

## LINEAR REGRESSION

Amazon\_cloths sells cloths online. Customers come in to the store, have meetings with a personal stylist, then they can go home and order either on a mobile app or website for the clothes they want.

The company is trying to decide whether to focus their efforts on their mobile app experience or their website. Following is predict is analysis for this company

Just follow the steps below to analyze the customer dat

import libaraies

```
In [9]: import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
%matplotlib inline
```

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

```
In [11]: customers = pd.read_csv('E:\Ecommerce Customers')
```

```
In [12]: customers.head()
```

```
Out[12]:
```

	Email	Address	Avatar	Avg. Session Length	Time on App	Ti W
0	mstephenson@fernandez.com	835 Frank Tunnel\nWrightmouth, MI 82180-9605	Violet	34.497268	12.655651	39.5
1	hduke@hotmail.com	4547 Archer Common\nDiazchester, CA 06566-8576	DarkGreen	31.926272	11.109461	37.2
2	pallen@yahoo.com	24645 Valerie Unions Suite 582\nCobbborough, D...	Bisque	33.000915	11.330278	37.1
3	riverarebecca@gmail.com	1414 David Throughway\nPort Jason, OH 22070-1220	SaddleBrown	34.305557	13.717514	36.7
4	mstephens@davidson-herman.com	14023 Rodriguez Passage\nPort Jacobville, PR 3...	MediumAquaMarine	33.330673	12.795189	37.5

```
In [13]: customers.describe()
```

```
Out[13]:
```

	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

	Avg. Session Length	Time on App	Time on Website	Length of Membership	Yearly Amount Spent
<b>std</b>	0.992563	0.994216	1.010489	0.999278	79.314782
<b>min</b>	29.532429	8.508152	33.913847	0.269901	256.670582
<b>25%</b>	32.341822	11.388153	36.349257	2.930450	445.038277
<b>50%</b>	33.082008	11.983231	37.069367	3.533975	498.887875
<b>75%</b>	33.711985	12.753850	37.716432	4.126502	549.313828
<b>max</b>	36.139662	15.126994	40.005182	6.922689	765.518462

In [16]:

```
customers.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 500 entries, 0 to 499
Data columns (total 8 columns):
#   Column                Non-Null Count  Dtype
---  -
0   Email                  500 non-null   object
1   Address                 500 non-null   object
2   Avatar                 500 non-null   object
3   Avg. Session Length    500 non-null   float64
4   Time on App            500 non-null   float64
5   Time on Website        500 non-null   float64
6   Length of Membership    500 non-null   float64
7   Yearly Amount Spent    500 non-null   float64
dtypes: float64(5), object(3)
memory usage: 31.4+ KB
```

## DATA ANALYSIS

In [18]:

