

Homework 5: Predicting Housing Prices

Due Date: 11:59pm Friday, July 26

Collaboration Policy

Data science is a collaborative activity. While you may talk with others about the homework, we ask that you **write your solutions individually**. If you do discuss the assignments with others please **include their names** in the collaborators cell below.

Collaborators: *list names here*

Introduction

In this homework, we will go through the iterative process of specifying, fitting, and analyzing the performance of a model.

In the first portion of the assignment, we will guide you through some basic exploratory data analysis (EDA), laying out the thought process that leads to certain modeling decisions. Next, you will add a new feature to the dataset, before specifying and fitting a linear model to a few features of the housing data to predict housing prices. Finally, we will analyze the error of the model and brainstorm ways to improve the model's performance.

After this homework, you should feel comfortable with the following:

- 1. Simple feature engineering
- 2. Using sklearn to build linear models
- 3. Building a data pipeline using pandas

Next week's homework will continue working with this dataset to address more advanced and subtle issues with modeling.

Score Breakdown

Question	Points
<u>Question 1</u>	3
<u>Question 2</u>	2
<u>Question 3</u>	1
<u>Question 4</u>	1
<u>Question 5</u>	2
<u>Question 6</u>	2
<u>Question 7a</u>	1
<u>Question 7b</u>	2
<u>Question 8a</u>	1
<u>Question 8b</u>	1
<u>Question 8c</u>	2
<u>Question 8d</u>	2
Total	20

```
import numpy as np

import pandas as pd
from pandas.api.types import CategoricalDtype

%matplotlib inline
import matplotlib.pyplot as plt
import seaborn as sns

# Plot settings
plt.rcParams['figure.figsize'] = (12, 9)
plt.rcParams['font.size'] = 12
```

The Ames dataset (<http://jse.amstat.org/v19n3/decock.pdf>) consists of 2930 records taken from the Ames, Iowa, Assessor's Office describing houses sold in Ames from 2006 to 2010. The data set has 23 nominal, 23 ordinal, 14 discrete, and 20 continuous variables (and 2 additional observation identifiers) --- 82 features in total. An explanation of each variable can be found in the included `codebook.txt` file. The information was used in computing assessed values for individual residential properties sold in Ames, Iowa from 2006 to 2010. **Some noise has been added to the actual sale price, so prices will not match official records.**

In [2]:

As a good sanity check, we should at least verify that the data shape matches the description.

[illegible]

The next order of business is getting a feel for the variables in our data. The Ames data set contains information that typical homebuyers would want to know. A more detailed description of each variable is included in `codebook.txt`. **You should take some time to familiarize yourself with the codebook before moving forward.**

In [4]:

```
training_data.columns.values
```

Out[4]:

```
array(['Order', 'PID', 'MS_SubClass', 'MS_Zoning', 'Lot_Frontage',  
      'Lot_Area', 'Street', 'Alley', 'Lot_Shape', 'Land_Contour',  
      'Utilities', 'Lot_Config', 'Land_Slope', 'Neighborhood',  
      'Condition_1', 'Condition_2', 'Bldg_Type', 'House_Style',  
      'Overall_Qual', 'Overall_Cond', 'Year_Built', 'Year_Remod/Add',  
,  
      'Roof_Style', 'Roof_Matl', 'Exterior_1st', 'Exterior_2nd',  
      'Mas_Vnr_Type', 'Mas_Vnr_Area', 'Exter_Qual', 'Exter_Cond',  
      'Foundation', 'Bsmt_Qual', 'Bsmt_Cond', 'Bsmt_Exposure',  
      'BsmtFin_Type_1', 'BsmtFin_SF_1', 'BsmtFin_Type_2', 'BsmtFin_SF_2',  
      'Bsmt_Unf_SF', 'Total_Bsmt_SF', 'Heating', 'Heating_QC',  
      'Central_Air', 'Electrical', '1st_Flr_SF', '2nd_Flr_SF',  
      'Low_Qual_Fin_SF', 'Gr_Liv_Area', 'Bsmt_Full_Bath',  
      'Bsmt_Half_Bath', 'Full_Bath', 'Half_Bath', 'Bedroom_AbvGr',  
      'Kitchen_AbvGr', 'Kitchen_Qual', 'TotRms_AbvGrd', 'Functional',  
,  
      'Fireplaces', 'Fireplace_Qu', 'Garage_Type', 'Garage_Yr_Blt',  
      'Garage_Finish', 'Garage_Cars', 'Garage_Area', 'Garage_Qual',  
      'Garage_Cond', 'Paved_Drive', 'Wood_Deck_SF', 'Open_Porch_SF',  
,  
      'Enclosed_Porch', '3Ssn_Porch', 'Screen_Porch', 'Pool_Area',  
      'Pool_QC', 'Fence', 'Misc_Feature', 'Misc_Val', 'Mo_Sold',  
      'Yr_Sold', 'Sale_Type', 'Sale_Condition', 'SalePrice'],  
      dtype=object)
```

Part 1: Exploratory Data Analysis

In this section, we will make a series of exploratory visualizations and interpret them.

