

Assignment #2- Git, GitHub, ML

Step 1:

Install git in your local machine (If it is already done, skip this step)

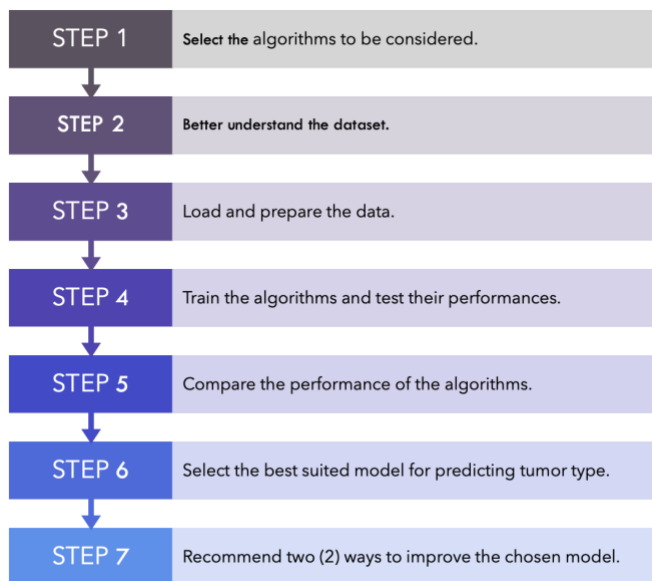
SKIPPED

Note: This is the third time I will be working with Wisconsin Breast Cancer Dataset. Therefore, the following steps were carried out using the code I wrote when I worked

Step 2:

Build a ML model for Breast Cancer Wisconsin (Diagnostic) Data set in Jupyter notebook.

[https://archive.ics.uci.edu/ml/datasets/breast+cancer+wisconsin+\(diagnostic\)](https://archive.ics.uci.edu/ml/datasets/breast+cancer+wisconsin+(diagnostic))



Step 6 and Step 7, in the above picture, will be carried at a later stage in this document.

Step 3:

Please provide screenshots for various stages of the design process (importing data, training, evaluation ...)

jupyter raj_dholakia_model_v1 Last Checkpoint: an hour ago (unsaved changes) Logout

File Edit View Insert Cell Kernel Widgets Help Trusted Python 3 (ipykernel)

Predicting Breast Tumor Type - V1

In this jupyter notebook, we will develop a classification algorithm to predict the diagnostic for the [UCI Machine Learning Repository's Breast Cancer Dataset](#). For the analysis, we will evaluate the performance following four (4) algorithms:

1. Support Vector Machines
2. K-Nearest Neighbors
3. Random Forest Classifier
4. Logistical Regression

Lecturer: Sajeewa Salgadoe
Module: AIDI-2004-02 - AI in Enterprise

STEP 1

Select the algorithms to be considered.

↓

STEP 2

Better understand the dataset.

↓

STEP 3

Load and prepare the data.

↓

STEP 4

Train the algorithms and test their performances.

↓

STEP 5

Compare the performance of the algorithms.

↓

STEP 6

Select the best suited model for predicting tumor type.

↓

STEP 7

Recommend two (2) ways to improve the chosen model.

Table of Contents

- [1. Business Problem Description](#)
- [2. Exploratory Data Analysis](#)
 - [2.1 Dataset Information](#)
 - [2.2 Loading Data](#)
 - [2.3 Exploratory Data Analysis](#)

2.2 Loading Data

```
In [1]: 1 # Load required libraries
2 import pandas as pd
3 import numpy as np
4 import matplotlib.pyplot as plt
5 %matplotlib inline
6 import seaborn as sns
7
8 from sklearn.model_selection import train_test_split
9 from sklearn.preprocessing import StandardScaler
10 from sklearn.linear_model import LogisticRegression
11 from sklearn.neighbors import KNeighborsClassifier
12 from sklearn.ensemble import RandomForestClassifier
13 from sklearn.tree import DecisionTreeClassifier
14 from sklearn.metrics import classification_report, confusion_matrix, auc
15
16 import os
17 import pathlib
18 from helper_functions import metric_evaluation, learning_curve, full_model_evaluation
19 # from helper_functions import plot_conf_mat, evaluate_model
```

```
In [2]: 1 # Define location of the data
2 data_dir = '../data'
3 filename = 'dataset.csv'
4 data_path = os.path.join(data_dir, filename)
5
6 if pathlib.Path(data_path).exists():
7     print(f"File {filename} found.")
8 else:
9     raise FileNotFoundError('No file found at the location defined.')
```



