

Linear Regression Project

CaseStudy - 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

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

```
In [31]: import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
import seaborn as sns
```

Get the Data

Well work with the Ecommerce Customers csv file from the company. It has Customer info, suchas 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.

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

```
In [32]: #Write your code here
df = read_csv('Ecommerce Customers')
df.head()
```

```
Out[32]:
```

	Email	Address	Avatar	Avg. Session Length	Time on App	Time on Website	Length of Membership	Yearly Amount Spent
0	mloehnerm@hmandes.com	835 Frank TurneyWinghamouth, MI 82180-9605	Vuot	34.497268	12.656651	39.577668	4.082621	587.951054
1	halak@hmad.com	4087 Asher CurrenWingBuckhause, CA 95969-8079	DuGOnen	31.982772	11.109461	37.269899	2.864034	382.204933
2	pat@h@yahoo.com	24645 Valens Unions Suite 525WCCoblesknogr, D	Rinep	33.009015	11.330278	37.110597	4.105453	487.547505
3	twiawebccad@gmail.com	3434 David ThroughwayWhef, Jason, OH 22070-1220	SadSkedrown	34.305557	13.715134	36.721263	3.120179	581.852344
4	mloehnerm@devotion-herman.com	14023 Rodriguez Passag@Port Jacobsville, PR 3	MedumAquaMaine	33.330673	12.795189	37.536653	4.446308	599.405092

#This is the output, DO NOT RUN THIS CELL OR ELSE OUTPUT WILL BE OVERRITTEN

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

```
In [4]: # WRITE YOUR CODE HERE
df.info()
```

```
Out[4]:
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 589 entries, 0 to 499
Data columns (total 8 columns):
 #   Column              Non-Null Count  Dtype
---  -
 0   Email               589 non-null    object
 1   Address             589 non-null    object
 2   Avatar              589 non-null    object
 3   Avg. Session Length 589 non-null    float64
 4   Time on App         589 non-null    float64
 5   Time on Website     589 non-null    float64
 6   Length of Membership 589 non-null    float64
 7   Yearly Amount Spent 589 non-null    float64
dtypes: float64(7), object(1)
memory usage: 31.4+ KB
```

#This is the output, DO NOT RUN THIS CELL OR ELSE OUTPUT WILL BE OVERRITTEN

```
In [6]: #WRITE YOUR CODE HERE
df.describe()
```

```
Out[6]:
```

	Avg. Session Length	Time on App	Time on Website	Length of Membership	Yearly Amount Spent
count	590.000000	590.000000	590.000000	590.000000	590.000000
mean	33.003384	12.030488	37.060545	3.532482	499.316028
std	0.992563	0.994216	1.015489	0.992778	79.314782
min	29.932429	8.508152	33.913847	0.296901	256.670682
25%	32.341822	11.388153	36.348575	2.830450	445.038277
50%	33.082008	11.983231	37.068367	3.538975	498.887875
75%	33.713195	12.733950	37.716432	4.126502	549.313828
max	36.139652	15.126994	40.005182	6.922089	765.818462

#This is the output, DO NOT RUN THIS CELL OR ELSE OUTPUT WILL BE OVERRITTEN

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.

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

```
In [8]: #WRITE YOUR CODE HERE
plt.figure(figsize=(8,6))
plt.xlabel('Time on Website')
plt.ylabel('Yearly amount spent')
plt.scatter(x=df['Time on Website'],y=df['Yearly Amount Spent'])
plt.show()
```

#This is the output, DO NOT RUN THIS CELL OR ELSE OUTPUT WILL BE OVERRITTEN

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

```
In [10]: #WRITE YOUR CODE HERE
plt.figure(figsize=(8,6))
plt.xlabel('Time on App')
plt.ylabel('Yearly amount spent')
plt.scatter(x=df['Time on App'],y=df['Yearly Amount Spent'])
plt.show()
```

