

E-commerce Store Analysis

Getting the data

```
In [2]: import pandas as pd
```

```
data = pd.read_csv("Ecommerce Customers")
data
```

```
Out[2]:
```

	Email	Address	Avatar	Avg. Session Length	Time on App	Time on Website	Length of Membership	Yearly Amount Spent
0	mstephenson@fernandez.com	835 Frank TunnelWrightmouth, MI 82180-9605	Violet	34.497268	12.656561	39.577668	4.082621	587.951054
1	hduke@hotmail.com	4547 Archer CommonrCazachester, CA 06566-8576	DarkGreen	31.926272	11.109461	37.268959	2.664034	392.204933
2	palen@yahoo.com	24645 Valerie Livons Suite 582rCobborough, D...	Bisque	33.000915	11.330278	37.110597	4.104543	487.547505
3	riverarebecca@gmail.com	1414 David ThoroughwayrPort Jason, OH 22070-1220	SaddleBrown	34.305557	13.717514	36.721283	3.120179	581.852344
4	mstephens@daivdon-heman.com	14023 Rodriguez PassagelrPort Jacobville, PR 3...	MediumAquaMarine	33.330673	12.795189	37.536653	4.446308	599.406092
...
495	lewissessica@craig-evans.com	4483 Jones Motorway Suite 872nLake Jamelurt...	Tan	33.237660	13.566160	36.417985	3.746573	573.847438
496	katriar56@gmail.com	172 Owen Divide Suite 497rWest Richard, CA 19320	PaleVioletRed	34.702529	11.695736	37.190288	3.576526	529.049004
497	daeb8@hotmail.com	0787 Andrews Ranch Apt. 633rSouth Chadburgh...	Cornsik	32.646777	11.499409	38.332576	4.958264	551.620145
498	owilson@hotmail.com	680 Jennifer Lodge Apt. 808rBrendachester, TX...	Teal	33.322501	12.391423	36.840086	2.336485	456.469510
499	hennahelison@daivdon.com	49791 Rachel Heights Apt. 889rEast Drewboroug...	DarkMagenta	33.715961	12.418808	35.771016	2.735160	497.778642

500 rows x 8 columns

```
In [3]: # check for any null values
data.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
```

```
In [4]: # delete the email address and avatar columns
data = data.drop(["Email", "Address", "Avatar"], axis=1)
```

```
In [7]: # Split the data into train and test data
(train_data = data[:int(500*0.8)])
(test_data = data[int(500*0.8):])
```

```
Out[7]:
```

	Avg. Session Length	Time on App	Time on Website	Length of Membership	Yearly Amount Spent
400	33.172331	13.078692	37.329819	5.405406	663.074818
401	33.247322	11.956426	36.517348	3.451751	506.375987
402	33.598913	13.252737	37.305961	2.935577	528.419330
403	33.085298	13.093537	38.315648	4.759360	632.123588
404	32.278443	12.527472	36.696367	3.531402	488.270298
...
495	33.237660	13.566160	36.417985	3.746573	573.847438
496	34.702529	11.695736	37.190288	3.576526	529.049004
497	32.646777	11.499409	38.332576	4.958264	551.620145
498	33.322501	12.391423	36.840086	2.336485	456.469510
499	33.715961	12.418808	35.771016	2.735160	497.778642