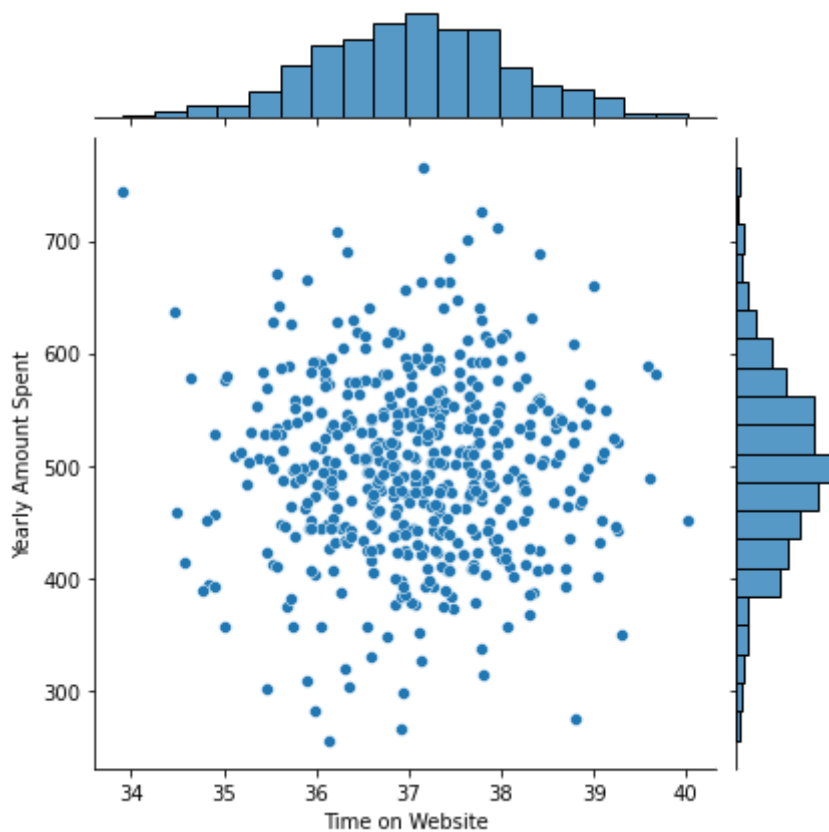
```
import seaborn as sns
```

In [19]:

```
sns.jointplot(customers['Time on Website'], customers['Yearly Amount Spent'])
```

```
C:\Users\apc\Anaconda3\lib\site-packages\seaborn\_decorators.py:36: FutureWarning: Pass the following variables as keyword args: x, y. From version 0.12, the only valid positional argument will be `data`, and passing other arguments without an explicit keyword will result in an error or misinterpretation.
  warnings.warn(
```

Out[19]: &lt;seaborn.axisgrid.JointGrid at 0x496e970d60&gt;



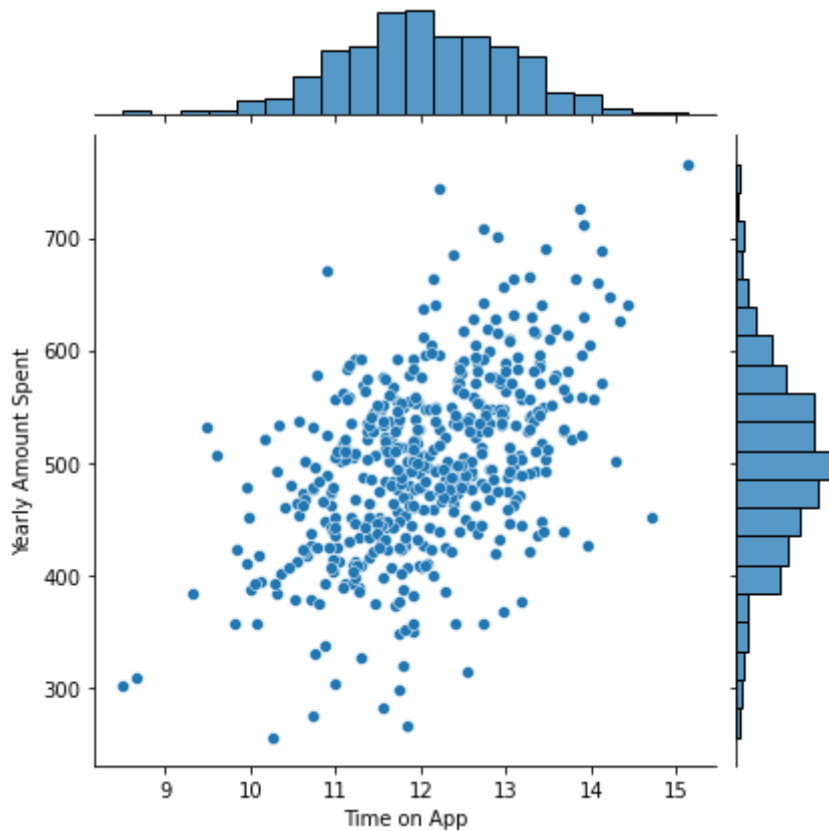
Do the same but with the Time on App column instead.

```
In [20]: sns.jointplot(customers['Time on App'],customers['Yearly Amount Spent'])
```

