COSC 3337 Dr. Rizk **About The Data**

Our goal for this lab is construct a model that can take a certain set of housing features and give us

back a price estimate. Since price is a continuous variable, linear regression may be a good place to

Week 4 Lab (Linear Regression)

start from. The dataset that we'll be using for this task comes from kaggle.com and contains the following attributes:

 'Avg. Area Income': Avg. income of residents of the city house is located in. 'Avg. Area House Age': Avg age of houses in same city

'Avg. Area Number of Rooms': Avg number of rooms for houses in same city

'Avg. Area Number of Bedrooms': Avg number of bedrooms for houses in same city

Area

Population

40173.072174

36882.159400

34310.242831

4.23 26354.109472 6.309435e+05

Avg. Area

Number of

Bedrooms

3.981330

1.234137

2.000000

3.140000

4.050000

4.490000

6.500000

float64

float64

float64

float64

float64

float64

object

5000.000000

3.26

Avg. Area

Rooms

Number of

5000.000000

6.987792

1.005833

3.236194

6.299250

7.002902

7.665871

10.759588

Non-Null Count Dtype

5000 non-null

Taking a closer look at price, we see that it's normally distributed with a peak around 1.232073e+06,

1.0

A scatterplot of Price vs. Avg. Area Income shows a strong positive linear relationship between the

Creating a boxplot of Avg. Area Number of Bedrooms lets us see that the median average area number of bedrooms is around 4, with a minimum of 2 and max of around 6.5. We can also so that there are no

Avg. Area Number of Bedrooms

Another important thing to look for while we're exploring our data is multicollinearity. Multicollinearity means that several variables are essentially measuring the same thing. Not only is there no point to having more than one measure of the same thing in a model, but doing so can actually cause our

model results to fluctuate. Luckily, checking for multicollinearity can be done easily with the help of a heatmap. Note: Depending on the situation, it may not be a problem for your model if only slight or moderate collinearity issue occur. However, it is strongly advised to solve the issue if severe

-0.0094

0.46

0.002

0.34

This dataset is quite clean, and so there's no severe collinearity issues. We'll later dive into some

We're now ready to begin creating and training our model. We first need to split our data into training and testing sets. This can be done using sklearn's train_test_split(X, y, test_size) function. This

function takes in your features (X), the target variable (y), and the test_size you'd like (Generally a test size of around 0.3 is good enough). It will then return a tuple of X_train, X_test, y_train, y_test sets for

X = housing_data[['Avg. Area Income', 'Avg. Area House Age', 'Avg. Area Number of Room

us. We will train our model on the training set and then use the test set to evaluate the model.

X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.3)

We'll now import sklearn's LinearRegression model and begin training it using the fit(train_data,

Now that we've finished training, we can make predictions off of the test data and evaluate our

1.0

Something that you may recall from MATH 3339 is that we'd like to see the residuals be normally

To get a rough idea of how well the model is predicting, we can make a scatterplot with the true test labels (y_test) on the x-axis, and our predictions on the y-axis. Ideally, we'd like a 45 degree line. The

train_data_labels) method. In a nutshell, fitting is equal to training. Then, after it is trained, the model can be used to make predictions, usually with a predict(test_data) method call. You can think of fit as

messier datasets which will require some type of feature engineering or PCA to resolve.

0.0061

-0.022

-0.016

-0.019

0.002

-0.022

0.45

0.34

0.17

Phice

Try plotting some of the other features for yourself to see if you can discover some interesting

findings. Refer back to the matplotlib lab if you're having trouble creating any graphs.

collinearity issue exists(e.g. correlation > 0.8 between 2 variables)

-0.0094

0.0061

-0.019

Age

sns.heatmap(housing data.corr(), annot=True)

-0.002

-0.011

0.02

-0.016

sns.scatterplot(x='Price', y='Avg. Area Income', data=housing data)

sns.boxplot(x='Avg. Area Number of Bedrooms', data=housing data)

