~	Presentation of the approach We will try to solve this task in two steps: 1. Because there is a strong correlation between micro-calcification and breast cancer, we can say that having cancer implies having malignant macro-calcification but because we assume all micros have the same label, having breast cancer implies that all the macros are malignant. And the same for the equivalence here: no-breast cancer <=> all micros are benign. So supervised learning, saying that label "breast cancer" is equivalent, in our problem, to "malignant/benign macro-calcification". 2. Given the learning algorithm we acquired in the first step, we can classify the micros and give them a label (even if the label is wrong), we do not have anymore the hypothesis of all micros having the same label. So we will add as a feature of the calcification, the column "malignant/benign macro-calcification", and based on that, we will do a supervised learning, with as label, the "breast cancer" column. Presentation of the technique
3	Presentation of the technique We will try here to test 4 different models to achieve the classification: 1. the Logistic Regression algorithm on a dimensionally reduced dataset (we use the Principal Component Analysis method to do so). PCA can be a good solution 2. the Support Vector Machine algorithm 3. a Neural Network 4. a Neural Network on a dimensionally reduced dataset (i.e applying PCA on the dataset) secause the dataset is unbalanced (see bellow), we will use as an error metric the F1 score that is adapted to unbalanced dataset by considering the harmonic mean of precision and recall error meatrics. Nevertheless, we need to take into account the false negatives case to choose the model. It corresponds to the case where the patient has/had cancer, but the model preasurer". Thus, we also calculate the number of false negatives case.
r	The Algorithm Importing some librairies Import numpy as np Import matplotlib.pyplot as plt Import numous del_selection import train_test_split Import pandas as pd Import pandas as pd Import pandas as pd
	<pre>print("Number of feature:",len(raw_data[:,0])) wumber of feature: 152</pre>
5	Number of micro-calcification: 3562 Reparating our data from the label needed for supervised learning. We also discard the feature "patient" that won't be usefull in our study X, y = raw_data[:,1:(len(raw_data[0])-1)], raw_data[:,len(raw_data[0])-1:len(raw_data[0])] Step 1 irst, we split the data in equal proportion between label 1 and 0 to not biais our training set
·	X_label_0, y_label_0 = X[:2020], y[:2020] X_label_1, y_label_1 = X[2020:], y[2020:] print(len(y_label_0)) print(len(y_label_1)) 2020 1542 We see that the labels are unbalanced between 0s and 1s. We notice this unbalance as a potential biais for our dataset so we choose as an error metric the f1 score We split into training-validating set and test set our data for label 0s and 1s separately
:	<pre>X_train_0, X_test_0, y_train_0, y_test_0 = train_test_split(X_label_0, y_label_0, test_size=0.025) print(len(X_train_0)) print(len(X_test_0)) 1969 51 X_train_1, X_test_1, y_train_1, y_test_1 = train_test_split(X_label_1, y_label_1, test_size=0.025) print(len(X_train_1)) print(len(X_test_1))</pre>
, ; V	## merge the labels to have same proportion of Os and 1s in training set and test set X_test, y_test = np.concatenate([X_test_0, X_test_1]), np.concatenate([y_test_0, y_test_1])
V	de reshuffle everything to have a good dataset for cross-validation train_indexes = np.arange(len(X_train)) np.random.shuffle(train_indexes) X_train = X_train[train_indexes] y_train = y_train[train_indexes] test_indexes = np.arange(len(X_test)) np.random.shuffle(test_indexes)