C:\Users\apc\Anaconda3\lib\site-packages\seaborn\\_decorators.py:36: FutureWarning: Pass the following variables as keyword args: x, y. From version 0.12, the only valid positional argument will be `data`, and passing other arguments without an explicit keyword will result in an error or misinterpretation.

```
warnings.warn(
```

```
Out[20]: <seaborn.axisgrid.JointGrid at 0x4971070c10>
```

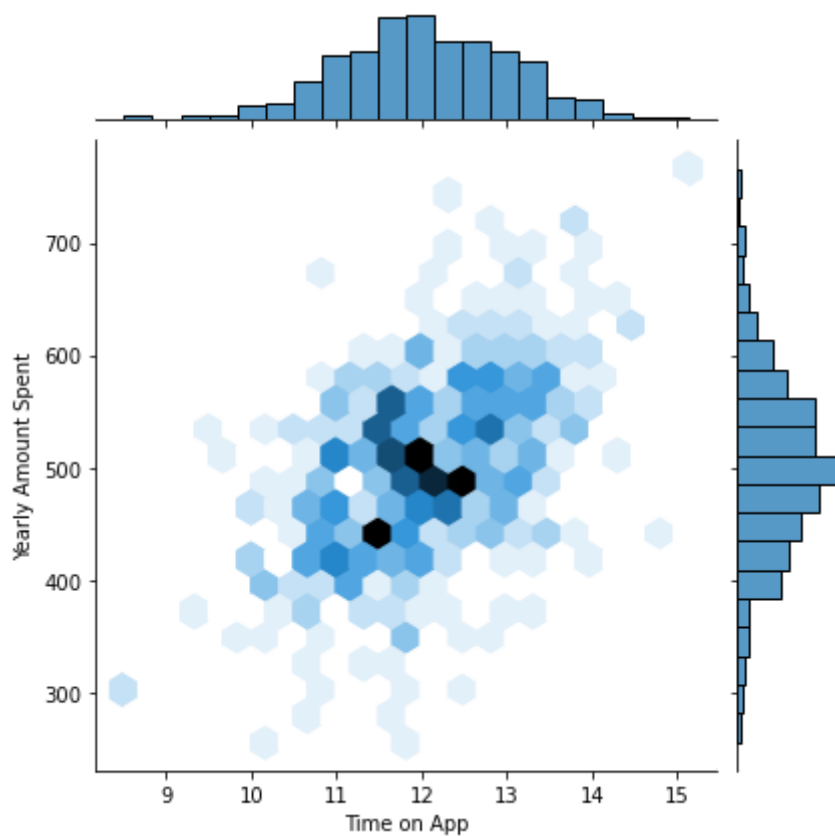


Use jointplot to create a 2D hex bin plot comparing Time on App and Length of Membership.

```
In [21]: sns.jointplot(customers['Time on App'],customers['Yearly Amount Spent'],kind='hex')
```

C:\Users\apc\Anaconda3\lib\site-packages\seaborn\\_decorators.py:36: FutureWarning: Pass the following variables as keyword args: x, y. From version 0.12, the only valid positional argument will be `data`, and passing other arguments without an explicit keyword will result in an error or misinterpretation.  
warnings.warn(

```
Out[21]: <seaborn.axisgrid.JointGrid at 0x496fb40fd0>
```

