



(<http://www.pieriandata.com>)

Linear Regression - Project Exercise

Congratulations! You just got some contract work with an Ecommerce company based in New York City that sells clothing online but they also have in-store style and clothing advice sessions. Customers come in to 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 want.

The company is trying to decide whether to focus their efforts on their mobile app experience or their website. They've hired you on contract to help them figure it out! Let's get started!

Just follow the steps below to analyze the customer data (it's fake, don't worry I didn't give you real credit card numbers or emails).

Imports

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

```
In [1]: import numpy as np
import pandas as pd
import matplotlib.pyplot as plt
import seaborn as sns
from matplotlib import rcParams
rcParams['patch.force_edgecolor'] = True
rcParams['patch.facecolor'] = 'b'
plt.style.use('seaborn')
%matplotlib inline
```

Get the Data

We'll work with the Ecommerce Customers csv file from the company. It has Customer info, such as Email, Address, and their color Avatar. Then it also has 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.

```
In [2]: df = pd.read_csv('Ecommerce Customers')
```

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

```
In [3]: df.head()
```

Out[3]:

	Email	Address	Avatar	Avg. Session Length	Time on App	Time on Website
0	mstephenson@fernandez.com	835 Frank Tunnel\nWrightmouth, MI 82180-9605	Violet	34.497268	12.655651	39.51
1	hduke@hotmail.com	4547 Archer Common\nDiazchester, CA 06566-8576	DarkGreen	31.926272	11.109461	37.26
2	pallen@yahoo.com	24645 Valerie Unions Suite 582\nCobbborough, D...	Bisque	33.000915	11.330278	37.11
3	riverarebecca@gmail.com	1414 David Throughway\nPort Jason, OH 22070-1220	SaddleBrown	34.305557	13.717514	36.74
4	mstephens@davidson-herman.com	14023 Rodriguez Passage\nPort Jacobville, PR 3...	MediumAquaMarine	33.330673	12.795189	37.56

In [277]:

Out[277]:

	Email	Address	Avatar	Avg. Session Length	Time on App	Time We
0	mstephenson@fernandez.com	835 Frank Tunnel\nWrightmouth, MI 82180-9605	Violet	34.497268	12.655651	39.57
1	hduke@hotmail.com	4547 Archer Common\nDiazchester, CA 06566-8576	DarkGreen	31.926272	11.109461	37.26
2	pallen@yahoo.com	24645 Valerie Unions Suite 582\nCobbborough, D...	Bisque	33.000915	11.330278	37.17
3	riverarebecca@gmail.com	1414 David Throughway\nPort Jason, OH 22070-1220	SaddleBrown	34.305557	13.717514	36.72
4	mstephens@davidson-herman.com	14023 Rodriguez Passage\nPort Jacobville, PR 3...	MediumAquaMarine	33.330673	12.795189	37.56

In [4]: df.describe()

Out[4]:

	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

In [278]:

Out[278]:

	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

In [5]: df.info()

```
<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
```

In [279]:

```
<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

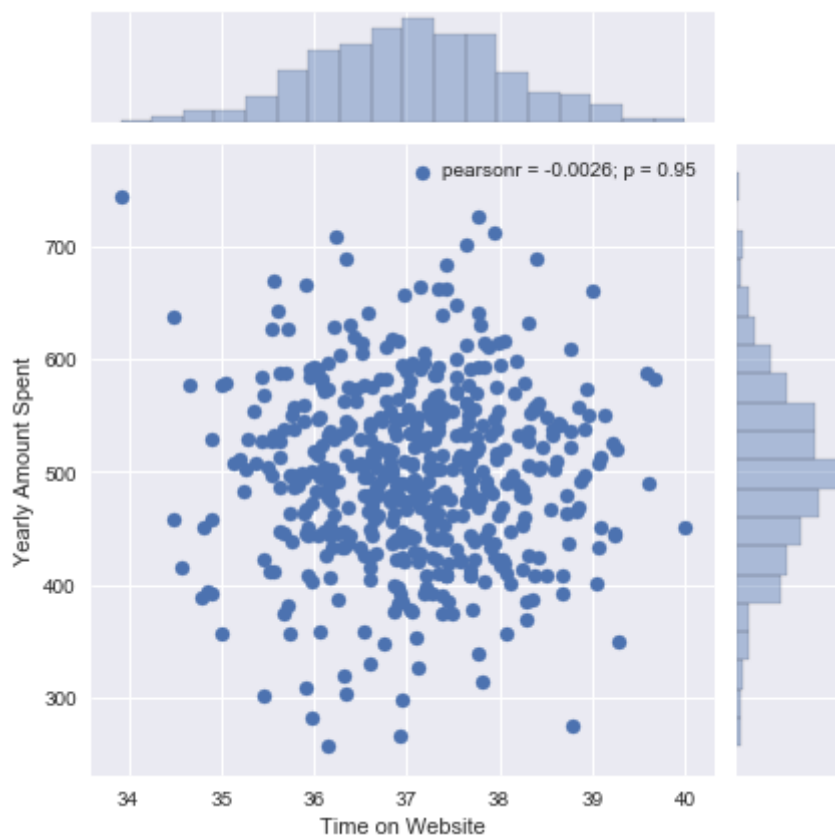
Let's explore the data!

For the rest of the exercise we'll only be using the numerical data of the csv file.

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

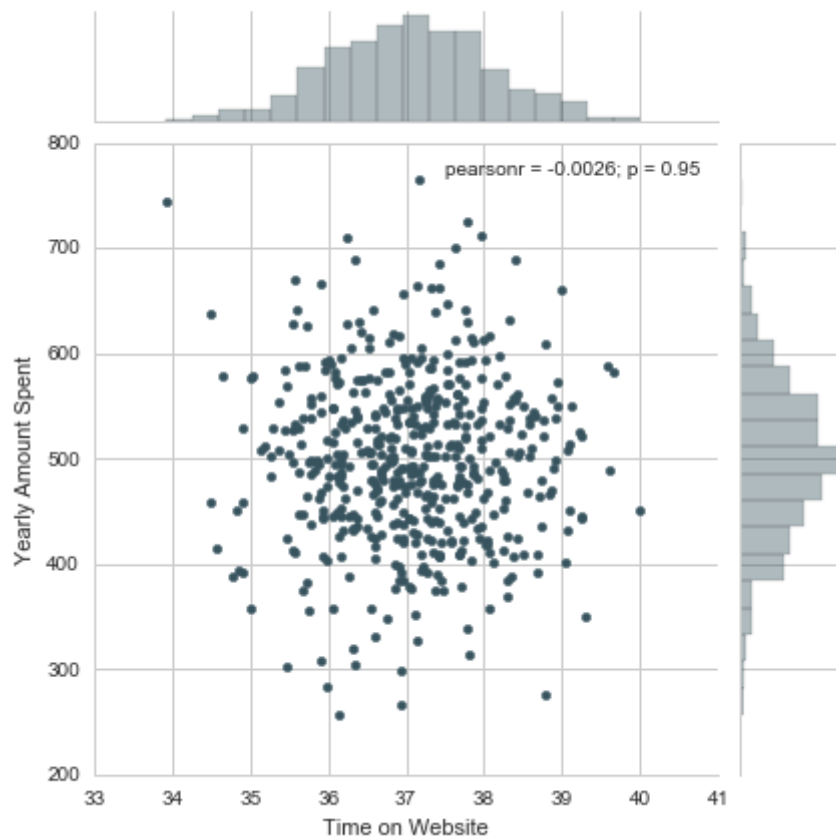
```
In [6]: sns.jointplot(df['Time on Website'], df['Yearly Amount Spent'])
```

