Breast Cancer Classification using Machine Learning This Data is derived from a test Called Fine needle aspiration. Fine needle aspiration is a type of biopsy procedure. In fine needle aspiration, a thin needle is inserted into an area of abnormal-appearing tissue or body fluid. As with other types of biopsies, the sample collected during fine needle aspiration can help make a diagnosis or rule out conditions such as cancer. Importing the Dependencies import numpy as np In [1]: import pandas as pd import sklearn.datasets from sklearn.model_selection import train_test_split from sklearn.linear_model import LogisticRegression from sklearn.metrics import accuracy_score Data Collection & Processing In [2]: # loading the data from sklearn breast_cancer_dataset = sklearn.datasets.load_breast_cancer() In [3]: print(breast_cancer_dataset) {'data': array([[1.799e+01, 1.038e+01, 1.228e+02, ..., 2.654e-01, 4.601e-01, [2.057e+01, 1.777e+01, 1.329e+02, ..., 1.860e-01, 2.750e-01, 8.902e-02], [1.969e+01, 2.125e+01, 1.300e+02, ..., 2.430e-01, 3.613e-01, 8.758e-02], [1.660e+01, 2.808e+01, 1.083e+02, ..., 1.418e-01, 2.218e-01, 7.820e-02], [2.060e+01, 2.933e+01, 1.401e+02, ..., 2.650e-01, 4.087e-01, 1.240e-01], [7.760e+00, 2.454e+01, 4.792e+01, ..., 0.000e+00, 2.871e-01, 0, 0, 1, 0, 1, 1, 1, 1, 1, 0, 0, 1, 0, 0, 1, 1, 1, 1, 0, 1, 0, 0, 1, 1, 1, 1, 0, 1, 0, 0, 1, 0, 1, 0, 0, 1, 1, 1, 0, 0, 1, 0, 0, 0, 1, 1, 1, 0, 1, 1, 0, 0, 1, 1, 1, 0, 0, 1, 1, 1, 1, 0, 1, 1, 0, 1, 1, 1, 1, 1, 1, 1, 1, 0, 0, 0, 1, 0, 0, 1, 1, 1, 0, 0, 1, 0, 1, 0, 1, 1, 0, 1, 1, 1, 1, 0, 0, 1, 0, 1, 1, 0, 0, 1, 1, 0, 0, 1, 1, 1, 1, 0, 1, 1, 0, 0, 0, 1, 0, 1, 0, 1, 1, 1, 0, 1, 1, 0, 0, 1, 0, 0, 0, 0, 1, 0, 0, 0, 1, 0, 1, 0, 1, 1, 0, 1, 0, 0, 0, 0, 1, 1, 0, 0, 1, 1, 1, 0, 1, 1, 1, 1, 1, 0, 0, 1, 1, 0, 1, 1, 0, 0, 1, 0, 1, 1, 1, 1, 0, 1, 1, 1, 1, 0, 1, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 1, 1, 1, 1, 1, 0, 1, 0, 1, 1, 0, 1, 1, 0, 1, 0, 0, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 0, 1, 0, 1, 0, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 0, 1, 1, 0, 1, 0, 1, 1, 1, 1, 1, 0, 0, 0, 1, 1, 1, 1, 0, 1, 0, 1, 0, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 0, 0, 0, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 0, 0, 1, 0, 0, 0, 1, 0, 0, 1, 1, 1, 1, 1, 0, 1, 1, 1, 1, 0, 1, 1, 1, 0, 1, 1, 0, 0, 1, 1, 1, 1, 1, 1, 0, 1, 1, 1, 1, 1, 1, 0, 1, 1, 1, 1, 0, 1, 1, 0, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 0, 1, 0, 0, 1, 0, 1, 1, 1, 1, 1, 0, 1, 1, 0, 1, 0, 1, 1, 0, 1, 0, 1, 1, 1, 1, 1, 1, 1, 1, 0, 0, 1, 1, 1, 1, 1, 0, 1, 1, 1, 1, 1, 1, 1, 1, 1, 0, 1, 1, 1, 1, 1, 1, 0, 1, 0, 1, 1, 0, 1, 1, 1, 1, 1, 0, 0, 1, 0, 1, 0, 1, 1, 1, 1, 1, 0, 1, 1, 0, 1, 0, 1, 0, 0, 1, 1, 1, 0, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 0, 0, 0, 0, 0, 0, 1]), 'frame': None, 'target_names': array(['malignant', 'benign'], dtype='<U9'), 'DESCR': '.._br east_cancer_dataset:\n\nBreast cancer wisconsin (diagnostic) dataset\n---------------------\n\n**Data Set Characteristics:**\n\n :Nu mber of Instances: 569\n\n :Number of Attributes: 30 numeric, predictive attributes and the class\n\n :Attribute Information:\n - radius (mean of distances from center to points on the perimeter)\n texture (standard deviation of gray-scale values)\n - perimeter\n - area∖n - compactness (perimeter^2 / area - 1.0)\n - concavity (severity of concave portions of the con smoothness (local variation in radius lengths)\n - concave points (number of concave portions of the contour)\n tour)\n - symmetry\n - fractal dimension ("coastline approximation" - 1)\n The mean, standard error, and "worst" or largest (mean of the three\n worst/largest values) of these features were computed for each image,\n resulting in 30 features. For instance, field 0 is Mean Radius, field\n 10 is Radius SE, field 20 is Worst Radius.