100 rows x 5 columns

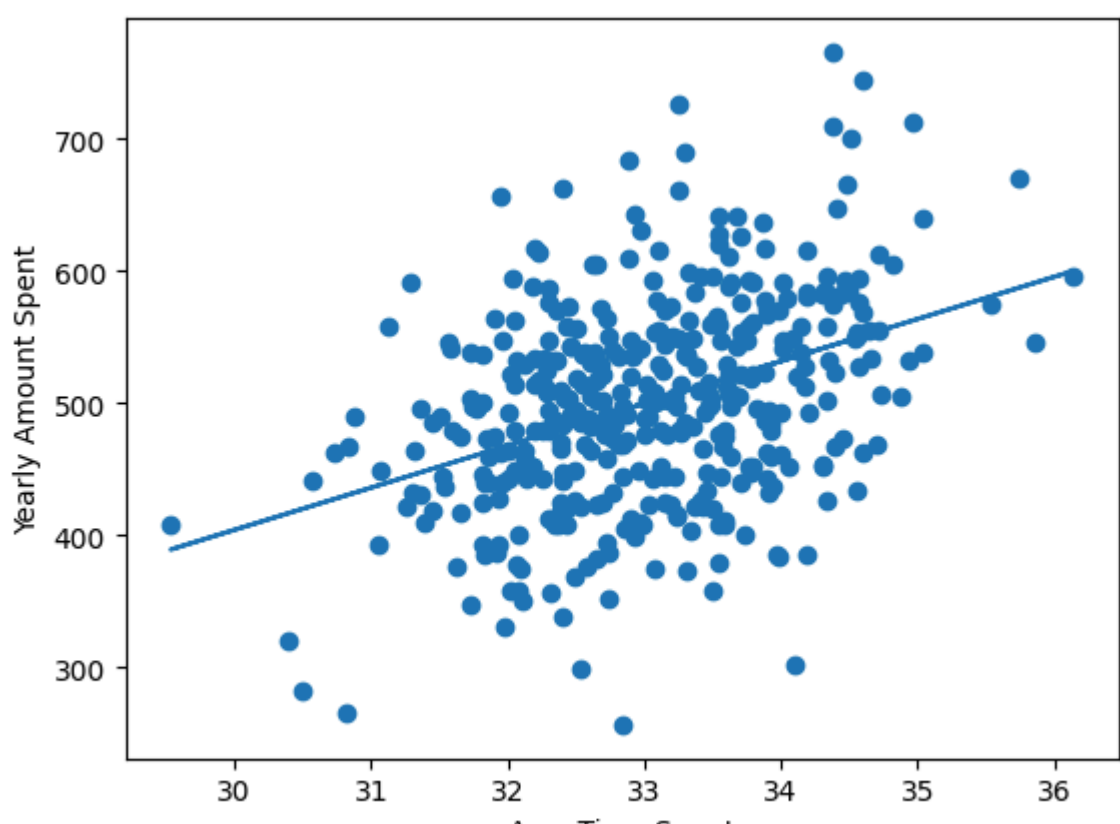
Check relationship between Avg. Session Length and Yearly Amount Spent

```
In [28]: import matplotlib.pyplot as plt
from scipy import stats

x = train_data["Avg. Session Length"]
y = train_data["Yearly Amount Spent"]
slope, intercept, r, p, std_err = stats.linregress(x,y)

def myfun(x):
    return slope*x + intercept

model = list(map(myfun, x))
plt.scatter(x, y)
plt.plot(x,model)
plt.xlabel("Avg. Time Spent")
plt.ylabel("Yearly Amount Spent")
plt.show()
print("The coefficient of correlation is: ", r)
```



The coefficient of correlation is: 0.39726338728882665

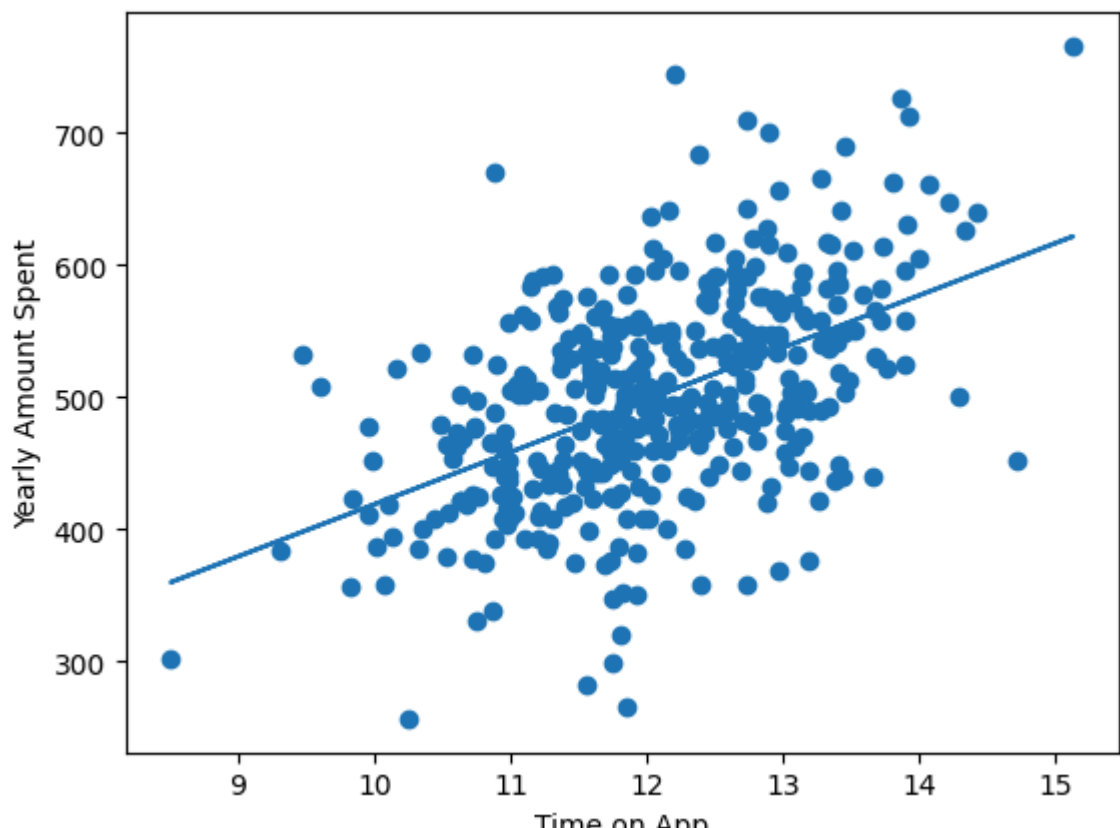
The coefficient is only 0.39 showing that there is not good relationship between avg. time spent and yearly amount spent. Therefore, it is not a good indicator

