

Assignment\_2\_data\_descriptions.ipynb\_jhwa426 ☆

파일 수정 보기 삽입 도구 도움말 모든 변경사항이 저장됨

+ 코드 + 텍스트

RAM 24.0GB / 24.0GB

수정 가능

## - Assignment 2

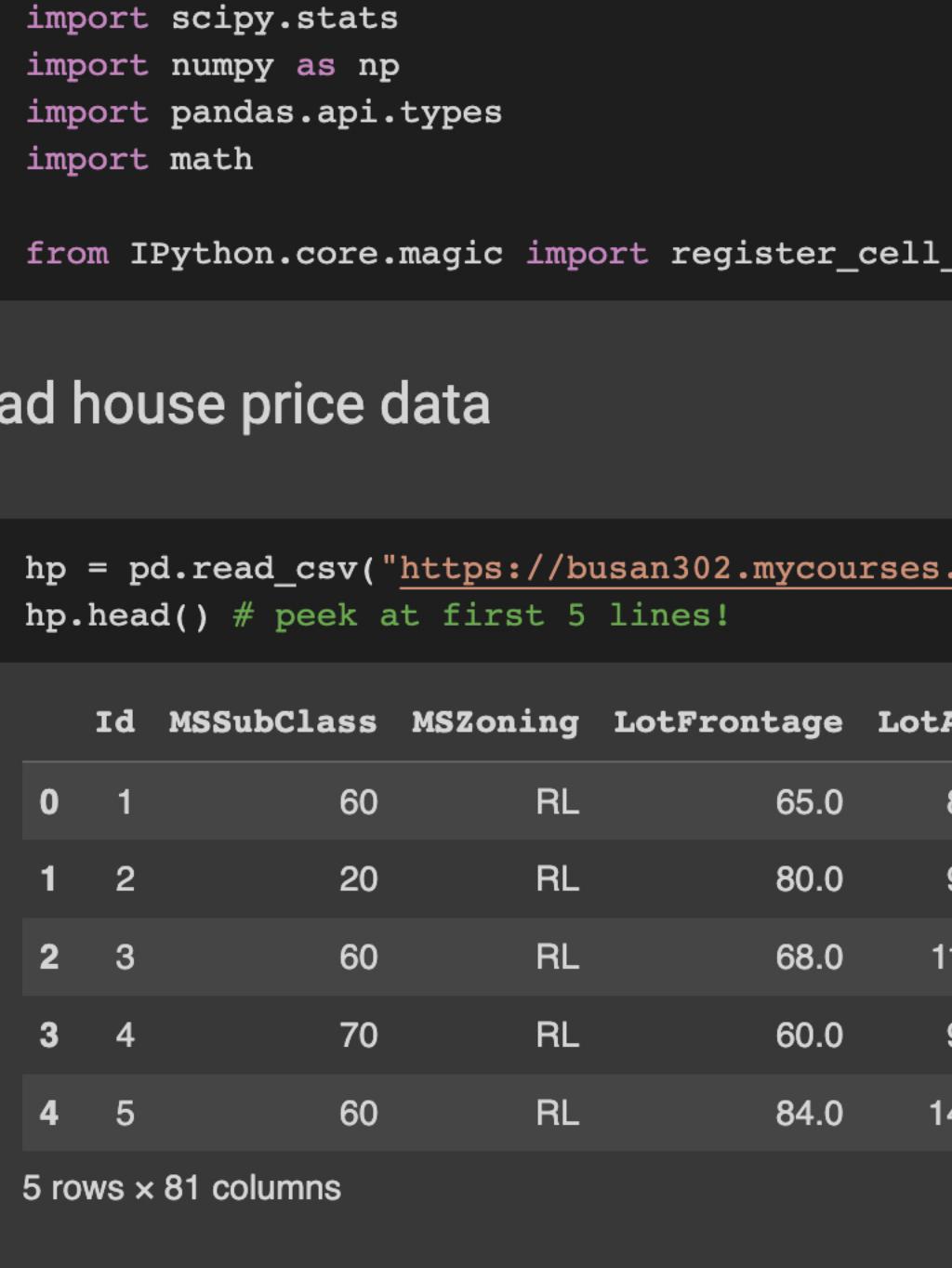
Welcome to assignment #2! In this assignment we're going to be doing some exploratory data analysis (EDA) on a dataset using pandas.

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

Make sure that you can save this notebook to your own Google Drive space! Otherwise, your answers will not be saved!

Before you start, I'd like you to watch the following video. This is a good first introduction and review of some of the topics we'll be covering in this course.

