

Diagnosing breast tumour using machine learning

Tomáš Chobola

Table of content

1. Overview	2
2. Dataset overview	2
2.1. Dataset description	2
2.2. Dataset example	2
3. Learning process and making predictions	3
4. Sample run of the program	4
5. Sources	6

1. Overview

The main purpose of the program is to predict breast cancer based on features computed from a digitised image of a fine needle aspirate (FNA) of a breast mass. For the prediction the program uses the k -nearest neighbours algorithm (k -NN) and dataset from the UCI Machine Learning Repository, specifically Breast Cancer Wisconsin (Diagnostic) Data Set.^[1]

2. Dataset overview

2.1. Dataset description

Dataset consists of 569 instances of cell nuclei from various breast tumours. Every instance describes features of the cell nuclei which are computed from a digitised image of a fine needle aspirate (FNA) of a breast mass.

The dataset has 32 columns - id number, diagnosis (M for malignant, B for benign), radius (mean of distances from centre to points on the perimeter), texture (deviation of grayscale values), perimeter (mean size of the core tumour), area, smoothness, compactness, concavity (mean of severity of concave portions of the contour), concave points (mean of severity of concave portions of the contour), symmetry, fractal dimension (mean for "coastline approximation" - 1), standard error of the radius, standard error of the texture, standard error of the smoothness, standard error of the compactness, standard error of the concavity, standard error of the concave points.^[2]

Application uses only the first 12 columns. First two columns are used merely for identifying the results of the cell analysis. Remaining 10 columns are used for training and validating.

2.2. Dataset example

853201	M	17.57	15.05	115	955.1	0.0984	0.1157	0.0988	0.0795	0.1739	0.06149
853401	M	18.63	25.11	124.8	1088	0.1064	0.1887	0.2319	0.1244	0.2183	0.06197
853612	M	11.84	18.7	77.93	440.6	0.1109	0.1516	0.1218	0.0518	0.2301	0.07799
854253	M	16.74	21.59	110.1	869.5	0.0961	0.1336	0.1348	0.0602	0.1896	0.05656
854002	M	19.27	26.47	127.9	1162	0.0940	0.1719	0.1657	0.0759	0.1853	0.06261

3. Learning process and making predictions

Application is using the k -nearest neighbours algorithm that stores all available cases and classifies new cases based on similarity (distance function).

Before the actual process of learning the whole dataset is split into two parts. The first part contains two thirds of the dataset and is used for learning and the remaining second part is used for validation. The training data are then fit into a classifier that is an instance of Python class `KNeighborsClassifier` from `scikit-learn` library. Application then finds the best k value for the inserted data by scoring the precision of the classifier using the validation part of the dataset. The classifier uses Euclidean distance as the metric function.

$$D = \sqrt{\sum_{i=1}^n (x_i - y_i)^2}$$

After picking the right k value the program is ready to make predictions based on values set by the user. The result is either that the breast tumour is malignant or benign. The program is able to predict that with precision around 92 % based on testing using the validation part of the original dataset, results may differ when the program is used for actual predictions.

The accuracy could be increased by using much larger dataset for learning and validating processes which would describe almost all possible situations, and taking measuring error into account. Another improvement in accuracy could be made by using different distance functions and comparing the results or using a simple heuristic that would assign weights to the values based on the importance of that particular value in recognising the malignancy of the breast tumour.

4. Sample run of the program

Entry fields for entering the values for which the user wants to make a prediction.

