# Advanced Concepts in Data Analytics

## Lab: Machine Learning

This page was intentionally left blank.

# Table of Contents

# Machine Learning

## Introduction

In the lab, you'll practice machine learning techniques to solve real-world problems. Using three different datasets, you'll perform classification, regression and clustering, respectively. You'll continue using pandas for data manipulation and seaborn for data visualization, and you'll also use scikit-learn, one of most widely used libraries in machine learning.

## Equipment and Materials

- BYOD laptop
- Python
- Visual Studio Code
- TXT file: requirements.txt
- Git file: gitignore
- Lab Activity 1 data and skeleton files:
    - titanic.zip
    - truth_titanic.csv
    - classficatioin.py
- Lab Activity 2 data and skeleton files:
    - house-prices-advanced-regression-techniques.zip
    - truth_house_prices.csv
    - regression.py
- Lab Activity 3 data and skeleton files:
    - Seeds_dataset.txt
    - clustering.py

## Instructions

## Lab Setup

1. Download all the data and skeleton files listed in the Equipment and Materials section above and create a folder for them on your computer (e.g., **lab3**).

2. Unzip each .zip file into your lab3 folder.

3. Open the folder in VS Code.

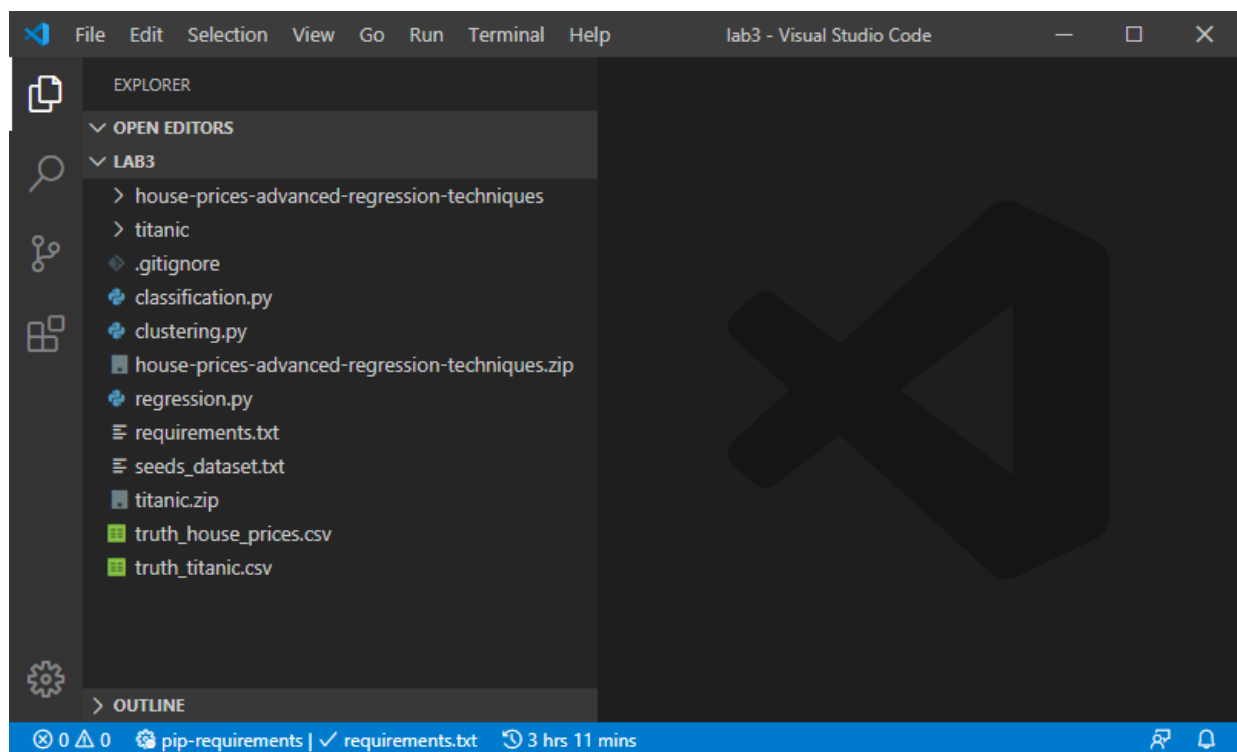   The initial structure should resemble the image below.