\n\n - class:\n - WDBC-Malignant\n 6.981 28.11\n texture (mean): Min Max\n radius (mean): 39.28\n perimeter (mean): 9.71 43.79 188.5\n area (mean): 143.5 2501.0\n smoothness (mean): 0.053 0.163\n compactness (mean): 0.019 0.345\n concavity (mean): 0.0 0.427\n concave points (mean): 0.201\n 0.106 0.304\n fractal dimension (mean): 0.0 symmetry (mean): 0.05 0.097\n radius (standard error): 0.112 2.873\n texture (standard error): 0.36 4.885\n perimeter (standard error): 0.757 21.98\n area (standard error): compactness (standard error): 6.802 542.2\n smoothness (standard error): 0.002 0.031\n 0.002 0.135\n concavity (standard error): concave points (standard error): 0.0 0.396\n 0.0 0.053\n symmetry (standard error): 0.008 0.079\n fractal dimension (standard radius (worst): 7.93 36.04\n texture (worst): 12.02 49.54\n error): 0.001 0.03\n perimeter (worst): compactness (worst): 185.2 4254.0\n 50.41 251.2\n 0.071 0.223\n area (worst): smoothness (worst): 0.027 1.058\n concavity (worst): 0.0 1.252\n concave points (worst): 0.291\n 0.0 symmetry (worst): 0.156 0.664\n fractal dimension (worst): 0.055 0.208\n :Missing Attribute Value s: None\n\n :Class Distribution: 212 - Malignant, 357 - Benign\n\n :Creator: Dr. William H. Wolberg, W. Nick Street, Olvi L. Mangasarian\n\n :Date: November, 1995\n\nThis is a copy of UCI ML Breast Cancer Wisconsin (Diagnostic) datasets.\nhttps://goo.gl/U2Uwz2\n\nFeatures are com puted from a digitized image of a fine needle\naspirate (FNA) of a breast mass. They describe\ncharacteristics of the cell nuclei present in the image.\n\nSe parating plane described above was obtained using\nMultisurface Method-Tree (MSM-T) [K. P. Bennett, "Decision Tree\nConstruction Via Linear Programming." Proc eedings of the 4th\nMidwest Artificial Intelligence and Cognitive Science Society,\npp. 97-101, 1992], a classification method which uses linear\nprogramming to construct a decision tree. Relevant features\nwere selected using an exhaustive search in the space of 1-4\nfeatures and 1-3 separating planes.\n\nThe act ual linear program used to obtain the separating plane\nin the 3-dimensional space is that described in:\n[K. P. Bennett and O. L. Mangasarian: "Robust Linear \nProgramming Discrimination of Two Linearly Inseparable Sets",\nOptimization Methods and Software 1, 1992, 23-34].\n\nThis database is also available through the UW CS ftp server:\n\nftp ftp.cs.wisc.edu\ncd math-prog/cpo-dataset/machine-learn/WDBC/\n\n.. topic:: References\n\n - W.N. Street, W.H. Wolberg and O.L. Mangasarian. Nuclear feature extraction \n for breast tumor diagnosis. IS&T/SPIE 1993 International Symposium on \n Electronic Imaging: Science and Te chnology, volume 1905, pages 861-870,\n San Jose, CA, 1993.\n - O.L. Mangasarian, W.N. Street and W.H. Wolberg. Breast cancer diagnosis and \n osis via linear programming. Operations Research, 43(4), pages 570-577, \n July-August 1995.\n - W.H. Wolberg, W.N. Street, and O.L. Mangasarian. Machin e learning techniques\n to diagnose breast cancer from fine-needle aspirates. Cancer Letters 77 (1994) \n 163-171.', 'feature_names': array(['mean rad ius', 'mean texture', 'mean perimeter', 'mean area', 'mean smoothness', 'mean compactness', 'mean concavity', 'mean concave points', 'mean symmetry', 'mean fractal dimension', 'radius error', 'texture error', 'perimeter error', 'area error', 'smoothness error', 'compactness error', 'concavity error', 'concave points error', 'symmetry error', 'fractal dimension error', 'worst radius', 'worst texture', 'worst