File dataset.csv found.

```
In [3]: 1 # Load data into a pandas DataFrame
2 data = pd.read_csv(data_path)
3 data.head()
```


Out[3]:

	id	diagnosis	radius_mean	texture_mean	perimeter_mean	area_mean	smoothness_mean
0	87139402	B	12.32	12.39	78.85	464.1	0.10280
1	8910251	B	10.60	18.95	69.28	346.4	0.09688
2	905520	B	11.04	16.83	70.92	373.2	0.10770
3	868871	B	11.28	13.39	73.00	384.8	0.11640
4	9012568	B	15.19	13.21	97.65	711.8	0.07963

5 rows × 32 columns

 jupyter Breast Cancer Diagnostic (unsaved changes)  Logout

File Edit View Insert Cell Kernel Widgets Help Not Trusted Python 3 (ipykernel)



3. Data Preparation

```
In [15]: 1 # Define x and y variables for CMS prediction
2 x = data.drop('diagnosis', axis=1).to_numpy()
3 y = data["diagnosis"].replace({'M': 1, "B": 0})
4
5 # Splitting data into train and test datasets
6 x_train, x_test, y_train, y_test = train_test_split(x, y, stratify=y, test_size=0.2,
7
8 # Scaling the data
9 sc = StandardScaler()
10 x_train_scaled = sc.fit_transform(x_train)
11 x_test_scaled = sc.transform(x_test)
12
13 # Created scaled x, to use cross-validated
14 x_scaled = sc.fit_transform(x)
```

4. Development of Models

In this section, we will train and test the performance of the model based on the **cross validated scores** of the following metrics:

1. accuracy
2. precision_weighted
3. recall_weighted
4. f1_weighted
5. roc_auc

4.1 Cross Validation

- A 5-fold cross validation score, first splits the data into five (5) parts.
- It then, considers one (1) part to be the test set and trains the model on the rest of the data. It repeats this for all the five (5) parts.

Note: in this analysis, an average of the cross validated scores is taken into consideration. One could also look at individual values to observe if the performance of a particular algorithm is consistent on all the five (5) parts/folds.

- Using this method will give us a better understanding of which algorithms perform better.
- The Brier Score uses probabilities and the expected values to calculate the mean squared error. The algorithms support probabilities.

jupyter Breast Cancer Diagnostic (unsaved changes) Logout

File Edit View Insert Cell Kernel Widgets Help Not Trusted | Python 3 (ipykernel) ○

+ ✂ 📄 📁 ⬆ ⬇ ▶ Run ■ ↺ ▶▶ Code 🔍

4.2 Model Training

```
In [16]: 1 models = {}
2
3 for name,method in [('Logistic Regreesion', LogisticRegression()),
4                     ('KNN', KNeighborsClassifier()),
5                     ('Random Forest', RandomForestClassifier()),
6                     ('Decision Tree', DecisionTreeClassifier())]:
7     # Train the model
8     method.fit(x_train_scaled, y_train)
9
10    # Make predictions
11    predictions = method.predict(x_test_scaled)
12
13    # Evaluate the performance
14    target_names=['benign','malignant']
15    print(f'\n----- {name} -----')
16    print(confusion_matrix(y_test, predictions))
17    print(classification_report(y_test, predictions, target_names=target_names))
18
19
20    # Save the models
21    models.update({name: method})
```

```
----- Logistic Regreesion -----
[[69  3]
 [ 0 42]]
              precision    recall  f1-score   support

   benign         1.00        0.96        0.98         72
  malignant         0.93        1.00        0.97         42

 accuracy                   0.97         114
 macro avg         0.97        0.98        0.97         114
weighted avg         0.98        0.97        0.97         114

----- KNN -----
[[71  1]
 [ 1 41]]
              precision    recall  f1-score   support

   benign         0.99        0.99        0.99         72
  malignant         0.98        0.98        0.98         42

 accuracy                   0.98         114
 macro avg         0.98        0.98        0.98         114
weighted avg         0.98        0.98        0.98         114
```

jupyter Breast Cancer Diagnostic (autosaved) Logout

File Edit View Insert Cell Kernel Widgets Help Not Trusted | Python 3 (ipykernel) ○