```
Out[6]: <seaborn.axisgrid.JointGrid at 0x108a95400>
```



```
In [281]:
```

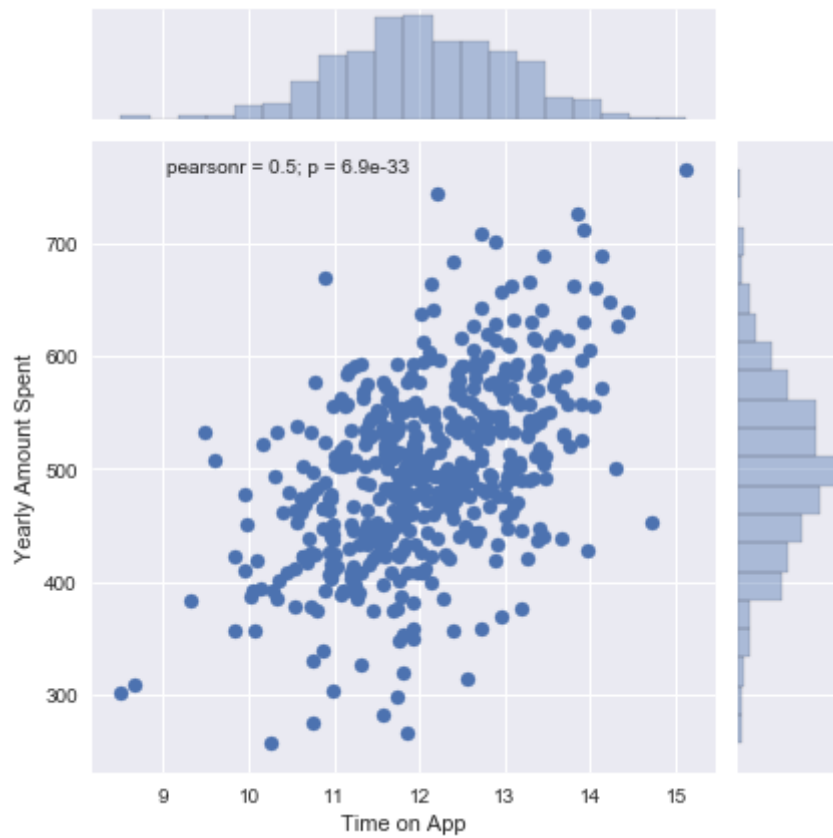
```
Out[281]: <seaborn.axisgrid.JointGrid at 0x120bfcc88>
```



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

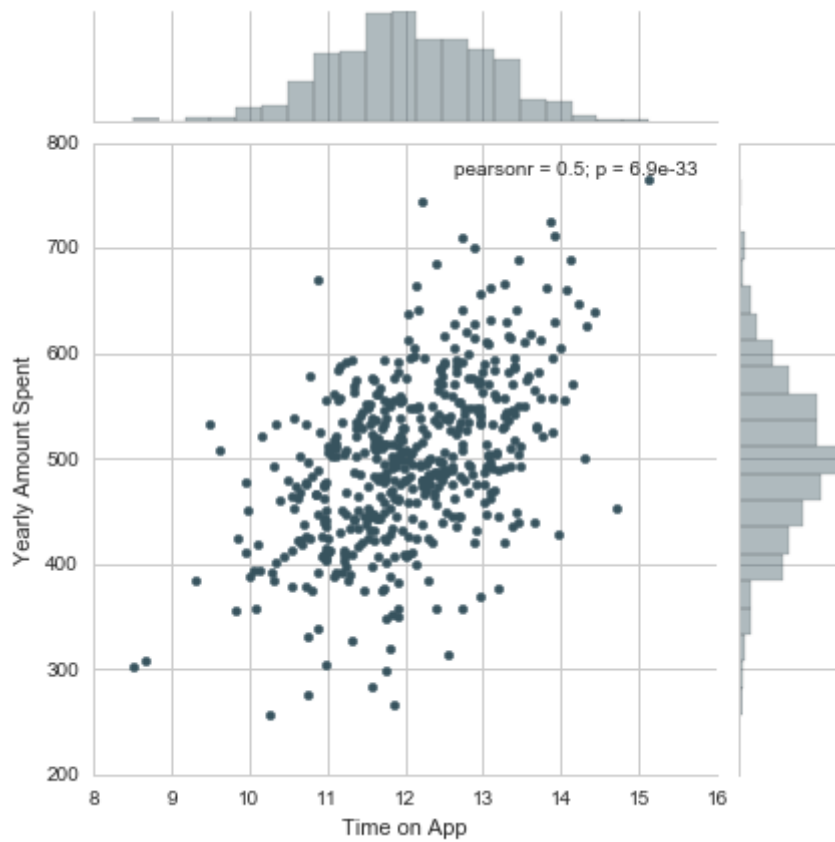
```
In [8]: sns.jointplot(df['Time on App'], df['Yearly Amount Spent'])
```

