

Πρώτη Σειρά Ασκήσεων

Στυλιανός Σύρρος: 4805

Διάβασμα Αρχείων

Στην παρακάτω εικόνα φαίνεται ο τρόπος με τον οποίο διαβάζουμε τα δεδομένα των δύο αρχείων.

```
24 airlines_delay = pd.read_csv("airlines_delay.csv")
25 sample_size = int(len(airlines_delay) * 0.01)
26 random_indices = random.sample(range(len(airlines_delay)), sample_size)
27 airlines_delay = airlines_delay.iloc[random_indices]
28 mobile_data = pd.read_csv("train.csv")
29
30 le = LabelEncoder()
31 airlines_delay['Airline'] = le.fit_transform(airlines_delay['Airline'])
32 airlines_delay['AirportFrom'] = le.fit_transform(airlines_delay['AirportFrom'])
33 airlines_delay['AirportTo'] = le.fit_transform(airlines_delay['AirportTo'])
34 airlines_delay = airlines_delay.drop('Flight', axis=1)
35
36 X_airline = airlines_delay.iloc[:, :-1].values
37 Y_airline = airlines_delay.iloc[:, -1].values
38 X_mobile = mobile_data.iloc[:, :-1].values
39 Y_mobile = mobile_data.iloc[:, -1].values
40
41 X_train_airline, X_test_airline, y_train_airline, y_test_airline = train_test_split(X_airline, Y_airline, test_size=0.3, random_state=42)
42 X_train_mobile, X_test_mobile, y_train_mobile, y_test_mobile = train_test_split(X_mobile, Y_mobile, test_size=0.3, random_state=42)
```

- Στις γραμμές 24-27 διαβάζουμε το αρχείο `airlines_delay.csv` και επιλέγουμε μόνο το 1% των δεδομένων τυχαία επειδή ο όγκος πληροφορία ήταν μεγάλος. Στην γραμμή 34 διαβάζουμε το άλλο αρχείο
- Στις γραμμές 31-33 μετατρέπουμε τα αλφαριθμητικά σε ακέραιες τιμές
- Στην γραμμή 34 διαγράφουμε την στήλη `Flight`
- Στις γραμμές 36-39 διαχωρίζουμε τις διαστάσεις από τις κατηγορίες
- Στις γραμμές 41-42 χωρίζουμε το σύνολο σε σύνολο `trainning` και σύνολο `test` και για τα δύο αρχεία.

Μέθοδος K – Nearest Neighbors

Αρχικά εφόσον θέλουμε να χρησιμοποιήσουμε ευκλείδεια απόσταση για τις συνεχείς μεταβλητές και Hamming απόσταση για τις διακριτές μεταβλητές, φτιάξαμε την συνάρτηση `custom_distance_airline` (για το αρχείο `airlines_delay.csv`) και `custom_distance_train` (για το

αρχείο train.csv) υπολογίζει την απόσταση με τον τρόπο αυτό. Παρακάτω φαίνονται αυτές οι συναρτήσεις.

```
dist_cont = DistanceMetric.get_metric('euclidean')
dist_disc = DistanceMetric.get_metric('hamming')

def custom_distance_train(x1, x2):
    x1_disc = x1[:, [1, 3, 5, 17, 18, 19]]
    x1_cont = x1[:, [i for i in range(x1.shape[1]) if i not in [1, 3, 5, 17, 18, 19]]]
    x2_disc = x2[:, [1, 3, 5, 17, 18, 19]]
    x2_cont = x2[:, [i for i in range(x2.shape[1]) if i not in [1, 3, 5, 17, 18, 19]]]

    dist_cont_mat = dist_cont.pairwise(x1_cont, x2_cont)
    dist_disc_mat = dist_disc.pairwise(x1_disc, x2_disc)

    combined_dist = dist_cont_mat + dist_disc_mat
    return combined_dist

def custom_distance_airlines(x1, x2):
    dist_cont_mat = dist_cont.pairwise(x1, x2)

    combined_dist = dist_cont_mat
    return combined_dist
```

