1DA2223 - Python programming and data analysis

Exercise 6 - Part 1 (Complete 2022L)    *Ostatnia modyfikacja: R. Szmurło 26.04.2022 07:59*

## Linear Regression - Project Exercise

### Data links:

- L06_Ecommerce_Customers.csv
- Salaries.csv
- titanic.csv

### Instruction

Imagine an Ecommerce company based in New York City that sells clothing online but they also have in-store style and clothing advice sessions. Cus the store, have sessions/meetings with a personal stylist, then they can go home and order either on a mobile app or website for the clothes they war

The company is trying to decide whether to focus their efforts on their mobile app experience or their website.

Just follow the steps below to analyze the customer data.

### Imports

\*\* Import pandas, numpy, matplotlib,and seaborn. Then set %matplotlib inline
(You'll import sklearn as you need it.)\*\*

### Get the Data

We'll work with the L06_Ecommerce_Customers.csv file attached to the exrecise. It has Customer info, such as Email, Address, and their color Avata numerical value columns:

- Avg. Session Length: Average session of in-store style advice sessions.
- Time on App: Average time spent on App in minutes
- Time on Website: Average time spent on Website in minutes
- Length of Membership: How many years the customer has been a member.

\*\* Read in the Ecommerce Customers csv file as a DataFrame called customers.\*\*

**Check the head of customers, and check out its info() and describe() methods.**

    . . .

| No. | Email | Address | Avatar | Avg. Session Length | Time on App | Time on Website | Length Members |
|-----|-------|---------|--------|---------------------|-------------|-----------------|----------------|
| 0 | mstephenson@fernandez.com | 835 Frank Tunnel\nWrightmouth, MI 82180-9605 | Violet | 34.497268 | 12.655651 | 39.577668 | 4.082621 |
| 1 | hduke@hotmail.com | 4547 Archer Common\nDiazchester, CA 06566-8576 | DarkGreen | 31.926272 | 11.109461 | 37.268959 | 2.664034 |
| 2 | pallen@yahoo.com | 24645 Valerie Unions Suite 582\nCobbborough, D... | Bisque | 33.000915 | 11.330278 | 37.110597 | 4.104543 |
| 3 | riverarebecca@gmail.com | 1414 David Throughway\nPort Jason, OH 22070-1220 | SaddleBrown | 34.305557 | 13.717514 | 36.721283 | 3.120179 |
| 4 | mstephens@davidson-herman.com | 14023 Rodriguez Passage\nPort Jacobville, PR 3... | MediumAquaMarine | 33.330673 | 12.795189 | 37.536653 | 4.446308 |

    . . .

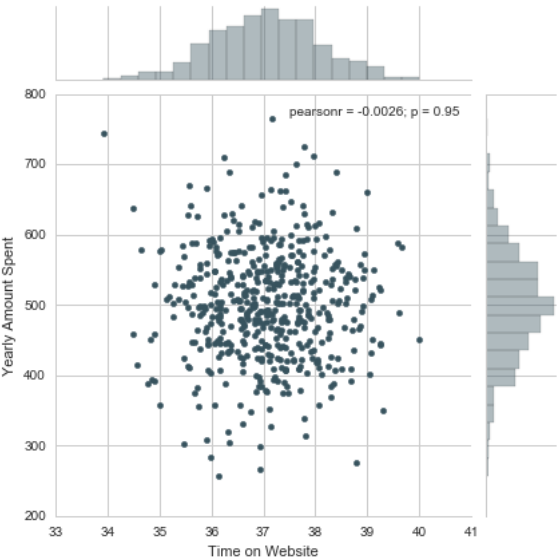| Measure | Avg. Session Length | Time on App | Time on Website | Length of Membership | Yearly Amount Spent |
|---|---|---|---|---|---|
| count | 500.000000 | 500.000000 | 500.000000 | 500.000000 | 500.000000 |
| mean | 33.053194 | 12.052488 | 37.060445 | 3.533462 | 499.314038 |
| std | 0.992563 | 0.994216 | 1.010489 | 0.999278 | 79.314782 |
| min | 29.532429 | 8.508152 | 33.913847 | 0.269901 | 256.670582 |
| 25% | 32.341822 | 11.388153 | 36.349257 | 2.930450 | 445.038277 |
| 50% | 33.082008 | 11.983231 | 37.069367 | 3.533975 | 498.887875 |
| 75% | 33.711985 | 12.753850 | 37.716432 | 4.126502 | 549.313828 |
| max | 36.139662 | 15.126994 | 40.005182 | 6.922689 | 765.518462 |