٧	X_test = X_test[test_indexes] y_test = y_test[test_indexes] y_test = y_test[test_indexes] We import the tools to center and reduce the data, plus the PCA for the first and the fourth techniques from sklearn.decomposition import PCA from sklearn.preprocessing import StandardScaler We center and reduce our data
	scaler = StandardScaler() scaler.fit(X_train) X_standardized = scaler.transform(X_train) X_standardized array([[-1.04504505e-01, -2.33066898e-01, -3.54485248e-01,,
	[6.84225334e-01, 1.21647855e-01, -3.36565768e-04,,
	<pre>X_transformed = pca.transform(X_standardized) X_transformed array([[-12.35865328, -2.38248564, 4.48287003,, 1.12712305,</pre>
	[-7. 60899869, 3.37293436, -0.455115 ,, -0.72704821, 0.39280078, -0.29810586], [19.866319 , 0.35838319, 9.64466459,, 0.70139385, -1.58196588, -0.27437916]]) We see how much information each axe found by the PCA has to summarize the data print(pca.explained_variance_ratio_) [0.44787719 0.11487575 0.0559034 0.04192724 0.02899031 0.02650838 0.02479933 0.01262256 0.01137507 0.09890236 0.09885644 0.00811045 0.00782032 0.004622578
٧	<pre>0.00614093 0.00553254 0.00535789 0.00526537 0.00488269 0.00455471 0.0044519 0.00429625] We choose to do a cross-validation on our training set, regarding the small number of data that we have from sklearn.model_selection import KFold from sklearn.metrics import f1_score kf = KFold(16) irst model : the Logistic Regression</pre>
	<pre>from sklearn.linear_model import LogisticRegression models_scores_val_log_reg = [] for train_indexes, val_indexes in kf.split(X_transformed): model_temp = LogisticRegression() model_temp = fit(X_transformed[train_indexes], y_train[train_indexes].ravel()) y_pred_temp = model_temp.predict(X_transformed[val_indexes]) f1_score_temp = f1_score(y_train[val_indexes], y_pred_temp)</pre>
	<pre>y_folds = y_train[val_indexes] score_false_negativ = 0 nb_positivs = 0 for i in range(len(y_train[val_indexes])): if y_folds[i][0]: nb_positivs+=1 if not(y_pred_temp[i]): score_false_negativ+=1 models_scores_val_log_reg.append([model_temp,f1_score_temp,score_false_negativ/nb_positivs])</pre>
	he different models given by the cross validation models_scores_val_log_reg [[LogisticRegression(), 0.7251461988304092, 0.33333333333333],
V	[LogisticRegression(), 0.7272727272727272, 0.2967032967032967], [LogisticRegression(), 0.71864406779661, 0.37], [LogisticRegression(), 0.7398333333333333, 0.297029702970297], [LogisticRegression(), 0.6144578313253013, 0.4069767441860465], [LogisticRegression(), 0.632183908045977, 0.47115384615384615], [LogisticRegression(), 0.5988394557823129, 0.488720930232558], [LogisticRegression(), 0.6198830409356725, 0.47], [LogisticRegression(), 0.7282051282951282, 0.3238095238095238], [LogisticRegression(), 0.7282051282051282, 0.3238095238095238], [LogisticRegression(), 0.666666666666666, 0.3829787234042553]] [Voc choose the model with the best f1 score given the lowest number of false negativ
I I	model_log_rec_selected = models_scores_val_log_reg[-7][0] print("F1 score:",models_scores_val_log_reg[-7][1]) print("Number of false negatives:",models_scores_val_log_reg[-7][2]) =1 score: 0.739583333333333 number of false negatives: 0.297029702970297 necond model, SVM from sklearn.svm import SVC
	<pre>for train_indexes, val_indexes in kf.split(X_standardized): model_temp = SVC(gamma='auto') model_temp.frit(X_standardized[train_indexes], y_train[train_indexes].ravel()) y_pred_temp = model_temp.predict(X_standardized[val_indexes]) f1_score_temp = f1_score(y_train[val_indexes], y_pred_temp) y_folds = y_train[val_indexes] score_false_negativ = 0 nb_positivs=0 for i in range(len(y_train[val_indexes])):</pre>