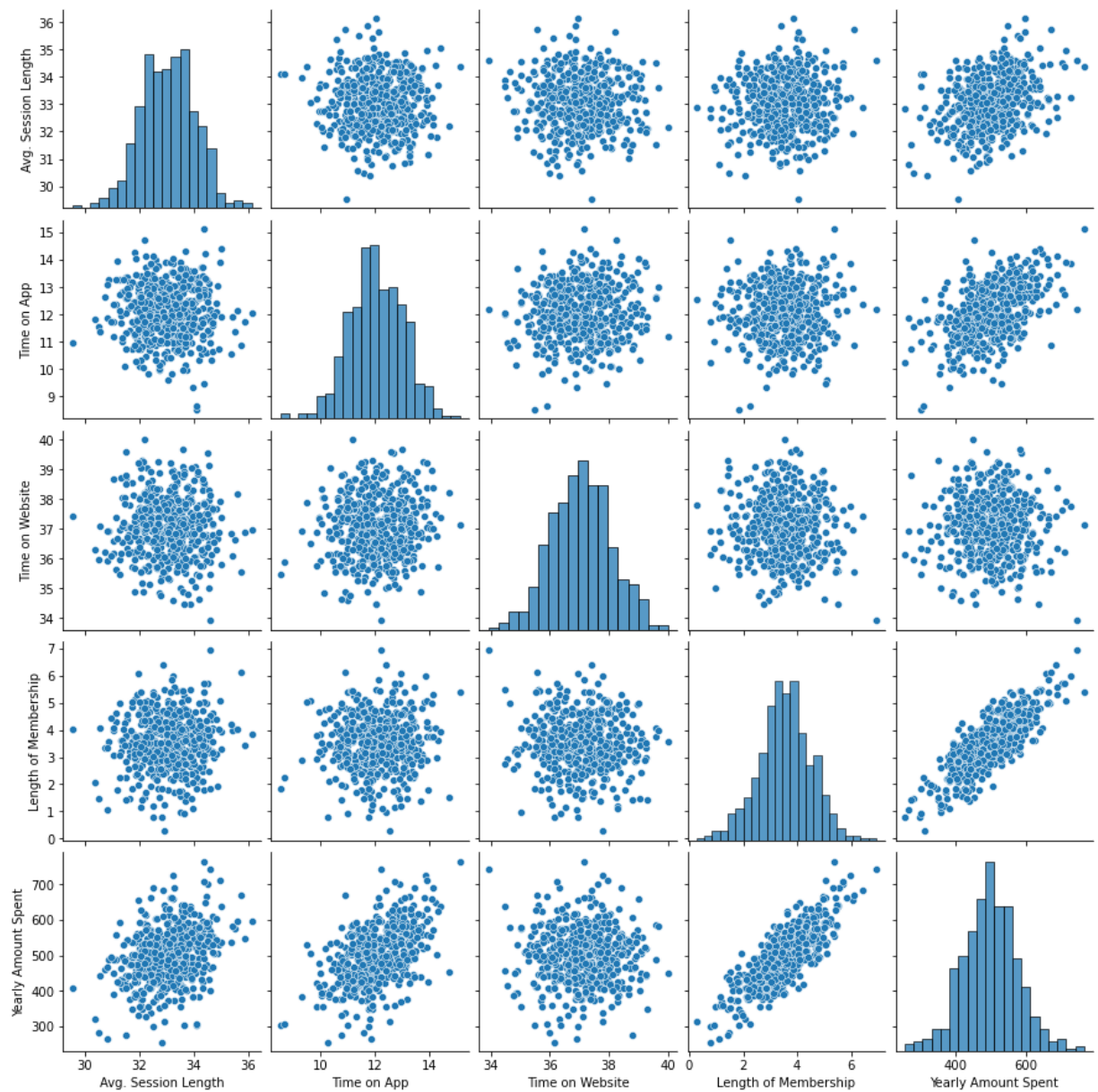


Let's explore these types of relationships across the entire data set

In [22]:

```
sns.pairplot(customers)
```

Out[22]: <seaborn.axisgrid.PairGrid at 0x497129fbb0>



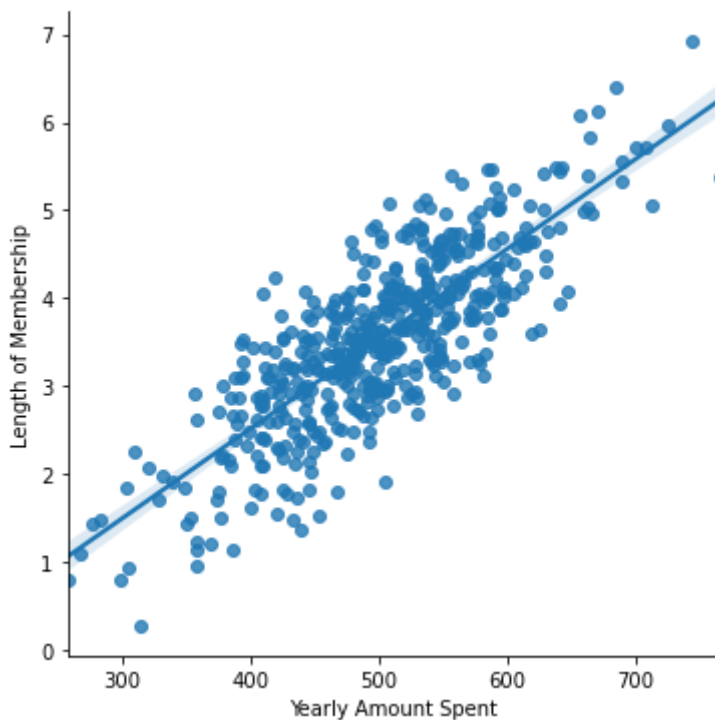
Based off this plot what looks to be the most correlated feature with Yearly Amount Spent?

## Length of Membership

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

```
In [23]: sns.lplot(x='Yearly Amount Spent', y='Length of Membership', data=customers)
```

```
Out[23]: <seaborn.axisgrid.FacetGrid at 0x49711c6820>
```



## Training and Testing Data

Training and Testing Data Now that we've explored the data a bit, let's go ahead and split the data into training and testing sets. Set a variable X equal to the numerical features of the customers and a variable y equal to the "Yearly Amount Spent" column.

```
In [24]: y = customers['Yearly Amount Spent']
```

```
In [25]: X = customers[['Avg. Session Length', 'Time on App', 'Time on Website', 'Length of Me
```

Use `model_selection.train_test_split` from `sklearn` to split the data into training and testing sets. Set `test_size=0.5` and `random_state=103`

```
In [26]: from sklearn.model_selection import train_test_split
```

```
In [27]: X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.5, random_stat
```

## Training the Model

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

Import `LinearRegression` from `sklearn.linear_model`

```
In [30]: from sklearn.linear_model import LinearRegression
```

Create an instance of a `LinearRegression()` model named `lm`

```
In [31]: lm = LinearRegression()
```

Train/fit lm on the training data.

```
In [32]: lm.fit(X_train,y_train)
```

```
Out[32]: LinearRegression()
```

Print out the coefficients of the model

```
In [33]: print('Coefficients: \n', lm.coef_)
```

```
Coefficients:  
[24.93886045  38.93370174 -0.3853061  61.85088684]
```