Note that we will perform EDA on the **training data** so that information from the test data does not influence our modeling decisions.

Sale Price

We begin by examining a raincloud plot (https://micahallen.org/2018/03/15/introducing-raincloud-plots/amp/?__twitter_impression=true) (a combination of a KDE, a histogram, a strip plot, and a box plot) of our target variable `SalePrice`. At the same time, we also take a look at some descriptive statistics of this variable.

In [5]:

```
fig, axs = plt.subplots(nrows=2)

sns.distplot(
    training_data['SalePrice'],
    ax=axs[0]
)
sns.stripplot(
    training_data['SalePrice'],
    jitter=0.4,
    size=3,
    ax=axs[1],
    alpha=0.3
)
sns.boxplot(
    training_data['SalePrice'],
    width=0.3,
    ax=axs[1],
    showfliers=False,
)

# Align axes
spacer = np.max(training_data['SalePrice']) * 0.05
xmin = np.min(training_data['SalePrice']) - spacer
xmax = np.max(training_data['SalePrice']) + spacer
axs[0].set_xlim((xmin, xmax))
axs[1].set_xlim((xmin, xmax))

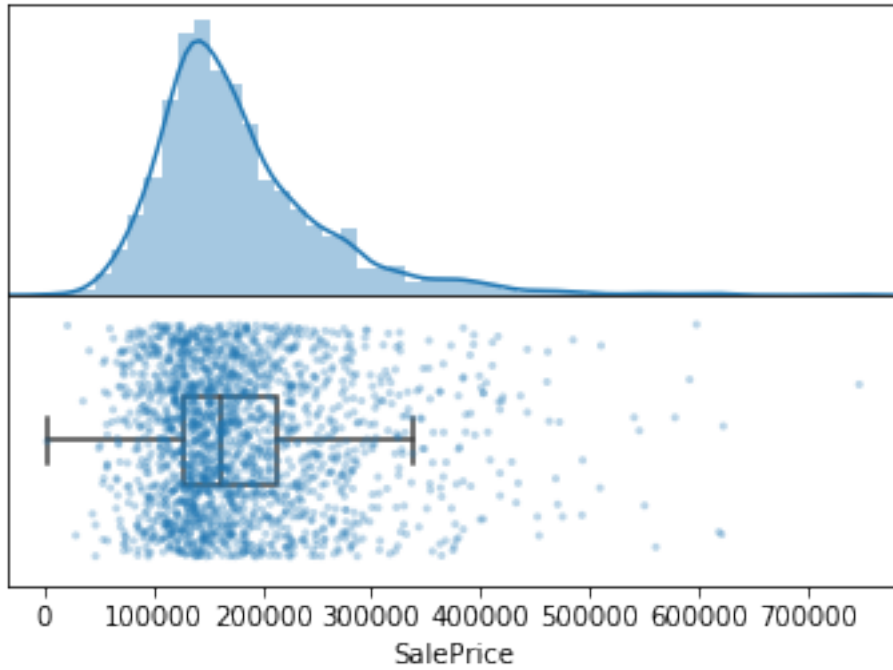
# Remove some axis text
axs[0].xaxis.set_visible(False)
axs[0].yaxis.set_visible(False)
axs[1].yaxis.set_visible(False)

# Put the two plots together
plt.subplots_adjust(hspace=0)

# Adjust boxplot fill to be white
axs[1].artists[0].set_facecolor('white')
```

```
C:\Users\Stephanie Djajadi\Anaconda3\lib\site-packages\matplotlib\axes\_axes.py:6462: UserWarning: The 'normed' kwarg is deprecated, and has been replaced by the 'density' kwarg.
```

```
warnings.warn("The 'normed' kwarg is deprecated, and has been "
```



```
In [6]:
```

```
training_data['SalePrice'].describe()
```

```
Out[6]:
```

```
count      2000.000000
mean      180775.897500
std       81581.671741
min        2489.000000
25%      128600.000000
50%      162000.000000
75%      213125.000000
max       747800.000000
Name: SalePrice, dtype: float64
```

Question 1

To check your understanding of the graph and summary statistics above, answer the following True or False questions:

1. The distribution of SalePrice in the training set is left-skew.
2. The mean of SalePrice in the training set is greater than the median.
3. At least 25% of the houses in the training set sold for more than \$200,000.00.

The provided tests for this question do not confirm that you have answered correctly; only that you have assigned each variable to True or False.

```
BEGIN QUESTION
```

```
name: q1
```

```
points: 3
```

In [7]:

```
# These should be True or False
q1statement1 = False # SOLUTION
q1statement2 = True # SOLUTION
q1statement3 = True # SOLUTION
```

In [8]:

```
# TEST
set([q1statement1, q1statement2, q1statement3]).issubset({False, True})
```

Out[8]:

True

In [9]:

```
# HIDDEN TEST
q1statement1
```

Out[9]:

False

In [10]:

```
# HIDDEN TEST  
q1statement1
```

Out[10]:

False

In [11]:

```
# HIDDEN TEST  
q1statement3
```

Out[11]:

True

SalePrice vs Gr_Liv_Area

Next, we visualize the association between SalePrice and Gr_Liv_Area. The codebook.txt file tells us that Gr_Liv_Area measures "above grade (ground) living area square feet."