	<pre>for i in range(len(y_train[val_indexes])): if y_folds[i][0]: nb_positivs+=1 if not(y_pred_temp[i]): score_false_negativ+=1 models_scores_val_SVC.append([model_temp, fi_score_temp, score_false_negativ/nb_positivs]) models_scores_val_SVC. [[SVC(gamma='auto'), 0.7560975609756097, 0.3333333333333], [SVC(gamma='auto'), 0.6857142857, 0.4],</pre>
	[SVC(gamma='auto'), 0.793478268695552, 0.27], [SVC(gamma='auto'), 0.691793233382766, 0.4025974025974026], [SVC(gamma='auto'), 0.6583859931677019, 0.39772727272777, [SVC(gamma='auto'), 0.6583859931677019, 0.397727272727777, [SVC(gamma='auto'), 0.658385931677019, 0.397272727277777, [SVC(gamma='auto'), 0.794191616766467, 0.27472527475], [SVC(gamma='auto'), 0.794191616766467, 0.27472527475], [SVC(gamma='auto'), 0.76757657656766, 0.29702970297], [SVC(gamma='auto'), 0.6712328767123289, 0.43023255813953487], [SVC(gamma='auto'), 0.6712328767123289, 0.43023255813953487], [SVC(gamma='auto'), 0.6619718309859154, 0.4534883720930232], [SVC(gamma='auto'), 0.66987573964497, 0.45], [SVC(gamma='auto'), 0.650887573964497, 0.45],
I	[SVC(gamma='auto'), 0.6927374301675978, 0.4055238095238095], [SVC(gamma='auto'), 0.6625, 0.43617021276595747]] model_sym_selected = models_scores_val_sVC[2][0] print("F1 score:", models_scores_val_sVC[2][1]) print("Number of false negatives:", models_scores_val_sVC[2][2]) F1 score: 0.7934782608695652 Number of false negatives: 0.27 hird model, Neural Network (without PCA)
	<pre>from sklearn.neural_network import MLPClassifier models_scores_val_NN = [] for train_indexes, val_indexes in kf.split(X_standardized): model_temp = MLPClassifier(hidden_layer_sizes=(6,10,6,4), activation="relu",max_iter=1000) model_temp.fit(X_standardized[train_indexes], y_train[train_indexes].ravel()) y_pred_temp = model_temp.predict(X_standardized[val_indexes]) f1_score_temp = f1_score(y_train[val_indexes],y_pred_temp)</pre>
	<pre>y_folds = y_train[val_indexes] score_false_negativ = 0 nb_positivs = 0 for i in range(len(y_train[val_indexes])): if y_folds[i][0]: nb_positivs+=1 if not(y_pred_temp[i]): score_false_negativ+=1 models_scores_val_NN.append([model_temp,f1_score_temp,score_false_negativ/nb_positivs])</pre>
	<pre>models_scores_val_NN [[MLPClassifier(hidden_layer_sizes=(6, 10, 6, 4), max_iter=1000),</pre>
	0.6853146853146853, 0.363636363636565], [MtPClassifier(hidden_layer_sizes=(6, 10, 6, 4), max_iter=1000), 0.625, 0.375], [MtPClassifier(hidden_layer_sizes=(6, 10, 6, 4), max_iter=1000), 0.6405228758169934, 0.3950617283950617), [MtPClassifier(hidden_layer_sizes=(6, 10, 6, 4), max_iter=1000), 0.7011494252873565, 0.3711340206185567], [MtPClassifier(hidden_layer_sizes=(6, 10, 6, 4), max_iter=1000), MtPClassifier(hidden_layer_sizes=(6, 10, 6, 4), max_iter=1000), MtPClassifier(hidden_layer_sizes=(6, 10, 6, 4), max_iter=1000), MtPClassifier(hidden_layer_sizes=(6, 10, 6, 4), max_iter=1000),
	0.7630057803468209, 0.2747252747257,475], [MLPCLassifier(hidden_layer_sizes=(6, 10, 6, 4), max_iter=1000), 0.6526315789473685, 0.38], [MLPCLassifier(hidden_layer_sizes=(6, 10, 6, 4), max_iter=1000), 0.749999999999, 0.749999999999, 0.8871287128712871], [MLPCLassifier(hidden_layer_sizes=(6, 10, 6, 4), max_iter=1000), 0.6225165562913907, 0.6225165562913907, 0.45348837209302323],