Check relationship between Time on App and Yearly Amount Spent

```
In [29]: x = train_data["Time on App"]
y = train_data["Yearly Amount Spent"]
slope, intercept, r, p, std_err = stats.linregress(x,y)

def myfun(x):
    return slope*x + intercept

model = list(map(myfun, x))
plt.scatter(x, y)
plt.plot(x,model)
plt.xlabel("Time on App")
plt.ylabel("Yearly Amount Spent")
plt.show()
print("The coefficient of correlation is: ", r)
```



The coefficient of correlation is: 0.5813178168518133

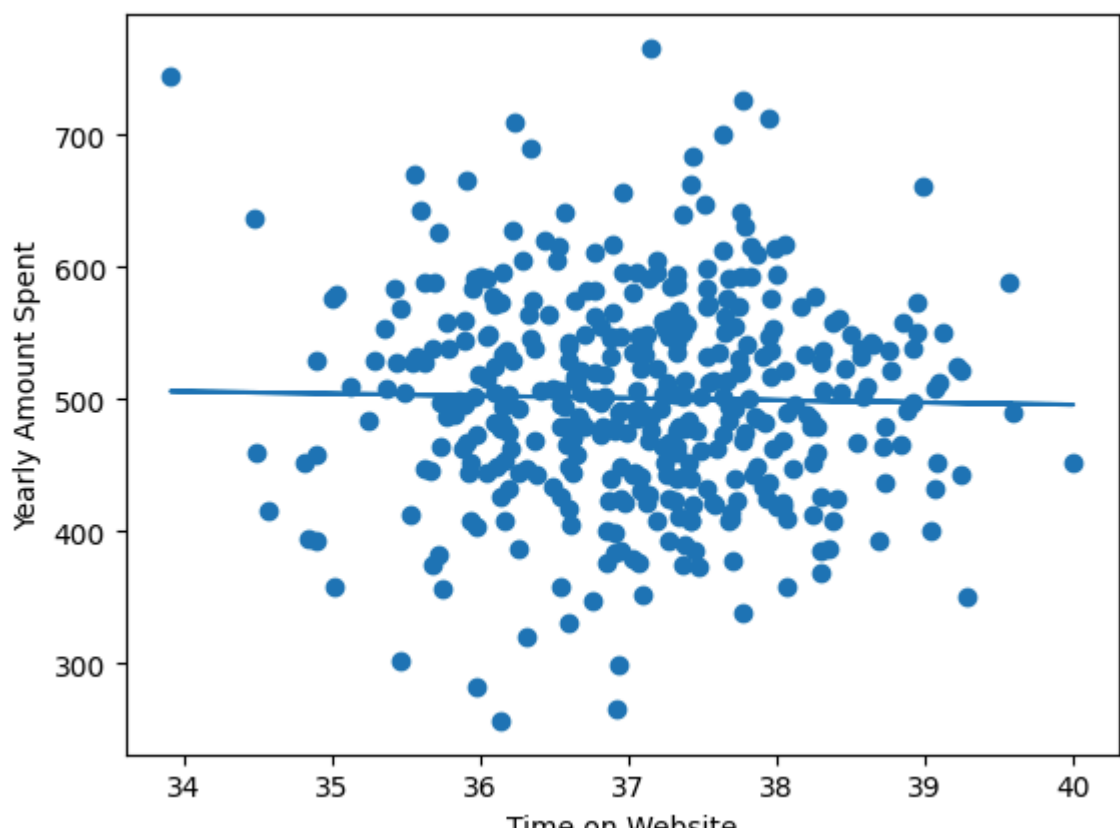
The coefficient is 0.5. Time on App is not a perfect indicator for predicting the Amount spent but it is certainly better indicator than Avg. Time Spent.