[22] from IPython.display import YouTubeVideo  
YouTubeVideo("11og6fjPax0")



[23] # First, load our libraries  
import pandas as pd  
import scipy.stats  
import numpy as np  
import pandas.api.types  
import math  
from IPython.core.magic import register\_cell\_magic

## - Load house price data

[24] hp = pd.read\_csv("https://busan302.mycourses.work/data/house\_price\_train.csv") #1460 rows  
hp.head() # peek at first 5 lines!

Id	MSSubClass	MSZoning	LotFrontage	LotArea	Street	Alley	LotShape	LandContour	Utilities	PoolArea	PoolQC	MiscFeature	MiscVal	NoSolds	YrSold	SaleType	SaleCondition	SalePrice		
0	1	60	RL	65.0	8450	Pave	NaN	Reg	Lvl	AllPub	...	0	Nan	Nan	0	2	2008	WD	Normal	208500
1	2	20	RL	80.0	9600	Pave	NaN	Reg	Lvl	AllPub	...	0	Nan	Nan	0	5	2007	WD	Normal	181500
2	3	60	RL	68.0	11250	Pave	NaN	IR1	Lvl	AllPub	...	0	Nan	Nan	0	9	2008	WD	Normal	223500
3	4	70	RL	60.0	9550	Pave	NaN	IR1	Lvl	AllPub	...	0	Nan	Nan	0	2	2006	WD	Abnorml	140000
4	5	60	RL	84.0	14260	Pave	NaN	IR1	Lvl	AllPub	...	0	Nan	Nan	0	12	2008	WD	Normal	250000

5 rows × 81 columns

## - 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!

[25] from IPython.display import IFrame  
IFrame(src="https://busan302.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
- MSZoning: The general zoning classification
- LotFrontage: Linear feet of street connected to property
- LotArea: Lot size in square foot
- Street: Type of road access
- Alley: Type of alley access
- LotConfig: Configuration 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

## - Exploratory data analysis

Let's see how pre-1945 houses differ from houses built post-1980s.

Q1. Complete the code in the cells below. (2 marks total)

Replace the ellipses '...' with appropriate goals to meet the defined goals.

[26] # Create subsets of old (pre-1945) and new houses (1980 onwards)  
#Old homes  
old\_homes = hp[hp["YearBuilt"] < 1945]  
  
#New homes  
new\_homes = hp[hp["YearBuilt"] >= 1980]

[27] # Find the means of the numeric variables in the new\_homes and old\_homes dataframes.  
#The means of New homes  
  
nh = np.mean(new\_homes)  
#The means of Old homes  
  
oh = np.mean(old\_homes)

[28] # Consider the data from this table  
df = pd.concat([oh, nh, axis=1]) # Take the Series created from the previous question. Concat will place them side by side in a dataframe  
df.columns = ['old\_homes', 'new\_homes'] # Give the columns descriptive names  
df # jupyter prints out the last line for us

	old_homes	new_homes
Id	739.682482	717.625949
MSSubClass	61.476102	60.751634
LotFrontage	62.796000	72.641326
LotArea	8977.890511	10258.459150
OverallQual	5.368613	7.163399
OverallCond	6.215320	5.116013
YearBuilt	1922.470803	2000.171569
YearRemodAdd	1971.463504	2000.928105
MasVnrArea	12.176832	157.677686
BsmtFinSF1	188.273723	626.539498
BsmtFinSF2	18.498350	22.266340
BsmtUnfSF	594.281022	694.290850
TotalBsmtSF	801.051098	1243.093137
1stFlrSF	953.708029	1268.500000
2ndFlrSF	479.817518	440.457516
LowQualFinSF	27.478102	0.392157
GrLivArea	1460.003650	1709.349873
BsmtFullBath	0.200730	0.508170
BsmtHalfBath	0.036466	0.027778
FullBath	1.299270	1.911765
HalfBath	0.222628	0.517974
BathroomAbvGr	2.930657	2.802288
KitchenAbvGr	1.087591	1.004902
TotRmsAbvGrd	6.554745	6.882745
Fireplaces	0.503650	0.697712
GarageCars	1.226277	2.233660
GarageArea	326.562044	586.947712
WoodDeckSF	50.529197	126.401961
OpenPorchSF	34.889781	65.282680
EnclosedPorch	69.726277	5.691176

A2. Write your answers in this cell

(a) Generally, new homes' floor space has bigger space than older homes.

(b) The most significant difference between newer and older homes is the average sale price of new homes are approximately \$100,000 higher than old homes and newer homes show that there is a high number of trading.

(c) The surprise is that the average floor space of the new houses is much bigger than the old ones, which means the new houses are supposed to be built with large bedrooms than other area. Also, there is no difference in annual home sales and monthly home sales between newer and older homes.

(d) The differences show that the construction trend has been developed. The floor space of the old homes is built with plenty of small bedrooms, however, newer homes tend to be built on a large scale of bedrooms.