```
Out[8]: <seaborn.axisgrid.JointGrid at 0x10eb10978>
```



In [282]:

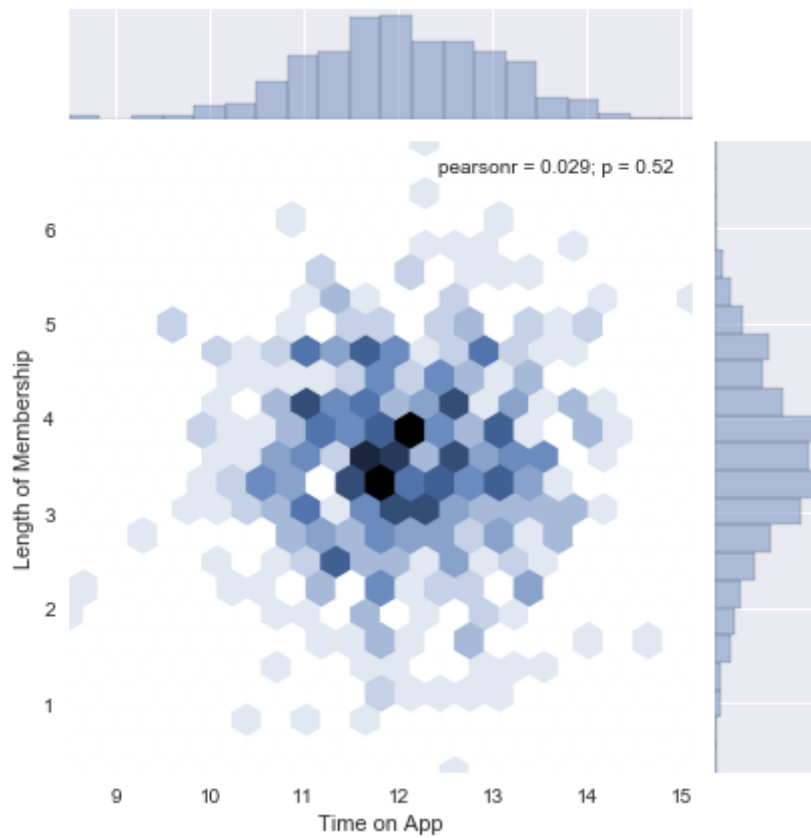
Out[282]: <seaborn.axisgrid.JointGrid at 0x132db5908>



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

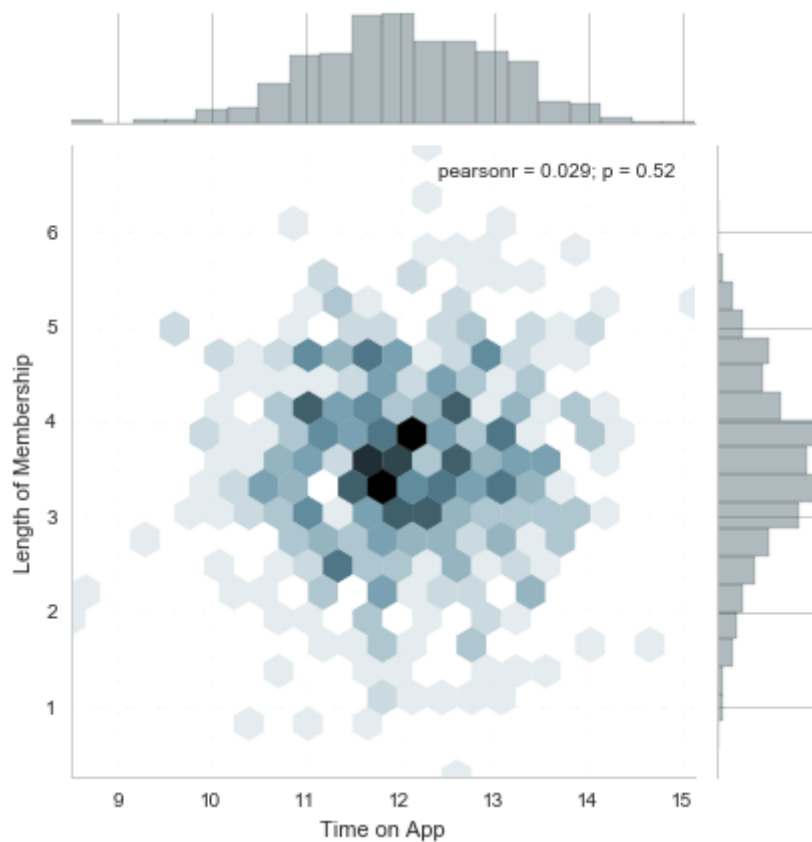