This variable represents the square footage of the house excluding anything underground. Some additional research (into real estate conventions) reveals that this value also excludes the garage space.

In [12]:

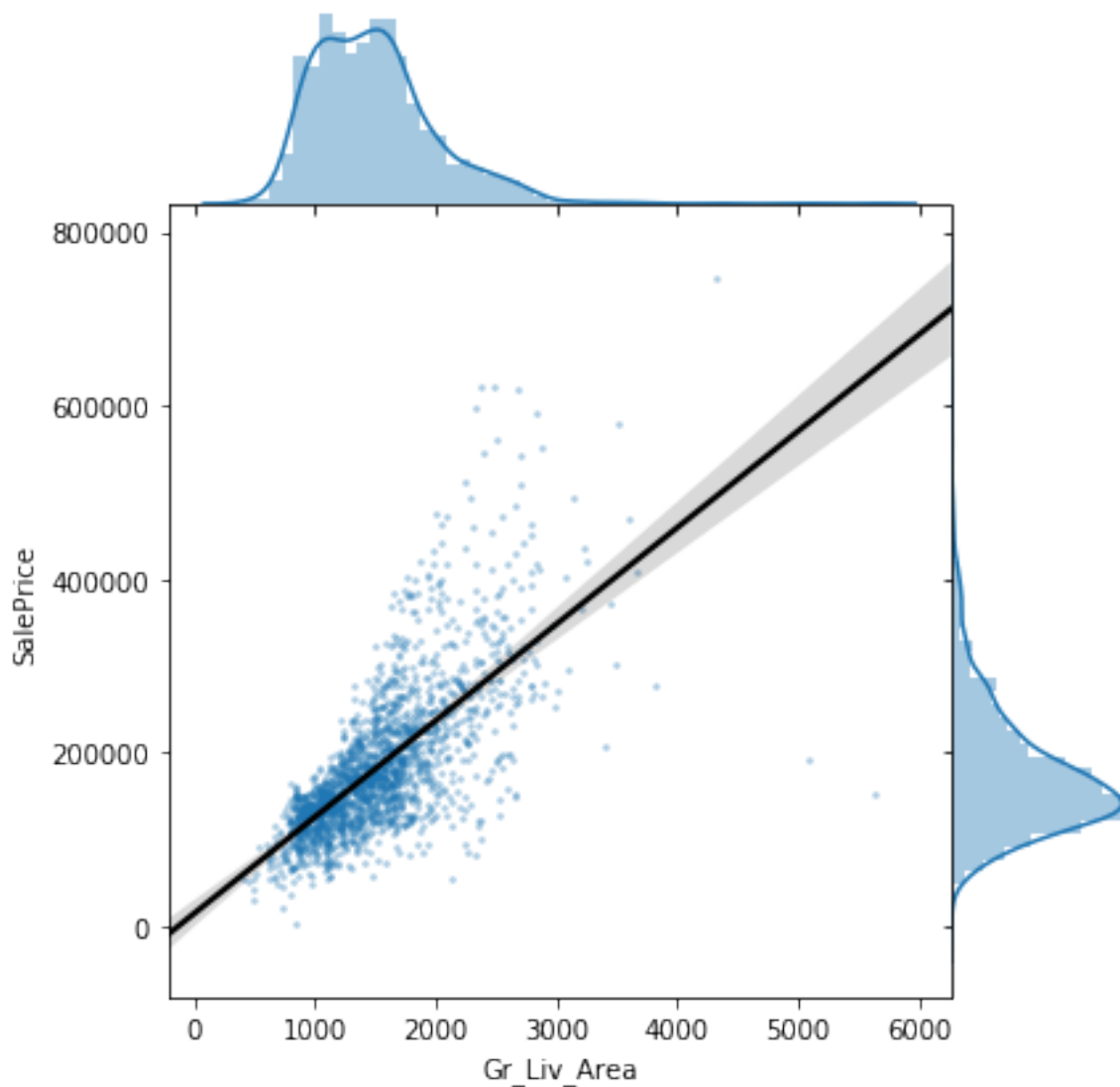
```
sns.jointplot(  
    x='Gr_Liv_Area',  
    y='SalePrice',  
    data=training_data,  
    stat_func=None,  
    kind="reg",  
    ratio=4,  
    space=0,  
    scatter_kws={  
        's': 3,  
        'alpha': 0.25  
    },  
    line_kws={  
        'color': 'black'  
    }  
);
```

```
C:\Users\Stephanie Djajadi\Anaconda3\lib\site-packages\matplotlib\axes\_axes.py:6462: UserWarning: The 'normed' kwarg is deprecated, and has been replaced by the 'density' kwarg.
```

```
warnings.warn("The 'normed' kwarg is deprecated, and has been "
```

```
C:\Users\Stephanie Djajadi\Anaconda3\lib\site-packages\matplotlib\axes\_axes.py:6462: UserWarning: The 'normed' kwarg is deprecated, and has been replaced by the 'density' kwarg.
```

```
warnings.warn("The 'normed' kwarg is deprecated, and has been "
```



There's certainly an association, and perhaps it's linear, but the spread is wider at larger values of both variables. Also, there are two particularly suspicious houses above 5000 square feet that look too inexpensive for their size.