**Figure 1: Opening the Lab Folder in VS Code**
Used with permission from Microsoft.

4. Create a virtual environment named **venv**, configure VS Code to use the newly created virtual environment as the default Python interpreter, and install the required dependencies. See Lab 2, Activity 1: Steps 3 to 6 for reference.

# Lab Activity 1: Classifying Titanic Survival

In this section, you'll investigate one of the most infamous disasters: the sinking of the Titanic. Although over 1,500 lives were lost, some people survived. You'll use a dataset containing a list of passengers and their characteristics (survival, sex, age, ticket class, cabin number, port of embarkation, etc.) and build a predictive model to classify whether a passenger would have survived based on their passenger data.

**Note:** Although you have already downloaded this dataset, it is also publicly available from Titanic: Machine Learning from Disaster (https://www.kaggle.com/c/titanic).

1.  Use the skeleton file **classification.py** and run through the nonempty code blocks.

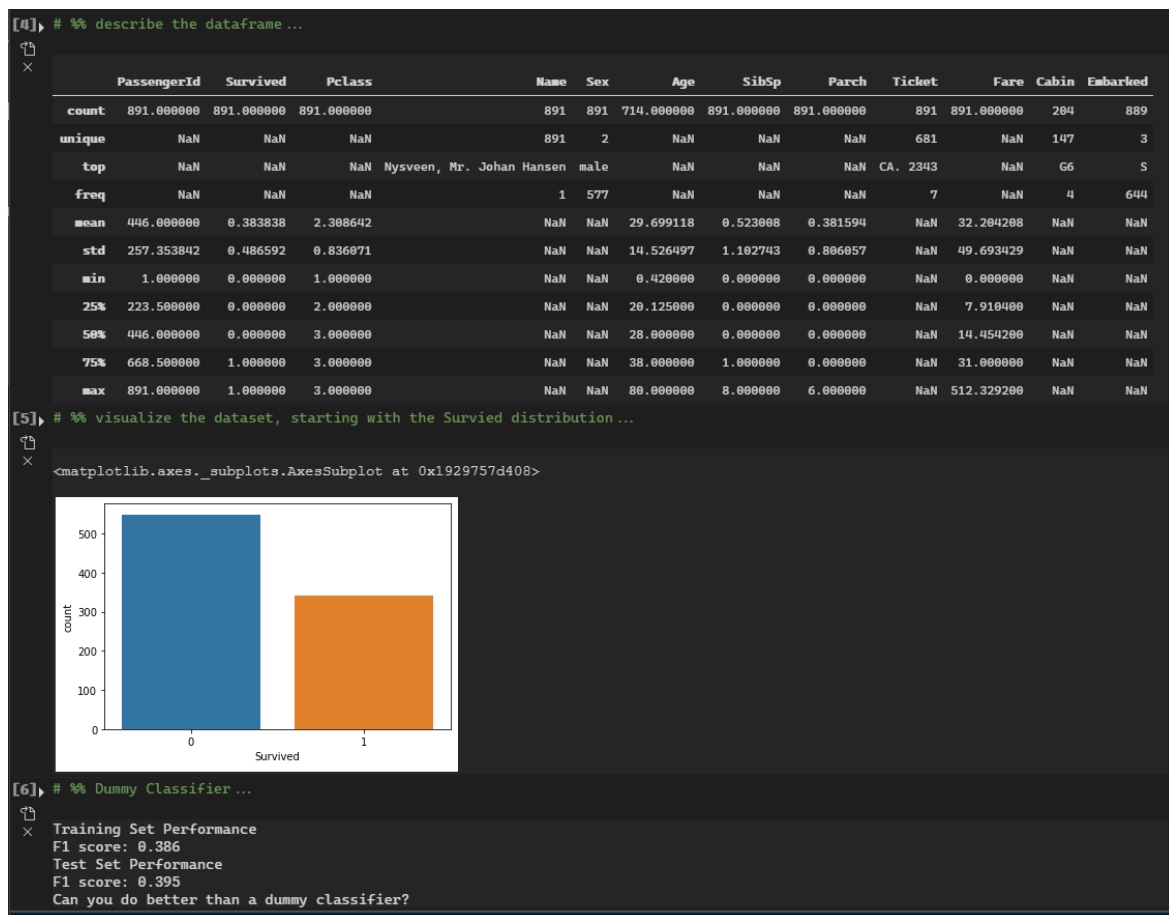    Your result should resemble the screen below.



