

# Logistic Regression Project

In this project we will be working with a fake advertising data set, indicating whether or not a particular internet user clicked on an Advertisement. We will try to create a model that will predict whether or not they will click on an ad based off the features of that user.

This data set contains the following features:

- 'Daily Time Spent on Site': consumer time on site in minutes
- 'Age': customer age in years
- 'Area Income': Avg. Income of geographical area of consumer
- 'Daily Internet Usage': Avg. minutes a day consumer is on the internet
- 'Ad Topic Line': Headline of the advertisement
- 'City': City of consumer
- 'Male': Whether or not consumer was male
- 'Country': Country of consumer
- 'Timestamp': Time at which consumer clicked on Ad or closed window
- 'Clicked on Ad': 0 or 1 indicated clicking on Ad

## Import Libraries

Import a few libraries you think you'll need (Or just import them as you go along!)

In [1]:

```
1
```

In [ ]:

```
1
```

## Get the Data

Read in the advertising.csv file and set it to a data frame called ad\_data.

In [2]:

```
1
```

Check the head of ad\_data

In [3]:

1	
---	--

Out[3]:

	Daily Time Spent on Site	Age	Area Income	Daily Internet Usage	Ad Topic Line	City	Male	Country	Timestamp	Clicked on Ad
0	68.95	35	61833.90	256.09	Cloned 5thgeneration orchestration	Wrightburgh	0	Tunisia	2016-03-27 00:53:11	
1	80.23	31	68441.85	193.77	Monitored national standardization	West Jodi	1	Nauru	2016-04-04 01:39:02	
2	69.47	26	59785.94	236.50	Organic bottom-line service-desk	Davidton	0	San Marino	2016-03-13 20:35:42	
3	74.15	29	54806.18	245.89	Triple-buffered reciprocal time-frame	West Terrifurt	1	Italy	2016-01-10 02:31:19	
4	68.37	35	73889.99	225.58	Robust logistical utilization	South Manuel	0	Iceland	2016-06-03 03:36:18	

Use info and describe() on ad\_data

In [4]:

1	
---	--

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 1000 entries, 0 to 999
Data columns (total 10 columns):
Daily Time Spent on Site    1000 non-null float64
Age                        1000 non-null int64
Area Income                 1000 non-null float64
Daily Internet Usage       1000 non-null float64
Ad Topic Line              1000 non-null object
City                       1000 non-null object
Male                       1000 non-null int64
Country                    1000 non-null object
Timestamp                  1000 non-null object
Clicked on Ad              1000 non-null int64
dtypes: float64(3), int64(3), object(4)
memory usage: 78.2+ KB
```

In [5]:

1

Out[5]:

	Daily Time Spent on Site	Age	Area Income	Daily Internet Usage	Male	Clicked on Ad
<b>count</b>	1000.000000	1000.000000	1000.000000	1000.000000	1000.000000	1000.000000
<b>mean</b>	65.000200	36.009000	55000.000080	180.000100	0.481000	0.500000
<b>std</b>	15.853615	8.785562	13414.634022	43.902339	0.499889	0.500250
<b>min</b>	32.600000	19.000000	13996.500000	104.780000	0.000000	0.000000
<b>25%</b>	51.360000	29.000000	47031.802500	138.830000	0.000000	0.000000
<b>50%</b>	68.215000	35.000000	57012.300000	183.130000	0.000000	0.500000
<b>75%</b>	78.547500	42.000000	65470.635000	218.792500	1.000000	1.000000
<b>max</b>	91.430000	61.000000	79484.800000	269.960000	1.000000	1.000000

## Exploratory Data Analysis

Let's use seaborn to explore the data!

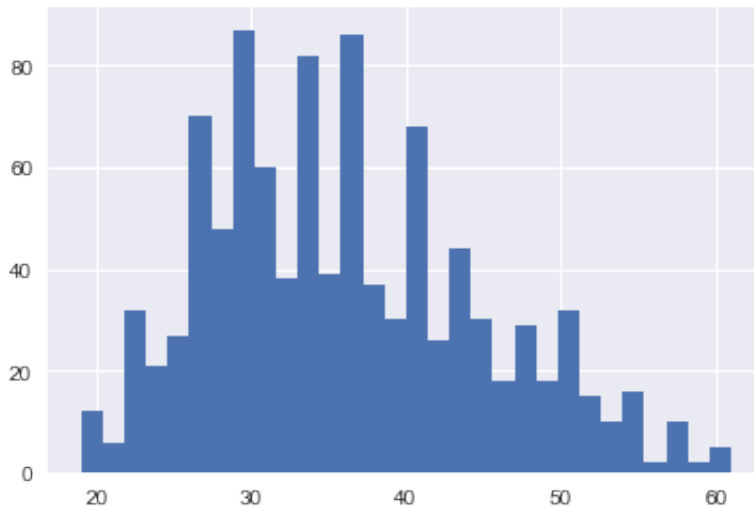
Try recreating the plots shown below!

**Create a histogram of the Age**

In [14]:

1