📁 + 🔍 📄 📄 ⬆ ⬇ ▶ Run ■ ↺ ▶▶ Code 🗨

```

[[69  3]
 [ 0 42]]

              precision    recall  f1-score   support

      benign       1.00        0.96        0.98         72
      malignant     0.93        1.00        0.97         42

   accuracy              0.97         114
  macro avg       0.97        0.98        0.97         114
 weighted avg       0.98        0.97        0.97         114

----- KNN -----
[[71  1]
 [ 1 41]]

              precision    recall  f1-score   support

      benign       0.99        0.99        0.99         72
      malignant     0.98        0.98        0.98         42

   accuracy              0.98         114
  macro avg       0.98        0.98        0.98         114
 weighted avg       0.98        0.98        0.98         114

----- Random Forest -----
[[70  2]
 [ 1 41]]

              precision    recall  f1-score   support

      benign       0.99        0.97        0.98         72
      malignant     0.95        0.98        0.96         42

   accuracy              0.97         114
  macro avg       0.97        0.97        0.97         114
 weighted avg       0.97        0.97        0.97         114

----- Decision Tree -----
[[69  3]
 [ 1 41]]

              precision    recall  f1-score   support

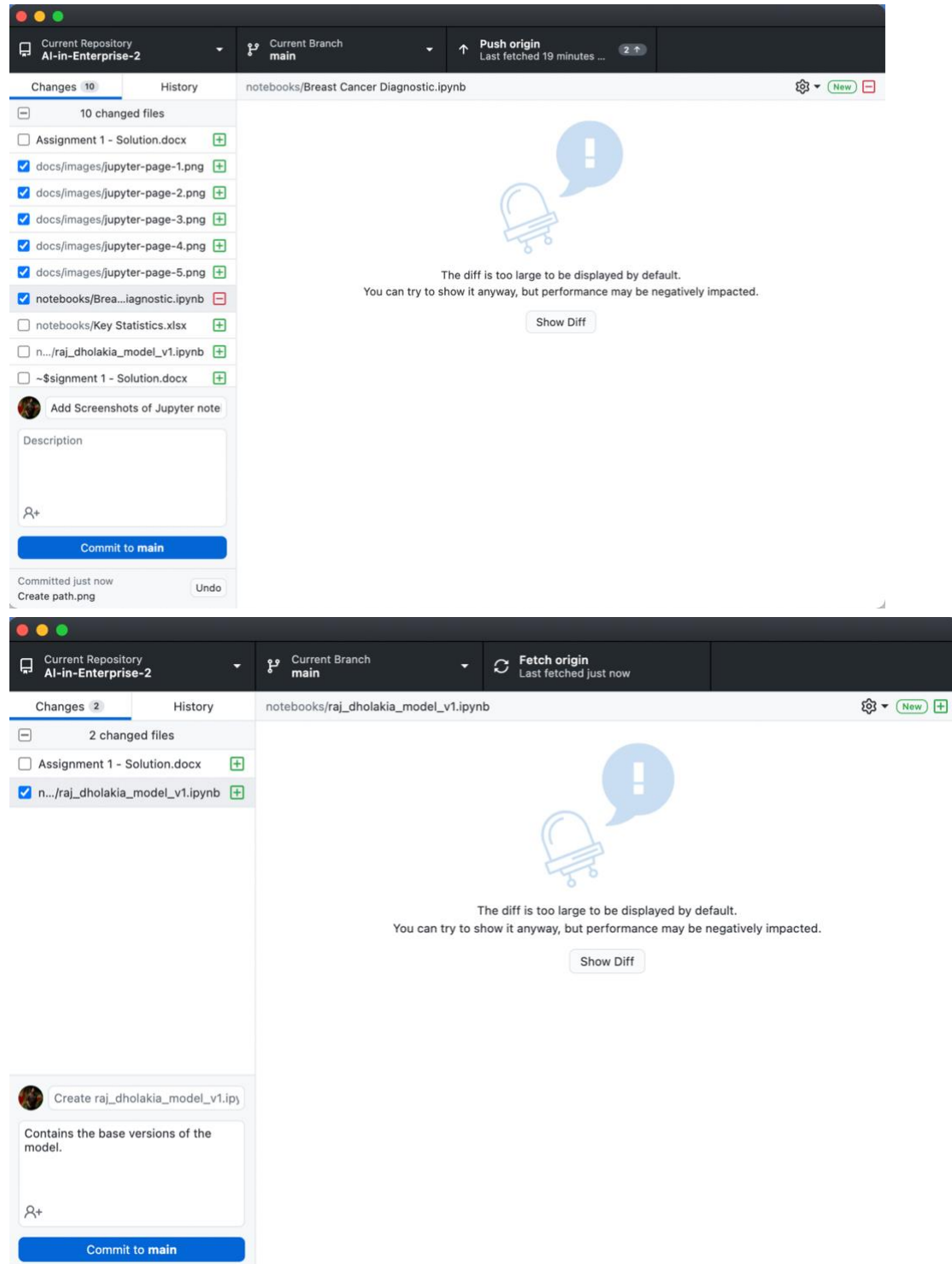
      benign       0.99        0.96        0.97         72
      malignant     0.93        0.98        0.95         42

   accuracy              0.96         114
  macro avg       0.96        0.97        0.96         114
 weighted avg       0.97        0.96        0.97         114