```
In [9]: sns.jointplot(df['Time on App'], df['Length of Membership'], kind='hexbin')
```

```
Out[9]: <seaborn.axisgrid.JointGrid at 0x10ef55c50>
```



```
In [283]:
```

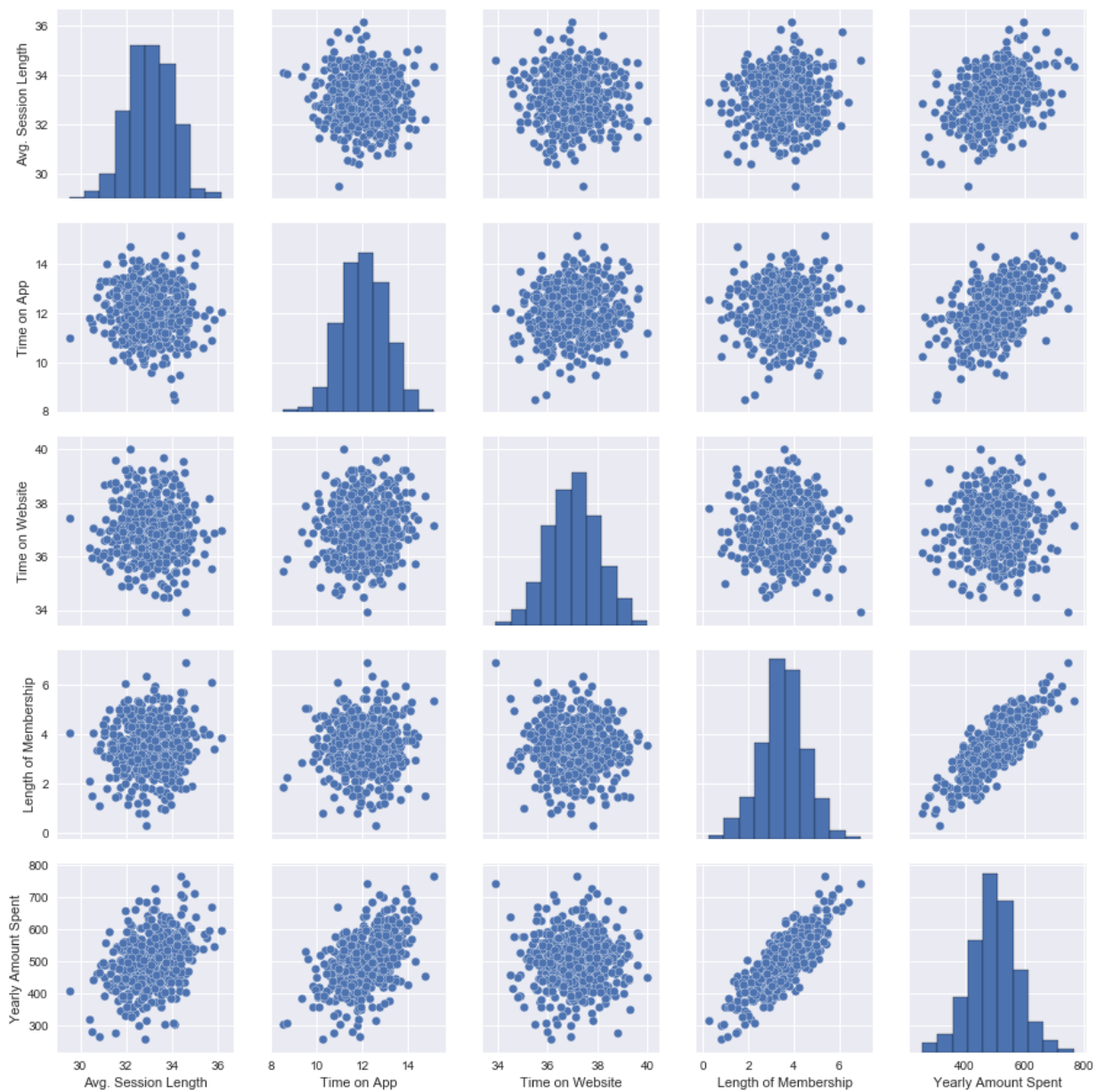
```
Out[283]: <seaborn.axisgrid.JointGrid at 0x130edac88>
```



Let's explore these types of relationships across the entire data set. Use `pairplot` (https://stanford.edu/~mwaskom/software/seaborn/tutorial/axis_grids.html#plotting-pairwise-relationships-with-pairgrid-and-pairplot) to recreate the plot below. (Don't worry about the the colors)

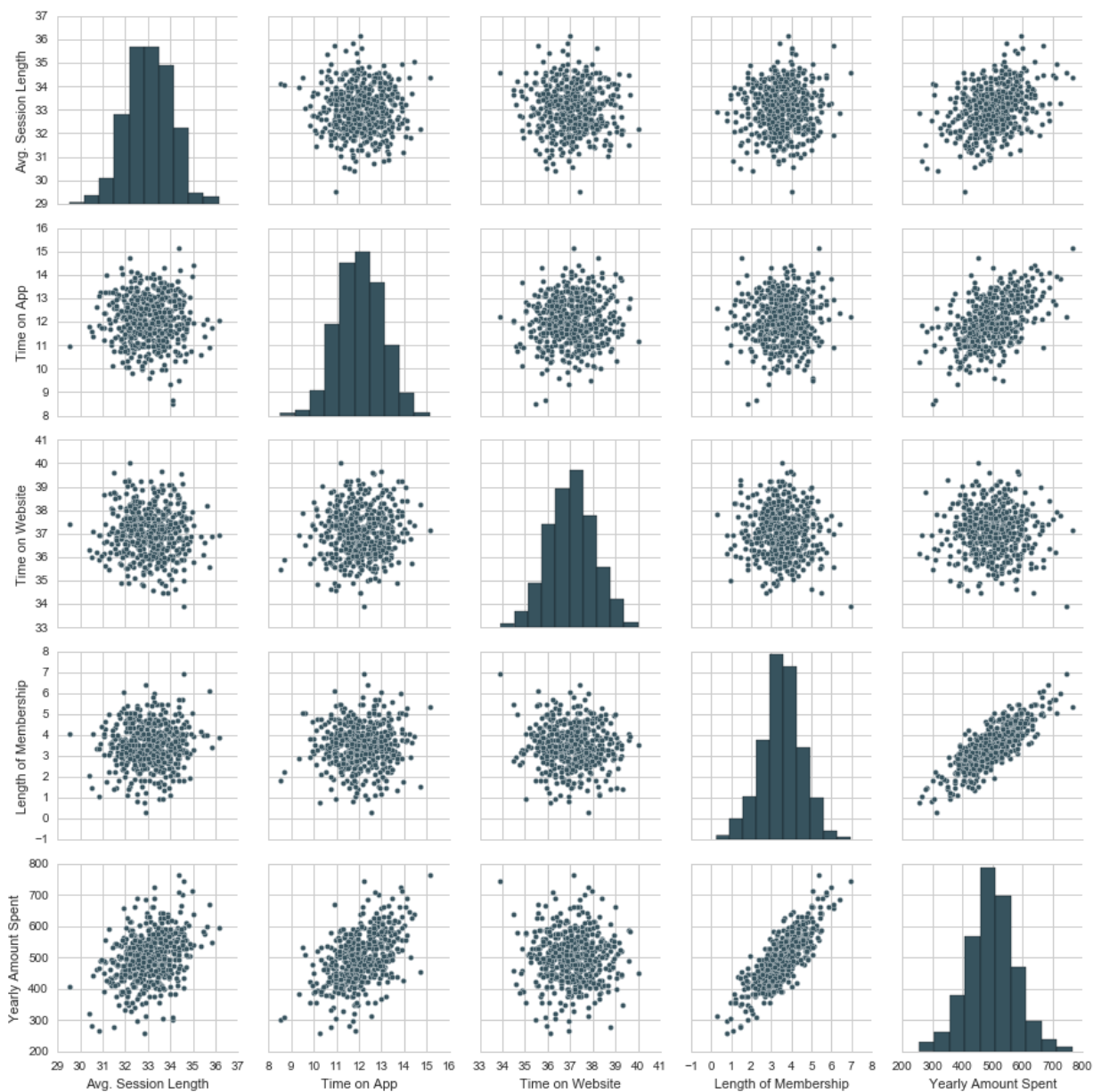
```
In [10]: sns.pairplot(df)
```