Price

and 75% of houses sold were at a price of 1.471210e+06 or lower.

sns.histplot(housing data['Price'])

5.000000e+03

1.232073e+06 3.531176e+05

1.593866e+04

9.975771e+05

1.232669e+06

1.471210e+06 2.469066e+06 Name: Price, dtype: float64

print(housing data['Price'].describe())

plt.show()

300

250

200

5 0 150

0

mean

std

min 25%

50%

75%

plt.show()

100000

80000

40000

20000

outliers present.

plt.show()

plt.show()

Avg. Area Income

Avg. Area House Age

Area Population

Price

Creating Our Linear Model

y = housing data['Price']

lm = LinearRegression() lm.fit(X train, y train)

Model Evaluation

Out[12]: LinearRegression()

plt.show()

2.0

1.5

1.0

0.0

In [14]:

In [11]: from sklearn.model selection import train test split

the step that finds the coefficients for the equation.

straighter the line, the better our predictions are.

0.5

residuals = y_test - predictions

sns.histplot(residuals)

Comparing these metrics:

the real world.

from sklearn import metrics

from sklearn.metrics import r2_score

Avg. Area Income

Area Population

Avg. Area Number of Rooms

Avg. Area Number of Bedrooms

What these coefficients mean:

increase of \$21.564645.

increase of \$15.309706.

use this notebook as a guide.

increase of \$166102.423648.

with an increase of \$122398.915857.

with an increase of \$887.665746.

Avg. Area House Age 166102.423648

R2 Score: 0.9150208174786678

coeff df

print('R2 Score: ', r2_score(y_test, predictions))

MAE: 83410.59496702514 MSE: 10608579825.136667 RMSE: 102997.96029600133

plt.show()

140

120

80 60 40

20

distributed in regression analysis. We can exam this as follows:

Here are the most common evaluation metrics for regression problems:

MAE is the easiest to understand, because it's the average error.

All of these are **loss functions**, because we want to minimize them.

Mean Squared Error (MSE) is the mean of the squared errors:

Mean Absolute Error (MAE) is the mean of the absolute value of the errors:

 $\frac{1}{n}\sum_{i=1}^n |y_i - \hat{y}_i|$

 $\frac{1}{n} \sum_{i=1}^{n} (y_i - \hat{y}_i)^2$

 $\sqrt{\frac{1}{n}\sum_{i=1}^n(y_i-\hat{y}_i)^2}$

• MSE is more popular than MAE, because MSE "punishes" larger errors, which tends to be useful in

RMSE is even more popular than MSE, because RMSE is interpretable in the "y" units.

Luckily, sklearn can calculate all of these metrics for us. All we need to do is pass the true labels

depends on your data. You can find a great example here, or refer back to the power points.

print('RMSE:', np.sqrt(metrics.mean_squared_error(y_test, predictions)))

Something we also like to look at is the coefficient of determination (R^2) , which is the percentage of

variation in y explained by all the x variables together. Usually an \mathbb{R}^2 of .70 is considered good.

Finally, let's see how we can interpret our model's coefficients. We can access the coefficients by calling coef_ on our linear model (Im in this case). We'll use this and put it in a nice pandas DataFrame

for visual purposes. Note: You can also call *intercept*_ if you'd like to get the intercept.

coeff_df = pd.DataFrame(lm.coef_, X.columns, columns=['Coefficient'])

Coefficient

21.564645

122398.915857

887.665746

15.309706

Holding all other features fixed, a 1 unit increase in Avg. Area Income is associated with an

Holding all other features fixed, a 1 unit increase in Avg. Area House Age is associated with an

Holding all other features fixed, a 1 unit increase in Avg. Area Number of Rooms is associated

• Holding all other features fixed, a 1 unit increase in Area Population is associated with an

Congratulations! You now know how to create and evaluate linear models using sklearn. As extra practice, I'd recommend now trying to find a used car or similar housing dataset on kaggle.com and