Question 2

What are the Parcel Identification Numbers for the two houses with Gr_Liv_Area greater than 5000 sqft?

The provided tests for this question do not confirm that you have answered correctly; only that you have assigned q2house1 and q2house2 to two integers that are in the range of PID values.

```
BEGIN QUESTION
```

```
name: q2
```

```
points: 2
```

In [14]:

```
# Hint: You can answer this question in one line  
q2house1, q2house2 = training_data.loc[training_data['Gr_Liv_Area'] > 5000, 'PID'  
' ] # SOLUTION
```

In [15]:

```
# TEST  
isinstance(q2house1, (int, np.integer)) # Make sure your answer is integer-value  
d
```

Out[15]:

True

In [16]:

```
# TEST  
isinstance(q2house2, (int, np.integer)) # Make sure your answer is integer-value  
d
```

Out[16]:

True

In [17]:

```
# TEST  
min(training_data['PID']) <= min(q2house1, q2house2)
```

Out[17]:

True

In [18]:

```
# TEST  
max(training_data['PID']) >= max(q2house1, q2house2)
```

Out[18]:

True

In [19]:

```
# HIDDEN TEST  
sorted(set(training_data.loc[training_data['Gr_Liv_Area'] > 5000, 'PID']))
```

Out[19]:

[908154195, 908154235]

Question 3

The codebook tells us how to manually inspect the houses using an online database called Beacon. These two houses are true outliers in this data set: they aren't the same time of entity as the rest. They were partial sales, priced far below market value. If you would like to inspect the valuations, follow the directions at the bottom of the codebook to access Beacon and look up houses by PID.

For this assignment, we will remove these outliers from the data. Write a function `remove_outliers` that removes outliers from a data set based off a threshold value of a variable. For example, `remove_outliers(training_data, 'Gr_Liv_Area', upper=5000)` should return a data frame with only observations that satisfy `Gr_Liv_Area` less than or equal to 5000.

The provided tests check that training_data was updated correctly, so that future analyses are not corrupted by a mistake. However, the provided tests do not check that you have implemented remove_outliers correctly so that it works with any data, variable, lower, and upper bound.

BEGIN QUESTION

name: q3

points: 1

In [20]:

```
def remove_outliers(data, variable, lower=-np.inf, upper=np.inf):
    """
    Input:
        data (data frame): the table to be filtered
        variable (string): the column with numerical outliers
        lower (numeric): observations with values lower than this will be removed
        upper (numeric): observations with values higher than this will be removed

    Output:
        a winsorized data frame with outliers removed

    Note: This function should not change mutate the contents of data.
    """
    # BEGIN SOLUTION
    return data.loc[(data[variable] > lower) & (data[variable] < upper), :]
    # END SOLUTION

training_data = remove_outliers(training_data, 'Gr_Liv_Area', upper=5000)
```

In [21]:

```
# TEST
training_data.shape[0] # Make sure that two observations were removed
```

Out[21]:

1998

In [22]:

```
# TEST  
# Make sure that the max Gr_Liv_Area is now below 5000  
max(training_data['Gr_Liv_Area']) < 5000
```

Out[22]:

True

In [23]:

```
# TEST  
# Make sure that remove_outliers doesn't mutate its input  
remove_outliers(training_data, 'Gr_Liv_Area', upper=2000).shape == training_data  
.shape
```

Out[23]:

False

In [24]:

```
# TEST  
# Make sure that the sum of Gr_Liv_Area values hasn't been altered somehow  
sum(training_data['Gr_Liv_Area'])
```

Out[24]:

2980752

In [25]:

```
# HIDDEN TEST  
# Make sure the function works for other upper/lower values  
sum(remove_outliers(training_data, 'SalePrice', lower=80000, upper=82000)['Gr_Li  
v_Area'])
```

Out[25]:

5245

Part 2: Feature Engineering

In this section we will create a new feature out of existing ones through a simple data transformation.

Bathrooms

Let's create a groundbreaking new feature. Due to recent advances in Universal WC Enumeration Theory, we now know that Total Bathrooms can be calculated as:

$$\text{TotalBathrooms} = (\text{BsmtFullBath} + \text{FullBath}) + \frac{1}{2}(\text{BsmtHalfBath} + \text{HalfBath})$$

The actual proof is beyond the scope of this class, but we will use the result in our model.

Question 4

Write a function `add_total_bathrooms(data)` that returns a copy of `data` with an additional column called `TotalBathrooms` computed by the formula above. **Treat missing values as zeros.** Remember that you can make use of vectorized code here; you shouldn't need any `for` statements.

The provided tests check that you answered correctly, so that future analyses are not corrupted by a mistake.