```
Out[10]: <seaborn.axisgrid.PairGrid at 0x10958fcf8>
```



```
In [284]:
```

```
Out[284]: <seaborn.axisgrid.PairGrid at 0x132fb3da0>
```



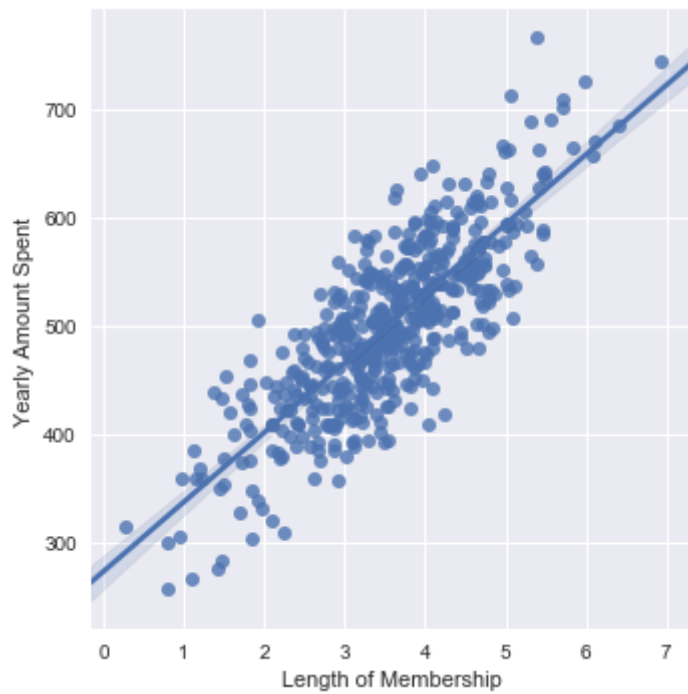
Based off this plot what looks to be the most correlated feature with Yearly Amount Spent?

```
In [285]:
```

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

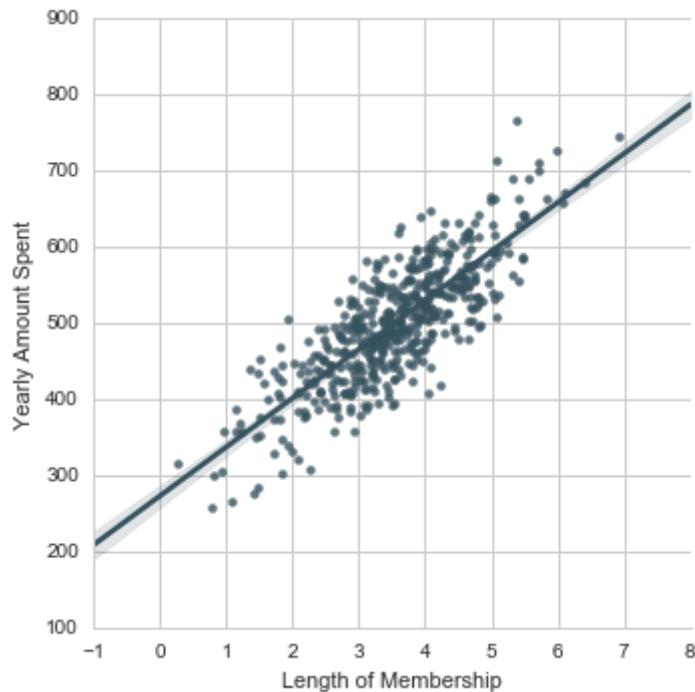
```
In [11]: sns.lmplot(x='Length of Membership', y='Yearly Amount Spent', data=df)
```