	[MLPClassifier(hidden_layer_sizes=(6, 10, 6, 4), max_iter=1000), 0.641304347826087, 0.4326923076923077], [MLPClassifier(hidden_layer_sizes=(6, 10, 6, 4), max_iter=1000), 0.6543209876543209, 0.38372093023255816], [MLPClassifier(hidden_layer_sizes=(6, 10, 6, 4), max_iter=1000), 0.7040816326530612, 0.31], [MLPClassifier(hidden_layer_sizes=(6, 10, 6, 4), max_iter=1000), 0.7027027027027027027027027027027027027027
V	[MLPClassifier (hidden_layer_sizes=(6, 10, 6, 4), max_iter=1000), 0.6844919786096257, 0.3191489361702128]] We could select the one doing the best in terms of f1 score but here we prefer selecting the one doing the less bad on numbers of false negatives model_NN_selected = models_scores_val_NN[0][0] print("F1 score:", models_scores_val_NN[0][1]) print("Number of false negatives:", models_scores_val_NN[0][2]) =1 score: 0.7474747474747475
F	<pre>Number of false negatives: 0.20430107526881722 ourth model, NN with PCA from sklearn.neural_network import MLPClassifier models_scores_val_NN_PCA = [] for train_indexes, val_indexes in kf.split(X_transformed): model_temp = MLPClassifier(hidden_layer_sizes=(6,10,6,4), activation="relu", max_iter=1000) model_temp.fit(X_transformed(train_indexes), y_train[train_indexes].ravel())</pre>
	<pre>model_temp.fit(X_transformed[vali_indexes], y_train[train_indexes], favel()) y_pred_temp = model_temp.predict(X_transformed[val_indexes]) f1_score_temp = f1_score(y_train[val_indexes], y_pred_temp) y_folds = y_train[val_indexes] score_false_negativ = 0 nb_positivs = 0 for i in range(len(y_train[val_indexes])): if y_folds[i][0]: nb_positivs+=1 if not(y_pred_temp[i]): score_false_negativ+=1 models_scores_val_NN_PCA.append([model_temp, f1_score_temp, score_false_negativ/nb_positivs])</pre>
	models_scores_val_NN_PCA [[MLPClassifier(hidden_layer_sizes=(6, 10, 6, 4), max_iter=1000), 0.7500000000000001, 0.2903225806451613], [MLPClassifier(hidden_layer_sizes=(6, 10, 6, 4), max_iter=1000), 0.6547619047619048, 0.45], [MLPClassifier(hidden_layer_sizes=(6, 10, 6, 4), max_iter=1000), 0.617619047619048, 0.45], [MLPClassifier(hidden_layer_sizes=(6, 10, 6, 4), max_iter=1000),
	0.768333333333333, 0.26], [MLPClassifier(hidden_layer_sizes=(6, 10, 6, 4), max_iter=1000), 0.6808510638297872, 0.37662337662376623, [MLPClassifier(hidden_layer_sizes=(6, 10, 6, 4), max_iter=1000), 0.6467065868263473, 0.386363636363635], [MLPClassifier(hidden_layer_sizes=(6, 10, 6, 4), max_iter=1000), 0.61038961033896104, 0.41975308641975306], [MLPClassifier(hidden_layer_sizes=(6, 10, 6, 4), max_iter=1000), 0.41975308641975306], [MLPClassifier(hidden_layer_sizes=(6, 10, 6, 4), max_iter=1000),
	0.6867469879518072, 0.41237113402061853], [MLPCLassifier(hidden_layer_sizes=(6, 10, 6, 4), max_iter=1000), 0.7701149425287357, 0.26373626373626374], [MLPClassifier(hidden_layer_sizes=(6, 10, 6, 4), max_iter=1000), 0.6885245901639344, 0.37], [MLPClassifier(hidden_layer_sizes=(6, 10, 6, 4), max_iter=1000), 0.6885245901639344, 0.37], [MLPClassifier(hidden_layer_sizes=(6, 10, 6, 4), max_iter=1000), 0.76237623762376237623, 0.23762376237623763,
	[MLPClassifier(hidden_layer_sizes=(6, 10, 6, 4), max_iter=1000), 0.6335403726708074, 0.40697674418604065], [MLPClassifier(hidden_layer_sizes=(6, 10, 6, 4), max_iter=1000), 0.6432748538011696, 0.47115384615334615], [MLPClassifier(hidden_layer_sizes=(6, 10, 6, 4), max_iter=1000), 0.649006625165563, 0.490032355813953487], [MLPClassifier(hidden_layer_sizes=(6, 10, 6, 4), max_iter=1000), 0.64000000000000001, 0.6400000000000001, 0.44000000000000001, 0.4400000000000001, 0.44000000000000001, 0.44000000000000001, 0.44000000000000001, 0.44000000000000001, 0.44000000000000001, 0.44000000000000001, 0.44000000000000001, 0.44000000000000001, 0.44000000000000001, 0.440000000000000001, 0.44000000000000001, 0.440000000000000001, 0.440000000000000001, 0.440000000000000001, 0.4400000000000000001, 0.440000000000000001, 0.4400000000000000001, 0.440000000000000000000000000001, 0.440000000000000000000000000000000000