Check relationship between Time on Website and Yearly Amount Spent

```
In [30]: x = train_data["Time on Website"]
y = train_data["Yearly Amount Spent"]
slope, intercept, r, p, std_err = stats.linregress(x,y)

def myfun(x):
    return slope*x + intercept

model = list(map(myfun, x))
plt.scatter(x, y)
plt.plot(x,model)
plt.xlabel("Time on Website")
plt.ylabel("Yearly Amount Spent")
plt.show()
print("The coefficient of correlation is: ", r)
```



The coefficient of correlation is: -0.8213893861369613256

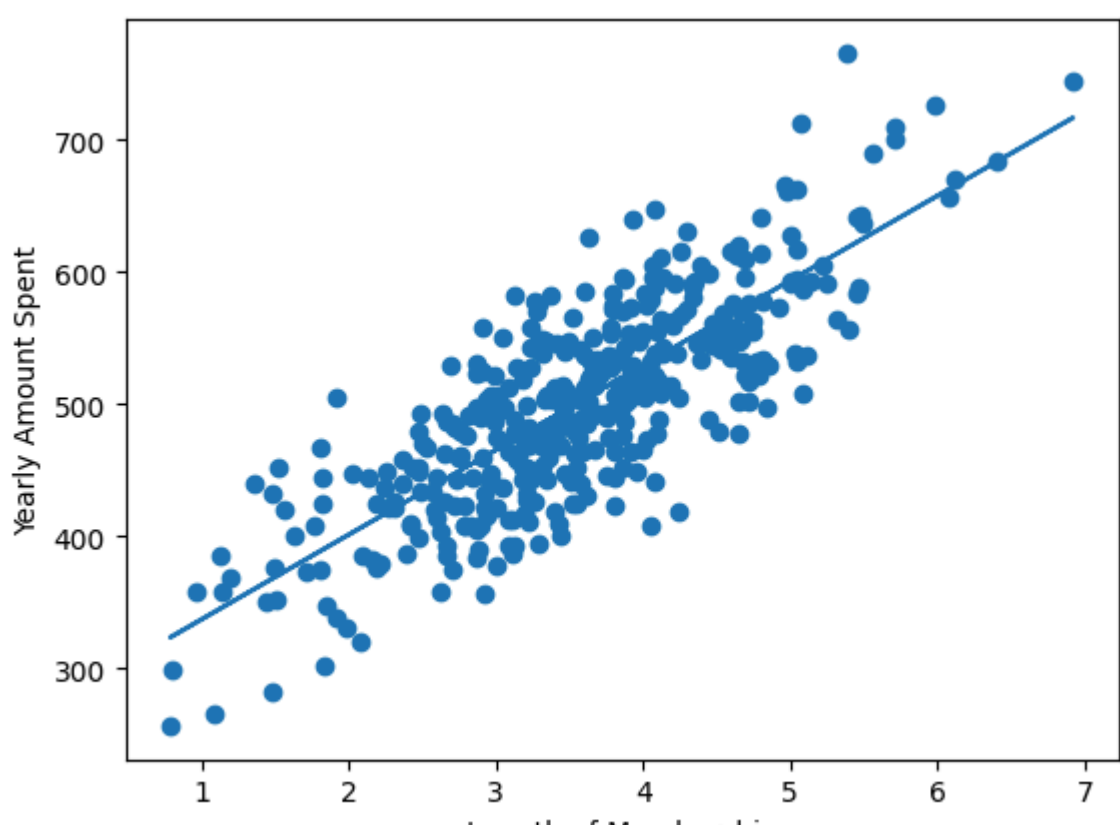
The coefficient of correlation is very low indicating that Time spent on websites by customers does not show any relation to how much the customers will spend money.