BEGIN QUESTION

name: q4

points: 1

In [26]:

```
def add_total_bathrooms(data):
    """
    Input:
        data (data frame): a data frame containing at least 4 numeric columns
                           Bsmt_Full_Bath, Full_Bath, Bsmt_Half_Bath, and Half_Bath
    """
    with_bathrooms = data.copy()
    bath_vars = ['Bsmt_Full_Bath', 'Full_Bath', 'Bsmt_Half_Bath', 'Half_Bath']
    weights = pd.Series([1, 1, 0.5, 0.5], index=bath_vars)
    # BEGIN SOLUTION
    with_bathrooms = with_bathrooms.fillna({var: 0 for var in bath_vars})
    with_bathrooms['TotalBathrooms'] = with_bathrooms[bath_vars].dot(weights)
    # END SOLUTION
    return with_bathrooms

training_data = add_total_bathrooms(training_data)
```

In [27]:

```
# TEST
not training_data['TotalBathrooms'].isnull().any() # Check that missing values are dealt with
```

Out[27]:

True

In [28]:

```
# TEST
training_data['TotalBathrooms'].sum() # Check that the values are as expected
```

Out[28]:

4421.5

In [29]:

```
# TEST
list(training_data.loc[training_data['PID'] == 903230120, 'TotalBathrooms'])
```

Out[29]:

[1.0]

In [30]:

```
# TEST
sum(training_data.loc[:, 'TotalBathrooms'] < 2.5)
```

Out[30]:

1124

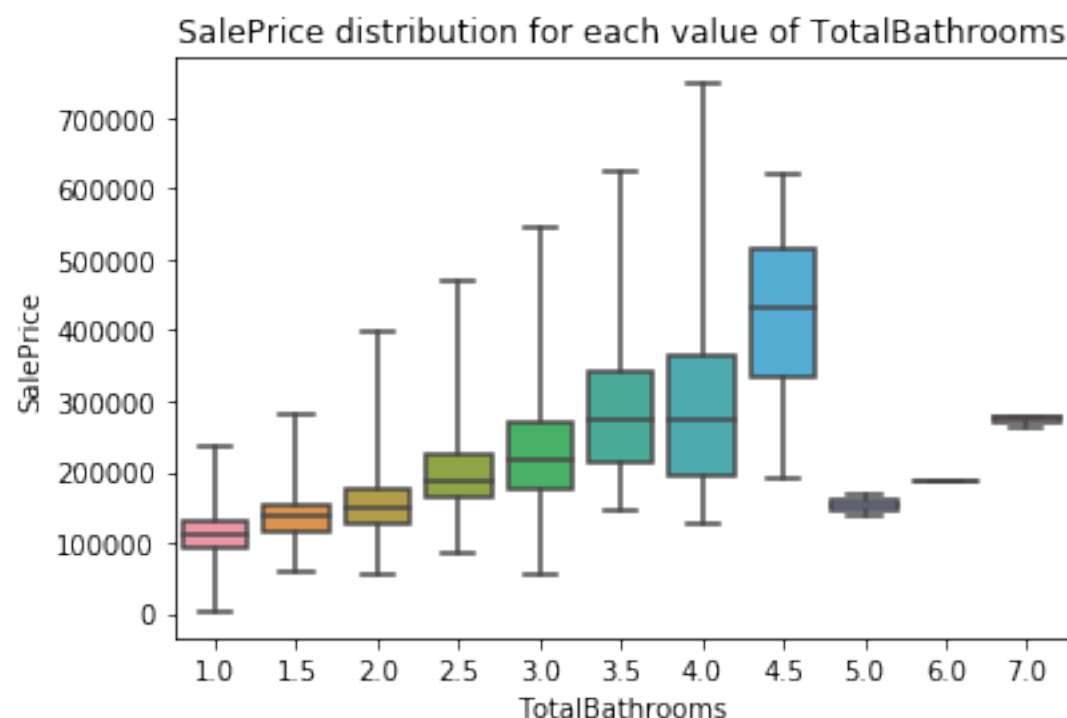
Question 5

Create a visualization that clearly and succinctly shows that TotalBathrooms is associated with SalePrice. Your visualization should avoid overplotting.

```
BEGIN QUESTION
name: q5
points: 2
manual: True
format: image
```

In [31]:

```
# BEGIN SOLUTION
sns.boxplot(x='TotalBathrooms', y='SalePrice', data=training_data, whis=5);
plt.title('SalePrice distribution for each value of TotalBathrooms');
# END SOLUTION
```



Part 3: Modeling

We've reached the point where we can specify a model. But first, we will load a fresh copy of the data, just in case our code above produced any undesired side-effects. Run the cell below to store a fresh copy of the data from `ames_train.csv` in a dataframe named `full_data`. We will also store the number of rows in `full_data` in the variable `full_data_len`.

