρισιμη κλιμακα. Αυτο ειναι σημαντικό για τον RBF πυρηνα, επειδη βασιζεται σε ευκλειδειες αποστασεις οπου χωρις καλη κλιμακωση, λιγα features με μεγαλες τιμες θα κυριαρχουσαν στον υπολογισμο της αποστασης και το SVM δεν θα μπορουσε να μαθει ενα σταθερο και καλα γενικευσιμο συνоро αποφασης.

Δοκιμαζουμε LinearSVC και στη συνεχεια RBF SVM με  $C = 10$  και  $\gamma = \text{"scale"}$ , στο ιδιο stratified subset. Το “scale” σημαινει οτι αν τα features εχουν μεγαλη διασπορα, το  $\gamma$  γινεται μικρο, ενω αν τα features εχουν μικρη διασπορα, το  $\gamma$  γινεται μεγαλυτερο. Ο συνδυασμος RBF SVM σε HOG χαρακτηριστικα, αποδειχθηκε ο πιο αποδοτικος και για αυτο εκπαιδευτηκε και σε ολοκληρο το CIFAR 10, αποτελωντας το τελικο SVM baseline για το συγκεκριμενο dataset.

Table 6 CIFAR 10 – SVM σε HOG χαρακτηριστικα (stratified subset, 10 000 train / 2 000 test)

Μοντελο	Χρονος train (sec)	Test accuracy	Macro F1
LinearSVC σε HOG	15.54	0.4780	0.4713
RBF SVC σε HOG	25.19	0.5590	0.5585

Τελος, τρεχουμε και στο full dataset το RBF αφου παρατηρουμε οτι τα results είναι promising με δυο διαφορετικες τιμες  $C$ .

Table 7 CIFAR 10 – SVM σε HOG χαρακτηριστικά (full dataset, 50 000 train / 10 000 test)

Μοντελο	Χρονος train (sec)	Test accuracy	Macro F1
RBF SVC σε HOG full (C = 10, gamma=scale)	1087.53	0.6362	0.6360
RBF SVC σε HOG full (C = 150, gamma=scale)	1203.82	0.6366	0.6364

Συμπερασματικά, σε raw pixels τα SVMs δίνουν μετρια αποδοση στο CIFAR 10. Ο πολυωνυμικός kernel σε full dataset βελτιώνει σημαντικά τα metrics, ενώ επιπλέον βελτίωση προκύπτει με την χρήση HOG χαρακτηριστικών. Με RBF SVM σε HOG (ιδίως στο full dataset με C = 150) αποτελεί το καλύτερο SVM baseline για CIFAR 10 στην πειραματική διαδικασία.

### Breast Cancer

Για αυτό το dataset για την επιλογή υπερπαραμετρών χρησιμοποιήθηκε 5-fold stratified cross validation (StratifiedKFold με n\_splits=5, shuffle=True, random\_state=0) σε συνδυασμό με Grid Search (GridSearchCV), με κριτήριο scoring="f1" πάνω στο training set. Δοκιμάστηκαν τρία βασικά SVM μοντέλα ώστε να συγκριθούν γραμμικά και μη γραμμικά συνόρα απόφασης στο ίδιο dataset.

- SVC με γραμμικό kernel (kernel="linear"),
- SVC με RBF kernel (kernel="rbf") και
- SVC με πολυωνυμικό kernel (kernel="poly") με βαθμό 3.

Table 8 Αποτελέσματα SVM στο breast cancer (70/30 split, 5 fold CV).

Μοντελο	Best params	Χρονος train (sec, wall)	Test accuracy	Test F1
SVM_poly	C=0.1, degree=3, coef0=0.0, gamma='scale'	0.42	0.9756	0.9660
SVM_linear	C=0.1	2.05	0.9610	0.9444
SVM_rbf	C=10, gamma=0.01	1.45	0.9610	0.9444

Ποιοτικά, όλα τα SVM μοντέλα αποδίδουν πολύ καλά στο breast cancer dataset. Οι τιμές test accuracy είναι πάνω από 0.96, ενώ το F1 score κινείται περίπου στο εύρος 0.94–0.97. Όπως φαίνεται και από τα classification reports, υπάρχει πολύ καλή ισορροπία μεταξύ precision και recall και για τις δύο κλάσεις (benign και malignant), χωρίς κάποια εμφάνιση μεροληψία υπέρ της μιας ή της άλλης.