	[MLPClassifier(hidden_layer_sizes=(6, 10, 6, 4), max_iter=1000),
A	Fils core: 0.77083333333333333333333333333333333333
	<pre>y_pred_X_test_NN = model_NN_selected.predict(X_test_standardized) f1_score_log_rec = f1_score(y_test.ravel(),y_pred_X_test_log_rec) score_false_negativ_log_rec = 0 for i in range(len(y_test)): if y_test[i] and not(y_pred_X_test_log_rec[i]): score_false_negativ_log_rec+=1 f1_score_SVM = f1_score(y_test.ravel(),y_pred_X_test_SVM) score_false_negativ_SVM = 0 for i in range(len(y_test)): if y_test[i] and not(y_pred_X_test_SVM[i]):</pre>
	<pre>score_false_negativ_SVM+=1 f1_score(y_test.ravel(),y_pred_X_test_NN) score_false_negativ_NN = 0 for i in range(len(y_test)); if y_test[i] and not(y_pred_X_test_NN[i]); score_false_negativ_NN+=1 f1_score(y_test.ravel(),y_pred_X_test_NN_PCA) score_false_negativ_NP_PA = 0 for i in range(len(y_test)); if y_test[i] and not(y_pred_X_test_NN_PCA]</pre>
	score_false_negativ_NN_PCA+=1 print(f"f1_score for Logistic Regression with PCA transformation on X_test: {f1_score_log_rec} and the number of false negativ: {score_false_negativ_log_rec}") print(f"f1_score for SVM on X_test: {f1_score_SVM} and the number of false negativ: {score_false_negativ_SVM}") print(f"f1_score for NN on X_test: {f1_score_NN} and the number of false negativ: {score_false_negativ_NN}") print(f"f1_score for NN with PCA transformation on X_test: {f1_score_NN_PCA} and the number of false negativ: {score_false_negativ_NN_PCA}") f1_score for Logistic Regression with PCA transformation on X_test: 0.6176470588235294 and the number of false negativ: 18 f1_score for SVM on X_test: 0.6451612903225806 and the number of false negativ: 19
N S V	filscore for NN on X_test: 0.6233766233766234 and the number of false negativ: 15 filscore for NN with PCA transformation on X_test: 0.6285714285714286 and the number of false negativ: 17 IN with PCA seems here to be a good model to deal with our database, with the second highest f1 score and the second smallest number of false negatives, which is a good tradeoff. I choose to focus on this model. Step 2 We go through the same process than above, using the 4 same models but trained on the new dataset made of the old X on which we added the column "benignant/malignant calcification". col_malignant = model_NN_PCA_selected_predict(pca_transform(scaler_transform(X))) X_new = np.hstack([X, np.atleast_2d(col_malignant).T]) X_new = np.hstack([X, np.atleast_2d(col_malignant).T])
	<pre>%_new_label_1 = X_new[2020:] %_new_label_1 = X_new[2020:] %_train_new_0, X_test_new_0, y_train_new_0, y_test_new_0 = train_test_split(X_new_label_0, y_label_0, test_size=0.025) %_train_new_1, X_test_new_1, y_train_new_1, y_test_new_1 = train_test_split(X_new_label_1, y_label_1, test_size=0.025) %_test_new, y_train_new = np.concatenate([X_test_new_0, X_test_new_1]), np.concatenate([y_test_new_0, y_test_new_1]) %_train_new = np.arange(len(X_train_new_0, X_train_new_1]), np.concatenate([y_train_new_0, y_train_new_1]) %_train_new = np.arange(len(X_train_new_0, X_train_new_1]), np.concatenate([y_train_new_0, y_train_new_1]) %_train_new = x_train_new[train_indexes_new] %_train_new = x_train_new[train_indexes_new] %_train_new = y_train_new[train_indexes_new] %_train_new = y_train_new[train_indexes_new] </pre>