In [32]:

```
# Load a fresh copy of the data and get its length
full_data = pd.read_csv("ames_train.csv")
full_data_len = len(full_data)
full_data.head()
```

Out[32]:

	Order	PID	MS_SubClass	MS_Zoning	Lot_Frontage	Lot_Area	Street	Alley
0	1	526301100	20	RL	141.0	31770	Pave	NaN
1	2	526350040	20	RH	80.0	11622	Pave	NaN
2	3	526351010	20	RL	81.0	14267	Pave	NaN
3	4	526353030	20	RL	93.0	11160	Pave	NaN
4	5	527105010	60	RL	74.0	13830	Pave	NaN

5 rows × 82 columns

Question 6

Now, let's split the data set into a training set and test set. We will use the training set to fit our model's parameters, and we will use the test set to estimate how well our model will perform on unseen data drawn from the same distribution. If we used all the data to fit our model, we would not have a way to estimate model performance on unseen data.

"Don't we already have a test set in `ames_test.csv`?" you might wonder. The sale prices for `ames_test.csv` aren't provided, so we're constructing our own test set for which we know the outputs.

In the cell below, split the data in `full_data` into two DataFrames named `train` and `test`. Let `train` contain 80% of the data, and let `test` contain the remaining 20% of the data.

To do this, first create two NumPy arrays named `train_indices` and `test_indices`. `train_indices` should contain a *random* 80% of the indices in `full_data`, and `test_indices` should contain the remaining 20% of the indices. Then, use these arrays to index into `full_data` to create your final `train` and `test` DataFrames.

The provided tests check that you not only answered correctly, but ended up with the exact same train/test split as our reference implementation. Later testing is easier this way.

BEGIN QUESTION
name: q6
points: 2

In [33]:

```
# This makes the train-test split in this section reproducible across different runs  
# of the notebook. You do not need this line to run train_test_split in general  
np.random.seed(1337)  
shuffled_indices = np.random.permutation(full_data_len)  
  
# Set train_indices to the first 80% of shuffled_indices and test_indices to the rest.  
train_indices = shuffled_indices[:int(full_data_len * 0.8)] # SOLUTION  
test_indices = shuffled_indices[int(full_data_len * 0.8):] # SOLUTION  
  
# Create train and test` by indexing into `full_data` using  
# `train_indices` and `test_indices`  
train = full_data.iloc[train_indices] # SOLUTION  
test = full_data.iloc[test_indices] # SOLUTION
```

In [34]:

```
# TEST  
train.shape == (1600, 82) # train should contain 80% of the data
```

Out[34]:

True

In [35]:

```
# TEST  
test.shape == (400, 82) # test should contain 20% of the data
```

Out[35]:

True

In [36]:

```
# TEST  
np.intersect1d(train_indices, test_indices).size # make sure train_indices and test_indices have no overlap
```

Out[36]:

0

In [37]:

```
# TEST  
np.intersect1d(train['PID'], test['PID']).size # make sure train and test have no overlapping houses
```

Out[37]:

0

In [38]:

```
# TEST  
sum(train['SalePrice']) # If this doesn't match, you might have still answered the question, but please adjust your code so that your split matches ours by following the implementation instructions about using shuffled_indices to split the data.
```

Out[38]:

287346111

In [39]:

```
# TEST  
sum(test['SalePrice']) # If this doesn't match, you might have still answered the question, but please adjust your code so that your split matches ours by following the implementation instructions about using shuffled_indices to split the data.
```

Out[39]:

74205684

Reusable Pipeline

Throughout this assignment, you should notice that your data flows through a single processing pipeline several times. From a software engineering perspective, it's best to define functions/methods that can apply the pipeline to any dataset. We will now encapsulate our entire pipeline into a single function `process_data_gm`. `gm` is shorthand for "guided model". We select a handful of features to use from the many that are available.

In [40]:

```
def select_columns(data, *columns):  
    """Select only columns passed as arguments."""  
    return data.loc[:, columns]  
  
def process_data_gm(data):  
    """Process the data for a guided model."""  
    data = remove_outliers(data, 'Gr_Liv_Area', upper=5000)  
  
    # Transform Data, Select Features  
    data = add_total_bathrooms(data)  
    data = select_columns(data,  
                           'SalePrice',  
                           'Gr_Liv_Area',  
                           'Garage_Area',  
                           'TotalBathrooms',  
                           )  
  
    # Return predictors and response variables separately  
    X = data.drop(['SalePrice'], axis = 1)  
    y = data.loc[:, 'SalePrice']  
  
    return X, y
```

Now, we can use `process_data_gm1` to clean our data, select features, and add our `TotalBathrooms` feature all in one step! This function also splits our data into `X`, a matrix of features, and `y`, a vector of sale prices.

Run the cell below to feed our training and test data through the pipeline, generating `X_train`, `y_train`, `X_test`, and `y_test`.

In [41]:

```
# Pre-process our training and test data in exactly the same way  
# Our functions make this very easy!  
X_train, y_train = process_data_gm(train)  
X_test, y_test = process_data_gm(test)
```

Fitting Our First Model

We are finally going to fit a model! The model we will fit can be written as follows:

$$\text{SalePrice} = \theta_0 + \theta_1 \cdot \text{Gr_Liv_Area} + \theta_2 \cdot \text{Garage_Area} + \theta_3 \cdot \text{TotalBathrooms}$$

In vector notation, the same equation would be written:

$$y = \vec{\theta} \cdot \vec{x}$$

where y is the SalePrice, $\vec{\theta}$ is a vector of all fitted weights, and \vec{x} contains a 1 for the bias followed by each of the feature values.

