

# Machine Learning Engineer Nanodegree

## Capstone Proposal

Jeyvison Nascimento

May 06, 2017

## Proposal

### Domain Background

Breast cancer is the most dangerous type of cancer, besides lung cancer, a woman may develop. It's death rate is high and it is expected that 12% of the american women will develop it during their lifetime.

### Problem Statement

The goal is to create a program to help the doctor choose a protocol to deal with breast cancer when it's malignant or benign. Many informations will be considered like radius, perimeter and smoothness. With those information the cancer will be classified as malignant or benign and the doctor can decide which steps to take next

The program will:

- Train the model based on a dataset
- Identify if the cancer is malignant or benign
- Inform the doctor about the result

### Datasets and Inputs

The [dataset](#) that will be used has 30 features and a binary target feature that tells if the cancer was malignant or benign. The input of the program will be a csv containing values for classification.

The dataset has 569 instances and is unbalanced so it will need to be stratified. For training purposes, the training will be splitted on 70% for training and 30% for testing using K-fold cross validation.

Here follows the list of the 30 features:

1. radius\_mean
2. texture\_mean
3. perimeter\_mean
4. area\_mean
5. smoothness\_mean
6. compactness\_mean
7. concavity\_mean
8. concave points\_mean
9. symmetry\_mean
10. fractal\_dimension\_mean
11. radius\_se
12. texture\_se
13. perimeter\_se
14. area\_se
15. smoothness\_se
16. compactness\_se
17. concavity\_se
18. concave points\_se
19. symmetry\_se
20. fractal\_dimension\_se
21. radius\_worst
22. texture\_worst
23. perimeter\_worst
24. area\_worst
25. smoothness\_worst
26. compactness\_worst
27. concavity\_worst
28. concave points\_worst
29. symmetry\_worst
30. fractal\_dimension\_worst

## Solution Statement

As we have many real-valued features in the dataset we will need to scale these values. SVC(SVM) will be used to perform the classification because we have a small dataset. Naive Bayes could be used for this one too but it does not provide a reliable probability prediction which can be useful for choice making.

## Benchmark Model

For benchmarking, a Scikit learn Dummy classifier will be used to compare the results with the SVC classifier.

## Evaluation Metrics

We have two metrics we may use to evaluate the model, Accuracy and F1 score. Both will be used to measure the model's performance.

## Project Design

The steps needed to develop the project are the following:

- Develop a simple python script
- Make the script load the dataset file(CSV)
- The dataset is given to the classifier for training
- The script loads a CSV evaluation file and uses its content for prediction
- The result showing the type of the breast cancer and the prediction probability will be shown in the screen.