Check relationship between Length of Membership and Yearly Amount Spent

```
In [31]: x = train_data["Length of Membership"]
y = train_data["Yearly Amount Spent"]
slope, intercept, r, p, std_err = stats.linregress(x,y)

def myfun(x):
    return slope*x + intercept

model = list(map(myfun, x))
plt.scatter(x, y)
plt.plot(x,model)
plt.xlabel("Length of Membership")
plt.ylabel("Yearly Amount Spent")
plt.show()
print("The coefficient of correlation is: ", r)
```



The coefficient of correlation is: 0.8823934262919785

Length of Membership definitely shows a significant impact on how much a customer might spend money on products. The people who have been a member for a longer time are likely to spend more money than other people. Length of Membership is by far the best indicator for predicting the customer's yearly expenditure.

```
In [32]: train_data
```

	Avg. Session Length	Time on App	Time on Website	Length of Membership	Yearly Amount Spent
0	34.497268	12.656561	39.577668	4.082621	587.951054
1	31.926272	11.109461	37.268959	2.664034	392.204933
2	33.000915	11.330278	37.110597	4.104543	487.547505
3	34.305557	13.717514	36.721283	3.120179	581.852344
4	33.330673	12.795189	37.536653	4.446308	599.406092
...
395	31.445972	12.846499	37.869217	3.420150	484.876965
396	35.742670	10.889828	35.565436	6.115199	669.987141
397	34.012619	12.914570	36.046204	3.488030	547.709989
398	34.140393	11.568527	38.918749	4.082855	537.825382
399	32.377990	11.971751	37.199388	2.829700	408.216902

400 rows x 5 columns

Now let's train our model using multi-variable regression and try to predict the yearly amount spent

```
In [46]: from sklearn import linear_model

x = train_data[["Avg. Session Length", "Time on App", "Time on Website", "Length of Membership"]]
y = train_data["Yearly Amount Spent"]

# train the model
reg = linear_model.LinearRegression()
reg.fit(x,y)

cdf = pd.DataFrame(reg.coef_ , x.columns , columns=["Coef" ])
cdf
```

```
Out[46]:
```

	Coef
Avg. Session Length	25.576288
Time on App	38.345474
Time on Website	0.593094
Length of Membership	61.182253

We can say that Time on Website is the worst predictor and Length of Membership is the best predictor. And as Time on App seems to have more impact in whether a customer buys a product, the company should put more focus on making the app more user-friendly and improve the customer experience. They should also look at why website sells are not performing that well and try to solve and improve customer experience on the website as well to improve their sells. They should definitely focus on making their loyalty and membership programs more attractive to the customers.

Let's predict the Yearly Amount Spent using our test data

```
In [49]: test = test_data[["Avg. Session Length", "Time on App", "Time on Website", "Length of Membership"]]
result = reg.predict(test)
```