```
Out[14]: (array([ 12.,   6.,  32.,  21.,  27.,  70.,  48.,  87.,  60.,  38.,  8
 2.,
           39.,  86.,  37.,  30.,  68.,  26.,  44.,  30.,  18.,  29.,  1
 8.,
           32.,  15.,  10.,  16.,   2.,  10.,   2.,   5.]),
 array([ 19. ,  20.4,  21.8,  23.2,  24.6,  26. ,  27.4,  28.8,  30.2,
        31.6,  33. ,  34.4,  35.8,  37.2,  38.6,  40. ,  41.4,  42.8,
        44.2,  45.6,  47. ,  48.4,  49.8,  51.2,  52.6,  54. ,  55.4,
        56.8,  58.2,  59.6,  61. ]),
<a list of 30 Patch objects>)
```

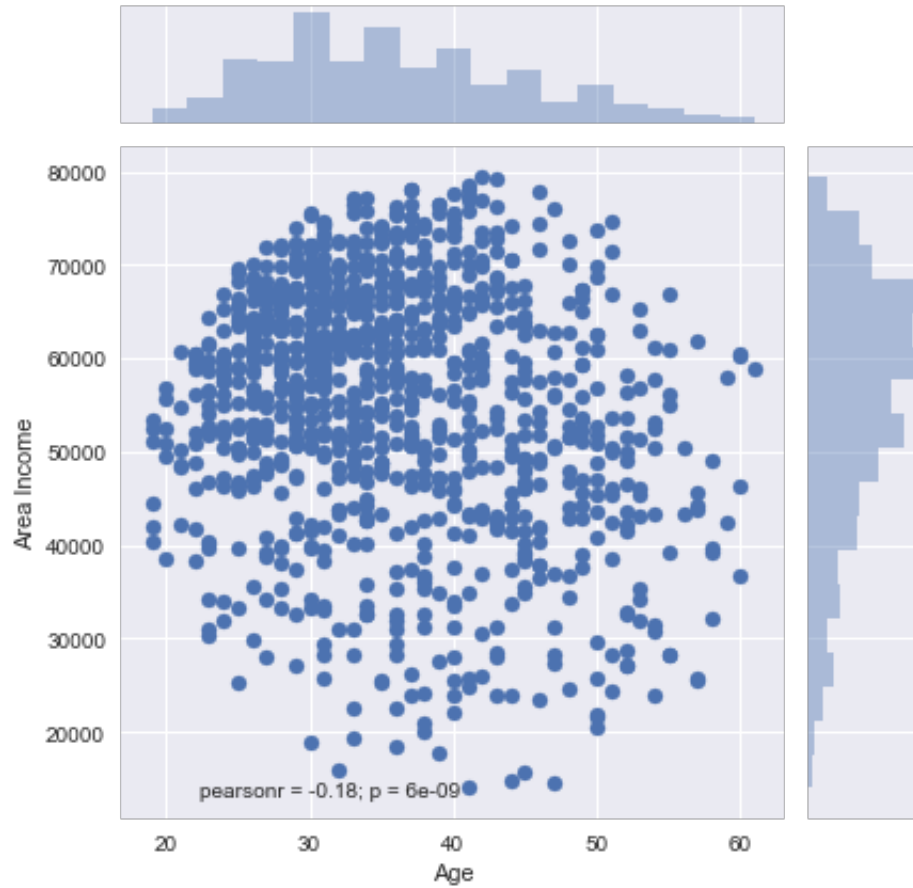


Create a jointplot showing Area Income versus Age.

In [7]:

1

Out[7]: <seaborn.axisgrid.JointGrid at 0x10cale450>



In [ ]:

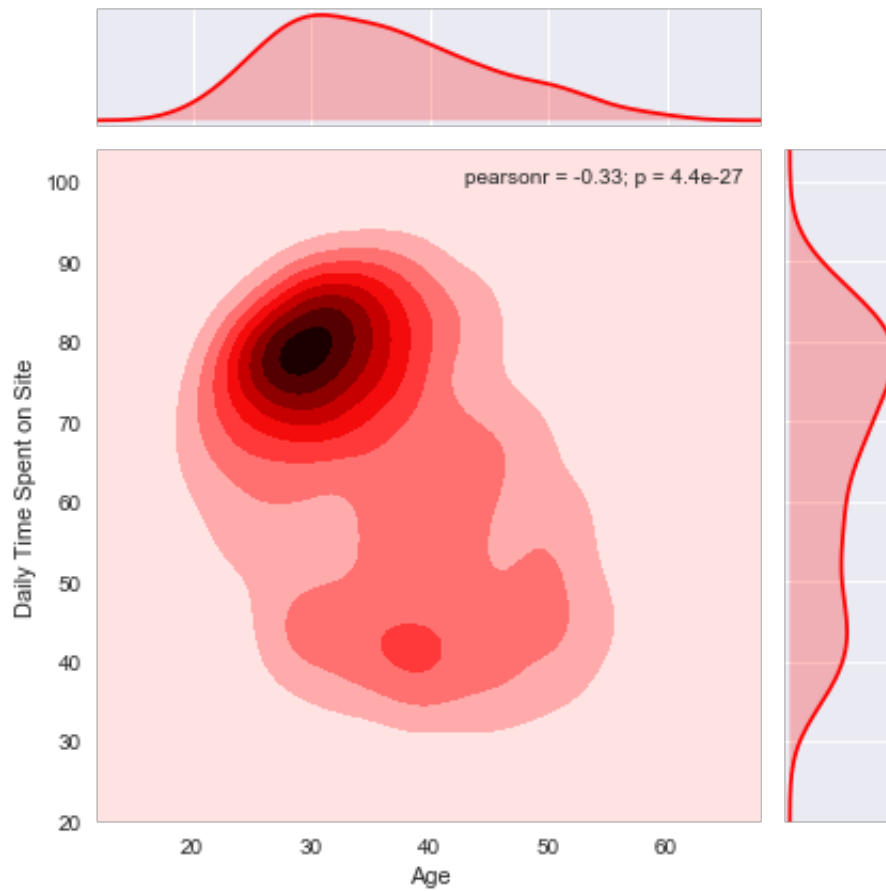
1