Note: Notice that all of our variables are continuous, except for TotalBathrooms, which takes on discrete ordered values (0, 0.5, 1, 1.5, ...). In this homework, we'll treat TotalBathrooms as a continuous quantitative variable in our model, but this might not be the best choice. The next homework may revisit the issue.

Question 7a

We will use a `sklearn.linear_model.LinearRegression` (https://scikit-learn.org/stable/modules/generated/sklearn.linear_model.LinearRegression.html) object as our linear model. In the cell below, create a `LinearRegression` object and name it `linear_model`.

Hint: See the `fit_intercept` parameter and make sure it is set appropriately. The intercept of our model corresponds to θ_0 in the equation above.

The provided tests check that you answered correctly, so that future analyses are not corrupted by a mistake.

```
BEGIN QUESTION
name: q7a
points: 1
```

In [42]:

```
from sklearn import linear_model as lm

linear_model = lm.LinearRegression(fit_intercept=True) # SOLUTION
```

In [43]:

```
# TEST
isinstance(linear_model, lm.LinearRegression)
```

Out[43]:

True

In [44]:

```
# TEST
linear_model.fit_intercept
```

Out[44]:

True

Question 7b

Now, remove the commenting and fill in the ellipses ... below with `X_train`, `y_train`, `X_test`, or `y_test`.

With the ellipses filled in correctly, the code below should fit our linear model to the training data and generate the predicted sale prices for both the training and test datasets.

The provided tests check that you answered correctly, so that future analyses are not corrupted by a mistake.

BEGIN QUESTION

name: q7b

points: 2

In [45]:

```
# Uncomment the lines below and fill in the ... with X_train, y_train, X_test, o
r y_test.
# linear_model.fit(..., ...)
# y_fitted = linear_model.predict(...)
# y_predicted = linear_model.predict(...)
# BEGIN SOLUTION NO PROMPT
linear_model.fit(X_train, y_train)
y_fitted = linear_model.predict(X_train)
y_predicted = linear_model.predict(X_test)
# END SOLUTION
```

In [46]:

```
# TEST
179000 <= y_fitted.mean() <= 180000
```

Out[46]:

True

In [47]:

```
# TEST
188000 <= y_predicted.mean() <= 190000
```

Out[47]:

True

Question 8a

Is our linear model any good at predicting house prices? Let's measure the quality of our model by calculating the Root-Mean-Square Error (RMSE) between our predicted house prices and the true prices stored in `SalePrice`.

$$\text{RMSE} = \sqrt{\frac{\sum_{\text{houses in test set}} (\text{actual price of house} - \text{predicted price of house})^2}{\text{\# of houses in data set}}}$$

In the cell below, write a function named `rmse` that calculates the RMSE of a model.

Hint: Make sure you are taking advantage of vectorized code. This question can be answered without any `for` statements.

The provided tests check that you answered correctly, so that future analyses are not corrupted by a mistake.

BEGIN QUESTION

name: q8a

points: 1

In [48]:

```
def rmse(actual, predicted):  
    """  
    Calculates RMSE from actual and predicted values  
    Input:  
        actual (1D array): vector of actual values  
        predicted (1D array): vector of predicted/fitted values  
    Output:  
        a float, the root-mean square error  
    """  
    return np.sqrt(np.mean((actual - predicted)**2)) # SOLUTION
```

In [49]:

```
# TEST  
rmse(np.array([77]), np.array([7])) # RMSE of one data point
```

Out[49]:

70.0

In [50]:

```
# TEST  
np.isclose(rmse(np.arange(0, 10), np.arange(10, 29, 2)), 14.7817, rtol=0.1) # RM  
SE of multiple data points
```

Out[50]:

True

Question 8b

Now use your `rmse` function to calculate the training error and test error in the cell below.

The provided tests for this question do not confirm that you have answered correctly; only that you have assigned each variable to a non-negative number.

```
BEGIN QUESTION  
name: q8b  
points: 1
```


In [51]:

```
training_error = rmse(y_fitted, y_train) # SOLUTION
test_error = rmse(y_predicted, y_test) # SOLUTION
(training_error, test_error)
```

Out[51]:

```
(46710.597505875914, 46146.64265682625)
```

In [52]:

```
# TEST
training_error > 0
```

Out[52]:

```
True
```

In [53]:

```
# TEST
test_error > 0
```

Out[53]:

```
True
```

In [54]:

```
# HIDDEN TEST
np.isclose(training_error, 46710.5975, atol=0.1) or np.isclose(training_error, 3
2163.67775, atol=0.1) # 1st case correct, 2nd case if switched root & mean in rm
se (our public tests didn't verify)
```

Out[54]:

```
True
```

In [55]:

```
# HIDDEN TEST
np.isclose(test_error, 46146.6426, atol=0.1) or np.isclose(test_error, 32788.747
34, atol=0.1) # 1st case correct, 2nd case if switched root & mean in rmse (our
public tests didn't verify)
```

Out[55]:

```
True
```

Question 8c

How much does including TotalBathrooms as a predictor reduce the RMSE of the model on the test set? That is, what's the difference between the RSME of a model that only includes Gr_Liv_Area and Garage_Area versus one that includes all three predictors?