#This is the output, DO NOT RUN THIS CELL OR ELSE OUTPUT WILL BE OVERRITTEN

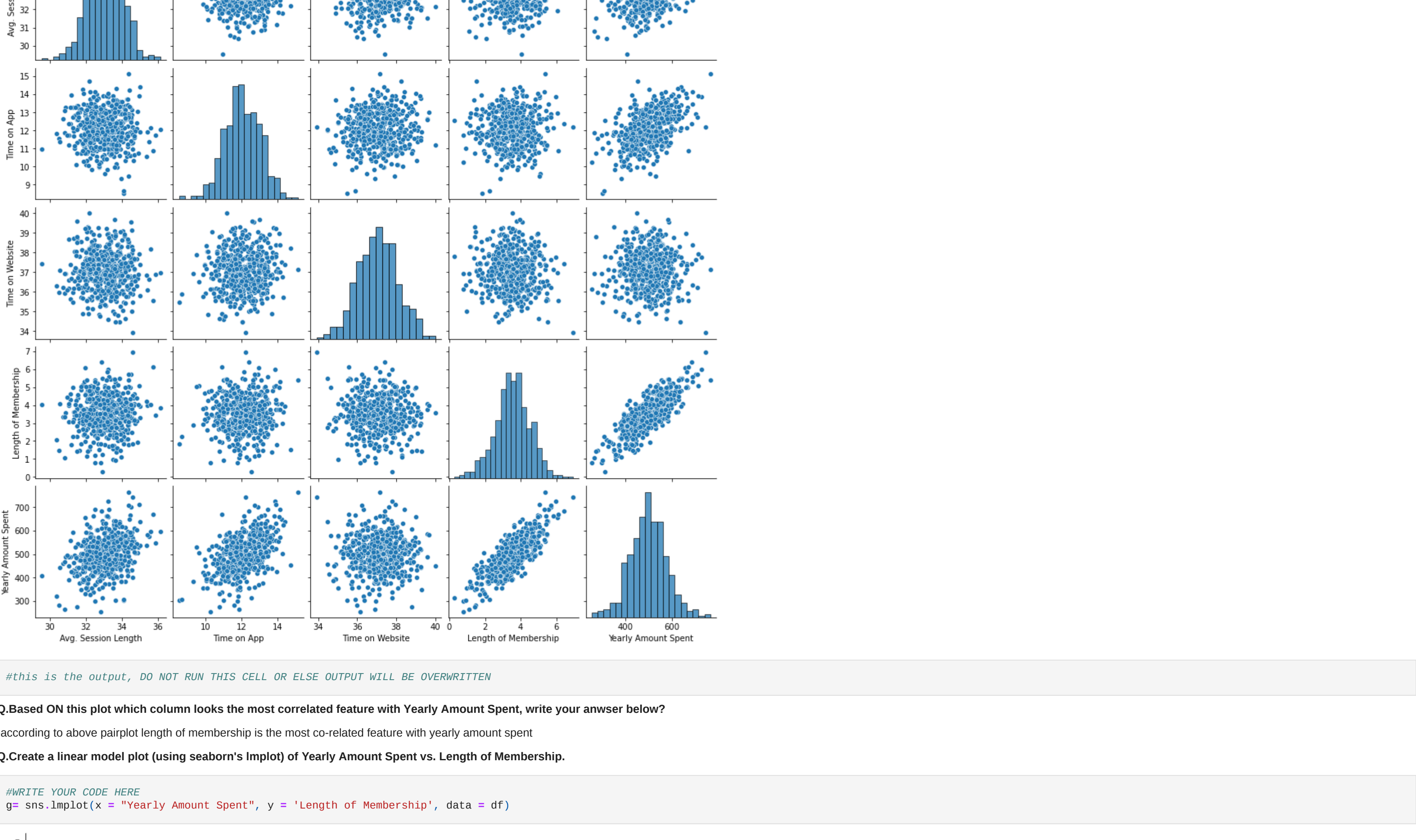
QUse jointplot to create a 2D hex bin plot comparing Time on App and Length of Membership.

```
In [12]: #WRITE YOUR CODE HERE
sns.jointplot(x = df['Time on App'], y = df['Length of Membership'],
              kind = "hex", figsize=(8,6))
plt.xlabel('Time on app')
plt.ylabel('Length of Membership')
plt.show()
```

```
In [13]: #This is the output, DO NOT RUN THIS CELL OR ELSE OUTPUT WILL BE OVERRITTEN
```

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

```
In [14]: #WRITE YOUR CODE HERE
sns.pairplot(df)
```



#This is the output, DO NOT RUN THIS CELL OR ELSE OUTPUT WILL BE OVERRITTEN

QBased ON this plot which column looks the most correlated feature with Yearly Amount Spent, write your answer below?

