

Assignment 3

Welcome to assignment #3! In this assignment, we will continue with exploratory data analysis (EDA) using pandas, examine correlations, and conduct a linear regression.

We've been working with pandas and python for a while now. This assignment is intended to be relatively straightforward. However, it is always possible that you run into problems. Part of this assignment is about seeing how you deal with problems.

If you do get stuck on a coding problem, you should abstract your problem before you ask the TAs or professor. (For example, we will gladly help you learn how to run a regression, but we will not help run the regression listed in this assignment.) Instead, you should show your problem example using the car data file. (You can download it from <https://busansbu.mycourses.work/data/mtcars.csv>) Make up a small example using that data, to reproduce the error you are getting.

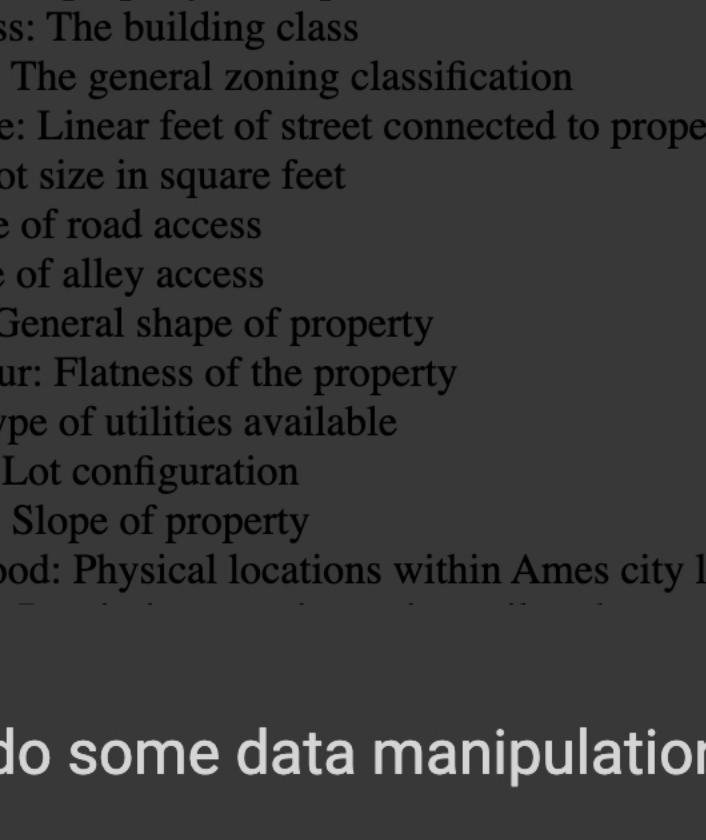
In other words, you can ask the TAs for help on specific technical issues but not for answers to problems. Also, "Is this right?" is not an acceptable use of TA or professor time. At a third year level you should be able to read a question written in plain English without needing excessive clarification.

DO NOT POST YOUR ATTEMPTED ANSWERS TO PIAZZA. This happened so much for Assignment 2 that as an exercise it is virtually meaningless. I don't want this to happen again.

Note: You need to answer all questions in complete sentences. Failure to do so will result in part or zero marks.

Before you start, I'd like you to watch the following video. This video is an excellent review of topics we'll cover in this assignment.

```
from IPython.display import YouTubeVideo
YouTubeVideo('PaFPbb6DxQ')
```



aka Linear Regression aka Least Squares

Let's look again at house price data!

Data dictionary

Here are brief explanations of our variables. You will rarely have a detailed dictionary in any organization. This means that your initial explorations of the data are absolutely crucial so that you can understand your data!

```
[30] from IPython.display import IFrame
IFrame(src="https://busansbu.mycourses.work/data/house_price_data_description_short.html", width=800, height=300)
```

Data fields

- SalePrice - the property's sale price in dollars. This is the target variable that you're trying to predict.
- MSSubClass: The building class
- MSCat: Major house zoning classification
- LotFrontage: Linear feet of street connected to property
- LotArea: Lot size in square feet
- Street: Type of road access
- Alley: Type of alley access
- LotShape: General shape of property
- LandContour: Flatness of the property
- Utilities: Type of utilities available
- LotConfig: Lot configuration
- LandSlope: Slope of property
- Neighborhood: Physical locations within Ames city limits

Download and do some data manipulation

```
[31] import pandas as pd
df = pd.read_csv('https://busansbu.mycourses.work/data/house_price_train.csv') # Use this for the assignment!
#df = pd.read_csv('https://busansbu.mycourses.work/data/mtcars.csv') # Use this to ask a question by making a separate example in another workbook.
```

Questions: EDA, Correlations and Regression

In the questions that follow, we're going to further explore the relationship between SalePrice and other explanatory variables.

In addition to LotFrontage, choose two other variables from the data frame you think are promising explanatory variables. Use these variables in the questions below.

Q1. Is exploratory data analysis the same thing as hypothesis testing. Explain your reasoning. (1 mark)

A1. No, exploratory data analysis(EDA) is an open-ended analysis in the process of transforming, visualising and summarising data while hypothesis testing is used to verify a priori hypotheses whether some observation seems unusual based on data.

Q2. Show a correlation matrix for SalePrice and the explanatory variables you chose.

- (a) Explain in plain English what correlation means.
(b) Plot and describe the variables with the highest correlation.

(2 marks)

A2(a) Correlation is the relation of being correlated that two variables are linearly related, such as height and weight.

```
[32] # A2(b)
correlation = df[['SalePrice', 'LotArea', 'YearBuilt']]
#correlation.head()
```

#Find the highest correlation.

print(correlation.corr())

print() #new line between values and plot

#Plot

from pandas.plotting import scatter_matrix

correlated_plot = scatter_matrix(correlation, alpha = .9, diagonal = 'hist')

print()

print(df)

print(df['SalePrice'].hist())

print(df['LotArea'].hist())

print(df['YearBuilt'].hist())

print(df['SalePrice'].hist())