```
Out[11]: <seaborn.axisgrid.FacetGrid at 0x10fea7a58>
```



```
In [286]:
```

```
Out[286]: <seaborn.axisgrid.FacetGrid at 0x13538d0b8>
```



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 [12]: from sklearn.model_selection import train_test_split
        from sklearn.linear_model import LinearRegression
        from sklearn import metrics
```

```
In [13]: df.columns
```

```
Out[13]: Index(['Email', 'Address', 'Avatar', 'Avg. Session Length', 'Time on App',
               'Time on Website', 'Length of Membership', 'Yearly Amount Spent'],
              dtype='object')
```

```
In [14]: X = df[['Avg. Session Length', 'Time on App', 'Time on Website', 'Length of
```

```
In [15]: y = df['Yearly Amount Spent']
```

```
In [ ]:
```

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`

```
In [17]: X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.3, ran
```

```
In [290]:
```

Training the Model

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

Import `LinearRegression` from `sklearn.linear_model`

```
In [291]:
```

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

```
In [18]: lm = LinearRegression()
```

Train/fit `lm` on the training data.

```
In [19]: lm.fit(X_train, y_train)
```

```
/usr/local/lib/python3.6/site-packages/scipy/linalg/basic.py:1226: RuntimeWarning: internal gelsd driver lwork query error, required iwork dimension not returned. This is likely the result of LAPACK bug 0038, fixed in LAPACK 3.2.2 (released July 21, 2010). Falling back to 'gelss' driver.
  warnings.warn(mesg, RuntimeWarning)
```

```
Out[19]: LinearRegression(copy_X=True, fit_intercept=True, n_jobs=1, normalize=False)
```

```
In [293]:
```

```
Out[293]: LinearRegression(copy_X=True, fit_intercept=True, n_jobs=1, normalize=False)
```

Print out the coefficients of the model

```
In [20]: print(lm.coef_)
```

```
[25.98154972 38.59015875  0.19040528 61.27909654]
```

```
In [294]:
```

```
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.

```
In [21]: predictions = lm.predict(X_test)
```

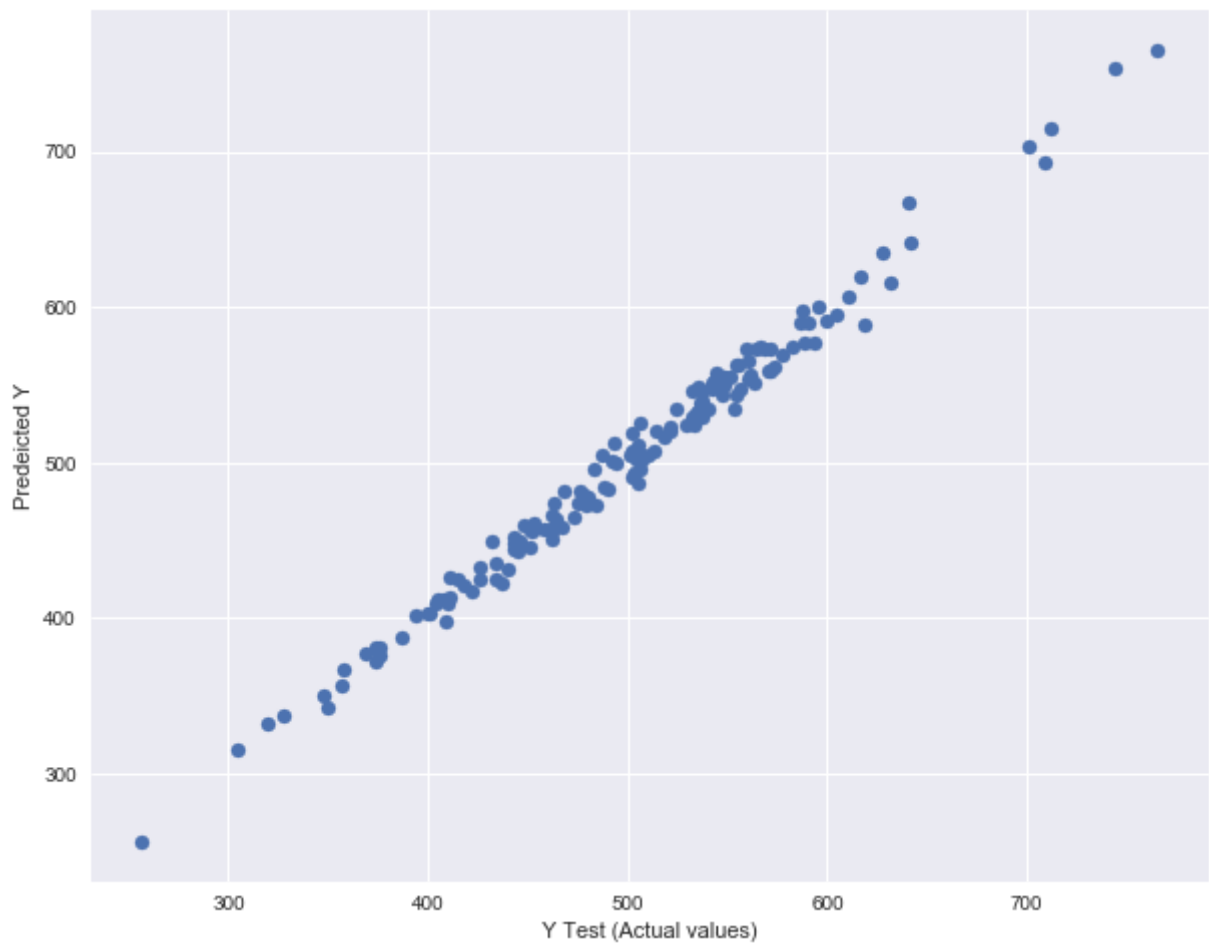
```
In [31]: predictions
```