Ακολούθως φτιάξαμε την μέθοδο KNN η οποία επιστρέφει το y_pred που είναι η εκτίμηση που θα συγκρίνουμε με τα πραγματικά δεδομένα για να βρούμε τα scores των μετρικών που χρησιμοποιούμε. Παρακάτω φαίνεται η συνάρτηση knn.

```
def knn(k, distance_fn, X_train, y_train, X_test):
    dist_mat = distance_fn(X_train, X_test)

    nn_indices = np.argpartition(dist_mat, k, axis=0)[:k, :]

    nn_labels = y_train[nn_indices]

    y_pred = np.empty(X_test.shape[0])
    for i in range(X_test.shape[0]):
        counter = Counter(nn_labels[:, i])
        y_pred[i] = counter.most_common(1)[0][0]

    return y_pred
```

Ακόμα φτιάξαμε την κλάση KNN_Custom για να μπορέσουμε να προσδιορίσουμε ποιες μετρικές να χρησιμοποιήσει για να υπολογίσει τις αποστάσεις , καθώς άμα χρησιμοποιήσουμε την έτοιμη KNeighborsClassifier δεν μπορούσαμε να προσδιορίσουμε ποια απόσταση να χρησιμοποιήσει(και πότε) οπότε χρησιμοποιούσε μόνο την default απόσταση που είναι η ευκλείδεια . Μετά ορίζουμε τις μετρικές του score οι οποίες είναι το accuracy και το f1_score και το αντικείμενο τύπου KNN_Custom το οποίο παίρνει ως παραμέτρους το k που είναι ο αριθμός των Nearest Neighbors και το distance_fn που είναι η συνάρτηση που θα χρησιμοποιήσει για τις αποστάσεις(π.χ. euclidean , hamming , custom_distance . Στην δικιά μας περίπτωση custom_distance) . Αν θέλουμε 10-fold cross validation χρησιμοποιούμε την έτοιμη συνάρτηση cross_val_score για κάθε μετρική ενώ αν

Θέλουμε 70 – 30 (train – test) split κάνουμε split τα δεδομένα μας , fit το KNN_Custom μοντέλο και υπολογίζουμε το y_pred . Παρακάτω φαίνεται η κλάση KNN_Custom.

```
class KNN_Custom:
    def __init__(self, k=5, distance_fn='euclidean'):
        self.k = k
        self.distance_fn = distance_fn

    def fit(self, X, y):
        self.X_train = X
        self.y_train = y
        return self

    def predict(self, X):
        return knn(self.k, self.distance_fn, self.X_train, self.y_train, X)

    def score(self, X, y, scoring='f1'):
        y_pred = self.predict(X)
        if scoring == 'f1':
            return f1_score(y, y_pred, average='weighted')
        elif scoring == 'accuracy':
            return accuracy_score(y, y_pred)
        else:
            raise ValueError("Invalid scoring metric")

    def get_params(self, deep=True):
        return {'k': self.k, 'distance_fn': self.distance_fn}
```

Τέλος υλοποιήσαμε την συνάρτηση NearestNeighbors η οποία χρησιμοποιεί ότι φτιάχαμε προηγουμένως και μας επιστρέφει το accuracy και το f1-score. Παρακάτω φαίνεται η συνάρτηση.

```
def NearestNeighbors(k,distance_fun,X_train=None,y_train=None,X_test=None,y_test=None,X=None,Y=None,type="test_set"):
    knn_custom = KNN_Custom(k=k,distance_fn=distance_fun)

    if type == 'test set':
        knn_custom.fit(X_train,y_train)
        y_pred = knn_custom.predict(X_test)
        f1 = f1_score(y_test, y_pred, average='weighted')
        accuracy = accuracy_score(y_test, y_pred)
    elif type == 'cross_val':
        scoring = {'accuracy': 'accuracy',
                   'f1 weighted': 'f1_weighted'}
        cv_results = cross_validate(estimator=knn_custom, X=X, y=Y, cv=10, n_jobs=-1, scoring=scoring)
        f1 = cross_val_score(knn_custom, X, Y, cv=10, scoring='f1 macro')
        accuracy = cross_val_score(knn_custom, X, Y, cv=10, scoring='accuracy')
        accuracy = accuracy.mean()

    return accuracy,f1
```

Εφαρμογή K – Nearest Neighbors στο αρχείο airlines_delay.csv

70-30 TRAIN TEST SPLIT			10-FOLD CROSS VALIDATION	
K	Accuracy	F1	Accuracy	F1
1	0.529307282415630	0.528489875017912	0.54483519918723	0.538376994383148
	5	9		3
3	0.564239194789816	0.562813776983123	0.547192706283615	0.538173691224678
	4	7		6
5	0.557134399052693	0.553593424633329	0.558769792425498	0.548276852711154
	9	5		8

1	0.569567791592658	0.565497213448483	0.566921773691599	0.553107137014814
0	4	6	6	8

Εφαρμογή K – Nearest Neighbors στο αρχείο train.csv

70-30 TRAIN TEST SPLIT			10-FOLD CROSS VALIDATION	
K	Accuracy	F1	Accuracy	F1
1	0.9066666666666666 6	0.906412515944712 2	0.9075	0.907246665013754 4
3	0.9133333333333333 3	0.913258928133431 5	0.9215	0.921887003564044 6
5	0.9183333333333333 3	0.918261100670334 3	0.9195	0.919598828730458 7
10	0.9283333333333333 3	0.928544091571694	0.9350000000000000 2	0.934752790358333 9

Μέθοδος Naïve Bayes

Για τη μέθοδο Naïve Bayes θέλουμε να χρησιμοποιήσουμε κανονική (Gaussian) κατανομή για κάθε ένα από τα χαρακτηριστικά συνεχούς τιμής και πολυωνυμική (Multinomial) κατανομή για τα διακριτά χαρακτηριστικά . Έτσι φτιάξαμε δύο μεταβλητές που κρατάν τους αριθμούς από τις στήλες του αρχείου που διαβάζουμε , όπου η μία έχει τους αριθμούς των στηλών με χαρακτηριστικά συνεχούς τιμής και η άλλη τους αριθμούς των στηλών με διακριτά χαρακτηριστικά . Ακολούθως φτιάχνουμε 2 pipelines , ένα για κάθε είδος χαρακτηριστικών , και χρησιμοποιούμε ένα preprocessor όπου συνδυάζει τα δυο pipelines . Μετατρέπουμε τα δεδομένα στην μορφή που θέλουμε με την βοήθεια του preprocessor , τα χωρίζουμε σε continuous_data(συνεχή χαρακτηριστικά) και discrete_data(διακριτά χαρακτηριστικά) και εκτελούμε Gaussian κατανομή στα continuous_data και Multinomial κατανομή στα discrete_data . Τέλος αν θέλουμε 10-fold cross validation παίρνουμε τα scores χρησιμοποιώντας την συνάρτηση cross_validate δύο φορές (μία για κάθε score ,accuracy και f1) , ενώ αν θέλουμε 70 – 30 (train – test) split ,κάνουμε split τα δεδομένα και χρησιμοποιούμε τις συναρτήσεις accuracy_score και f1_score για να υπολογίσουμε τα scores. Παρακάτω φαίνεται η συνάρτηση NaiveBayes η οποία υλοποιεί ότι αναφέραμε προηγουμένως.

```

def NaiveBayes(X_train=None,y_train=None,X_test=None,y_test=None,X=None,Y=None,file=None,type="test_set"):
    if file == 'train.csv':
        continuous_indices = [2,7]
        discrete_indices = [i for i in range(X.shape[1]) if i not in continuous_indices]

        continuous_pipeline = Pipeline([
            ('imputer', SimpleImputer(strategy='median')),
            ('scaler', StandardScaler())
        ])

        discrete_pipeline = Pipeline(steps=[
            ('imputer', SimpleImputer(strategy='most_frequent')),
            ('encoder', OrdinalEncoder())
        ])

        preprocessor = ColumnTransformer([
            ('continuous', continuous_pipeline, continuous_indices),
            ('discrete', discrete_pipeline, discrete_indices),
        ])

        preprocessor.fit(X)

        X_transformed = preprocessor.transform(X)

        continuous_data = X_transformed[:, :len(continuous_indices)]
        discrete_data = X_transformed[:, len(continuous_indices):]

        gaussian_nb = GaussianNB()

        multinomial_nb = MultinomialNB()

        gaussian_nb.fit(continuous_data, Y)
        multinomial_nb.fit(discrete_data, Y)
        if type == "test_set":
            accuracyG = accuracy_score(Y, gaussian_nb.predict(continuous_data))
            f1G = f1_score(Y, gaussian_nb.predict(continuous_data), average='weighted')

            accuracyM = accuracy_score(Y, multinomial_nb.predict(discrete_data))
            f1M = f1_score(Y, multinomial_nb.predict(discrete_data), average='weighted')

        elif type == "cross_val":
            scoring = {'accuracy': 'accuracy', 'f1_weighted': 'f1_weighted'}
            gaussian_nb_scores = cross_validate(gaussian_nb, continuous_data, Y, cv=10, scoring=scoring)
            accuracyG = np.mean(gaussian_nb_scores['test_accuracy'])
            f1G = np.mean(gaussian_nb_scores['test_f1_weighted'])

            multinomial_nb_scores = cross_validate(multinomial_nb, discrete_data, Y, cv=10, scoring=scoring)
            accuracyM = np.mean(multinomial_nb_scores['test_accuracy'])
            f1M = np.mean(multinomial_nb_scores['test_f1_weighted'])

    return accuracyG,accuracyM,f1G,f1M

```

```

elif file == "airlines_delay.csv":
    continuous_indices = []
    discrete_indices = [i for i in range(X.shape[1]) if i not in continuous_indices]

    continuous_pipeline = Pipeline([
        ('imputer', SimpleImputer(strategy='median')),
        ('scaler', StandardScaler())
    ])

    discrete_pipeline = Pipeline(steps=[
        ('imputer', SimpleImputer(strategy='most_frequent')),
        ('encoder', OrdinalEncoder())
    ])

    preprocessor = ColumnTransformer([
        ('continuous', continuous_pipeline, continuous_indices),
        ('discrete', discrete_pipeline, discrete_indices),
    ])

    preprocessor.fit(X)
    X_transformed = preprocessor.transform(X)

    continuous_data = X_transformed[:, :len(continuous_indices)]
    discrete_data = X_transformed[:, len(continuous_indices):]

    multinomial_nb = MultinomialNB()
    multinomial_nb.fit(discrete_data, Y)
    if type == "test_set":
        accuracy = accuracy_score(Y, multinomial_nb.predict(discrete_data))
        f1 = f1_score(Y, multinomial_nb.predict(discrete_data), average='weighted')
    elif type == "cross_val":
        scoring = {'accuracy': 'accuracy', 'f1_weighted': 'f1_weighted'}
        multinomial_nb_scores = cross_validate(multinomial_nb, discrete_data, Y, cv=10, scoring=scoring)
        accuracy = np.mean(multinomial_nb_scores['test_accuracy'])
        f1 = np.mean(multinomial_nb_scores['test_f1_weighted'])

    return accuracy,f1

```

Εφαρμογή Naïve Bayes στο αρχείο airlines_delay.csv

70-30 TRAIN TEST SPLIT		10-FOLD CROSS VALIDATION	
Accuracy	F1	Accuracy	F1

0.5499289267945985	0.5512481475761554	0.5493988735991099	0.550482182756145
--------------------	--------------------	--------------------	-------------------

Εφαρμογή Naïve Bayes στο αρχείο train.csv

70-30 TRAIN TEST SPLIT		10-FOLD CROSS VALIDATION	
Accuracy	F1	Accuracy	F1
0.4567	0.4419826850041815	0.4557	0.4427529950384546

Μέθοδος Neural Network

Στην παρακάτω εικόνα φαίνεται η συνάρτηση την οποία φτιάξαμε και υλοποιεί το νευρωνικό δίκτυο.

```
def NeuralNetwork(hidden_layers, n_neurons1, n_neurons2=None,X_train=None,y_train=None,X_test=None,y_test=None,X=None,Y=None, activation_hidden=None,type="test_set"):
    if hidden_layers == 1:
        hidden_layer_sizes = (n_neurons1)
    elif hidden_layers == 2:
        if n_neurons2 is not None:
            hidden_layer_sizes = (n_neurons1, n_neurons2)
        else:
            raise ValueError("Invalid number of neurons for the second hidden layer.")
    else:
        raise ValueError("Invalid number of hidden layers. Must be 1 or 2.")

    nn = MLPClassifier(hidden_layer_sizes=hidden_layer_sizes,
                       activation=activation_hidden,
                       solver='sgd',
                       learning_rate='constant',
                       learning_rate_init=0.01,
                       max_iter=10000,
                       tol=0.001)

    nn.out_activation_ = 'softmax'

    if(type == "test_set"):
        nn.fit(X_train, y_train)
        mean_accuracy = nn.score(X_test, y_test)
    elif type == "cross_val":
        accuracy = cross_val_score(nn, X, Y, cv=10)
        mean_accuracy = accuracy.mean()

    return mean_accuracy
```

Τα ορίσματα που χρησιμοποιούμε στην συνάρτηση είναι κατανοητά. Ανάλογα με το όρισμα type θα εφαρμόζεται είτε η μέθοδος 10-fold cross_validation είτε ο τυχαίος διαχωρισμός δεδομένων σε 70% για trainning και 30% για test. Όπως φαίνεται στην συνάρτηση χρησιμοποιούμε την κλάση MLPClassifier της βιβλιοθήκης scikit-learn Θέσαμε τα παρακάτω ορίσματα της κλάσης ως εξής:

- Hidden_layer_sizes: Μία tuple όπου περιέχει τον αριθμό νευρώνων κάθε κρυμμένου επιπέδου
- Activation: Συνάρτηση ενεργοποίησης την οποία δίνει ο χρήστης
- Solver: Θέσαμε sgd για να εφαρμοστεί η μέθοδος Gradient Descent
- Learning_rate_init: Θέσαμε στο 0.01 το οποίο είναι για το πόσο θα αλλάζουν τα βάρη κατά την εκπαίδευση
- Max_iter: Αριθμός εποχών που πρέπει να γίνουν για να τελειώσει η εκπαίδευση
- Tol: Θέσαμε στο 0.001 προκειμένου να σταματήσει η εκπαίδευση όταν το τετραγωνικό σφάλμα είναι μικρότερο από 0.001

Ο αριθμός νευρώνων στο επίπεδο εξόδου θέτεται αυτόματα από την βιβλιοθήκη όταν εκτελείται ο μέθοδο fit.

Εφαρμογή Neural Network στο αρχείο `airlines_delay.csv`

LOGISTIC ACTIVATION FUNCTION					
70-30 TRAIN TEST				10-FOLD CROSS VAL	
K1	K2	ACCURACY	F1-Score	MEAN ACCURACY	F1-Score
50	-	0.5530945833204586	0.3940187385535 758	0.5546032312250949	0.39613355832765634
100	-	0.5530822235268671	0.3942746401582 7424	0.5490061677038968	0.39628996666387567
200	-	0.5139140376355715	0.5531069431140 5	0.5546810979590988	0.3943572111427122
50	25	0.5529400859005654	0.3937598641075 988	0.5545717136254044	0.3957441829208744
100	50	0.5529400859005654	0.3937598641075 988	0.5545828374052657	0.39591876039749974
200	100	0.552995704971727	0.3939798133415 7545	0.5546236246324625	0.39588738339186685

TANH ACTIVATION FUNCTION					
70-30 TRAIN TEST				10-FOLD CROSS VAL	
K1	K2	ACCURACY	F1-Score	MEAN ACCURACY	F1-Score
50	-	0.553038964249297	0.3938790072142997	0.5545828377833555	0.39634361077354574
100	-	0.553026604455705 6	0.39443527622434343	0.54370388168521	0.45372012051971133
200	-	0.553032784352501 4	0.3947608040579752	0.5110515765000174	0.5030395731325239
50	25	0.552940085900565 4	0.3944371778494296	0.5546551424040119	0.39622206113602
100	50	0.553119302907641 4	0.3940019811490781	0.5546644121002618	0.39911356118101327
200	100	0.553125482804437 2	0.3944924641465437	0.5546644122033773	0.3963123410355083

LOGISTIC ACTIVATION FUNCTION					
70-30 TRAIN TEST				10-FOLD CROSS VAL	
K1	K2	ACCURACY	F1-Score	MEAN ACCURACY	F1-Score
50	-	0.5183333333333333	0.50658797200 58647	0.5065	0.488910340848 5215
100	-	0.5066666666666667	0.54925160662 32078	0.512	0.478652680089 99906
200	-	0.5383333333333333	0.44775638681 314767	0.5010000000000001	0.489601189104 6367
50	25	0.5183333333333333	0.47823243457 240494	0.5105000000000001	0.495477887855 9075
100	50	0.5316666666666666	0.49901194986 46232	0.5255000000000001	0.510022385158 4835
200	100	0.5466666666666666	0.48330691674 667214	0.52	0.506589045279 4952

Εφαρμογή Neural Network στο αρχείο train.csv

TANH ACTIVATION FUNCTION					
70-30 TRAIN TEST				10-FOLD CROSS VAL	
K1	K2	ACCURACY	F1-Score	MEAN ACCURACY	F1-Score
50	-	0.5116666666666666	0.28064734907507 294	0.446	0.408295913644573 74
10 0	-	0.4833333333333333	0.32757448076055 823	0.4955	0.434397703948531 67
20 0	-	0.475	0.46231674030545 566	0.4615	0.456590589293326 74
50	25	0.3183333333333333	0.38460239405624 597	0.39349999999999	0.387704984164880 1
10 0	50	0.495	0.33199122350113 536	0.439500000000000	0.415247106107944 5
20 0	10 0	0.5066666666666666	0.13240069728507 57	0.45149999999999	0.436754358002402

Μέθοδος SVM (Support Vector Machine)

Στην παρακάτω εικόνα φαίνεται η συνάρτηση SVM η οποία υλοποιεί της λειτουργίες που ζητήθηκαν.

```
24 airlines_delay = pd.read_csv("airlines_delay.csv")
25 sample_size = int(len(airlines_delay) * 0.8)
26 random_indices = random.sample(range(len(airlines_delay)), sample_size)
27 airlines_delay = airlines_delay.iloc[random_indices]
28 mobile_data = pd.read_csv("train.csv")
29
30 le = LabelEncoder()
31 airlines_delay['Airline'] = le.fit_transform(airlines_delay['Airline'])
32 airlines_delay['AirportFrom'] = le.fit_transform(airlines_delay['AirportFrom'])
33 airlines_delay['AirportTo'] = le.fit_transform(airlines_delay['AirportTo'])
34 airlines_delay = airlines_delay.drop('Flight', axis=1)
35
36 X_airline = airlines_delay.iloc[:, :-1].values
37 Y_airline = airlines_delay.iloc[:, -1].values
38 X_mobile = mobile_data.iloc[:, :-1].values
39 Y_mobile = mobile_data.iloc[:, -1].values
40
41 X_train_airline, X_test_airline, y_train_airline, y_test_airline = train_test_split(X_airline, Y_airline, test_size=0.3, random_state=42)
42 X_train_mobile, X_test_mobile, y_train_mobile, y_test_mobile = train_test_split(X_mobile, Y_mobile, test_size=0.3, random_state=42)
```

Όταν έχουμε 2 κατηγορίες (αρχείο airlines_delay) χρησιμοποιούμε την την στρατηγική 'ovo' και για πυρήνα rbf και για πυρήνα linear, διότι ο default είναι ο ovr. Η κλάση SVC που χρησιμοποιούμε έχει τα παρακάτω πεδία

- Kernel: ο πυρήνας που θα χρησιμοποιήσουμε (είτε linear είτε rbf)
 - Σε περίπτωση που χρησιμοποιήσουμε rbf έχουμε και την παράμετρο gamma
- Decision_function_shape: για να θέσουμε την αρχιτεκτονική one-versus-all
- C: επηρεάζει το margin

Επίσης προτού χρησιμοποιήσουμε την συνάρτηση που φτιάξαμε πρέπει να κάνουμε Scale τα δεδομένα επειδή υπολογίζονται αποστάσεις.

Στις γραμμές 305-314 ανάλογα ποια μέθοδο χρησιμοποιούμε για εκπαίδευση και test υπολογίζουμε επιπλέον και το accuracy και f1-score.

Εφαρμογή SVM στο αρχείο airlines_delay.csv

70-30 TRAIN TEST			10-FOLD CROSS VALIDATION	
C	ACCURACY	F1-SCORE	MEAN ACCURACY	F1-SCORE
0.01	0.5777276519482125	0.5584483725741 101	0.5804161091781331	0.561300878484 3746
0.1	0.5777338318450082	0.5584458475878 042	0.5804086932561485	0.561299316989 0511
1	0.5777276519482125	0.5584535721683 915	0.5804105472366448	0.561294632498 7503
4	0.5777214720514168	0.5584381229250 552	0.5804161091437614	0.561300124593 8432

10	0.5783580014213763	0.5567893700565 896	0.5802492511053409	0.560341563992 5611
----	--------------------	------------------------	--------------------	------------------------

RBF KERNEL					
70-30 TRAIN TEST				10-FOLD CROSS VAL	
C	gamma	ACCURACY	F1-Scope	MEAN ACCURACY	F1-Scope
0.01	0.01	0.2433333333333335	0.09524575513851 653	0.7084999999999999	0.7124918722279729
0.01	0.1	0.2433333333333335	0.09524575513851 653	0.6825	0.6924207876406397
0.01	1	0.09524575513851653	0.09524575513851 653	0.253	0.1165471716860094 1
0.01	5	0.2433333333333335	0.09524575513851 653	0.2514999999999999 5	0.1056997736188991 1
0.01	10	0.2433333333333335	0.09524575513851 653	0.296	0.2469459128953125
0.1	0.01	0.4933333333333335	0.48035402885908 707	0.7194999999999999	0.7240657586871548
0.1	0.1	0.2433333333333335	0.09537360178970 918	0.6825	0.6922505922634586
0.1	1	0.2433333333333335	0.09524575513851 653	0.253	0.1165471716860094 1
0.1	5	0.2433333333333335	0.09524575513851 653	0.2514999999999999 5	0.1056997736188991 1
0.1	10	0.2433333333333335	0.09524575513851 653	0.2784999999999999 7	0.2135639749748295 4
1	0.01	0.903333333333333	0.90406797339141 01	0.931	0.931229334525816
1	0.1	0.8066666666666666	0.80885333662452 11	0.8404999999999999	0.8419432768070635
1	1	0.245	0.10459418921183 626	0.312	0.2886486901359700 5
1	5	0.2433333333333335	0.09524575513851 653	0.2505	0.1699222457993575
1	10	0.2433333333333335	0.09524575513851 653	0.251	0.1038157582173359
4	0.01	0.925	0.92508078678742 32	0.9215	0.9212349628515061
4	0.1	0.818333333333334	0.81981843538641 77	0.8370000000000001	0.838081800437288

4	1	0.24666666666666667	0.10796795304760 423	0.3259999999999999 6	0.3089042775538753 4
4	5	0.24333333333333335	0.09524575513851 653	0.253	0.1430324585131191 3
4	10	0.24333333333333335	0.09524575513851 653	0.2505	0.1036873707139740 2
10	0.01	0.9083333333333333	0.90842695001207 44	0.9319999999999998	0.9319455992868556

Εφαρμογή SVM στο αρχείο train.csv

10	0.1	0.8183333333333334	0.81981843538641 77	0.8375	0.8385303466874743
10	1	0.246666666666666667	0.10796795304760 423	0.3259999999999999 6	0.3089042775538753 4
10	5	0.2433333333333335	0.09524575513851 653	0.253	0.1430324585131191 3
10	10	0.2433333333333335	0.09524575513851 653	0.2505	0.1036873707139740 2

LINEAR KERNEL					
70-30 TRAIN TEST			10-FOLD CROSS VAL		
C	ACCURACY	F1-Score	MEAN ACCURACY	F1-Score	
0.01	0.905	0.9061022053037766	0.9259999999999999	0.926372341 5292042	
0.1	0.9383333333333334	0.9384500710825442	0.9514999999999999	0.951423382 1942482	
1	0.9533333333333334	0.9534558763219786	0.9604999999999999	0.960431991 8792089	
4	0.9583333333333334	0.9583630559238058	0.9649999999999999	0.964987563 8272258	
10	0.96666666666666667	0.9667887792152499	0.9700000000000001	0.969987063 2199403	

Βέλτιστη Μέθοδος

Όπως φαίνεται και από τις μετρήσεις που κάναμε η βέλτιστη μέθοδο είναι SVM με πυρήνα linear, όπου μας έδινε accuracy > 0,90 και f1-score > 0,90 σε αντίθεση με τις άλλες μεθόδους.