According to above pairplot length of membership is the most co-related feature with yearly amount spent

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

```
In [16]: #WRITE YOUR CODE HERE
g = sns.lmplot(x = "Yearly Amount Spent", y = "Length of Membership", data = df)
```

#This is the output, DO NOT RUN THIS CELL OR ELSE OUTPUT WILL BE OVERRITTEN

Training and Testing Data

Now that we've explored the data a bit, lets 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 [18]: #WRITE YOUR CODE HERE
df=df.select_dtypes(float)
df
```

```
Out[18]:
```

	Avg. Session Length	Time on App	Time on Website	Length of Membership	Yearly Amount Spent
0	34.497268	12.656651	39.577668	4.082621	587.951054
1	31.982772	11.109461	37.269899	2.864034	382.204933
2	33.009015	11.330278	37.110597	4.105453	487.547505
3	34.305557	13.715134	36.721263	3.120179	581.852344
4	33.330673	12.795189	37.536653	4.446308	599.405092
...
495	33.327660	13.566160	36.417895	3.746573	573.947439
496	34.702529	11.695736	37.190268	3.576526	529.049004
497	32.646777	11.499409	38.332576	4.968204	551.620145
498	33.322501	12.393423	36.840086	2.336485	456.469510
499	33.715981	12.418908	35.771016	2.735100	497.778642

500 rows x 5 columns

#This is the output, DO NOT RUN THIS CELL OR ELSE OUTPUT WILL BE OVERRITTEN

```
In [20]: # DIVIDE THE DATASET INTO X & Y
x=df.iloc[:,1:]
y=df.iloc[:,0]
x
```

```
Out[20]:
```

	Avg. Session Length	Time on App	Time on Website	Length of Membership
0	34.497268	12.656651	39.577668	4.082621
1	31.982772	11.109461	37.269899	2.864034
2	33.009015	11.330278	37.110597	4.105453
3	34.305557	13.715134	36.721263	3.120179
4	33.330673	12.795189	37.536653	4.446308
...
495	33.327660	13.566160	36.417895	3.746573
496	34.702529	11.695736	37.190268	3.576526
497	32.646777	11.499409	38.332576	4.968204
498	33.322501	12.393423	36.840086	2.336485
499	33.715981	12.418908	35.771016	2.735100

500 rows x 4 columns

```
In [21]: y
```

```
Out[21]:
```

0	587.951054
1	382.204933
2	487.547505
3	581.852344
4	599.405092
...	...
495	573.947439
496	529.049004
497	551.620145
498	456.469510
499	497.778642

Names: Yearly Amount Spent, Length: 500, dtype: float64

Q Split the data into training and testing sets. Set test_size=0.3 and random_state=101

```
In [22]: #import here
from sklearn.model_selection import train_test_split as t
```

```
Out[22]:
```

	Avg. Session Length	Time on App	Time on Website	Length of Membership
0	34.497268	12.656651	39.577668	4.082621
1	31.982772	11.109461	37.269899	2.864034
2	33.009015	11.330278	37.110597	4.105453
3	34.305557	13.715134	36.721263	3.120179
4	33.330673	12.795189	37.536653	4.446308
...
495	33.327660	13.566160	36.417895	3.746573
496	34.702529	11.695736	37.190268	3.576526
497	32.646777	11.499409	38.332576	4.968204
498	33.322501	12.393423	36.840086	2.336485
499	33.715981	12.418908	35.771016	2.735100

500 rows x 4 columns

```
In [23]: y
```

```
Out[23]:
```

0	587.951054
1	382.204933
2	487.547505
3	581.852344
4	599.405092
...	...
495	573.947439
496	529.049004
497	551.620145
498	456.469510
499	497.778642

Names: Yearly Amount Spent, Length: 500, dtype: float64

Q Split the data into training and testing sets. Set test_size=0.3 and random_state=101

