

**ΑΛΕΞΑΝΔΡΕΙΟ ΤΕΙ ΘΕΣΣΑΛΟΝΙΚΗΣ**

**ΤΜΗΜΑ ΜΗΧΑΝΙΚΩΝ ΠΛΗΡΟΦΟΡΙΚΗΣ**

**ΜΑΘΗΜΑ: ΜΗΧΑΝΙΚΗ ΜΑΘΗΣΗ**

**ΚΑΘΗΓΗΤΕΣ : ΚΩΣΤΑΣ ΔΙΑΜΑΝΤΑΡΑΣ, ΚΩΣΤΑΣ ΓΟΥΛΙΑΝΑΣ**

## **ΕΡΓΑΣΤΗΡΙΑΚΗ ΑΣΚΗΣΗ 2**

### **ΓΡΑΜΜΙΚΟΣ ΤΑΞΙΝΟΜΗΤΗΣ**

**Σκοπός της άσκησης:** Η εκτίμηση της επίδοσης ενός γραμμικού ταξινομητή δύο κλάσεων σε άγνωστα δεδομένα. Θα γίνει χρήση της μεθόδου διασταύρωσης (Cross-Validation).

Θα χρησιμοποιηθούν τα εξής κριτήρια επίδοσης:

1. Ακρίβεια (accuracy)
2. Ευστοχία (precision)
3. Ανάκληση (recall)
4. F-Measure
5. Ευαισθησία (Sensitivity)
6. Προσδιοριστικότητα (Specificity)

Δείτε τη σημασία αυτών των κριτηρίων στο επισυναπτόμενο κείμενο [Κριτήρια επίδοσης ταξινομητών.pdf](#)

#### **Βήματα υλοποίησης:**

1. Χρησιμοποιήστε το σύνολο δεδομένων IRIS από το προηγούμενο εργαστήριο, καθώς και τον κώδικα από το εργαστήριο αυτό. Θυμίζουμε ότι τα πρότυπα χωρίστηκαν σε δύο κλάσεις ως εξής:
  - Κλάση 0 (στόχος  $t=0$ ): αποτελείται από τα πρότυπα των κατηγοριών "Iris-setosa" + "Iris-virginica",
  - Κλάση 1 (στόχος  $t=1$ ): αποτελείται από τα πρότυπα της κατηγορίας "Iris-versicolor".
2. Θα εφαρμοστεί η μέθοδος `train_test_split(...)` για  $K=9$  folds. Στο αντίστοιχο loop θα πρέπει να κάνετε τα εξής:
  - Επαυξήστε τον πίνακα των προτύπων προσθέτοντας σε κάθε πρότυπο τον αριθμό 1, δηλαδή προσθέστε μια γραμμή με 1 στον πίνακα  $x$ . Χρησιμοποιήστε τη συνάρτηση `hstack()` του `numpy`.

**Για κάθε fold:**

- Έχετε ήδη δημιουργήσει τους πίνακες `xtrain`, `xtest` καθώς και τα διανύσματα στόχων `ttrain`, `ttest`. Φροντίστε να είναι `numpy arrays` τύπου `float`.
- Βρείτε το πλήθος των προτύπων στο train set ( $P_{train}$ ) και στο test set ( $P_{test}$ ) για παράδειγμα χρησιμοποιώντας τη συνάρτηση `len()`.
- Μετατρέψτε τα διανύσματα στόχων `ttrain[]`, `ttest[]` έτσι ώστε
  - `Av ttrain[pattern] == 1` → `ttrain1[pattern] = 1`
  - `Av ttrain[pattern] == 0` → `ttrain1[pattern] = -1`
  - `Av ttest[pattern] == 1` → `ttest1[pattern] = 1`
  - `Av ttest[pattern] == 0` → `ttest1[pattern] = -1`

- Βρείτε το διάνυσμα βαρών  $\tilde{\mathbf{w}}$  του γραμμικού ταξινομητή

$$y = \tilde{\mathbf{x}} \tilde{\mathbf{w}}$$

$$\tilde{\mathbf{w}} = \begin{bmatrix} \mathbf{w} \\ w_0 \end{bmatrix}, \quad \tilde{\mathbf{x}} = [\mathbf{x} \ 1]$$

όπου  $\mathbf{w}$  είναι το διάνυσμα βαρών,  $w_0$  είναι η πόλωση, και  $\mathbf{x}$  είναι το πρότυπο εισόδου.