	test_indexes_new = np.arange(len(X_test_new)) np.random.shuffle(test_indexes_new) X_test_new = X_test_new[test_indexes_new] y_test_new = y_test_new[test_indexes_new] scaler_2 = StandardScaler() scaler_2.fit(X_train_new) X_new_standardized = scaler_2.transform(X_train_new) pca_2 = PCA(0.9)
	<pre>pca_2.fit(X_new_standardized) X_new_transformed = pca_2.transform(X_new_standardized) models_scores_val_log_reg_new = [] for train_indexes, val_indexes in kf.split(X_new_transformed): model_temp = LogisticRegression() model_temp.fit(X_new_transformed[train_indexes], y_train_new[train_indexes].ravel()) y_pred_temp = model_temp.predict(X_new_transformed[val_indexes])</pre>
	<pre>f1_score_temp = f1_score(y_train_new[val_indexes],y_pred_temp) y_folds = y_train_new[val_indexes] score_false_negativ = 0 nb_positivs = 0 for i in range(len(y_train[val_indexes])): if y_folds[i][0]: nb_positivs+1 if not(y_pred_temp[i]): score_false_negativ+1 models_scores_val_log_reg_new.append([model_temp, f1_score_temp, score_false_negativ/nb_positivs])</pre>
	models_scores_val_log_reg_new [[LogisticRegression(), 0.776470588235294, 0.2826086956521739], [LogisticRegression(), 0.6628571428571428, 0.414141414141414], [LogisticRegression(), 0.738095238095238, 0.30337078651685395], [LogisticRegression(), 0.6716417910447761, 0.4155844155844156], [LogisticRegression(), 0.75, 0.3030303030303034], [LogisticRegression(), 0.7675675675675675, 0.2828282828282828], [LogisticRegression(), 0.7362637362637363, 0.33], [LogisticRegression(), 0.7329192546583853, 0.2716049382716049],
	[LogisticRegression(), 0.7692307692307692, 0.3010752688172043], [LogisticRegression(), 0.756097560975609756098, 0.29545454545454547], [LogisticRegression(), 0.756097560975609756098, 0.29545454545454547], [LogisticRegression(), 0.7560976767676767, 0.2952962962962963], [LogisticRegression(), 0.76767676767676, 0.2962962962963], [LogisticRegression(), 0.77484662576687118, 0.3146067415730337], [LogisticRegression(), 0.7771428571428571, 0.313131313131313], [LogisticRegression(), 0.77712020725388601, 0.29245283018867924]] model_log_rec_selected_new = models_scores_val_log_reg_new[0][0] print("F1 score:", models_scores_val_log_reg_new[0][1])
I	<pre>print("Number of false negatives:",models_scores_val_log_reg_new[0][2]) =1 score: 0.776470588235294 vumber of false negatives: 0.2826086956521739 models_scores_val_SVC_new = [] for train_indexes, val_indexes in kf.split(X_new_standardized): model_temp = SVC(gamma='auto') model_temp = SVC(gamma='auto') model_temp.fit(X_new_standardized[train_indexes], v_train_new[train_indexes], v_train_indexes], v_train_new[train_indexes], v_train_new[train_indexes], v_train_indexes].</pre>
	<pre>y_pred_temp = model_temp.predict(X_new_standardized[val_indexes]) f1_score_temp = f1_score(y_train_new[val_indexes],y_pred_temp) y_folds = y_train_new[val_indexes] score_false_negativ = 0 nb_positivs=0 for i nrange(len(y_train[val_indexes])): if y_folds[i][0]: nb_positivs+=1 if not(y_pred_temp[i]): score_false_negativ+=1 models_scores_val_SVC_new.append([model_temp,f1_score_temp,score_false_negativ/nb_positivs])</pre>
	models_scores_val_SVC_new [[SVC(gamma='auto'), 0.788235294117647, 0.2717391304347826], [SVC(gamma='auto'), 0.6666666666667, 0.4141414141414], [SVC(gamma='auto'), 0.771084337349376, 0.2808988764044944], [SVC(gamma='auto'), 0.661647058823529, 0.4155844155844156], [SVC(gamma='auto'), 0.7734866629834253, 0.292929292929293], [SVC(gamma='auto'), 0.7734966629834253, 0.2929292929293], [SVC(gamma='auto'), 0.77617391304347825, 0.29], [SVC(gamma='auto'), 0.77617391304347825, 0.29], [SVC(gamma='auto'), 0.75, 0.25925925925924],