perimeter', 'worst area', 'worst smoothness', 'worst compactness', 'worst concavity', 'worst concave points', 'worst symmetry', 'worst fractal dimension'], dtype='<U23'), 'filename': 'breast_cancer.csv', 'data_module': 'sklearn.datasets.data'} # loading the data to a data frame data_frame = pd.DataFrame(breast_cancer_dataset.data, columns = breast_cancer_dataset.feature_names) # print the first 5 rows of the dataframe data_frame.head() Out[5]: mean mean worst mean mean mean mean mean mean worst worst worst worst worst worst worst mean mean fractal concave concave radius texture perimeter area smoothness compactness concavity symmetry radius texture perimeter area smoothness compactness concavity points dimension points 0.7119 17.99 10.38 122.80 1001.0 0.11840 0.27760 0.3001 0.14710 0.2419 0.07871 ... 25.38 17.33 184.60 2019.0 0.1622 0.6656 0.2654 0.07017 1 20.57 17.77 132.90 1326.0 0.08474 0.07864 0.0869 0.1812 24.99 23.41 158.80 1956.0 0.1238 0.1866 0.2416 0.1860 0.05667 ... 19.69 130.00 1203.0 0.10960 0.15990 0.12790 0.05999 ... 0.4245 0.4504 21.25 0.1974 0.2069 23.57 25.53 152.50 1709.0 0.1444 0.2430 0.14250 0.28390 0.10520 0.2597 26.50 0.6869 11.42 20.38 77.58 386.1 0.2414 0.09744 ... 14.91 98.87 567.7 0.2098 0.8663 0.2575 20.29 14.34 135.10 1297.0 0.10030 0.13280 0.1980 0.10430 0.1809 0.05883 ... 22.54 16.67 152.20 1575.0 0.1374 0.2050 0.4000 0.1625 5 rows × 30 columns In [6]: # adding the 'target' column to the data frame data_frame['label'] = breast_cancer_dataset.target # print last 5 rows of the dataframe data_frame.tail() Out[7]: mean mean worst mean mean worst worst worst worst worst mean mean mean mean mean mean worst wors fractal concave concave radius texture perimeter area smoothness compactness concavity symmetry texture perimeter area smoothness compactness concavity symmetr dimension points points 21.56 22.39 142.00 1479.0 0.11100 0.11590 0.24390 0.13890 166.10 2027.0 0.21130 0.4107 0.2216 564 0.1726 0.05623 26.40 0.14100 0.206 20.13 28.25 131.20 1261.0 0.09780 0.10340 0.14400 0.09791 0.1752 0.05533 38.25 155.00 1731.0 0.11660 0.19220 0.3215 0.1628 0.257 565 108.30 858.1 0.10230 0.09251 0.05302 0.05648 ... 0.11390 0.30940 0.3403 566 16.60 28.08 0.08455 0.1590 34.12 126.70 1124.0 0.1418 0.221 140.10 1265.0 0.11780 0.27700 0.15200 0.07016 184.60 1821.0 0.16500 0.86810 0.9387 0.2650 567 20.60 29.33 0.35140 0.2397 39.42 0.408 0.04362 0.00000 0.05884 ... 0.08996 0.06444 0.0000 0.0000 568 7.76 24.54 47.92 181.0 0.05263 0.00000 0.1587 30.37 59.16 268.6 0.287 5 rows × 31 columns In [8]: # number of rows and columns in the dataset data_frame.shape (569, 31) Out[8]: In [9]: # getting some information about the data data_frame.info() <class 'pandas.core.frame.DataFrame'> RangeIndex: 569 entries, 0 to 568 Data columns (total 31 columns): Column Non-Null Count Dtype -----569 non-null 0 mean radius float64 1 mean texture 569 non-null float64 mean perimeter 569 non-null float64 mean area 569 non-null float64 mean smoothness 569 non-null float64 569 non-null 5 mean compactness float64 569 non-null float64 6 mean concavity 7 mean concave points 569 non-null float64 mean symmetry 569 non-null 8 float64 mean fractal dimension 569 non-null float64 9 569 non-null 10 radius error float64 569 non-null float64 11 texture error 12 perimeter error 569 non-null float64 13 area error 569 non-null float64 14 smoothness error 569 non-null float64 569 non-null 15 compactness error float64 16 concavity error 569 non-null float64 17 concave points error 569 non-null float64 18 symmetry error 569 non-null float64 19 fractal dimension error 569 non-null float64 20 worst radius 569 non-null float64 21 worst texture 569 non-null