```
Out[31]: array([456.44186104, 402.72005312, 409.2531539 , 591.4310343 ,
590.01437275, 548.82396607, 577.59737969, 715.44428115,
473.7893446 , 545.9211364 , 337.8580314 , 500.38506697,
552.93478041, 409.6038964 , 765.52590754, 545.83973731,
693.25969124, 507.32416226, 573.10533175, 573.2076631 ,
397.44989709, 555.0985107 , 458.19868141, 482.66899911,
559.2655959 , 413.00946082, 532.25727408, 377.65464817,
535.0209653 , 447.80070905, 595.54339577, 667.14347072,
511.96042791, 573.30433971, 505.02260887, 565.30254655,
460.38785393, 449.74727868, 422.87193429, 456.55615271,
598.10493696, 449.64517443, 615.34948995, 511.88078685,
504.37568058, 515.95249276, 568.64597718, 551.61444684,
356.5552241 , 464.9759817 , 481.66007708, 534.2220025 ,
256.28674001, 505.30810714, 520.01844434, 315.0298707 ,
501.98080155, 387.03842642, 472.97419543, 432.8704675 ,
539.79082198, 590.03070739, 752.86997652, 558.27858232,
523.71988382, 431.77690078, 425.38411902, 518.75571466,
641.9667215 , 481.84855126, 549.69830187, 380.93738919,
555.18178277, 403.43054276, 472.52458887, 501.82927633,
473.5561656 , 456.76720365, 554.74980563, 702.96835044,
534.68884588, 619.18843136, 500.11974127, 559.43899225,
574.8730604 , 505.09183544, 529.9537559 , 479.20749452,
424.78407899, 452.20986599, 525.74178343, 556.60674724,
425.7142882 , 588.8473985 , 490.77053065, 562.56866231,
495.75782933, 445.17937217, 456.64011682, 537.98437395,
367.06451757, 421.12767301, 551.59651363, 528.26019754,
493.47639211, 495.28105313, 519.81827269, 461.15666582,
528.8711677 , 442.89818166, 543.20201646, 350.07871481,
401.49148567, 606.87291134, 577.04816561, 524.50431281,
554.11225704, 507.93347015, 505.35674292, 371.65146821,
342.37232987, 634.43998975, 523.46931378, 532.7831345 ,
574.59948331, 435.57455636, 599.92586678, 487.24017405,
457.66383406, 425.25959495, 331.81731213, 443.70458331,
563.47279005, 466.14764208, 463.51837671, 381.29445432,
411.88795623, 473.48087683, 573.31745784, 417.55430913,
543.50149858, 547.81091537, 547.62977348, 450.99057409,
561.50896321, 478.30076589, 484.41029555, 457.59099941,
411.52657592, 375.47900638])
```

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

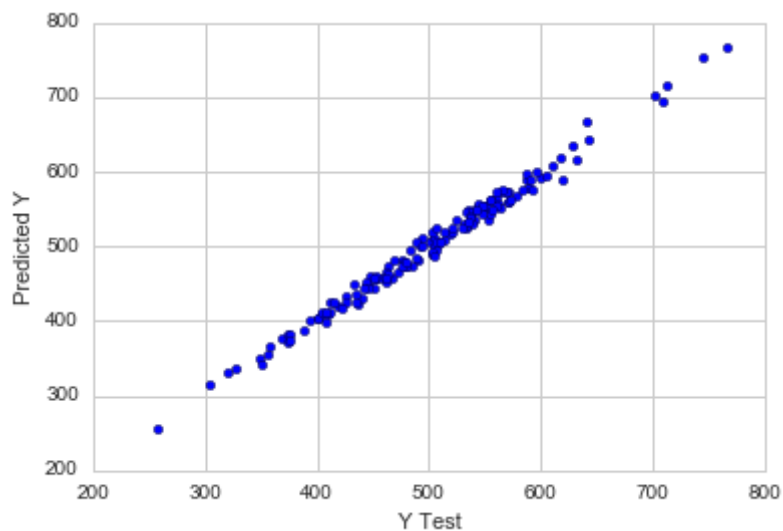

```
In [22]: plt.figure(figsize=(10,8))  
plt.scatter(y_test, predictions)  
plt.xlabel('Y Test (Actual values)')  
plt.ylabel('Predeicted Y')
```

```
Out[22]: Text(0,0.5,'Predeicted Y')
```



```
In [296]:
```

```
Out[296]: <matplotlib.text.Text at 0x135546320>
```



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

```
In [23]: print("MAE: " + str(metrics.mean_absolute_error(y_test, predictions)))
print("MSE: " + str(metrics.mean_squared_error(y_test, predictions)))
print("RMSE: " + str(np.sqrt(metrics.mean_squared_error(y_test, predictions)))
print("R squared (explained variance): " + str(metrics.explained_variance_score(y_test, predictions)))
### See below R-squared is almost 99% so we can explain almost 99% of the variance
```

```
MAE: 7.2281486534308215
MSE: 79.81305165097424
RMSE: 8.933815066978623
R squared (explained variance): 0.9890771231889607
```