**Figure 2: Executing the Skeleton File**
Used with permission from Microsoft.

This reads the data files into dataframes, shows some simple statistics, and visualizes the survival distribution. It also provides an evaluation function and a dummy classifier.

However, the dummy classifier only achieves a score of approximately 40% F1 on the test data.

2. Use your data mining skills to explore and understand the dataset. Some sample questions are provided in the skeleton file.

   For example, if you choose to answer **Survived w.r.t Age distribution**, your result should resemble the screen below.
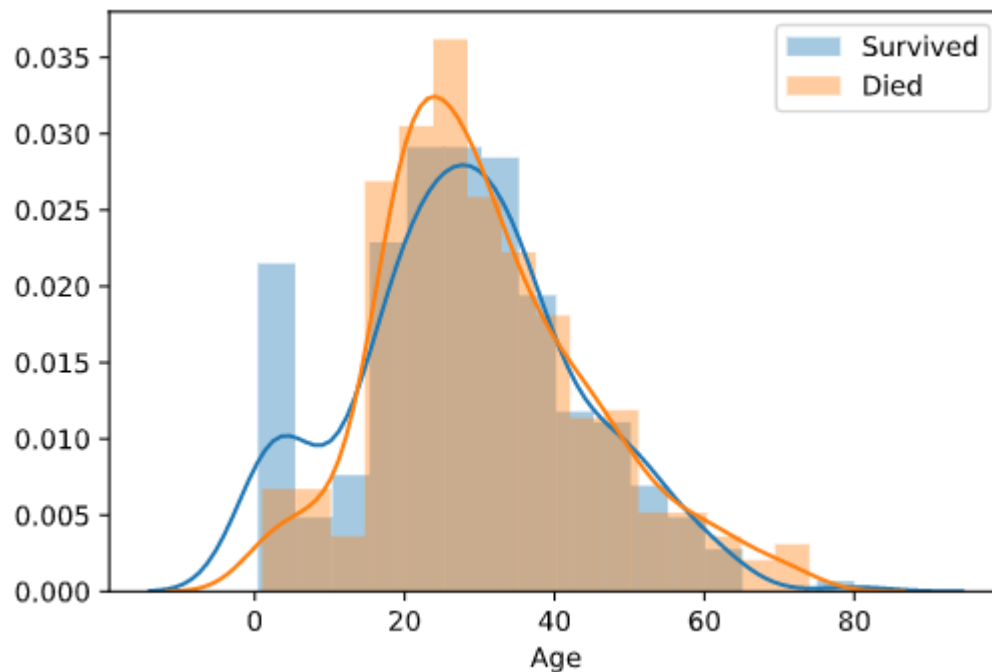


**Figure 3: Age Distribution Survival Rate**
© 2020, Southern Alberta Institute of Technology

3. In the last code block, provide your solution to the classification problem.

   Your goal is to provide a better solution than the dummy classifier. In other words, your F1 score on the test set should be higher than 0.395.

# Lab Activity 2: Predicting House Prices

In this section, you will work with the Ames Housing dataset. You'll use the characteristics of a house and build a predictive model to forecast its sale price.

**Note:** Although you have already downloaded this dataset, it is also publicly available from [House Prices: Advanced Regression Techniques](https://www.kaggle.com/c/house-prices-advanced-regression-techniques) (https://www.kaggle.com/c/house-prices-advanced-regression-techniques).