	SVC(gamma='auto'), 0.783132530120482, 0.3010752688172043], SVC(gamma='auto'), 0.7534246575342465, 0.32926829268292684], SVC(gamma='auto'), 0.7439024390243902, 0.3068181818181818], SVC(gamma='auto'), 0.743684210526315, 0.30392156862745996], SVC(gamma='auto'), 0.78172588324873, 0.28703703703703703], SVC(gamma='auto'), 0.8, 0.2808988764044944], SVC(gamma='auto'), 0.7909664519774011, 0.2929292929293], SVC(gamma='auto'), 0.7748691099476439, 0.3018867924528302]] model_svm_selected_new = models_scores_val_SVC_new[0][0] print("F1 score:", models_scores_val_SVC_new[0][1])
1	print("Number of false negatives:", models_scores_val_SVC_new[0][2]) =1 score: 0.788235294117647 Number of false negatives: 0.2717391304347826 models_scores_val_NN_new = [] for train_indexes, val_indexes in kf.split(X_new_standardized):
	<pre>y_pred_temp = model_temp.predict(X_new_standardized[val_indexes]) f1_score_temp = f1_score(y_train_new[val_indexes], y_pred_temp) y_folds = y_train_new[val_indexes] score_false_negativ = 0 nb_positivs=0 for i in range(len(y_train[val_indexes])): if y_folds[i][0]: nb_positivs+=1 if not(y_pred_temp[i]): score_false_negativ+=1 models_scores_val_NN_new.append([model_temp,f1_score_temp,score_false_negativ/nb_positivs])</pre>
	models_scores_val_NN_new [[MLPClassifier(hidden_layer_sizes=(6, 10, 6, 4), max_iter=1000),
	0.3146067415730337], [MLPClassifier(hidden_layer_sizes=(6, 10, 6, 4), max_iter=1000), 0.6715328467153284, 0.4025974025974026], [MLPClassifier(hidden_layer_sizes=(6, 10, 6, 4), max_iter=1000), 0.7046632124352332, 0.313131313131315], [MLPClassifier(hidden_layer_sizes=(6, 10, 6, 4), max_iter=1000), 0.700630475935829, 0.2727272727272727, [MLPClassifier(hidden_layer_sizes=(6, 10, 6, 4), max_iter=1000), 0.700530475935829, 0.2727272727272727, [MLPClassifier(hidden_layer_sizes=(6, 10, 6, 4), max_iter=1000), 0.7643979057591622,
	0.27], [MLPClassifier(hidden_layer_sizes=(6, 10, 6, 4), max_iter=1000), 0.7450980392156864, 0.2962962962963], [MLPClassifier(hidden_layer_sizes=(6, 10, 6, 4), max_iter=1000), 0.689655172413793, 0.3543387096774194], [MLPClassifier(hidden_layer_sizes=(6, 10, 6, 4), max_iter=1000), 0.7189542483660131, 0.32926829268293, [MLPClassifier(hidden_layer_sizes=(6, 10, 6, 4), max_iter=1000), 0.7189542483660131, 0.32926829268292684], [MLPClassifier(hidden_layer_sizes=(6, 10, 6, 4), max_iter=1000), 0.7389535031847134,
	0.738853631847134, 0.34090909090909], [MLPClassifier(hidden_layer_sizes=(6, 19, 6, 4), max_iter=1000), 0.6984126984126984, 0.535294117647058826], [MLPClassifier(hidden_layer_sizes=(6, 10, 6, 4), max_iter=1000), 0.7282051282051282, 0.342525925925926], [MLPClassifier(hidden_layer_sizes=(6, 10, 6, 4), max_iter=1000), 0.76845454545454, 0.30337078651685395], [MLPClassifier(hidden_layer_sizes=(6, 10, 6, 4), max_iter=1000), 0.7045454545454, 0.30337078651685395], [MLPClassifier(hidden_layer_sizes=(6, 10, 6, 4), max_iter=1000), 0.7513227514,
	0.7513227513227514, 0.28282828282828], [MLPClassifier(hidden_layer_sizes=(6, 10, 6, 4), max_iter=1000), 0.7192118226609984, 0.3113207547169811]] print("F1 score:", models_scores_val_NN_new[5][0] print("F1 score:", models_scores_val_NN_new[5][1]) print("Number of false negatives:", models_scores_val_NN_new[5][2]) F1 score: 0.770053475935829 Number of false negatives: 0.272727272727272727
	<pre>Number of false negatives: 0.27272727272727</pre> models_scores_val_NN_PCA_new = [] for train_indexes, val_indexes in kf.split(X_new_transformed): model_temp = MLPClassifier(hidden_layer_sizes=(6,10,6,4), activation="relu",max_iter=1000) model_temp, fit(X_new_transformed[train_indexes], y_train_new[train_indexes].ravel()) y_pred_temp = model_temp.predict(X_new_transformed[val_indexes]) fl_score_temp = fl_score(y_train_new[val_indexes]) y_folds = y_train_new[val_indexes], y_pred_temp) y_folds = y_train_new[val_indexes] score_false_negativ = 0