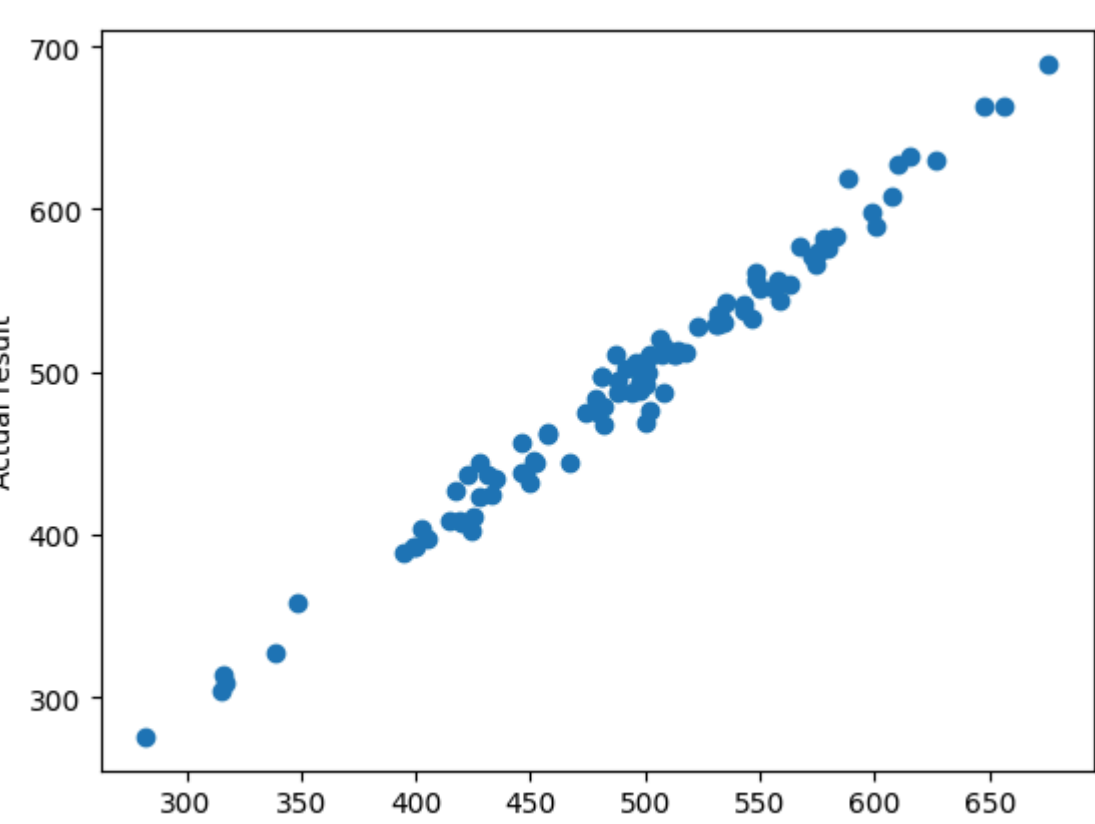
```
Out[49]: array([556.53182746, 495.4856683 , 522.99217893, 615.38255148,
        497.4964398 , 598.76291281, 425.28748928, 419.81697395,
        481.74527484, 518.09632796, 607.98569662, 600.85228474,
        467.37139755, 497.76514136, 546.16417418, 282.28888956,
        597.2568919 , 446.14773139, 502.25551065, 478.21939216,
        647.96177913, 558.33921848, 626.93993935, 457.69487285,
        497.5862481 , 576.1633185 , 534.55273655, 578.18656094,
        597.78976655, 491.18624097, 547.92338666, 479.58629829,
        598.25897821, 454.88185925, 514.92189784, 572.89624926,
        585.96376383, 451.01628418, 599.60248079, 574.38348986,
        599.71533089, 486.93836324, 516.97439683, 548.08992539,
        427.87789982, 513.78459653, 531.20891959, 315.73224449,
        478.909599 , 452.07878091, 473.8909771 , 431.08979947,
        597.78939966, 486.19878795, 449.67649253, 447.94855879,
        594.0864566 , 511.91616161, 542.55778388, 419.60224666,
        587.88989154, 495.69260639, 405.4429639 , 399.05999346,
        671.86293757, 535.06617611, 567.29218779, 422.36637356,
        562.92595229, 417.08859786, 433.48601325, 542.96189902,
        598.18376077, 427.85248271, 499.88484978, 532.2166441 ,
        413.31678871, 493.68629989, 487.64388829, 422.87121686,
        549.48083273, 488.78388887, 487.93878276, 481.44727885,
        457.88527012, 508.34325912, 579.67437459, 348.68989291,
        598.54885678, 338.48788318, 502.15572929, 513.24377263,
        482.28686828, 618.08871328, 512.6449285 , 574.86288351,
        538.60078222, 555.76538226, 445.96493687, 489.83964621])
```

Let's compare our predicted data with the actual given data

```
In [50]: actual_result = test_data["Yearly Amount Spent"]
```

```
plt.scatter(result , actual_result)
plt.xlabel("Predicted result")
plt.ylabel("Actual result")
```

```
Out[50]: Text(0, 0.5, 'Actual result')
```



```
In [49]: from sklearn.metrics import mean_absolute_error , mean_squared_error
import math

print("Mean absolute error: " , mean_absolute_error(result , actual_result))
print("Mean squared error: " , mean_squared_error(result , actual_result))
print("Root mean squared error: " , math.sqrt(mean_squared_error(result , actual_result)))

Mean absolute error: 0.22257352553893
Mean squared error: 112.5193969810785
Root mean squared error: 10.60751749698788

Residuals analysis
```

```
In [50]: residual = actual_result - result
plt.hist(residual , bins=10)
```

```
Out[50]: array([-30.8681242 , -24.54458389, -16.4488555 , -12.2726330 ,
        -6.89744297,  0.87797734,  6.25359765, 12.42921786,
        18.80483827, 24.78845858, 30.95687889]),
<BarContainer object of 38 artists>)
```