...

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 500 entries, 0 to 499
Data columns (total 8 columns):
Email                   500 non-null object
Address                 500 non-null object
Avatar                  500 non-null object
Avg. Session Length     500 non-null float64
Time on App             500 non-null float64
Time on Website         500 non-null float64
Length of Membership    500 non-null float64
Yearly Amount Spent     500 non-null float64
dtypes: float64(5), object(3)
memory usage: 31.3+ KB
```
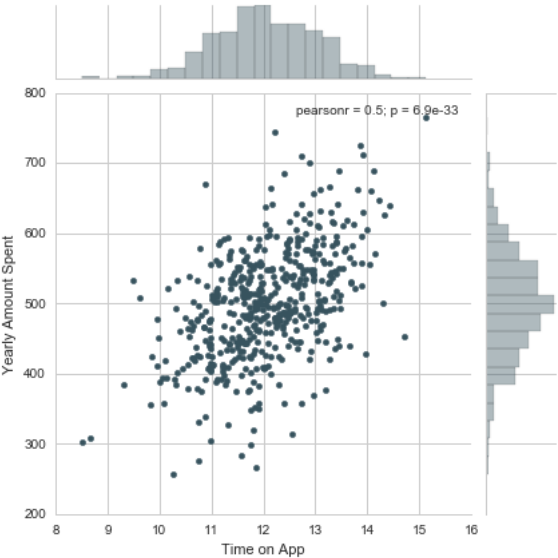
## Exploratory Data Analysis

**Use seaborn to create a jointplot to compare the Time on Website and Yearly Amount Spent columns. Does the correlation make sense?**
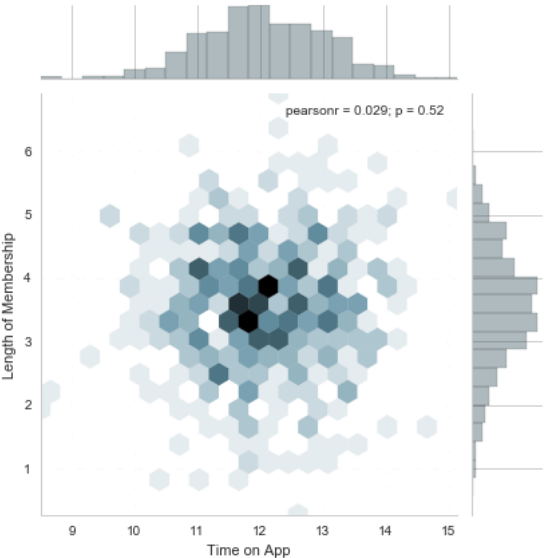
...



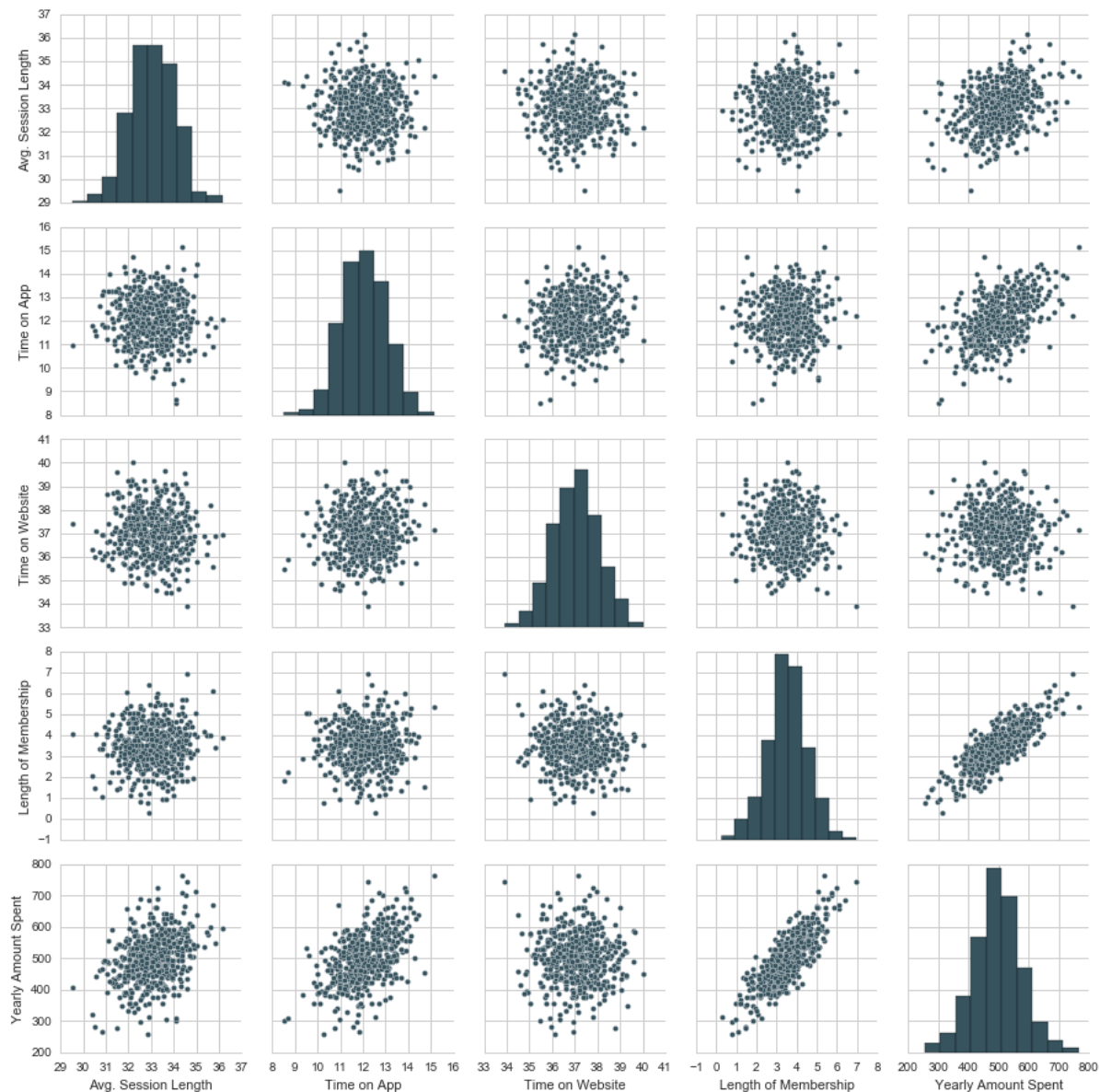** Do the same but with the Time on App column instead. **

...

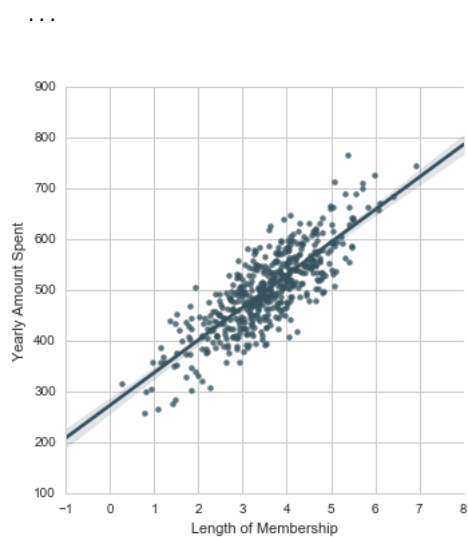** Use jointplot to create a 2D hex bin plot comparing Time on App and Length of Membership.**

. . .



**Let's explore these types of relationships across the entire data set. Use pairplot to recreate the plot below.(Don't worry about the the colo**

. . .

**Based off this plot what looks to be the most correlated feature with Yearly Amount Spent?**

**Create a linear model plot (using seaborn's lmplot) of Yearly Amount Spent vs. Length of Membership. **

. . .



## Training and Testing Data

** Set a variable X equal to the numerical features of the customers and a variable y equal to the "Yearly Amount Spent" column. **

** Use model_selection.train_test_split from sklearn to split the data into training and testing sets. Set test_size=0.3 and random_state=101**

## Training the Model

Now its time to train our model on our training data!

** Import LinearRegression from sklearn.linear_model **

**Create an instance of a LinearRegression() model named lm.**

** Train/fit lm on the training data.**