- Υπολογίστε το βέλτιστο διάνυσμα βαρών  $\tilde{\mathbf{w}}$ . Σύμφωνα με τη θεωρία των γραμμικών ταξινομητών, το βέλτιστο διάνυσμα είναι

$$\tilde{\mathbf{w}} = \tilde{\mathbf{X}}_{train}^+ \mathbf{t}_{train}$$

όπου

$\mathbf{t}_{train} = [t_1 \ t_2 \ \dots \ t_p]^T$  είναι το διάνυσμα των τροποποιημένων στόχων (-1/1),

$$\tilde{\mathbf{X}}_{train} = \begin{bmatrix} \tilde{\mathbf{x}}_1 \\ \tilde{\mathbf{x}}_2 \\ \vdots \\ \tilde{\mathbf{x}}_p \end{bmatrix} = \begin{bmatrix} \mathbf{x}_1 & 1 \\ \mathbf{x}_2 & 1 \\ \vdots & \vdots \\ \mathbf{x}_p & 1 \end{bmatrix} \text{ είναι ο πίνακας των επαυξημένων προτύπων του train set,}$$

$\mathbf{X}^+$  είναι ο τελεστής του ψευδο-αντίστροφου του πίνακα  $\mathbf{X}$ . Στην python ο ψευδο-αντίστροφος υλοποιείται με τη συνάρτηση `numpy.linalg.pinv()`.

- Υπολογίστε την έξοδο του ταξινομητή για όλα τα πρότυπα του test set:

$$\mathbf{y}_{test} = \tilde{\mathbf{X}}_{test} \tilde{\mathbf{w}}$$

- Υπολογίστε την εκτίμηση που κάνει ο ταξινομητής για την κλάση στην οποία ανήκουν τα πρότυπα του test set:

$$predict_{test}(i) = \begin{cases} 0, & y_{test}(i) < 0 \\ 1, & y_{test}(i) \geq 0 \end{cases}$$

- Υλοποιήστε τη συνάρτηση `evaluate()` με τρεις εισόδους και μια έξοδο ως εξής:

```
def evaluate( t, predict, criterion ):
    """
    Είσοδος t : διάνυσμα με τους πραγματικούς στόχους (0/1)
    Είσοδος predict : διάνυσμα με τους εκτιμώμενους στόχους (0/1)
    Είσοδος criterion : text-string με τις εξής πιθανές τιμές:
        'accuracy'
        'precision'
        'recall'
        'fmeasure'
        'sensitivity'
        'specificity'
    Έξοδος value : η τιμή του κριτηρίου που επιλέξαμε.
    """
```

- Πρώτα υπολογίστε τα true-negatives, false-negatives, true-positives, false positives, τα οποία ορίζονται ως εξής:

(α) **true negatives** (πραγματικά αρνητικά) οι περιπτώσεις όπου:  
 το πρότυπο βγήκε αρνητικό, δηλ.  $y_i = 0$   
 και όντως ανήκει στην κλάση 0, δηλ.  $t_i = 0$

(β) **false negatives** (εσφαλμένα αρνητικά) οι περιπτώσεις όπου:

το πρότυπο βγήκε αρνητικό, δηλ.  $y_i = 0$   
αλλά ανήκει στην κλάση 1, δηλ.  $t_i = 1$

(γ) **true positives** (πραγματικά θετικά) οι περιπτώσεις όπου:

το πρότυπο βγήκε θετικό, δηλ.  $y_i = 1$   
και όντως ανήκει στην κλάση 1, δηλ.  $t_i = 1$

(δ) **false positives** (εσφαλμένα θετικά) οι περιπτώσεις όπου:

το πρότυπο βγήκε θετικό, δηλ.  $y_i = 1$   
αλλά ανήκει στην κλάση 0, δηλ.  $t_i = 0$

Σημαντικό! Φροντίστε τα  $tp$ ,  $tn$ ,  $fp$ ,  $fn$  να είναι τύπου **float**. Χρησιμοποιήστε την συνάρτηση **float()**