1. Use the skeleton file **regression.py** and run through the nonempty code blocks.

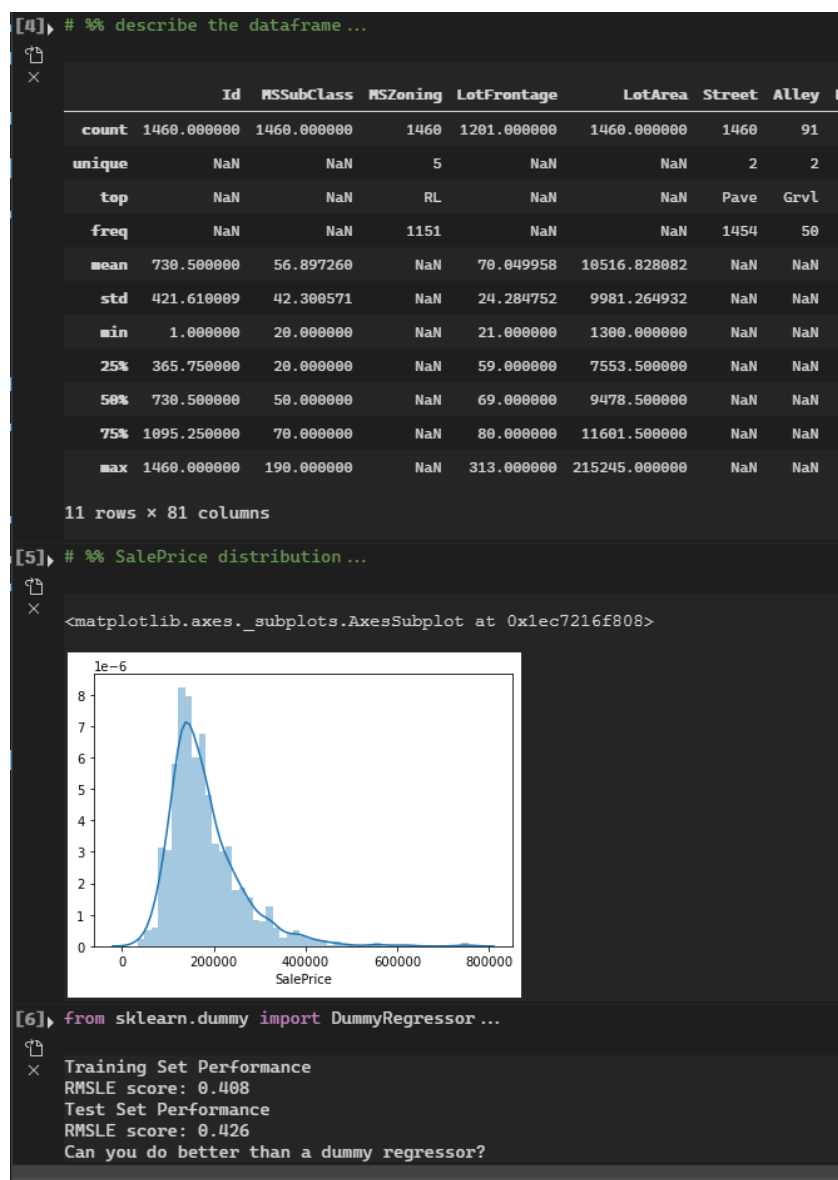   Your result should resemble the screen below.



**Figure 4: Executing the Skeleton File**
Used with permission from Microsoft.

2. Use your data mining skills to explore and understand the dataset. Some sample questions are provided in the skeleton file.

   For example, if you choose to answer **SalePrice distribution w.r.t YearBuilt**, your result should resemble the screen below.
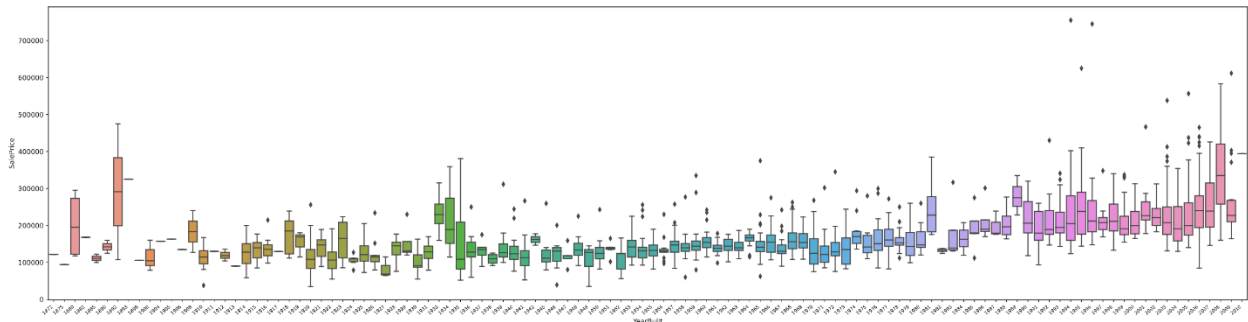