Το SVM με polynomial kernel βαθμού 3 (SVM\_poly) προκύπτει ως το καλύτερο μοντέλο στο notebook και πετυχαίνει accuracy γύρω στο 0.976 και F1 περίπου 0.966. Στην confusion matrix βλέπουμε ότι η benign κλάση ταξινομείται σχεδόν τέλεια (129 σωστά, 4 λάθος), ενώ και για την malignant κλάση οι επιδόσεις είναι αντίστοιχα πολύ καλές (71 σωστά, 1 λάθος). Επιπλέον, όλα αυτά επιτυγχάνονται με πολύ μικρό χρόνο εκπαίδευσης καθώς προκειται για ένα μικρό dataset.

Τα SVMs με γραμμικό kernel (SVM\_linear) και με RBF kernel (SVM\_rbf) εμφανίζουν σχεδόν ταυτιζόμενα αποτελέσματα: accuracy γύρω στο 0.961 και F1 περίπου 0.944, με confusion matrix  $\begin{bmatrix} 129 & 4 \\ 4 & 68 \end{bmatrix}$ . Αυτά τα μοντέλα προσφέρουν ήδη πολύ υψηλή απόδοση και λειτουργούν ως ισχυρά baselines, δείχνοντας ότι ακόμα και σχετικά απλά SVMs, πάνω σε καλά χαρακτηριστικά, είναι ικανά να λύσουν αποτελεσματικά το συγκεκριμένο πρόβλημα ταξινόμησης.

## k-Nearest Neighbors (kNN) και Nearest Class Centroid (NCC)

### CIFAR 10

Ουσιαστικά, και τα δύο μοντέλα (kNN και Nearest Class Centroid) βασίζονται στις αποστάσεις μεταξύ δειγμάτων στον χώρο των χαρακτηριστικών για να προβλέψουν την κλάση. Δεν κάνουν πολυπλοκή μαθηση παραμετρών όπως τα SVMs, αλλά χρησιμοποιούν την γεωμετρία των δεδομένων για να πάρουν αποφάσεις ταξινόμησης.

Στο **kNN**, για κάθε νέο δείγμα αναζητούνται οι  $k$  κοντινότεροι γείτονες (π.χ. με ευκλείδεια ή manhattan απόσταση) και η προβλεψή βασίζεται στην πλειοψηφία των κλάσεων τους. Όταν χρησιμοποιούμε `weights="distance"`, οι πιο κοντινοί γείτονες έχουν μεγαλύτερη επιδραση στην απόφαση, κάτι που μπορεί να βελτιώσει την απόδοση σε περιπτώσεις όπου η τοπική γειτονία είναι πυκνή.

Στο Nearest Class Centroid (**NCC**), η λογική είναι πιο απλή. Για κάθε κλάση υπολογίζεται ένα centroid, δηλαδή το μέσο διάνυσμα όλων των δειγμάτων της. Ένα νέο δείγμα ταξινομείται στην κλάση της οποίας το centroid βρίσκεται πιο κοντά. Έτσι, κάθε κλάση αναπαριστάται από ένα μόνο σημείο στον χώρο των χαρακτηριστικών, κάτι που κάνει το μοντέλο ιδιαίτερα γρήγορο και εύκολα ερμηνεύσιμο, αλλά λιγότερο αποδοτικό σε σχέση με πιο πολυπλοκά μοντέλα.

Όπως και στα SVMs αρχικά δουλεύουμε με ένα subset του CIFAR-10, ώστε να βρούμε γρήγορα τις καλύτερες παραμέτρους. Επιλέγουμε ένα stratified subset με 10,000 δείγματα για train και 2,000 για test, διατηρώντας την ίδια αναλογία κλάσεων όπως στο αρχικό σύνολο. Για το **kNN** χρησιμοποιήθηκε Grid Search με ένα πλέγμα τιμών για τον αριθμό γειτόνων `n_neighbors` (1, 3, 5, 7, 9, 11), δύο μετρικές απόστασης (euclidean, manhattan) και δύο σχήματα βαρών (uniform, distance), ενώ ο αλγόριθμος αναζήτησης γειτόνων είναι στο "auto" ώστε να επιλεγεται αυτομάτως η καταλληλότερη στρατηγική. Αντιστοίχα, για το *Nearest Class Centroid* ορίστηκε grid πάνω στη metric (euclidean, manhattan) και στην παράμετρο `shrink_threshold`, η οποία μπορεί να εφαρμόσει shrinkage στα centroids (τιμές [None, 0.05, 0.1, 0.5, 1.0]), με στόχο να βελτιωθεί η σταθερότητα σε θορυβώδη δεδομένα υψηλής διαστάσης.

Και στις δύο περιπτώσεις η επιλογή των υπερ-παραμετρών έγινε με GridSearchCV, χρησιμοποιώντας `scoring="accuracy"`, 3-fold cross validation (`cv=3`) και `n_jobs=-1` για παράλληλη εκτέλεση, ώστε να βρεθούν οι συνδυασμοί που δίνουν την καλύτερη απόδοση στο stratified subset του CIFAR 10.

Table 9 Αποτελέσματα subset του CIFAR 10 για  $k$  Nearest Neighbors (kNN) και Nearest Class Centroid (NCC).

Μοντέλο	Καλύτερα params	Χρόνος εκπαίδευσης (GridSearchCV, sec)	Test accuracy	Test macro F1
<b>kNN</b>	algorithm = auto, metric = manhattan, <code>n_neighbors</code> = 9, weights = distance	268.50	0.3225	0.3100

Μοντελο	Καλυτερα params	Χρονος εκπαιδευσης (GridSearchCV, sec)	Test accuracy	Test macro F1
NCC	metric = manhattan, shrink_threshold = None	7.96	0.2615	0.2416

Επειτα, με τις καλυτερες υπερ-παραμετρους που βρισκουµε το τρεχουµε στο full dataset (50 000 train / 10 000 test), οπου παρατηρειται βελτιωση στον kNN.

Table 10 Αποτελεσµατα στο full dataset CIFAR 10, για k Nearest Neighbors (kNN) και Nearest Class Centroid (NCC).

Μοντελο	Params (απο subset)	Χρονος fit (sec)	Test accuracy	Test macro F1
kNN	algorithm = auto, metric = manhattan, n_neighbors = 9, weights = distance	0.17	0.3952	0.3895
NCC	metric = manhattan, shrink_threshold = None	8.23	0.2734	0.2528

Στην πειραµατικη µας διαδικασια για να εχουµε σωστη συγκριση µε τα SVMs εφαρµοστηκε η ιδια µεθοδος µε τα HOG features οπως περιγραφηκε στη ενοτητα SVMs και βρηκαµε αυτα τα αποτελεσµατα:

Table 11 CIFAR 10 – kNN / NCC σε HOG + StandardScaler (stratified subset, 10 000 train / 2 000 test)

Μοντελο	Best params (απο grid σε raw pixels)	Χρονος fit (sec)	Test accuracy	Test macro F1
kNN σε HOG	algorithm = auto, metric = manhattan, n_neighbors = 9, weights = distance	0.00	0.4780	0.4681
NCC σε HOG	metric = manhattan, shrink_threshold = None	0.08	0.4070	0.4031

Table 12 CIFAR 10 – kNN / NCC σε HOG + StandardScaler (FULL dataset, 50 000 train / 10 000 test)

Μοντελο	Best params (απο grid σε raw pixels)	Χρονος fit (sec)	Test accuracy	Test macro F1
kNN σε HOG full	algorithm = auto, metric = manhattan, n_neighbors = 9, weights = distance	0.01	0.5435	0.5356

Μοντελο	Best params (απο grid σε raw pixels)	Χρονος fit (sec)	Test accuracy	Test macro F1
NCC σε HOG full	metric = manhattan, shrink_threshold = None	0.40	0.4084	0.4034

### Breast Cancer

Για τον ταξινομητή **kNN** εφαρμόστηκε διαδικασία grid search προκειμένου να επιλεγουν οι βελτιστές υπερπαραμετροί. Συγκεκριμένα, δοκιμάστηκαν διαφορές τιμές για τον αριθμό γειτόνων  $n\_neighbors$  ([1, 3, 5, 7, 9, 11, 15]), δύο μετρικά απόστασης (euclidean και manhattan), καθώς και δύο σχήματα βαρών (uniform και distance). Η εκπαίδευση έγινε με χρήση Grid Search πάνω στο train set, το οποίο περιλαμβάνει 478 δείγματα, χρησιμοποιώντας 5 fold stratified cross validation ώστε να διατηρείται η αναλογία των κλάσεων σε κάθε fold. Αφού ολοκληρωθεί το grid search, το μοντέλο επανεκπαίδευεται (refit) στα συνολικά δεδομένα εκπαίδευσης χρησιμοποιώντας τον συνδυασμό υπερπαραμετρών που έδωσε τη μέγιστη τιμή στην cross-validation accuracy.

Για την αξιολόγηση χρησιμοποιείται το test set, το οποίο αποτελείται από 205 δείγματα. Ως μετρικές επίδοσης υπολογίζονται η accuracy και η F1, σε δυαδική μορφή, με τη θετική κλάση να αντιστοιχεί στα malignant περιστατικά.

Table 13 Αποτελέσματα kNN

Μοντελο	Best params	Χρονος train (sec)	Test accuracy	Test F1 (binary)
kNN	metric = euclidean, $n\_neighbors = 7$ , weights = uniform	2.7610	0.9610	0.9437

Για τον ταξινομητή Nearest Class Centroid (**NCC**) εφαρμόστηκε αντιστοιχη διαδικασία grid search ώστε να επιλεγουν οι βελτιστές υπερπαραμετροί. Συγκεκριμένα, δοκιμάστηκαν δύο μετρικά απόστασης (euclidean και manhattan), καθώς και διαφορετικές τιμές για την παραμετρο shrink\_threshold ([None, 0.05, 0.1, 0.5, 1.0]), η οποία ελέγχει τον βαθμό συρρικνώσης των κεντροειδών. Έτσι, πάλι με Grid Search έχουμε τα εξής αποτελέσματα:

Table 14 Αποτελέσματα NCC.

Μοντελο	Best params	Χρονος train (sec)	Test accuracy	Test F1 (binary)
NCC	metric = euclidean, shrink_threshold = None	0.1149	0.9659	0.9510

Συμπερασματικά, και ο kNN και ο NCC πετυχαίνουν πολύ υψηλή επίδοση στο breast cancer dataset. Ο απλός γραμμικός NCC φτάνει και ξεπερνά ελαφρώς τον kNN, ενώ εκπαιδεύεται σημαντικά πιο γρήγορα.



## Συζήτηση

Τα αποτελέσματα της μελέτης δείχνουν καθαρά ότι ο ρόλος των χαρακτηριστικών είναι καθοριστικός για την απόδοση των ταξινομητών. Όταν χρησιμοποιούμε raw pixels, η ευκλείδεια ή η manhattan απόσταση μεταξύ των εικόνων δεν αντανάκλα πάντα την πραγματική οπτική ομοιότητα, με αποτέλεσμα τα μοντέλα βασισμένα σε αποστάσεις, όπως kNN και NCC, να παρουσιάζουν περιορισμένη επίδοση. Ακόμη και τα SVMs, σε αυτή τη μορφή δεδομένων, χρειάζονται ισχυρά kernels για να μπορέσουν να μοντελοποιήσουν σωστά τα όρια απόφασης. Αντιθέτως, με HOG χαρακτηριστικά στο CIFAR-10, η απόδοση των SVMs βελτιώνεται θεαματικά, δείχνοντας ότι η σωστή αναπαράσταση των δεδομένων είναι κρίσιμη προϋπόθεση για καλά αποτελέσματα.

Τα SVMs αναδεικνύονται ως πολύ καλή γενική λύση, ιδίως σε προβλήματα υψηλής διαστάσης. Σε τέτοιες συνθήκες μπορούν να διαχειριστούν καλύτερα το πρόβλημα σε σχέση με τον kNN/NCC, καθώς δεν βασίζονται στον υπολογισμό πλήρων αποστάσεων από όλα τα δείγματα του train set για κάθε νέο δείγμα. Με την χρήση του kernel trick (π.χ. RBF kernel) μπορούν να προσαρμόσουν πολύπλοκα, μη-γραμμικά όρια απόφασης στον χώρο των χαρακτηριστικών, διατηρώντας παράλληλα καλή ικανότητα γενίκευσης.

Όμως, και τα μοντέλα απόστασης, όπως kNN και NCC είναι ιδιαίτερα χρήσιμα ως baselines, ιδίως σε μικρά, καλά κλιμακωμένα datasets. Στο breast cancer dataset, για παράδειγμα, kNN και NCC παρέχουν γρήγορα και εύκολα στην υλοποίηση μοντέλα με απόδοση που συχνά πλησιάζει αυτή των SVMs, ιδιαίτερα όταν η σχέση μεταξύ χαρακτηριστικών και κλάσης είναι σχετικά απλή. Ο NCC, συγκεκριμένα, προσφέρει ένα εξαιρετικά ελαφρύ υπολογιστικό μοντέλο, το οποίο, παρά την απλότητά του, μπορεί σε ορισμένες περιπτώσεις να πλησιάσει ή και να ξεπεράσει πιο σύνθετους ταξινομητές.