

### Description:

- Implement the **naïve Bayes classifier (NBC)** and **logistic regression with gradient descent** algorithms in Python or MATLAB to solve a **binary classification problem** for a dataset of your choosing with continuous input and binary output.
- You may use all built-in functions. The built-in NBC and logistic regression models in Python and MATLAB that can be used for this assignment are:
  - Python - Gaussian naïve Bayes (for continuous input): [https://scikit-learn.org/stable/modules/generated/sklearn.naive\\_bayes.GaussianNB.html](https://scikit-learn.org/stable/modules/generated/sklearn.naive_bayes.GaussianNB.html)
  - Python – logistic regression: [https://scikit-learn.org/stable/modules/generated/sklearn.linear\\_model.SGDClassifier.html](https://scikit-learn.org/stable/modules/generated/sklearn.linear_model.SGDClassifier.html) with the “loss” parameter set to ‘log\_loss’
  - MATLAB - Gaussian naïve Bayes (for continuous input): <https://www.mathworks.com/help/stats/classificationnaivebayes.html> (default is Gaussian/normal)
  - MATLAB – logistic regression: <https://www.mathworks.com/help/stats/incrementalclassificationlinear.html> with the “learner” parameter set to ‘logistic’

Train and test your models with a dataset of your choosing that meets the following criteria:

- Number of input features: 3+
- Input characteristics: Continuous real-valued
- Output characteristics: Binary

**Note:** If you find a dataset that has 3 or more output classes, you can use it by removing the data examples from other classes

- Use 80% of the dataset for training and 20% for testing your model.
- Use one of the following sources to find a dataset:
  - University of California, Irvine Machine Learning Repository <https://archive.ics.uci.edu/ml/index.php>
  - Kaggle <https://www.kaggle.com/>
  - Awesome Public Datasets <https://github.com/awesomedata/awesome-public-datasets>
  - Google Dataset Search Engine <https://datasetsearch.research.google.com/>
  - Microsoft Research Open Data <https://msrpendata.com/>
  - U.S. Government’s Open Data <https://www.data.gov/>
  - Registry of Research Data Repositories <https://www.re3data.org/>
  - CMU Libraries <https://guides.library.cmu.edu/machine-learning/datasets>

Summarize your approach and results in a report that includes at least the following:

- The dataset you used, its source and characteristics.
- The data preprocessing steps you took (if any).
- The solution  $w$  (parameter vector) for logistic regression. This vector should include the intercept (bias term).

- Relevant evaluation metrics for NBC (accuracy, sensitivity, specificity, f1 score, log loss) for BOTH the training and test datasets.
- Relevant evaluation metrics for logistic regression (accuracy, sensitivity, specificity, f1 score, log loss) for BOTH the training and test datasets.