```
In [24]: #import here
from sklearn.model_selection import train_test_split as t
```

```
Out[24]:
```

	Avg. Session Length	Time on App	Time on Website	Length of Membership
0	34.497268	12.656651	39.577668	4.082621
1	31.982772	11.109461	37.269899	2.864034
2	33.009015	11.330278	37.110597	4.105453
3	34.305557	13.715134	36.721263	3.120179
4	33.330673	12.795189	37.536653	4.446308
...
495	33.327660	13.566160	36.417895	3.746573
496	34.702529	11.695736	37.190268	3.576526
497	32.646777	11.499409	38.332576	4.968204
498	33.322501	12.393423	36.840086	2.336485
499	33.715981	12.418908	35.771016	2.735100

500 rows x 4 columns

```
In [25]: # DIVIDE THE DATASET INTO X & Y
x=df.iloc[:,1:]
y=df.iloc[:,0]
x
```

```
Out[25]:
```

	Avg. Session Length	Time on App	Time on Website	Length of Membership
0	34.497268	12.656651	39.577668	4.082621
1	31.982772	11.109461	37.269899	2.864034
2	33.009015	11.330278	37.110597	4.105453
3	34.305557	13.715134	36.721263	3.120179
4	33.330673	12.795189	37.536653	4.446308
...
495	33.327660	13.566160	36.417895	3.746573
496	34.702529	11.695736	37.190268	3.576526
497	32.646777	11.499409	38.332576	4.968204
498	33.322501	12.393423	36.840086	2.336485
499	33.715981	12.418908	35.771016	2.735100

500 rows x 4 columns

```
In [26]: # DIVIDE THE DATASET INTO X & Y
x=df.iloc[:,1:]
y=df.iloc[:,0]
x
```

```
Out[26]:
```

	Avg. Session Length	Time on App	Time on Website	Length of Membership
0	34.497268	12.656651	39.577668	4.082621
1	31.982772	11.109461	37.269899	2.864034
2	33.009015	11.330278	37.110597	4.105453
3	34.305557	13.715134	36.721263	3.120179
4	33.330673	12.795189	37.536653	4.446308
...
495	33.327660	13.566160	36.417895	3.746573
496	34.702529	11.695736	37.190268	3.576526
497	32.646777	11.499409	38.332576	4.968204
498	33.322501	12.393423	36.840086	2.336485
499	33.715981	12.418908	35.771016	2.735100

500 rows x 4 columns

```
In [27]: # DIVIDE THE DATASET INTO X & Y
x=df.iloc[:,1:]
y=df.iloc[:,0]
x
```

```
Out[27]:
```

	Avg. Session Length	Time on App	Time on Website	Length of Membership
0	34.497268	12.656651	39.577668	4.082621
1	31.982772	11.109461	37.269899	2.864034
2	33.009015	11.330278	37.110597	4.105453
3	34.305557	13.715134	36.721263	3.120179
4	33.330673	12.795189	37.536653	4.446308
...
495	33.327660	13.566160	36.417895	3.746573
496	34.702529	11.695736	37.190268	3.576526
497	32.646777	11.499409	38.332576	4.968204
498	33.322501	12.393423	36.840086	2.336485
499	33.715981	12.418908	35.771016	2.735100

500 rows x 4 columns

```
In [28]: # DIVIDE THE DATASET INTO X & Y
x=df.iloc[:,1:]
y=df.iloc[:,0]
x
```

```
Out[28]:
```

	Avg. Session Length	Time on App	Time on Website	Length of Membership
0	34.497268	12.656651	39.577668	4.082621
1	31.982772	11.109461	37.269899	2.864034
2	33.009015	11.330278	37.110597	4.105453
3	34.305557	13.715134	36.721263	3.120179
4	33.330673	12.795189	37.536653	4.446308
...
495	33.327660	13.566160	36.417895	3.746573
496	34.702529	11.695736	37.190268	3.576526
497	32.646777	11.499409	38.332576	4.968204
498	33.322501	12.393423	36.840086	2.336485
499	33.715981	12.418908	35.771016	2.735100

500 rows x 4 columns

```
In [29]: # DIVIDE THE DAT
```