```


Step 4:

Upload your model (Python script, let's called it <yourname>_model_v1) to GitHub. Provide screenshot of all your git commands and your command prompt showing success of commit of your model files in the remote host.



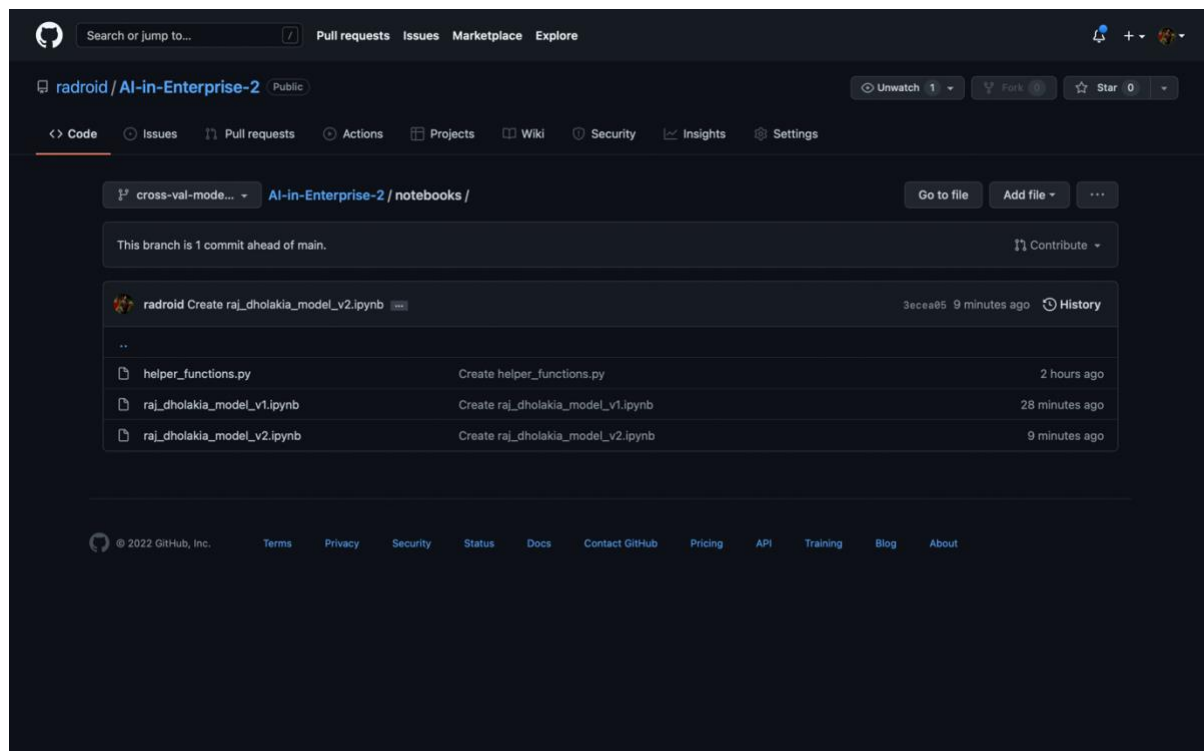
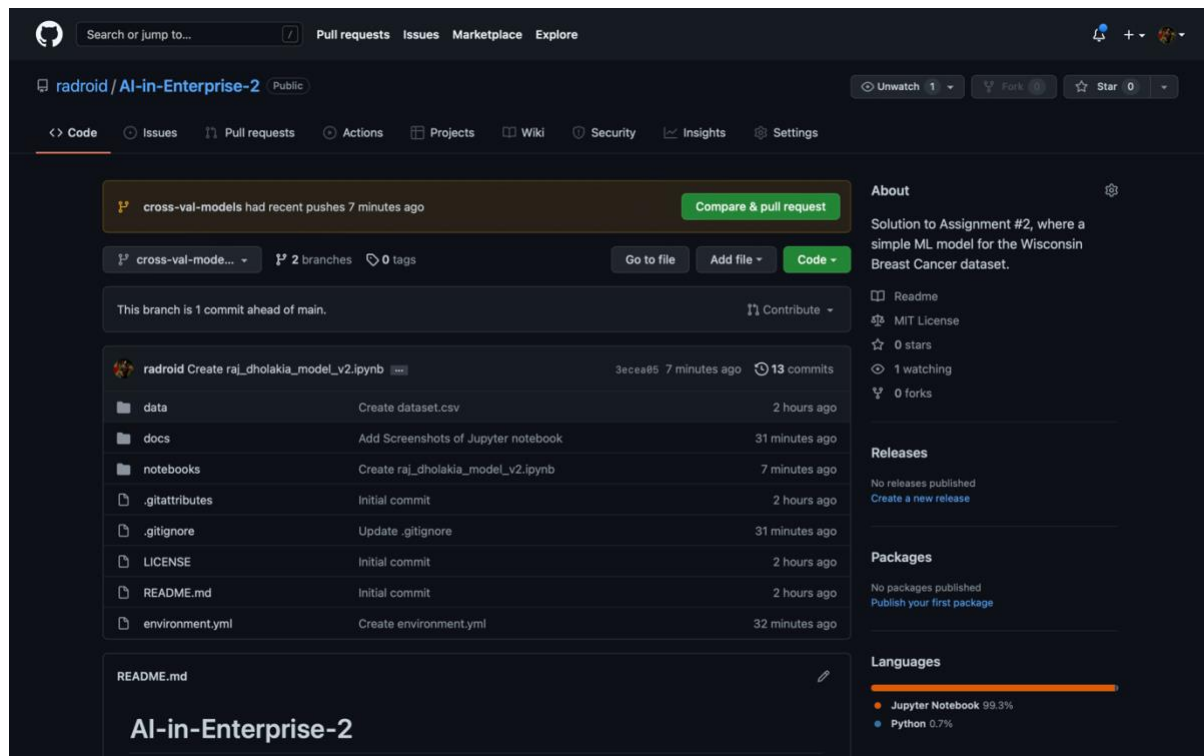
Step 5:

Create a branch in your repo and upload another ML model (may be using a different algorithm and named the file: <yourname>_model_v2) of your choice for the same dataset into that branch.



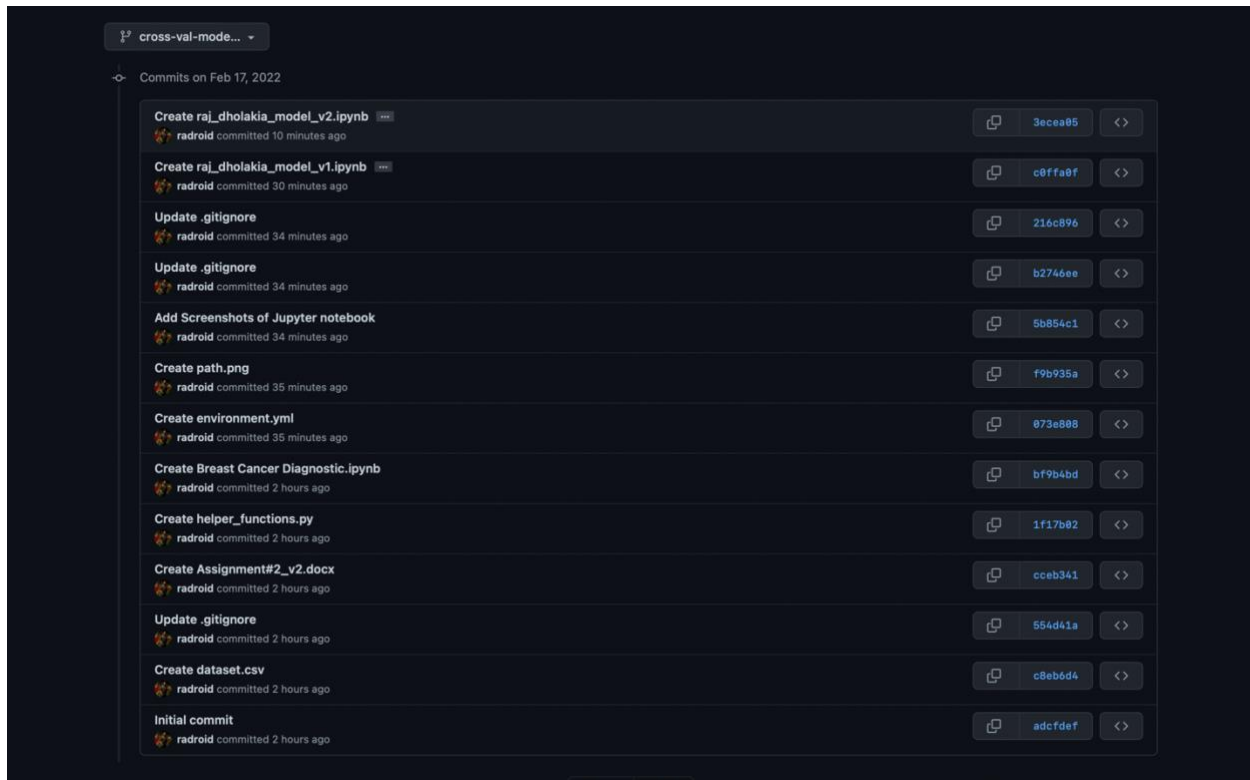
Step 6:

Navigate to your newly created branch and provide screenshot showing status of your repo.



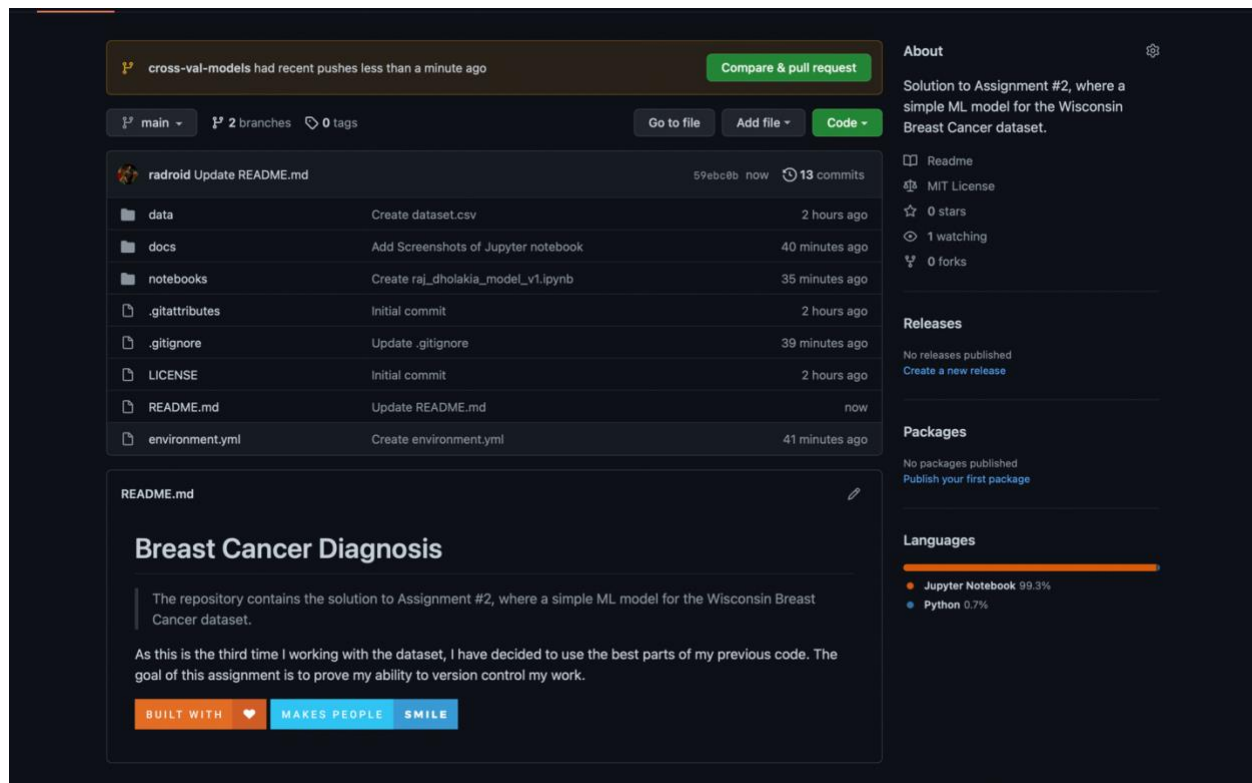
Step 7:

Provide a screenshot showing your log of activities and perform your final commit.



Step 8:

Provide a description of your program in the README.md file.

**Step 9:**

Make your repo public and share the link of your repo for check.

Link to the repository: [AI-in-Enterprise-2](#)

THANK YOU FOR READING!