float64 22 worst perimeter 569 non-null float64 23 worst area 569 non-null float64 24 worst smoothness 569 non-null float64 25 worst compactness 569 non-null float64 569 non-null float64 26 worst concavity 27 worst concave points 569 non-null float64 float64 28 worst symmetry 569 non-null 29 worst fractal dimension 569 non-null float64 30 label 569 non-null int32 dtypes: float64(30), int32(1) memory usage: 135.7 KB In [10]: # checking for missing values data_frame.isnull().sum() 0 mean radius Out[10]: mean texture mean perimeter mean area mean smoothness mean compactness mean concavity mean concave points mean symmetry mean fractal dimension radius error texture error perimeter error area error smoothness error compactness error concavity error concave points error symmetry error fractal dimension error worst radius worst texture worst perimeter worst area worst smoothness worst compactness worst concavity worst concave points worst symmetry worst fractal dimension 0 label dtype: int64 In [11]: # statistical measures about the data data_frame.describe() Out[11]: mean mean mean mean mean mean mean mean mean worst worst worst mean area concave fractal worst area radius smoothness compactness smoothness compactne texture perimeter concavity symmetry texture perimeter points dimension **count** 569.000000 569.000000 569.000000 569.000000 569.000000 569.000000 569.000000 569.000000 569.000000 569.000000 569.000000 569.000000 569.000000 569.000000 569.0000 14.127292 19.289649 91.969033 654.889104 0.096360 0.104341 0.088799 0.048919 0.181162 0.062798 25.677223 107.261213 880.583128 0.132369 0.2542 mean 3.524049 4.301036 24.298981 351.914129 0.052813 0.079720 0.038803 0.027414 0.007060 6.146258 33.602542 569.356993 0.022832 0.1573 std 0.014064 6.981000 9.710000 43.790000 0.019380 0.000000 0.000000 0.106000 0.049960 12.020000 0.071170 0.0272 143.500000 0.052630 50.410000 185.200000 min 11.700000 16.170000 75.170000 420.300000 0.086370 0.064920 0.029560 0.020310 0.161900 0.057700 21.080000 84.110000 515.300000 0.116600 0.1472 25% 25.410000 50% 13.370000 18.840000 86.240000 551.100000 0.095870 0.092630 0.061540 0.033500 0.179200 0.061540 97.660000 686.500000 0.131300 0.2119 **75**% 15.780000 21.800000 104.100000 782.700000 0.105300 0.130400 0.130700 0.074000 0.195700 0.066120 29.720000 125.400000 1084.000000 0.146000 0.3391 28.110000 39.280000 188.500000 0.163400 0.345400 0.304000 0.222600 1.0580 max 2501.000000 0.426800 0.201200 0.097440 ... 49.540000 251.200000 4254.000000 8 rows × 31 columns # checking the distribution of Target Varibale data_frame['label'].value_counts() 357 Out[12]: 212 Name: label, dtype: int64 1 --> Benign 0 --> Malignant data_frame.groupby('label').mean() Out[13]: mean mean mean mean mean mean mean mean worst worst worst mean area concave fractal worst area smoothness compactne radius texture perimeter smoothness compactness concavity symmetry radius texture perimeter points dimension label **0** 17.462830 21.604906 115.365377 978.376415 0.160775 0.087990 0.144845 0.102898 0.145188 0.192909 0.062680 ... 21.134811 29.318208 141.370330 1422.286321 0.3748 558.899440 **1** 12.146524 17.914762 78.075406 462.790196 0.092478 0.080085 0.046058 0.025717 0.174186 0.062867 ... 13.379801 23.515070 87.005938 0.124959 0.1826 2 rows × 30 columns Separating the features and target In [14]: X = data_frame.drop(columns='label', axis=1) Y = data_frame['label'] In [15]: Out[15]: mean mean worst mean mean mean mean mean mean mean mean worst worst worst worst worst worst worst concave fractal concave area smoothness compactness area smoothness compactness concavity concavity radius texture perimeter symmetry radius texture perimeter dimension points points 17.99 10.38 122.80 1001.0 0.11840 0.27760 0.30010 0.14710 0.2419 0.07871 2019.0 0.16220 0.66560 0.7119 0.2654 ... 