**Figure 5: SalePrice distribution w.r.t YearBuilt**
© 2020, Southern Alberta Institute of Technology

3. Plot as many figures as possible to visually comprehend the relationship between variables and sale prices.

4. In the last code block, provide your solution to the regression problem. You goal is to provide a better solution than the dummy regressor. In other words, your RMSLE score on the test set should be lower than **0.426**.

## Lab Activity 3: Clustering Seeds

In this section, you will work with a wheat seed dataset with three varieties of wheat. You'll use the characteristics of the seed's internal kernel structure to build a predictive model to group the seeds into clusters.

**Note:** Although you have already downloaded this dataset, it is also publicly available from UCI Machine Learning Repository: Seeds Data Set (https://archive.ics.uci.edu/ml/datasets/seeds).

1.  Use the skeleton file **clustering.py** and run through the first three code blocks.

    Your result should resemble the image below.



**Figure 6: Executing Skeleton File**
Used with permission from Microsoft.

2. Use your data mining skills to explore and understand the dataset. Plot as many figures as need to visually comprehend the relationship between different variables.

For example, the relationship between **perimeter** and **compactness** can be visualized as shown in the image below.
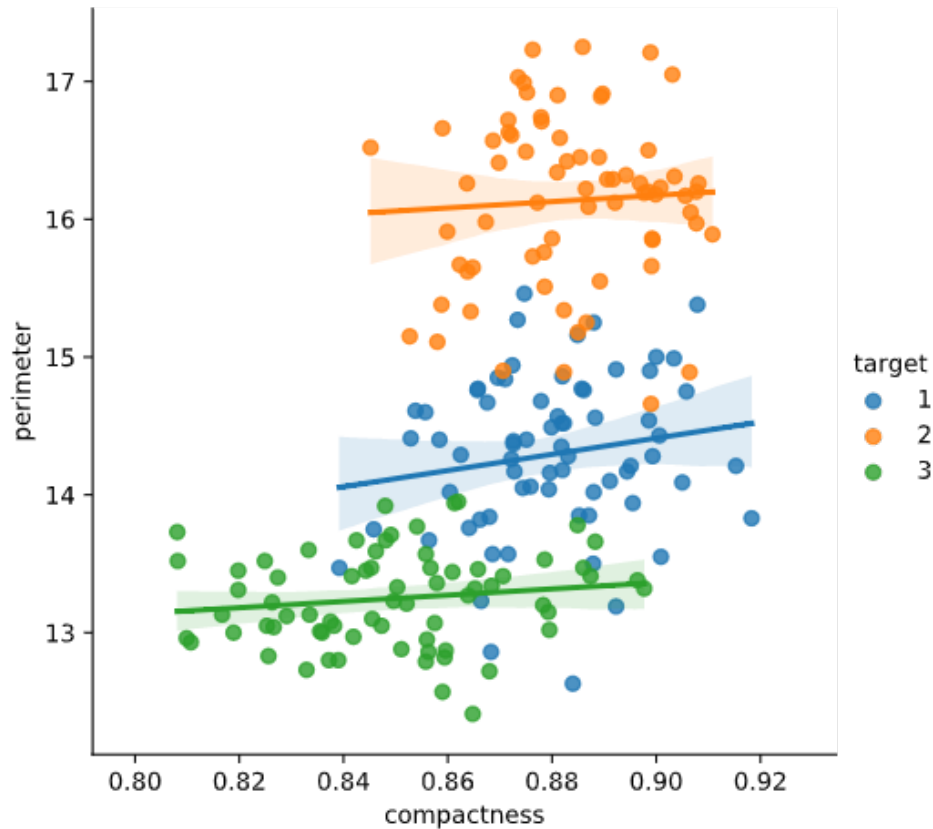


**Figure 7. Exploring Perimeter vs. Compactness**
© 2020, Southern Alberta Institute of Technology

3. Suppose you didn't know the exact number of wheat varieties. Determine the best number of clusters by iterating a list of candidate numbers and finding the "elbow" point for the change in inertia and/or homogeneity.

