# U D A C I T Y

# AI and EHR

| REVIEW |
| :---: |
| CODE REVIEW |
| HISTORY |

## Meets Specifications

Udacity student,

your project shows how committed you are to the course. Working with AI and Healthcare data is challenging and also important for future researches and companies to enhance this area.

Keep the great work on the next sections of this amazing Udacity course!!! I would like to share some nice resources:

A guide to start your path in Data Science and Machine Learning:

How can Artificial Intelligence Help Health Care?

AI and machine learning for healthcare

Finally, I would like to share interview questions tips:

Data Science Interview

😄

## Exploratory Data Analysis

The project correctly identified which field(s) has/have a high amount of missing/zero values.

The project correctly identified which field(s) has/have a Gaussian distribution shape based on the histogram analysis.

The project correctly identified fields with high cardinality.

The project justified why these fields had high cardinality.

The project correctly describes the distributions for age and gender.

Optional: The project uses Tensorflow Data Validation Visualizations to analyze which fields have a high amount of missing/null values or high cardinality.

## Exploratory Data Analysis

Great job! 😄

In this part of the project, you are asked to:

[✔] - The project correctly identified which field(s) has/have a high amount of missing/zero values.
[✔] - The project correctly identified which field(s) has/have a Gaussian distribution shape based on the histogram analysis.
[✔] - The project correctly identified fields with high cardinality.
[✔] - The project justified why these fields had high cardinality.
[✔] - The project correctly describes the distributions for age and gender.

The project correctly identifies whether to include/exclude payer_code and weight fields.

The project justified why these fields should be included/excluded by using supporting data analysis.

## Exploratory Data Analysis

Awesome! 😄

The project correctly identifies whether to include/exclude payer_code and weight fields. You also justified why these fields should be included/excluded by using supporting data analysis.

# Data Preparation

The project uses the correct level(s) for the given EHR dataset (line, encounter, patient) and transforms, aggregates and filters appropriately.
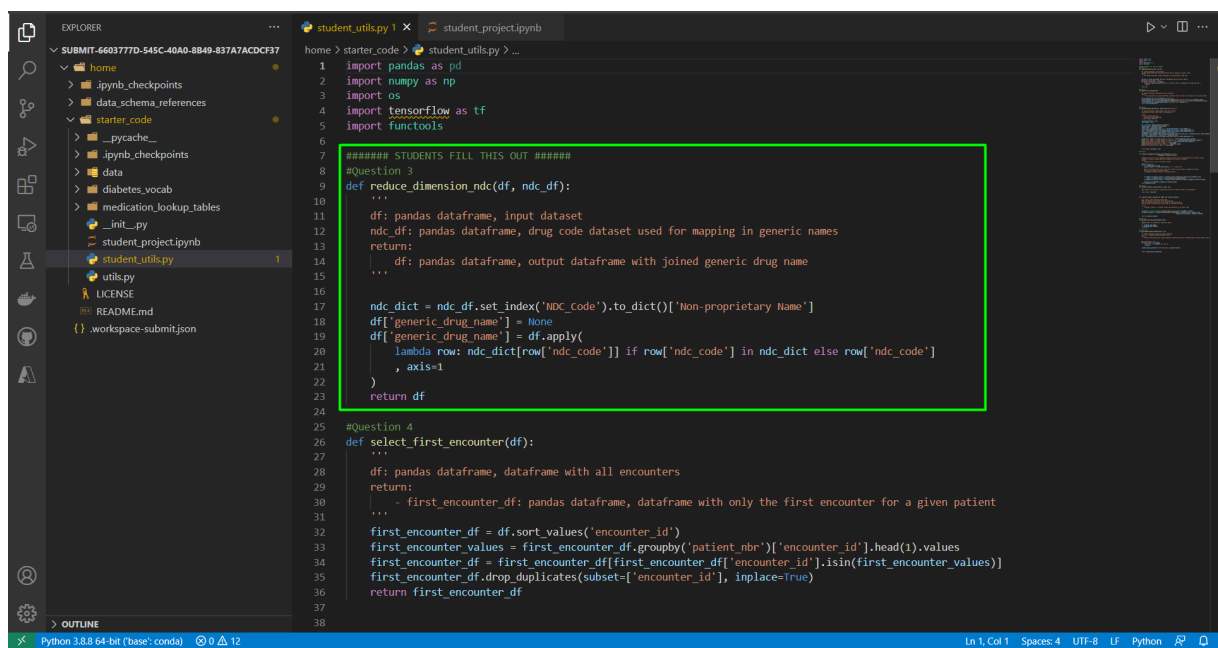
## Data Preparation

Good job! 😄

You project uses the correct level(s) for the EHR dataset. One way to validate is the **passed test** in your notebook! 😉

The project correctly maps NDC codes to generic drug names and prints out the correct mappings in the notebook.

## Data Preparation

Perfect! 😄

Your file `student_utils.py` is completed and have all necessary implementations. The function is working correctly and you have printed the output.



The project has correctly split the original dataset into train, validation, and test datasets.

The projects dataset splits do not contain patient or encounter data leakage.

The Projects code passes the Encounter Test.

## Data Preparation

All set here! 😄

This part of the project you have successfully completed! The dataset has a train, validation and test sets.

### Notebook

## Function



# Feature Engineering

The project correctly completes the categorical feature transformer boilerplate function.

The project successfully uses this function to transform the demo dataset with at least one new categorical feature.

## Feature Engineering

Great! 😉

The project successfully uses this function to transform the demo dataset with at least one new categorical feature.

The project correctly completes the numerical feature transformer boilerplate function.

The project successfully uses this function to transform the demo dataset with at least one new numerical feature.

The project's transformer function correctly incorporates the provided z-score normalizer function for normalization or another custom normalizer.

## Feature Engineering

Good job here as well! 😉

The project successfully uses functions to transform the demo dataset with at least one new numerical feature.

# Model Building and Analysis

The project has prepared the regression model predictions for TF Probability and binary classification outputs by doing the following:
Correctly utilized TF Probability to provide mean and standard deviation prediction outputs
Created an output prediction dataset that has the labels correctly mapped to a binary prediction and actual value.

## Model Building and Analysis

Great job! 😄

Once again your project set the regression model for TF Probability and binary classification outputs. You also created the dataset with the necessary outputs.

## Feature Engineering

The model has been evaluated across the following classification metrics: AUC, F1, precision, and recall.

Students have completed both questions for the model summary and address bias-variance tradeoff in regard to this problem.

# Model Building and Analysis

Nice job! 😄

Your model has been evaluated and you have presented the necessary metrics like AUC, F1, precision and recall.

As a suggestion, you can also use the ROC curve to improve your project

## Notebook



The project contains a bias report with the following:

- A visualization of at least two key metric(s) for patient selection
- A visualization showing at least one reference group fairness example and its comparison on at least one metric (e.g. TPR).
- Justification for analysis made about at least one visualization

# Model Building and Analysis

Good job on this last part of your project! 😄

You have provided the key metrics for patient selection and justified your point of view!

[⤓] **DOWNLOAD PROJECT**

RETURN TO PATH

Rate this review

START