25.380 17.33 184.60 20.57 17.77 132.90 1326.0 0.08474 0.07864 0.08690 0.07017 0.1812 0.05667 ... 24.990 23.41 158.80 1956.0 0.12380 0.18660 0.2416 0.1860 11.42 77.58 386.1 0.14250 0.28390 0.24140 0.10520 0.2597 0.09744 ... 14.910 98.87 567.7 0.20980 0.86630 0.6869 0.2575 135.10 1297.0 0.10030 0.13280 0.19800 0.10430 0.1809 152.20 1575.0 0.13740 0.20500 0.4000 0.1625 20.29 14.34 0.05883 ... 22.540 16.67 142.00 1479.0 0.05623 ... 25.450 21.56 22.39 0.11100 0.11590 0.24390 0.13890 0.1726 26.40 166.10 2027.0 0.14100 0.21130 0.4107 0.2216 564 20.13 28.25 131.20 1261.0 0.09780 0.10340 0.14400 0.09791 0.1752 0.05533 ... 23.690 38.25 155.00 1731.0 0.11660 0.19220 0.3215 0.1628 565 0.05648 ... 18.980 0.11390 0.30940 0.3403 566 16.60 28.08 108.30 858.1 0.08455 0.10230 0.09251 0.05302 0.1590 34.12 126.70 1124.0 0.1418 29.33 140.10 1265.0 0.11780 0.27700 0.35140 0.15200 0.2397 0.07016 ... 25.740 39.42 184.60 1821.0 0.16500 0.86810 0.9387 0.2650 20.60 0.05263 0.04362 0.05884 ... 9.456 0.08996 0.06444 0.0000 7.76 24.54 47.92 181.0 0.00000 0.00000 0.1587 30.37 268.6 0.0000 568 59.16 569 rows × 30 columns In [16]: 0 Out[16]: 0 2 0 3 0 0 564 0 565 0 566 0 567 0 568 Name: label, Length: 569, dtype: int32 Splitting the data into training data & Testing data X_train, X_test, Y_train, Y_test = train_test_split(X, Y, test_size=0.2, random_state=2) print(X.shape, X_train.shape, X_test.shape) In [18]: (569, 30) (455, 30) (114, 30) **Model Training** Logistic Regression model = LogisticRegression() In [19]: # training the Logistic Regression model using Training data model.fit(X_train, Y_train) C:\Users\Azlan\anaconda3\envs\tensorflow_env\lib\site-packages\sklearn\linear_model_logistic.py:444: ConvergenceWarning: lbfgs failed to converge (status=1): STOP: TOTAL NO. of ITERATIONS REACHED LIMIT. Increase the number of iterations (max_iter) or scale the data as shown in: https://scikit-learn.org/stable/modules/preprocessing.html Please also refer to the documentation for alternative solver options: https://scikit-learn.org/stable/modules/linear_model.html#logistic-regression n_iter_i = _check_optimize_result(Out[20]: ▼ LogisticRegression LogisticRegression() Model Evaluation Accuracy Score In [21]: # accuracy on training data X_train_prediction = model.predict(X_train) training_data_accuracy = accuracy_score(Y_train, X_train_prediction) print('Accuracy on training data = ', training_data_accuracy) Accuracy on training data = 0.9494505494505494In [23]: # accuracy on test data X_test_prediction = model.predict(X_test) test_data_accuracy = accuracy_score(Y_test, X_test_prediction) In [24]: print('Accuracy on test data = ', test_data_accuracy) Accuracy on test data = 0.9298245614035088 Building a Predictive System In [25]: input_data = (13.54,14.36,87.46,566.3,0.09779,0.08129,0.06664,0.04781,0.1885,0.05766,0.2699,0.7886,2.058,23.56,0.008462,0.0146,0.02387,0.01315,0.0198,0.0023,15 # change the input data to a numpy array input_data_as_numpy_array = np.asarray(input_data) # reshape the numpy array as we are predicting for one datapoint input_data_reshaped = input_data_as_numpy_array.reshape(1,-1) prediction = model.predict(input_data_reshaped) print(prediction) if (prediction[0] == 0): print('The Breast cancer is Malignant') print('The Breast Cancer is Benign') [1] The Breast Cancer is Benign C:\Users\Azlan\anaconda3\envs\tensorflow_env\lib\site-packages\sklearn\base.py:450: UserWarning: X does not have valid feature names, but LogisticRegression w as fitted with feature names warnings.warn(