4. Complete the second-last code block, and then run the final code block.

   a. Since you know that the actual number of clusters is 3, use a list from 1 to 10 as the candidate.

   b. In each iteration, fit the data to KMeans with a different candidate number of clusters.

   c. Get the inertia score and homogeneity score and store them in a dictionary.
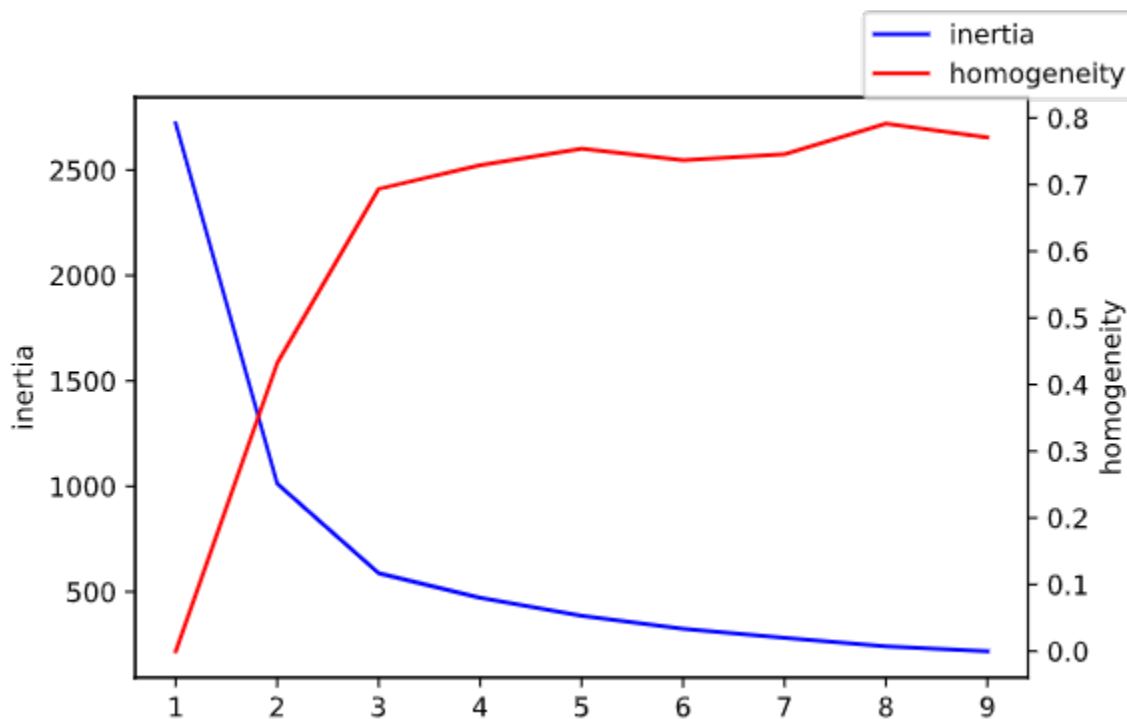
Your results should be similar to the image below.



**Figure 8. Determining Best Number of Clusters**
© 2020, Southern Alberta Institute of Technology

# References

Charytanowicz, M., Niewczas, J., Kulczycki, P., Kowalski, P. A., Lukasik, S., & Zak, S. (2010). A complete gradient clustering algorithm for features analysis of X-ray images. In Pietka, E. and Kawa, J. (Eds.), *Information technologies in biomedicine*, (pp. 15–24). Berlin: Springer-Verlag.

De Cock, D. (2011). Ames, Iowa: Alternative to the Boston housing data as an end of semester regression project. *Journal of Statistics Education, 19*(3). doi: 10.1080/10691898.2011.11889627