The provided tests for this question do not confirm that you have answered correctly; only that you have assigned the answer variable to a non-negative number.

```
BEGIN QUESTION
name: q8c
points: 2
```

In [56]:

```
# BEGIN SOLUTION NO PROMPT
X_train_no_bath = X_train.loc[:, : 'Garage_Area']
linear_model.fit(X_train_no_bath, y_train)
y_predicted_no_bath = linear_model.predict(X_test.loc[:, : 'Garage_Area'])
test_error_no_bath = rmse(y_predicted_no_bath, y_test)
# END SOLUTION
test_error_difference = test_error_no_bath - test_error
test_error_difference
```

Out[56]:

```
2477.0084636470347
```

In [57]:

```
# TEST
test_error_difference >= 0
```

Out[57]:

```
True
```

In [58]:

```
# HIDDEN TEST
np.isclose(test_error_difference, 2477.00846, atol=0.1) or np.isclose(test_error_difference, 2039.10204, atol=0.1) # 1st case correct, 2nd case if switched root & mean in rmse (our public tests didn't verify)
```

Out[58]:

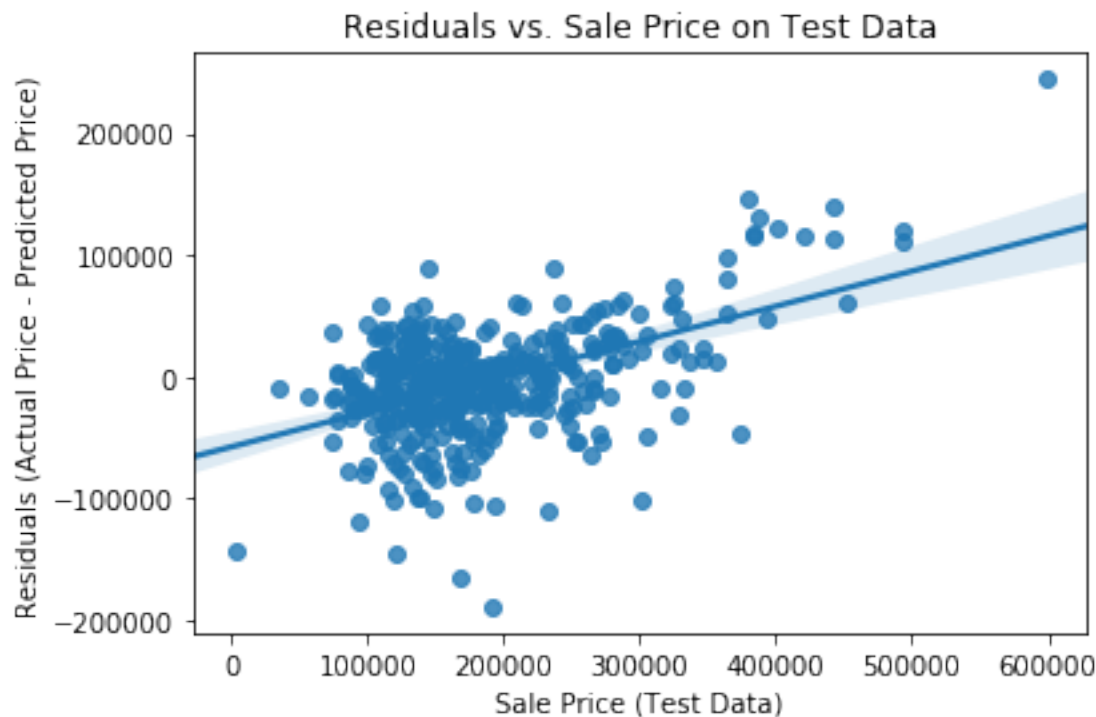
```
True
```

Residual Plots

One way of understanding the performance (and appropriateness) of a model is through a residual plot. Run the cell below to plot the actual sale prices against the residuals of the model for the test data.

In [59]:

```
residuals = y_test - y_predicted
ax = sns.regplot(y_test, residuals)
ax.set_xlabel('Sale Price (Test Data)')
ax.set_ylabel('Residuals (Actual Price - Predicted Price)')
ax.set_title("Residuals vs. Sale Price on Test Data");
```



Ideally, we would see a horizontal line of points at 0 (perfect prediction!). The next best thing would be a homogenous set of points centered at 0.

But alas, our simple model is probably too simple. The most expensive homes are systematically more expensive than our prediction.

Question 8d

What changes could you make to your linear model to improve its accuracy and lower the test error? Suggest at least two things you could try in the cell below, and carefully explain how each change could potentially improve your model's accuracy.

```
BEGIN QUESTION
name: q8d
points: 2
manual: True
```

SOLUTION: We could add more features that are likely to correlate with expensive houses. We could also encode the neighborhoods as features, as houses from the same neighborhood might be closer in price. We could also introduce regularization to improve our model.