Q3. The presenter in the video talked about two types of data to conduct a t-test: categorical and numeric variables (2 marks).

(a) What is the null hypothesis of the t-test?

(b) Give an example of the categorical and numeric variables we are considering when we conduct the t-test in the code below.

A3. Place your answers below.

(a) The null hypothesis is the p-value equivalent to 0.

(b) Example variables are:

- categorical: Age group, Gender
- numeric: Weight, Height

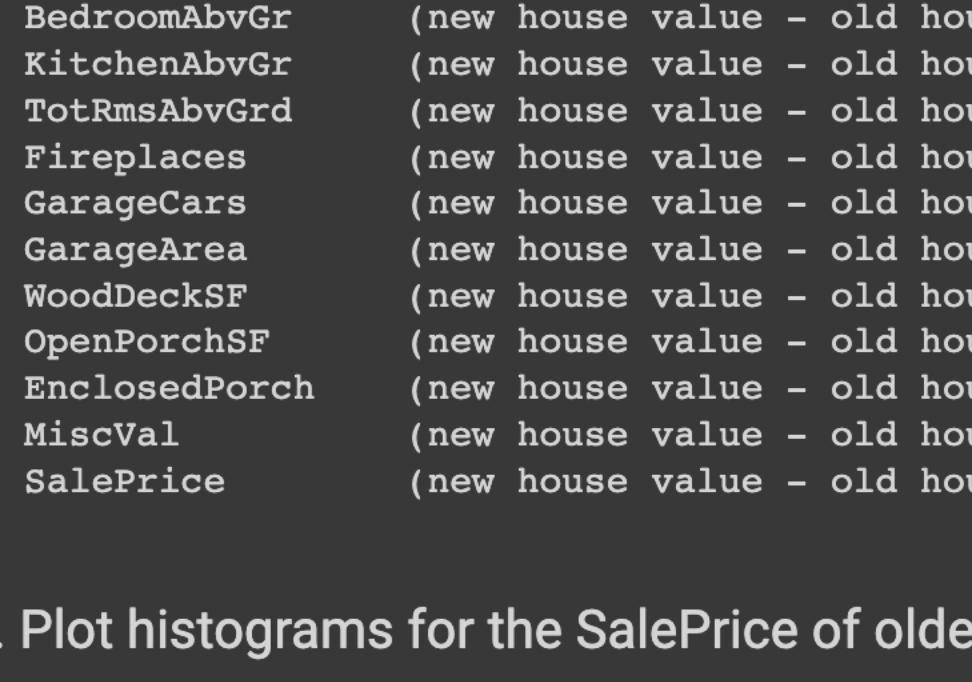
[29] for var in old\_homes.columns:  
 # Next line checks that the dataframe column is numeric.  
 if not pandas.api.types.is\_numeric\_dtype(old\_homes[var].dtype):  
 continue # skip all non-numeric columns  
  
 # t-test: do a t-test to compare means for the column in the old and new house dataframes  
 t = scipy.stats.ttest\_ind(old\_homes[var], new\_homes[var])  
  
 if t[1] < 0.05: # Only print cases where the p-value is less than 5%. We'll say for our purposes this is significant.  
 # The \*\* syntax is a little cryptic if you're new to python!  
 print(f"\*\*{var}\*\* ({new\_homes[var].mean() - old\_homes[var].mean():.2f}) \t p-value={round(t[1],4)}")

LotArea	(new house value - old house value)	p-value
OverallQual	1.280.57	0.0009
OverallCond	-1.17	0.0001
YearBuilt	-7.7	0.0001
YearRemodAdd	-2.46	0.0001
DensUnifSF	336.26	0.0001
BsmtUnfSF	100.01	0.0001
TotalBsmtSF	442.04	0.0001
1stFlrSF	315.79	0.0001
2ndFlrSF	249.55	0.0001
LowQualFinSF	0.31	0.0001
GrLivArea	0.6	0.0001
BsmtFullBath	0.3	0.0001
BathroomAbvGr	0.13	0.0001
KitchenAbvGr	-0.08	0.0001
TotRmsAbvGrd	0.13	0.0001
Fireplaces	0.13	0.0001
GarageCars	0.13	0.0001
GarageArea	260.39	0.0001
WoodDeckSF	75.87	0.0001
OpenPorchSF	30.59	0.0001
EnclosedPorch	-1.75	0.0001
MiscVal	100085.04	0.0001

Q4. Plot histograms for the SalePrice of older homes and newer homes. Describe the shape of each histogram (3 marks).

You can find out more about [how to plot histograms in pandas here](#).

[30] # older homes histogram  
oh\_hist = old\_homes.hist(column="SalePrice")



[31] # newer homes histogram  
nh\_hist = new\_homes.hist(column="SalePrice")