```
...
```

**Print out the coefficients of the model**

```
...
```

```
Coefficients:
 [ 25.98154972  38.59015875   0.19040528  61.27909654]
```
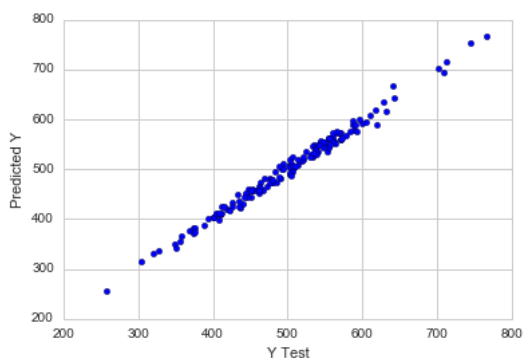
## Predicting Test Data

Now that we have fit our model, let's evaluate its performance by predicting off the test values!

** Use lm.predict() to predict off the X_test set of the data.**

** Create a scatterplot of the real test values versus the predicted values. **

```
...
```



## Evaluating the Model

Let's evaluate our model performance by calculating the residual sum of squares and the explained variance score (R^2).

** Calculate the Mean Absolute Error, Mean Squared Error, and the Root Mean Squared Error. Refer to the lecture or to Wikipedia for the formulas**
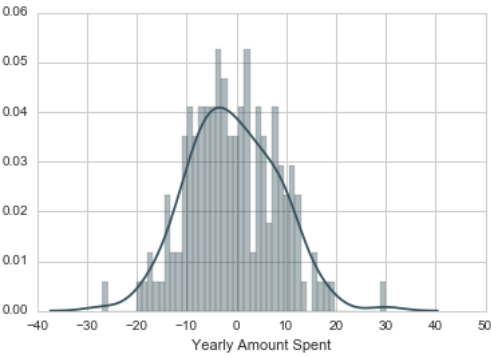
```
...
```

```
MAE: 7.22814865343
MSE: 79.813051651
RMSE: 8.93381506698
```

## Residuals

**Plot a histogram of the residuals and make sure it looks normally distributed. Use either seaborn distplot, or just plt.hist().**

```
...
```

## Conclusion

We still want to figure out the answer to the original question, do we focus our efforts on mobile app or website development? Or maybe that doesn't and Membership Time is what is really important. Let's see if we can interpret the coefficients at all to get an idea.

** Recreate the dataframe below. **

. . .

| Coeffecient | |
| --- | --- |
| Avg. Session Length | 25.981550 |
| Time on App | 38.590159 |
| Time on Website | 0.190405 |
| Length of Membership | 61.279097 |

** How can you interpret these coefficients? **

**Do you think the company should focus more on their mobile app or on their website?**

*Answer?*