```
In [303]:
```

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

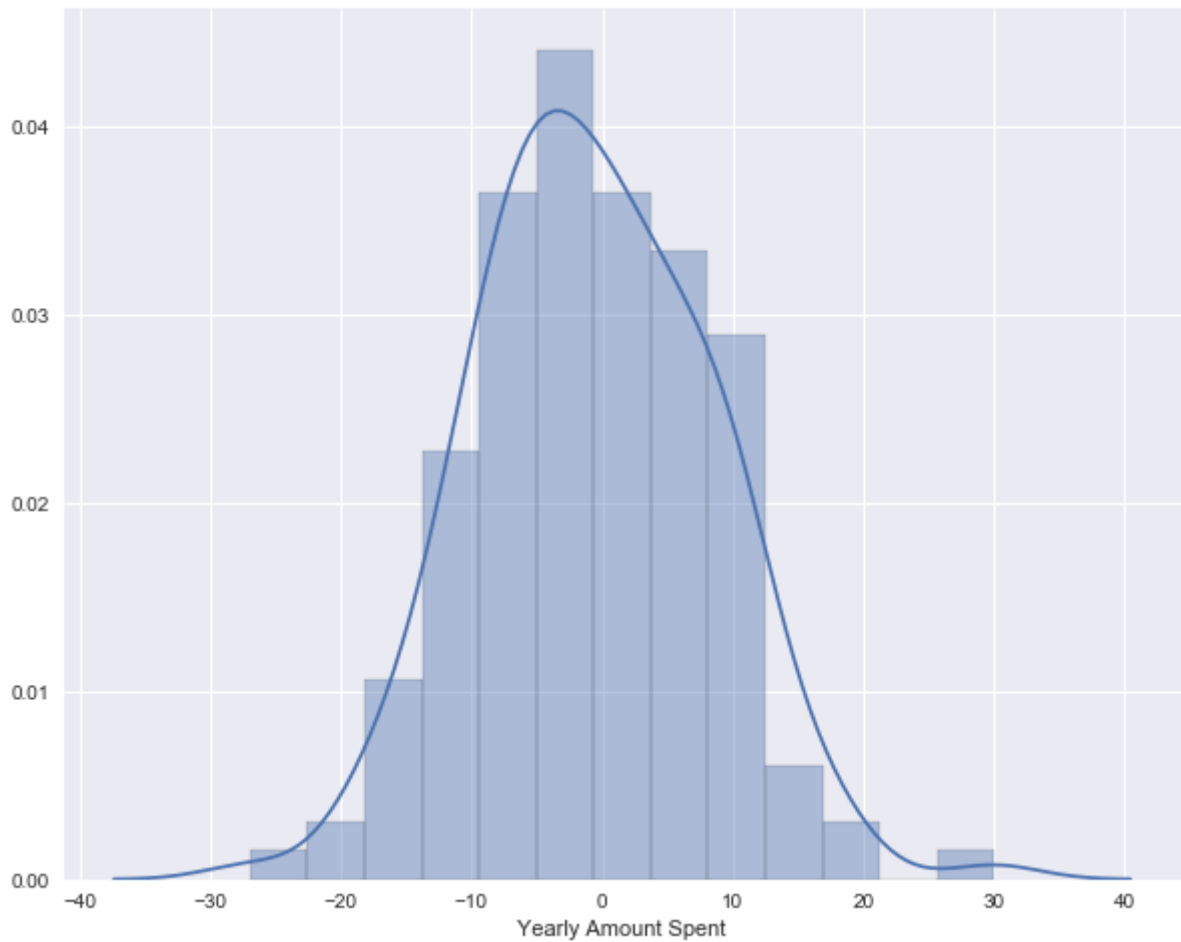
Residuals

You should have gotten a very good model with a good fit. Let's quickly explore the residuals to make sure everything was okay with our data.

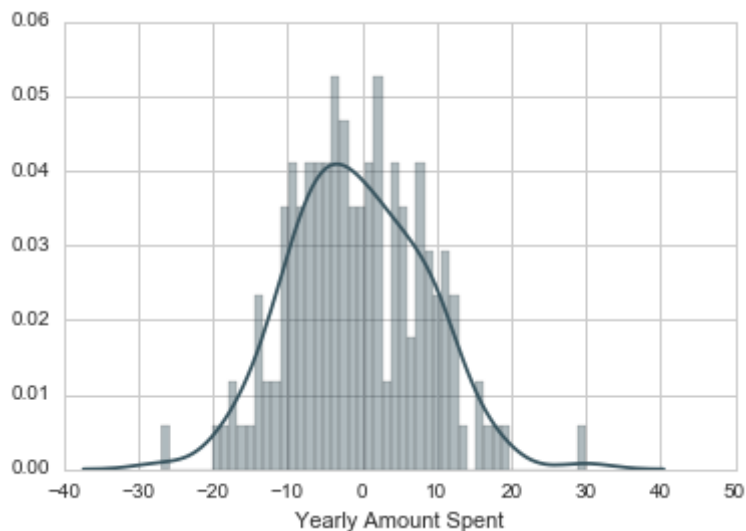
Plot a histogram of the residuals and make sure it looks normally distributed. Use either seaborn distplot, or just plt.hist().

```
In [24]: # Residuals
plt.figure(figsize=(10,8))
sns.distplot((y_test - predictions))
```

Out[24]: <matplotlib.axes._subplots.AxesSubplot at 0x111fc3160>



In [317]:



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 [25]: cdf = pd.DataFrame(lm.coef_, index=X.columns, columns=['Coeff'])
cdf
```

Out[25]:

	Coeff
Avg. Session Length	25.981550
Time on App	38.590159
Time on Website	0.190405
Length of Membership	61.279097

In [298]:

Out[298]:

	Coefficient
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?

Type *Markdown* and LaTeX: α^2

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

Answer here

Great Job!

Congrats on your contract work! The company loved the insights! Let's move on.