	<pre>score_false_negativ = 0 nb_positivs=0 for i in range(len(y_train[val_indexes])): if y_folds[i][0]: nb_positivs+=1 if not(y_pred_temp[i]): score_false_negativ+=1 models_scores_val_NN_PCA_new.append([model_temp, f1_score_temp, score_false_negativ/nb_positivs])</pre> models_scores_val_NN_PCA_new
	[[MLPClassifier(hidden_layer_sizes=(6, 10, 6, 4), max_iter=1000),
	0.4875324675324675], [MLPClassifier(hidden_layer_sizes=(6, 10, 6, 4), max_iter=1000), 0.7085714285714285, 0.373737373737376], [MLPClassifier(hidden_layer_sizes=(6, 10, 6, 4), max_iter=1000), 0.7789473684210526, 0.25525252525254], [MLPClassifier(hidden_layer_sizes=(6, 10, 6, 4), max_iter=1000), 0.7407407407407409, 0.31, [MLPClassifier(hidden_layer_sizes=(6, 10, 6, 4), max_iter=1000), 0.7407407407407409, 0.31, [MLPClassifier(hidden_layer_sizes=(6, 10, 6, 4), max_iter=1000), 0.7407407407407409, 0.31,
	[MLPClassifier(hidden_layer_sizes=(6, 10, 6, 4), max_iter=1000), 0.7044025157232704, 0.308641975308641961, [MLPClassifier(hidden_layer_sizes=(6, 10, 6, 4), max_iter=1000), 0.7399393939394, 0.34408602150537637], [MLPClassifier(hidden_layer_sizes=(6, 10, 6, 4), max_iter=1000), 0.837980891719746, 0.3414634146341637], [MLPClassifier(hidden_layer_sizes=(6, 10, 6, 4), max_iter=1000), 0.7388535031847134, 0.3409090909099], [MLPClassifier(hidden_layer_sizes=(6, 10, 6, 4), max_iter=1000), 0.4388535031847134, 0.4090909090999], [MLPClassifier(hidden_layer_sizes=(6, 10, 6, 4), max_iter=1000),
	[MLPClassifier(hidden_layer_sizes=(6, 10, 6, 4), max_iter=1000), 0.71875, 0.3235294117647059], [MLPClassifier(hidden_layer_sizes=(6, 10, 6, 4), max_iter=1000), 0.78, 0.77, 0.77777777777778], [MLPClassifier(hidden_layer_sizes=(6, 10, 6, 4), max_iter=1000), 0.70303030303030, 0.70303030303030, 0.34831460674157305], [MLPClassifier(hidden_layer_sizes=(6, 10, 6, 4), max_iter=1000), 0.7292817679558011, 0.7292817679558011, 0.33333333333333,
	<pre>X_test_standardized_new = scaler_2.transform(X_test_new) X_test_transformed_new = pca_2.transform(X_test_standardized_new) y_pred_X_test_log_rec_new = model_log_rec_selected_new.predict(X_test_transformed_new) y_pred_X_test_SVM_new = model_svm_selected_new.predict(X_test_standardized_new) y_pred_X_test_NN_PCA_new = model_NN_PCA_selected_new.predict(X_test_transformed_new) y_pred_X_test_NN_new = model_NN_pca_selected_new.predict(X_test_standardized_new) y_pred_X_test_NN_new = model_NN_selected_new.predict(X_test_standardized_new) f1_score_log_rec_new = f1_score(y_test_new.ravel(),y_pred_X_test_log_rec_new) score_false_negativ_log_rec_new = 0 for in range(len(y_test_new)):</pre>
	<pre>score_false_negativ_log_rec_new = 0 for i in range(len(y_test_new)): if y_test_new[i] and not(y_pred_X_test_log_rec_new[i]): score_false_negativ_log_rec_new+=1 f1_score_SVM_new = f1_score(y_test_new.ravel(),y_pred_X_test_SVM_new) score_false_negativ_SVM_new = 0 for i in range(len(y_test_new)): if y_test_new[i] and not(y_pred_X_test_SVM_new[i]): score_false_negativ_sVM_new+=1 f1_score(y_test_new.ravel(),y_pred_X_test_NN_new) score_false_negativ_sVM_new+=0</pre>
	<pre>score_false_negativ_NN_new = 0 for i in range(len(y_test_new)): if y_test_new[i] and not(y_pred_X_test_NN_new[i]): score_false_negativ_NN_new+=1 f1_score_NN_PCA_new = f1_score(y_test_new.ravel(),y_pred_X_test_NN_PCA_new) score_false_negativ_NN_PCA_new = 0 for i in range(len(y_test_new)):</pre>
	<pre>if y_test_new[i] and not(y_pred_X_test_NN_PCA_new[i]): score_false_negativ_NN_PCA_new+=1</pre>