**Create a jointplot showing the kde distributions of Daily Time spent on site vs. Age.**

In [13]:

1

Out[13]: <seaborn.axisgrid.JointGrid at 0x110b75d50>



In [ ]:

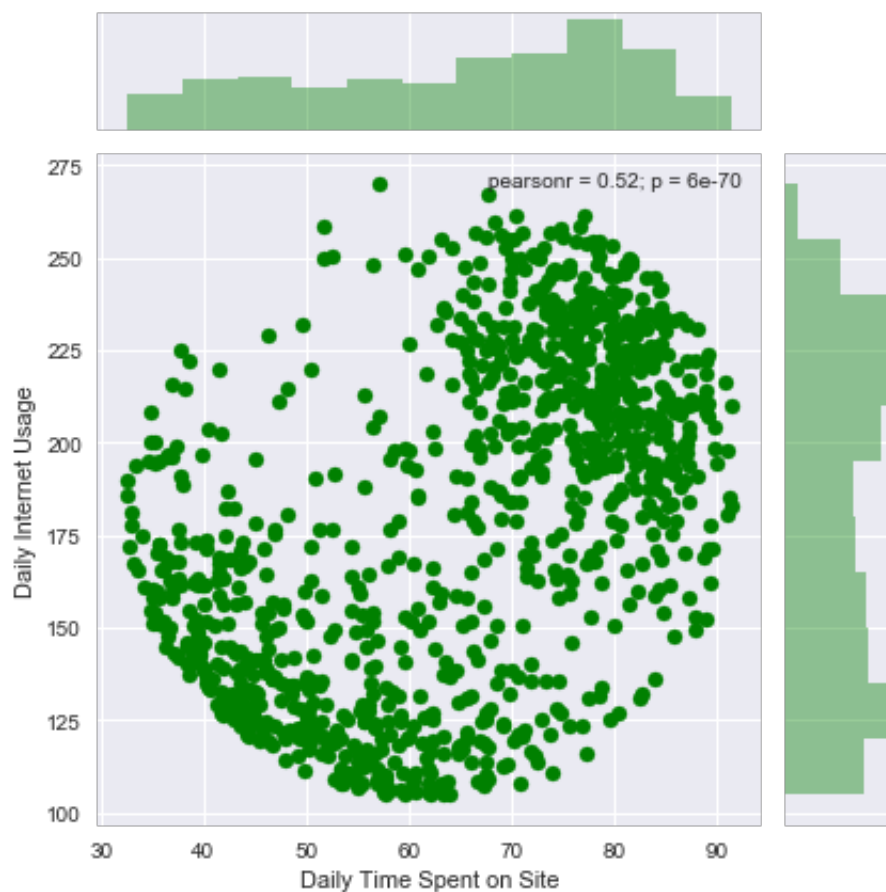
1

Create a jointplot of 'Daily Time Spent on Site' vs. 'Daily Internet Usage'

In [15]:

1

Out[15]: <seaborn.axisgrid.JointGrid at 0x10ffd3f50>



Finally, create a pairplot with the hue defined by the 'Clicked on Ad' column feature.

In [17]:

1

Out[17]: <seaborn.axisgrid.PairGrid at 0x113787090>



## Logistic Regression

Now it's time to do a train test split, and train our model!

You'll have the freedom here to choose columns that you want to train on!

**Split the data into training set and testing set using `train_test_split`**

In [38]:

1

In [39]:

1



In [47]:

1

**Train and fit a logistic regression model on the training set.**

In [48]:

1

In [49]:

1

In [50]:

1

```
Out[50]: LogisticRegression(C=1.0, class_weight=None, dual=False, fit_intercept
= True,
                        intercept_scaling=1, max_iter=100, multi_class='ovr', n_jobs
= 1,
                        penalty='l2', random_state=None, solver='liblinear', tol=0.0
001,
                        verbose=0, warm_start=False)
```

## Predictions and Evaluations

**Now predict values for the testing data.**

In [51]:

1

**Create a classification report for the model.**

In [52]:

1

In [53]:

1

	precision	recall	f1-score	support
0	0.87	0.96	0.91	162
1	0.96	0.86	0.91	168
avg / total	0.91	0.91	0.91	330

**Great Job!**