Κατόπιν

- Αν επιλεγεί **criterion = 'accuracy'** τότε η συνάρτηση επιστρέφει

$$value = Accuracy = \frac{tp + tn}{tp + tn + fp + fn}$$

- Αν επιλεγεί **criterion = 'precision'** τότε η συνάρτηση επιστρέφει

$$value = Precision = \frac{tp}{tp + fp}$$

- Αν επιλεγεί **criterion = 'recall'** τότε η συνάρτηση επιστρέφει

$$value = Recall = \frac{tp}{tp + fn}$$

- Αν επιλεγεί **criterion = 'fmeasure'** τότε η συνάρτηση επιστρέφει

$$value = Fmeasure = \frac{Precision \cdot Recall}{(Precision + Recall)/2}$$

- Αν επιλεγεί **criterion = 'sensitivity'** τότε η συνάρτηση επιστρέφει

$$value = Sensitivity = \frac{tp}{tp + fn}$$

- Αν επιλεγεί **criterion = 'specificity'** τότε η συνάρτηση επιστρέφει

$$value = Specificity = \frac{tn}{tn + fp}$$

- Καλέστε τη συνάρτηση **evaluate()** όσες φορές χρειάζεται έτσι ώστε για το συγκεκριμένο fold να υπολογίσετε το Accuracy, Precision, Recall, F-measure, Sensitivity και Specificity.
- Χρησιμοποιώντας κατάλληλο subplot σε grid 3x3 στο figure(1) τυπώστε το εξής γράφημα:

- δείξτε με μπλε τελείες τους πραγματικούς στόχους  $t_{test}[i]$  για όλα τα πρότυπα του test set
- δείξτε με κόκκινους κύκλους τους εκτιμώμενους στόχους  $predict_{test}[i]$  για όλα τα πρότυπα του test set

# end for

Μετά το τέλος του loop υπολογίστε και τυπώστε στην οθόνη τα εξής:

1. τη μέση τιμή του Accuracy για όλα τα folds
2. τη μέση τιμή του Precision για όλα τα folds
3. τη μέση τιμή του Recall για όλα τα folds
4. τη μέση τιμή του F-Measure για όλα τα folds
5. τη μέση τιμή του Sensitivity για όλα τα folds
6. τη μέση τιμή του Specificity για όλα τα folds

Θα χρησιμοποιήσετε τις παρακάτω εντολές ή συναρτήσεις:

- Υλοποιήστε τη λειτουργία switch με dictionary.
- Χρησιμοποιήστε τη συνάρτηση pinv() του numpy.linalg για τον υπολογισμό ψευδο-αντίστροφου ενός πίνακα

# Παράδειγμα:

```
>>> import numpy as np
>>> A = np.array([[0.3, 0.1, 0.2],
                  [0.5, 1, 0.2]])
>>> A1 = np.linalg.pinv(A)
>>> A1
array([[ 2.50777202, -0.1761658 ],
       [-1.66839378,  1.15025907],
       [ 2.07253886, -0.31088083]])
```

- Πολλαπλασιασμός πίνακα επί διάνυσμα

```
>>> w = np.array([3, -2, 0.5])
>>> A.dot(w)
array([ 0.8, -0.4])
>>> np.matmul(A, w)
array([ 0.8, -0.4]) # Ακριβώς ίδιο με το .dot()
```

- Λογικοί τελεστές με διανύσματα:

# Παράδειγμα:

```
>>> z = np.array([0.1, -0.5, 0, 0.5, -0.1, 0.3, 0.2])
>>> y1 = (z>0)
>>> y1
array([ True, False, False,  True, False,  True,  True], dtype=bool)
>>> y2 = (z==0)
>>> y2
```

```
array([False, False,  True, False, False, False, False], dtype=bool)
>>> y3
array([False,  True, False, False,  True, False, False], dtype=bool)
```

Στην Python υπάρχουν οι Boolean σταθερές *True* και *False*. Συνεπώς μπορούμε να κάνουμε λογικές πράξεις μεταξύ binary διανυσμάτων ως εξής:

```
# Παράδειγμα λογικού OR:
>>> y1|y2
array([ True, False,  True,  True, False,  True,  True], dtype=bool)
# Παράδειγμα λογικού NOT:
>>> y1
array([ True, False, False,  True, False,  True,  True], dtype=bool)
>>> ~y1
array([False,  True,  True, False,  True, False, False], dtype=bool)
```

- Η συνάρτηση `sum()`: αν το όρισμά της είναι διάνυσμα, επιστρέφει ένα αριθμό που είναι το άθροισμα των στοιχείων του διανύσματος. Αν το όρισμα είναι πίνακας υπολογίζει το άθροισμα των στοιχείων κάθε στήλης και επιστρέφει διάνυσμα.

```
# Παράδειγμα:
>>> z = np.array([0.1, -0.5, 0, 0.5, -0.1, 0.3, 0.2])
>>> y
array([ True, False, False,  True, False,  True,  True], dtype=bool)
>>> sum(y)
4
```