Holding all other features fixed, a 1 unit increase in Avg. Area Number of Bedrooms is associated

print('MAE:', metrics.mean_absolute_error(y_test, predictions)) print('MSE:', metrics.mean_squared_error(y_test, predictions))

(y_test) and our predictions to the functions below. What's more important is that we understand what each of these means. Root Mean Square Error (RMSE) is what we'll most commonly use, which is the standard deviation of the residuals (prediction errors). Residuals are a measure of how far from the regression line data points are; RMSE is a measure of how spread out these residuals are. In other words, it tells us how concentrated the data is around the line of best fit. Determining a good RMSE

Root Mean Squared Error (RMSE) is the square root of the mean of the squared errors:

predictions = lm.predict(X_test) plt.scatter(y_test,predictions)

from sklearn.linear model import LinearRegression

model's performance using the corresponding test data labels (y_test).

Avg. Area Number of Rooms

Avg. Area Number of Bedrooms

Area Income 60000 **Price**

1.505891e+06

1.260617e+06

Address

3701...

208 Michael Ferry Apt.

674\nLaurabury, NE

188 Johnson Views

USS Barnett\nFPO AP

USNS Raymond\nFPO

1.058988e+06 Stravenue\nDanieltown,

Area

Population

5000.000000

36163.516039

9925.650114

172.610686

29403.928702

36199.406689

42861.290769

69621.713378

Suite 079\nLake Kathleen, CA... 9127 Elizabeth

WI 06482...

44820

AE 09386

Price

5.000000e+03

1.232073e+06

3.531176e+05

1.593866e+04

9.975771e+05

1.232669e+06

1.471210e+06

2.469066e+06

Exploratory Data Analysis Let's begin by importing some necessary libraries that we'll be using to explore the data. import numpy as np In [1]:

'Area Population': Population of city house is located in

'Price': Price that the house sold at (target)

· 'Address': Address for the house

import pandas as pd import matplotlib.pyplot as plt import seaborn as sns

from matplotlib import rcParams rcParams['figure.figsize'] = 15, 5 sns.set_style('darkgrid')

housing_data = pd.read_csv('USA_Housing.csv') housing_data.head() Avg.

Avg. Avg. Area Area Number Avg. Area Area Number

Our first step is to load the data into a pandas DataFrame

Income House of

of

Bedrooms

Rooms

79545.458574 5.682861 7.009188 4.09 23086.800503 1.059034e+06

3.09

79248.642455 6.002900 6.730821

61287.067179 5.865890

8.512727

3 63345.240046 7.188236 5.586729

59982.197226 5.040555 7.839388

see if we have any missing values. housing data.describe()

From here, it's always a good step to use describe() and info() to get a better sense of the data and Avg. Area Avg. Area Income **House Age**

In [4]: count 5000.000000 mean

68583.108984 10657.991214 std

Out[4]:

5000.000000 5.977222 0.991456 17796.631190 2.644304 min 5.322283 25% 61480.562388 68804.286404 50% 5.970429

75783.338666 75% 107701.748378

6.650808 9.519088 The info below lets us know that we have 5,000 entries and 5,000 non-null values in each feature/column. Therefore, there are no missing values in this dataset.

housing_data.info() <class 'pandas.core.frame.DataFrame'> RangeIndex: 5000 entries, 0 to 4999 Data columns (total 7 columns): Column

Avg. Area Income Avg. Area House Age Avg. Area Number of Rooms Avg. Area Number of Bedrooms Area Population Price Address dtypes: float64(6), object(1)

0

1

5

Awg

memory usage: 273.6+ KB A quick pairplot lets us get an idea of the distributions and relationships in our dataset. From here, we could choose any interesting features that we'd like to later explore in greater depth. Warning: The more features in our dataset, the harder our pairplot will be to interpret. sns.pairplot(housing data) plt.show()

40000 20000

20000 2.5 1.5 1.0