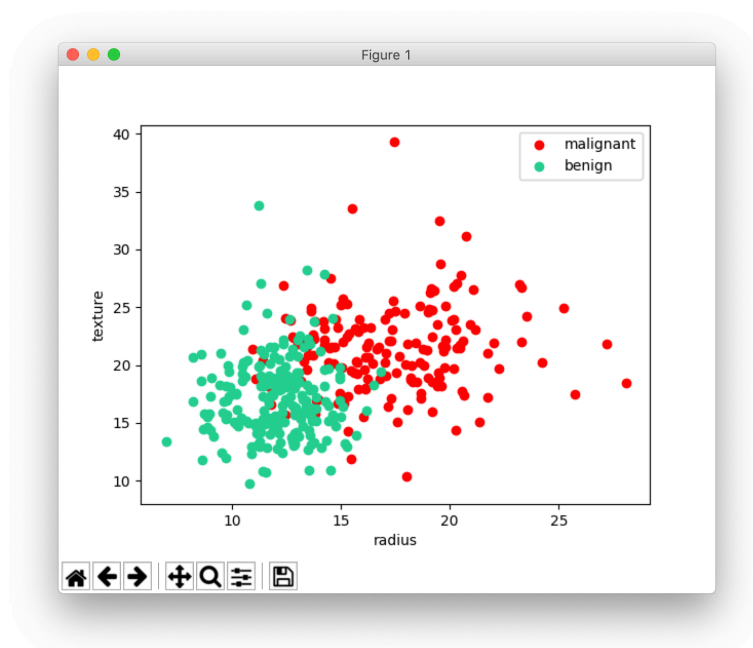
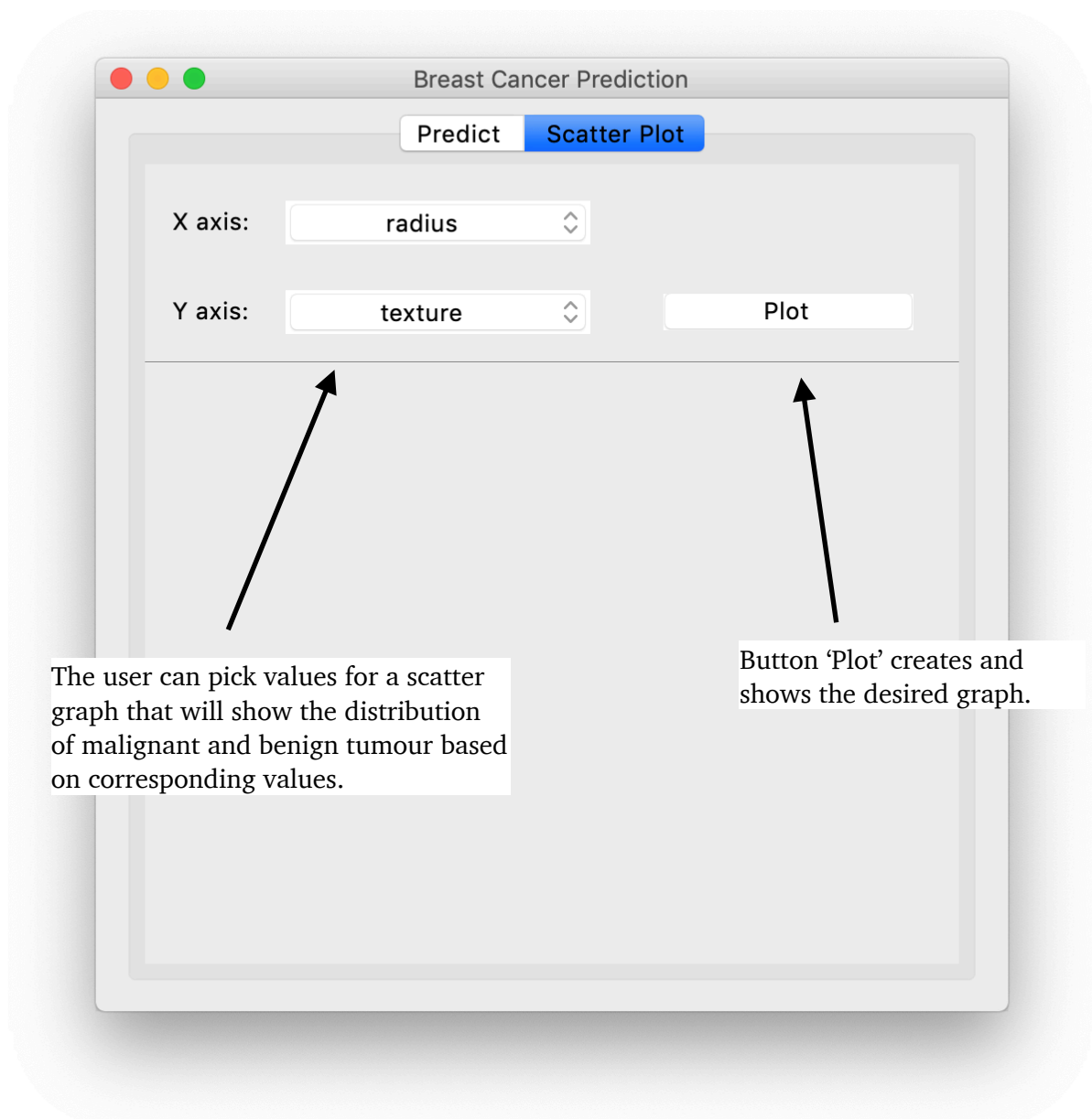
The screenshot shows a window titled "Breast Cancer Prediction". It has two tabs: "Predict" (which is selected and highlighted in blue) and "Scatter Plot". Below the tabs, there are ten input fields arranged in two columns. The left column contains fields for "radius", "texture", "perimeter", "area", and "smoothness". The right column contains fields for "compactness", "concavity", "concave points", "symetry", and "fractal dimension". Below these input fields is a large white rectangular area labeled "result". To the right of the "result" area is a button labeled "Calculate".

Annotations with arrows point to the following elements:

- An arrow points to the "radius" input field, with the text "Entry fields for entering the values for which the user wants to make a prediction." above it.
- An arrow points to the "compactness" input field.
- An arrow points to the "result" text field.
- An arrow points to the "Calculate" button.

Text field for showing the result of the prediction.
If not all of the entry fields are filled or do not contain numbers it is going to show an error message based on the value that is incorrect.

By pressing the button 'Calculate' the program will preform prediction based on the values in the entry fields.



Example of plotted graph.

5. Sources

- [1] Wolberg, William H. “Breast Cancer Wisconsin (Diagnostic) Data Set.” Edited by Nick Street and Olvi Mangasarian, *UCI Machine Learning Repository: Flags Data Set*, [archive.ics.uci.edu/ml/datasets/Breast+Cancer+Wisconsin+\(Diagnostic\)](http://archive.ics.uci.edu/ml/datasets/Breast+Cancer+Wisconsin+(Diagnostic)).
- [2] Uci. “Breast Cancer Wisconsin (Diagnostic) Data Set.” *RSNA Pneumonia Detection Challenge | Kaggle*, 25 Sept. 2016, www.kaggle.com/uciml/breast-cancer-wisconsin-data.