## Predicting Test Data

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

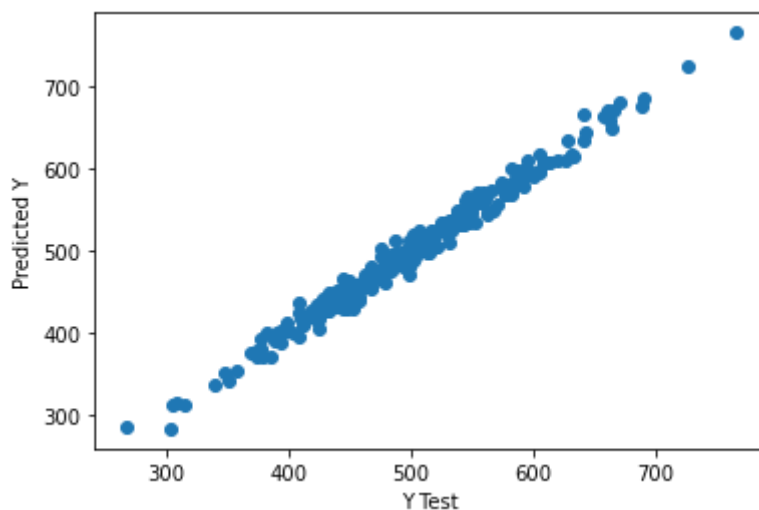
Use `lm.predict()` to predict off the `X_test` set of the data.

```
In [34]: predictions = lm.predict(X_test)
```

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

```
In [35]: plt.scatter(y_test,predictions)  
plt.xlabel('Y Test')  
plt.ylabel('Predicted Y')
```

```
Out[35]: Text(0, 0.5, 'Predicted Y')
```



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.

```
In [37]: from sklearn import metrics  
  
print('MAE:', metrics.mean_absolute_error(y_test, predictions))  
print('MSE:', metrics.mean_squared_error(y_test, predictions))  
print('RMSE:', np.sqrt(metrics.mean_squared_error(y_test, predictions)))
```



MAE: 8.084327052696858  
 MSE: 100.83296403823485  
 RMSE: 10.041561832615226

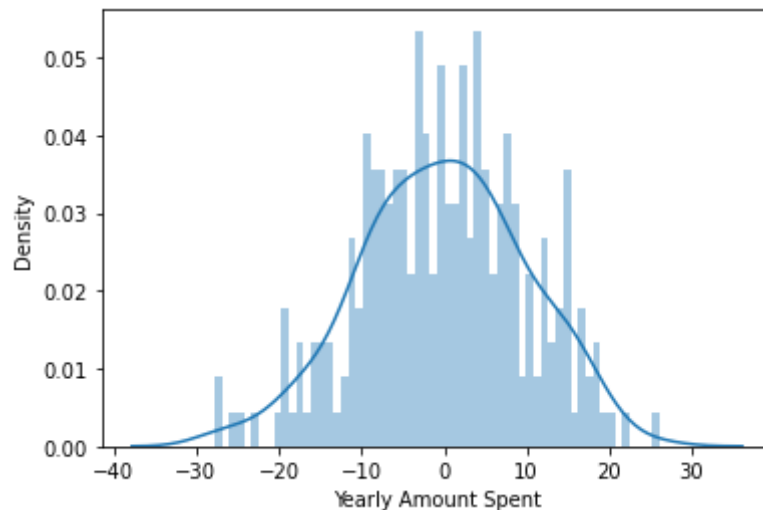
## Residuals

Let's quickly explore the residuals to make sure everything was okay with our data.

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

```
In [42]: sns.distplot((y_test-predictions),bins=60);
```

C:\Users\apc\Anaconda3\lib\site-packages\seaborn\distributions.py:2557: FutureWarning: `distplot` is a deprecated function and will be removed in a future version. Please adapt your code to use either `displot` (a figure-level function with similar flexibility) or `histplot` (an axes-level function for histograms).  
 warnings.warn(msg, FutureWarning)



## Conclusion

We still want to figure out the answer to the original question, do we focus our effort on mobile app or website development? Or maybe that doesn't even really matter, 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.

```
In [43]: coefficients = pd.DataFrame(lm.coef_,X.columns)
coefficients.columns = ['Coefficient']
coefficients
```

```
Out[43]:
```

	Coefficient
<b>Avg. Session Length</b>	24.938860
<b>Time on App</b>	38.933702
<b>Time on Website</b>	-0.385306
<b>Length of Membership</b>	61.850887

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

## **MOBILE APP**

In [ ]: