

# **Employee Absenteeism Prediction**

**Gursimran Singh**

**May 2019**

## Introduction

Employee absenteeism is a significant problem for most organizations. In 2017 the U.S. Department of Labour (DOL) estimated that almost 3 per cent of an employer's workforce was absent on any given day. The high price of absenteeism affects organizations even more when lost productivity, morale and temporary labour costs are considered. The cost of absence is often misunderstood, seen as not easily measured or dismissed as a negligible amount because the costs are largely included in payroll expenses. In addition, employers often fail to carefully track absenteeism.

Many firms still underestimate the magnitude of the problem of employee absenteeism. Instead, they consider such absences to be part of "the cost of doing business." From that perspective, they fail to quantify the full impact of employee absenteeism. In addition, they also fail to appreciate the value of solutions that can reduce the costs and lost time that result from employees' being off work.

But with the emergence of Machine Learning and Data Sciences technologies, most of the corporate giants started maintaining data of its employees and using that data to predict the absenteeism rate. Besides this companies are also able to predict the resigning rate of their employees from the number of their working hours or on other parameters of the collected data.

## Problem Statement

XYZ is a courier company. As we appreciate that human capital plays an important role in collection, transportation and delivery. The company is passing through genuine issue of Absenteeism. The company has shared it dataset and requested to have an answer on the following areas:

1. What changes company should bring to reduce the number of absenteeism?
2. How many losses every month can we project in 2011 if same trend of absenteeism continues?

Our task is to build a machine learning model which will predict the absenteeism hours of the employees depending on various other factors. Besides this we have to assist the top business leaders in making decision by exploring data and answering their required questions. Sample from the whole dataset is shown below:-

Id	AbsentReason	AbsentMonth	WeekDay	Seasons	Expenses	ResidentDistance
11	26	7	3	1	289	36
36	0	7	3	1	118	13
3	23	7	4	1	179	51
7	7	7	5	1	279	5
11	23	7	5	1	289	36

Table 1.1: Sample Data

Service time	Age	AverageWorkLoad	HitTarget	DisciplineFailure	Education	Son
13	33	2,39,554	97	0	1	2
18	50	2,39,554	97	1	1	1
18	38	2,39,554	97	0	1	0
14	39	2,39,554	97	0	1	2
13	33	2,39,554	97	0	1	2

**Table 1.2 : Sample Data**

SocialDrinker	Social smoker	Pet	Weight	Height	Body mass index	AbsentTime
1	0	1	90	172	30	4
1	0	0	98	178	31	0
1	0	0	89	170	31	2
1	1	0	68	168	24	4
1	0	1	90	172	30	2

**Table 1.3: Sample Data**

Descriptions of the attributes are given below:-

<b>Attribute</b>	<b>Description</b>
<b>Id</b>	Individual identification
<b>AbsentReason</b>	<p>Reason for absence</p> <p>Absences attested by the International Code of Diseases (ICD) stratified into 21 categories (I to XXI) as follows:</p> <ul style="list-style-type: none"> <li>I Certain infectious and parasitic diseases</li> <li>II Neoplasms</li> <li>III Diseases of the blood and blood-forming organs and certain disorders involving the immune mechanism</li> <li>IV Endocrine, nutritional and metabolic diseases</li> <li>V Mental and behavioural disorders</li> <li>VI Diseases of the nervous system</li> <li>VII Diseases of the eye and adnexa</li> <li>VIII Diseases of the ear and mastoid process</li> <li>IX Diseases of the circulatory system</li> <li>X Diseases of the respiratory system</li> <li>XI Diseases of the digestive system</li> <li>XII Diseases of the skin and subcutaneous tissue</li> <li>XIII Diseases of the musculoskeletal system and connective tissue</li> <li>XIV Diseases of the genitourinary system</li> <li>XV Pregnancy, childbirth and the puerperium</li> <li>XVI Certain conditions originating in the perinatal period</li> <li>XVII Congenital malformations, deformations and chromosomal abnormalities</li> <li>XVIII Symptoms, signs and abnormal clinical and laboratory findings, not elsewhere classified</li> <li>XIX Injury, poisoning and certain other consequences of external causes</li> </ul>

	XX External causes of morbidity and mortality
	XXI Factors influencing health status and contact with health services.
	And 7 categories without (CID) patient follow-up (22), medical consultation (23), blood donation (24), laboratory examination (25), unjustified absence (26), physiotherapy (27), dental consultation (28).
<b>AbsentMonth</b>	Month of absence
<b>WeekDay</b>	Day of the week (Monday (2), Tuesday (3), Wednesday (4), Thursday (5), Friday (6))
<b>Seasons</b>	Seasons (summer (1), autumn (2), winter (3), spring (4))
<b>Expenses</b>	Transportation expense
<b>ResidentDistance</b>	Distance from Residence to Work (kilometres)
<b>ServiceTIme</b>	Service time
<b>Age</b>	Age
<b>WorkLoad</b>	Work load Average/day
<b>HitTarget</b>	Hit target
<b>DiscplineFailure</b>	Disciplinary failure (yes=1; no=0)
<b>Education</b>	Education (high school (1), graduate (2), postgraduate (3), master and doctor (4))
<b>Son</b>	Son (number of children)
<b>SocialDrinker</b>	Social drinker (yes=1; no=0)
<b>SocialSmoker</b>	Social smoker (yes=1; no=0)
<b>Pet</b>	Pet (number of pet)
<b>Weight</b>	Weight
<b>Height</b>	Height
<b>BodyMassIndex</b>	Body mass index
<b>AbsentTime</b>	Absenteeism time in hours (target)

**Table 1.3: Attribute Descriptions**

## **Methodology**

Any predictive modelling requires to look at the data before start modelling. However, in data mining terms *looking at data* refers to so much more than just looking. Looking at data refers to exploring the data, cleaning the data as well as visualizing the data through graphs and plots. This is often called as Exploratory Data Analysis (EDA).

Exploratory data analysis (EDA) is a very important step which takes place after feature engineering and acquiring data and it should be done before any modelling. This is because it is very important for a data scientist to be able to understand the nature of the data without making assumptions.

The purpose of EDA is to use summary statistics and visualizations to better understand data, and find clues about the tendencies of the data, its quality and to formulate assumptions and the hypothesis of our analysis. EDA is not about making fancy visualizations or even aesthetically pleasing ones, the goal is to try and answer questions with data. A goal should be to be able to create a figure which someone can look at in a couple of seconds and understand what is going on. If not, the visualization is too complicated (or fancy) and something similar should be used.

EDA is also very iterative since we first make assumptions based on our first exploratory visualizations, and then build some models. We then make visualizations of the model results and tune our models.

Remember the quality of our inputs decide the quality of our output. So, once we have got our business hypothesis ready, it makes sense to spend lot of time and efforts here. Estimating, data exploration, cleaning and preparation can take up to 70% of our total project time.

## **Variable Identification**

Variable identification is the first step in the exploratory data analysis. Identification of the variables in the dataset is totally dependent on the business requirements and need of the client. The main task here is to identify the Target variable on which future decision has to be made. Besides these identifying the independent variables are equally important because end result of the target variable are totally dependent on the independent/predictor variables. Brief introduction of predictor and target variables are given below:-

### **Predictor Variable**

Predictor variables are those variables or attributes in the dataset on which the result of the target variable is totally dependent. These variables are those on which decisions are made by the clients to get the maximum profit from the business.

### **Target Variable**

It is the variable or attribute in the whole dataset on which client is mostly interested. Based on the business requirements category of the target variable is identified by the data sciences experts.

Once the target variable and predictor variables are identified, our next task to identify the data types and categories of the variable. From analysing the dataset of bike renting company, the detailed description of all the variables are given below in the diagram on the different parameters.

In the further stages of exploratory data analysis process, we have to dive deep into the understanding of the each variable present in the dataset. From Business point of view each and every variable is crucial and even a minute mistake here can cause a loss of millions to your client. So to get the detailed summary of all the variables in the dataset on statistical parameters will help to better understanding of data to a data sciences expert.

<b><u>Predictor Variables</u></b>	<b><u>Target Variable</u></b>
Id	AbsentTime
AbsentReason	
AbsentMonth	
WeekDay	
Seasons	
Expenses	
ResidentDistance	
ServiceTlme	
Age	
WorkLoad	
HitTarget	
DisciplineFailure	

---

Education  
 Son  
 SocialDrinker  
 SocialSmoker  
 Pet  
 Weight  
 Height  
 BodyMassIndex

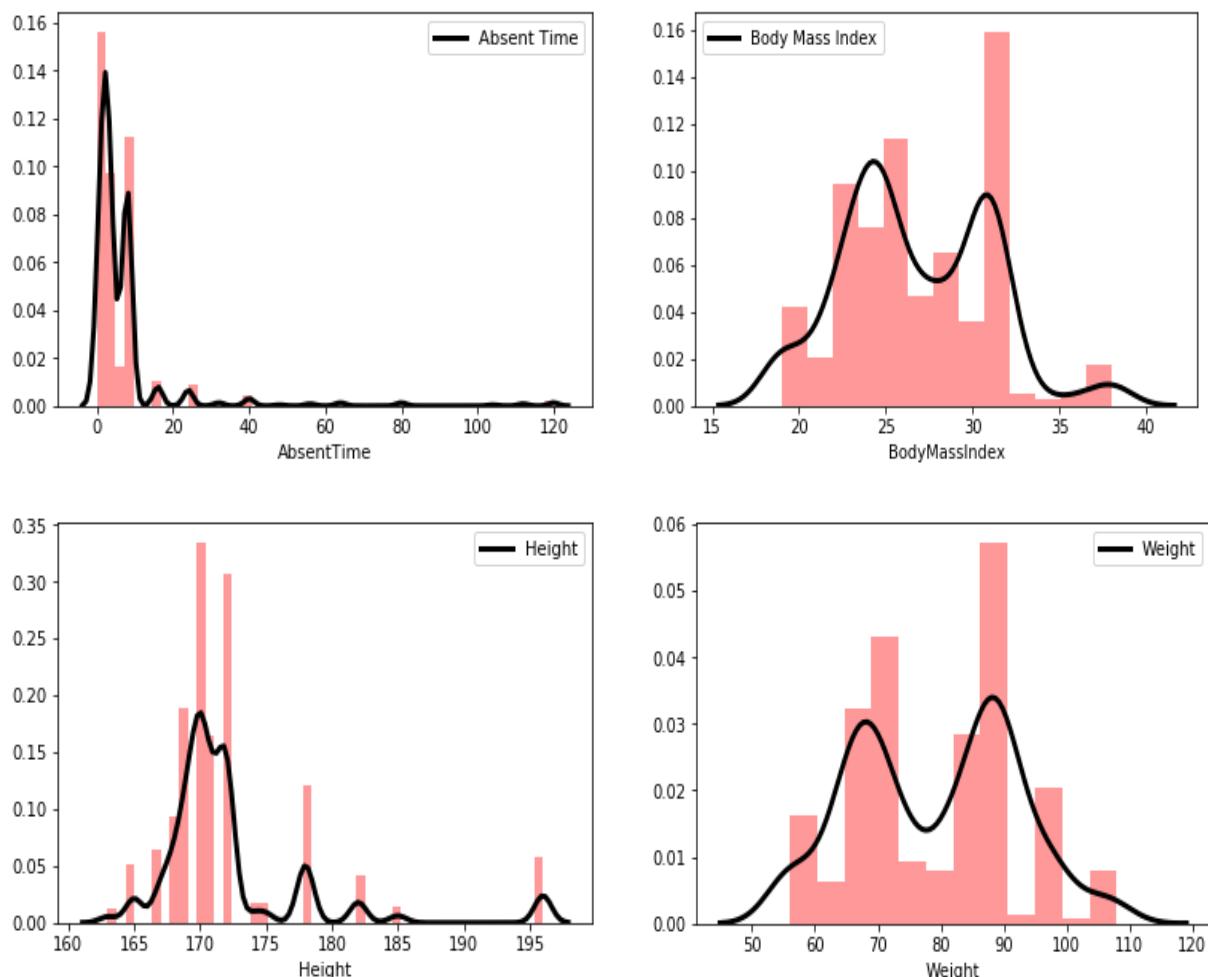
---

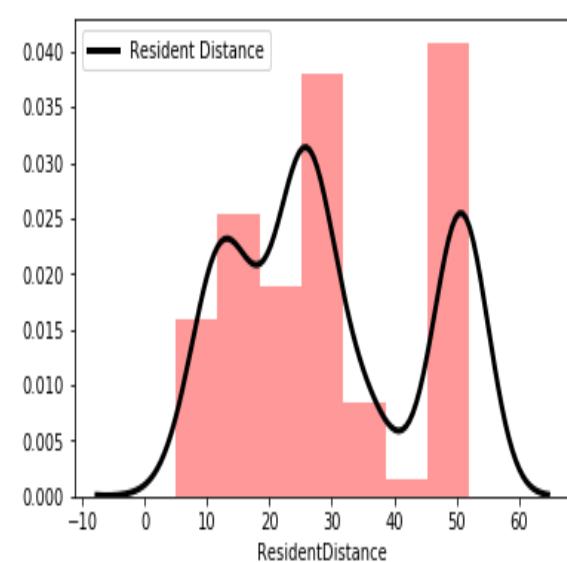
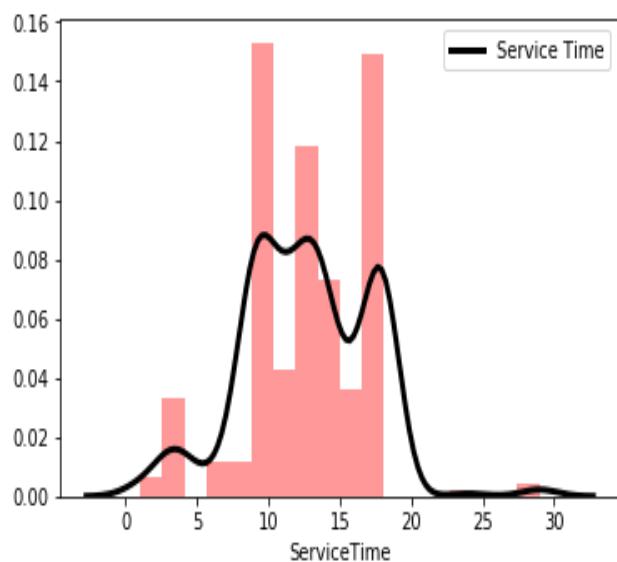
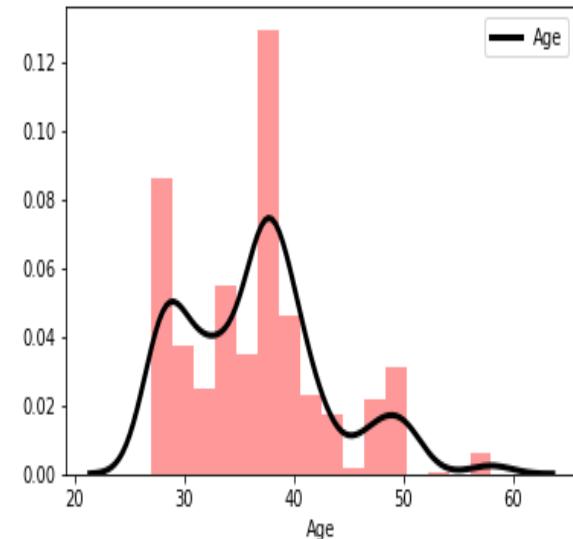
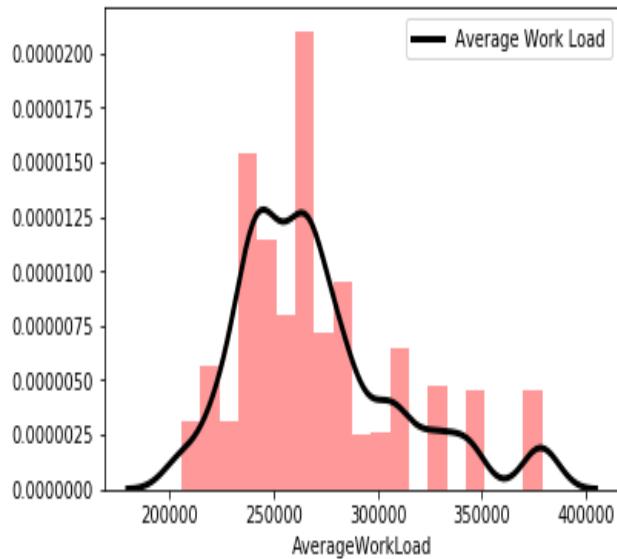
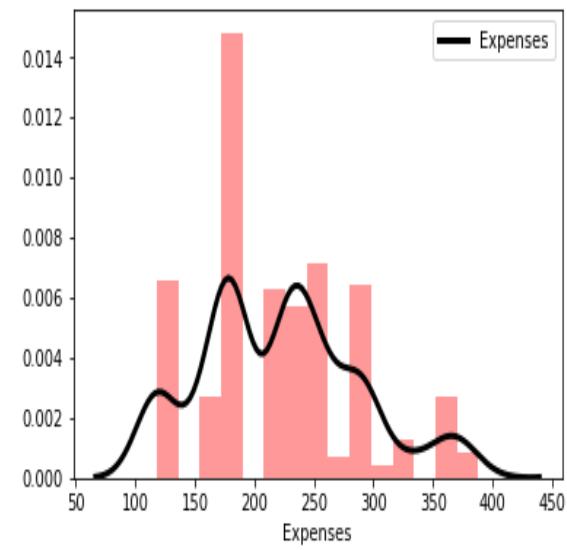
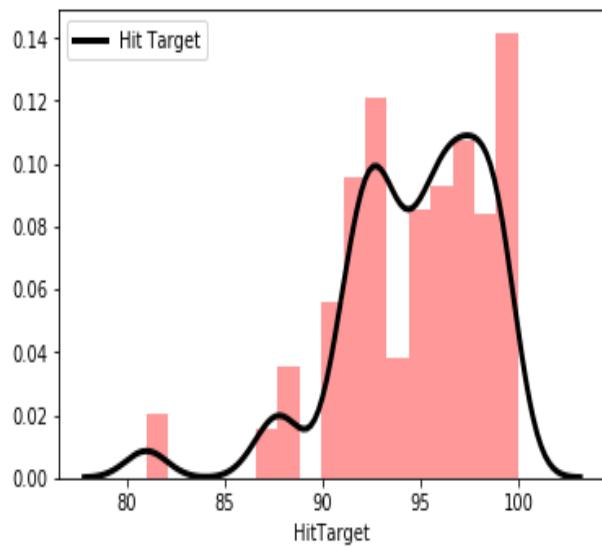
**NOTE:** For my better understanding I have changed the names of the dataset variables.

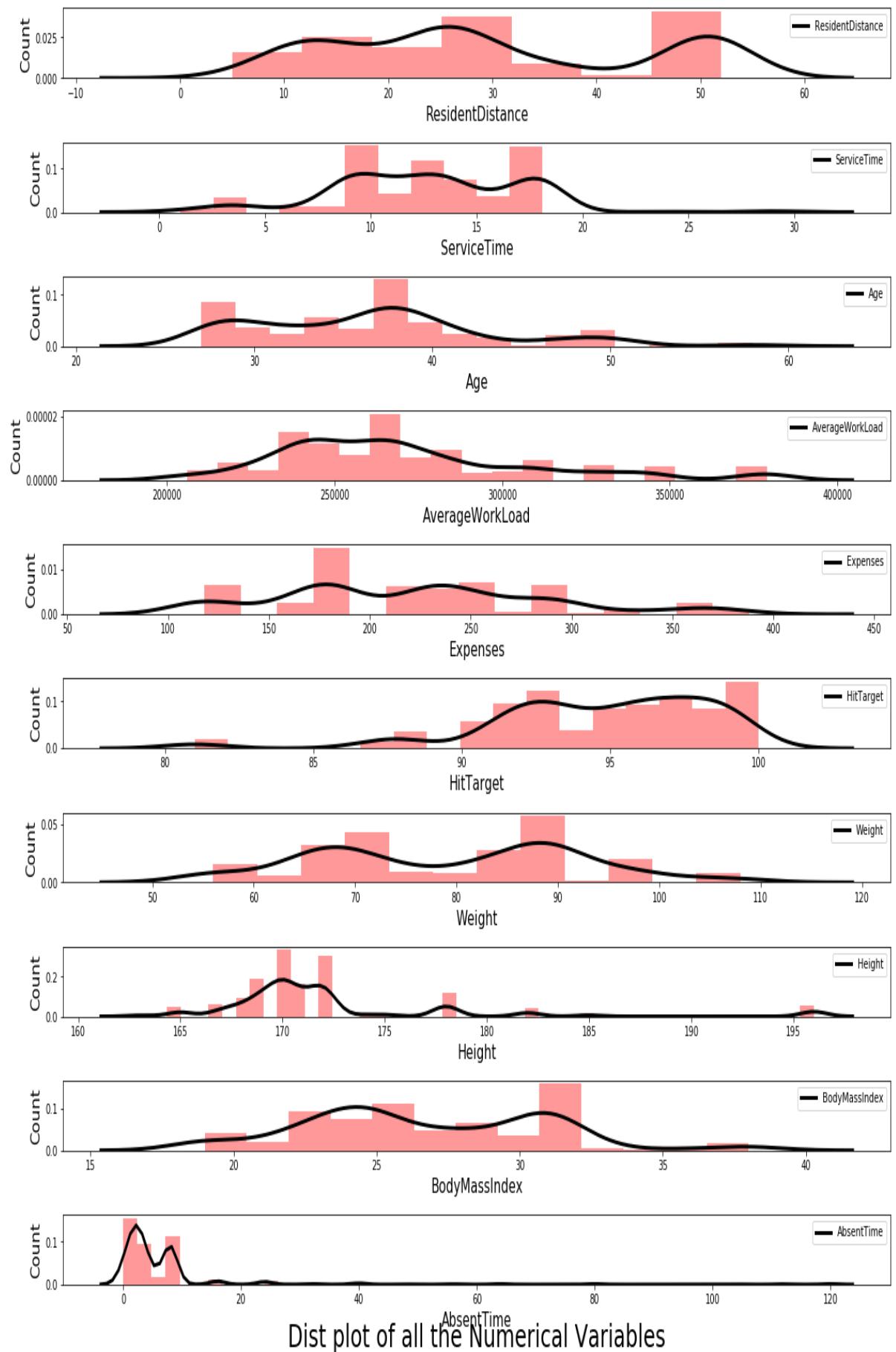
### Univariate Analysis

Univariate analysis is the simplest form of analysing data. “Uni” means “one”, so in other words data has only one variable. It doesn’t deal with causes or relationships (unlike regression) and its major purpose is to describe. It takes data, summarizes that data and finds patterns in the data.

Method to perform univariate analysis will depend on whether the variable type is categorical or continuous. Plots of all continuous variables are shown below:-

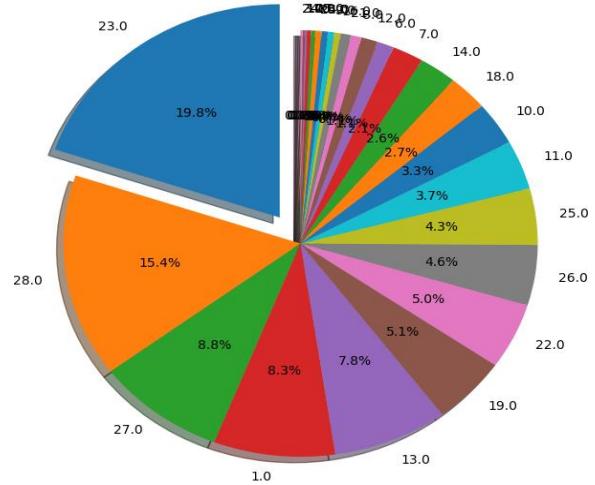
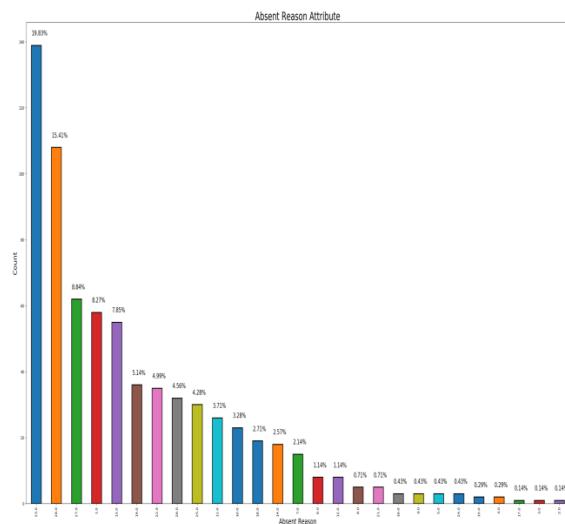
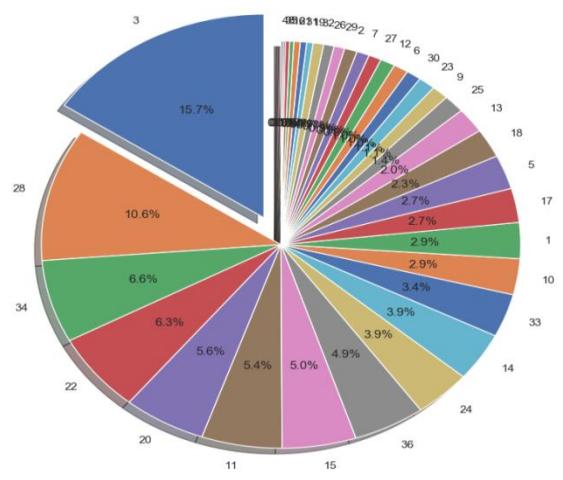
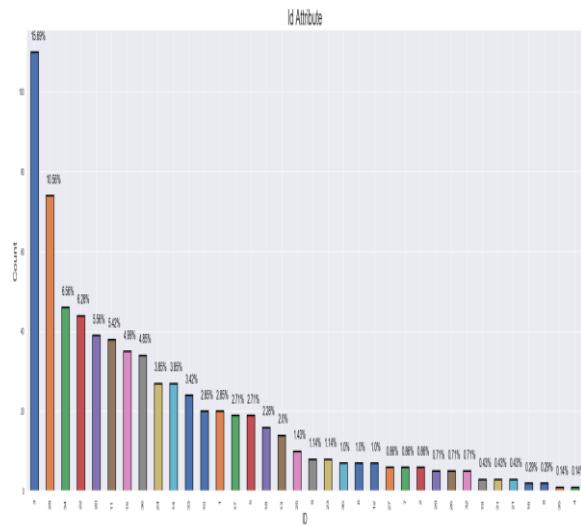


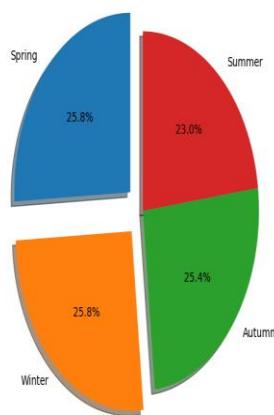
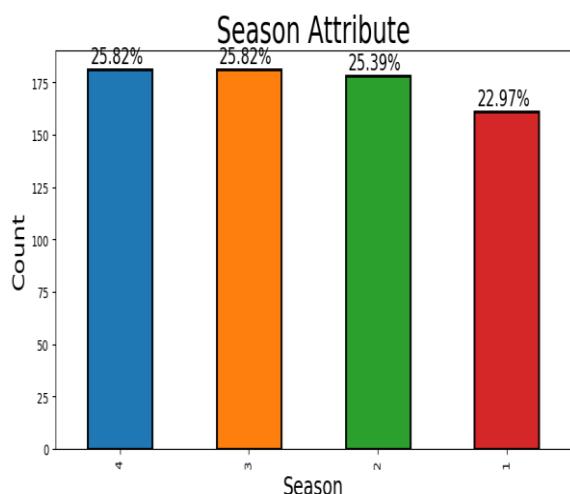
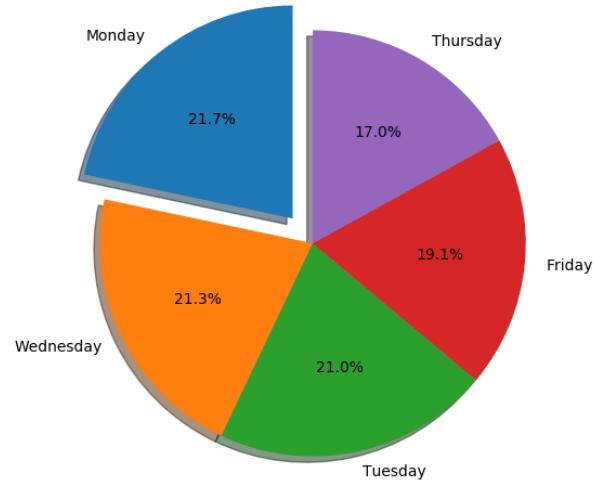
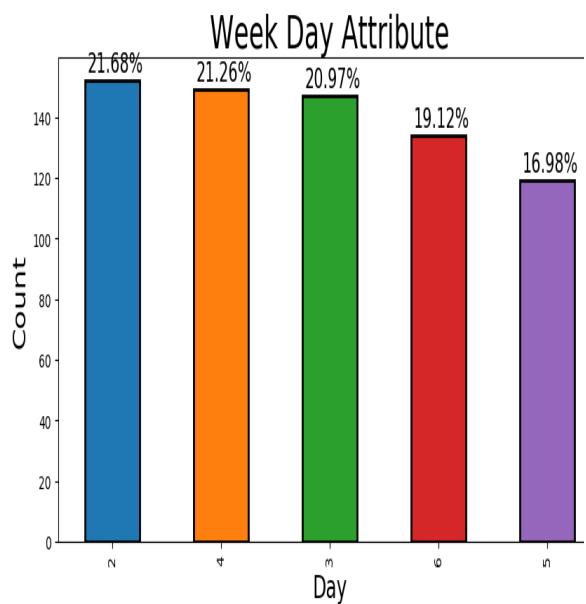
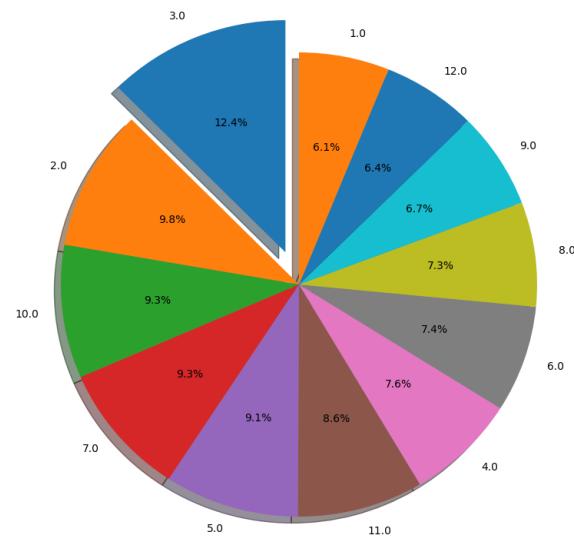
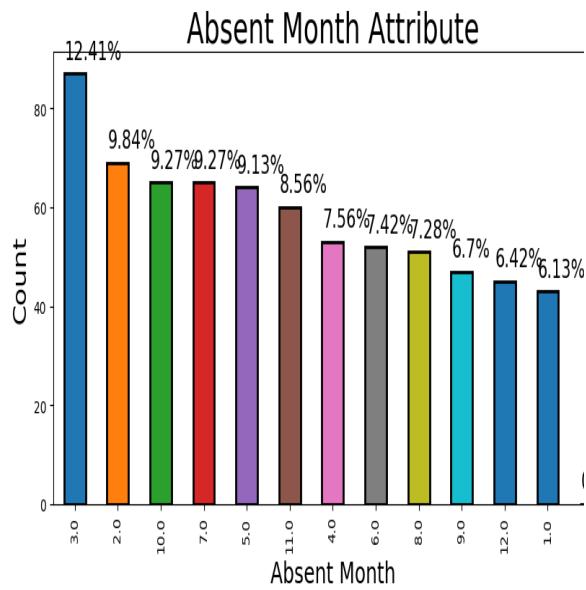


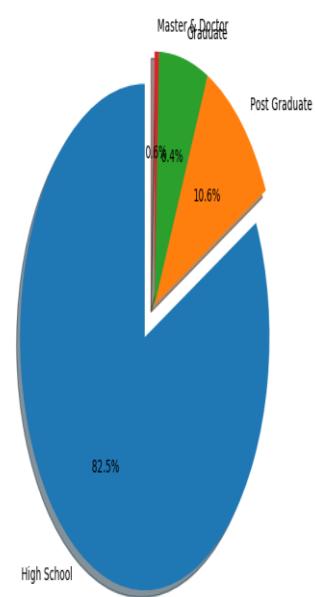
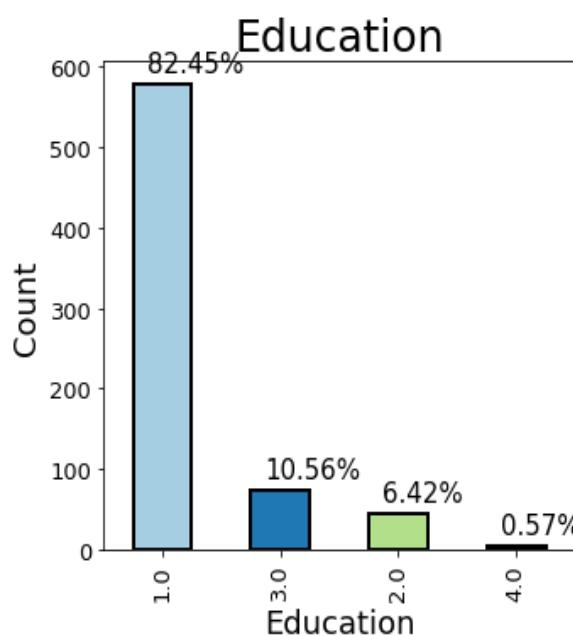
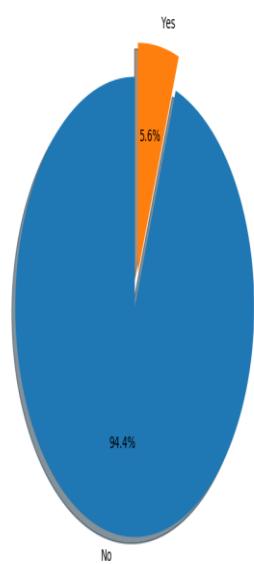
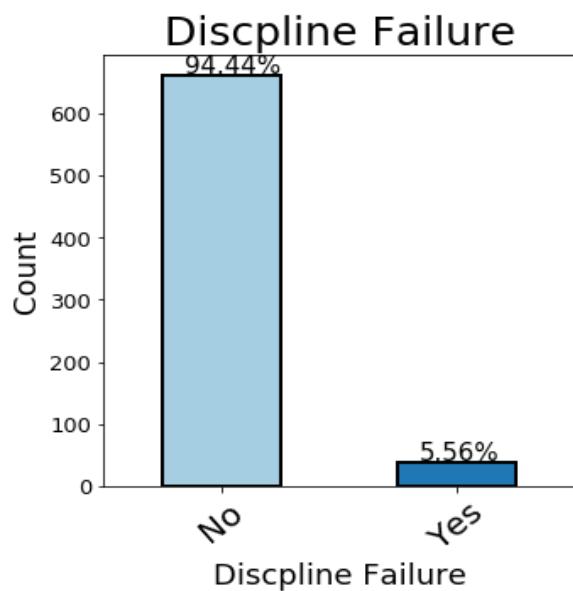


## Categorical Variables

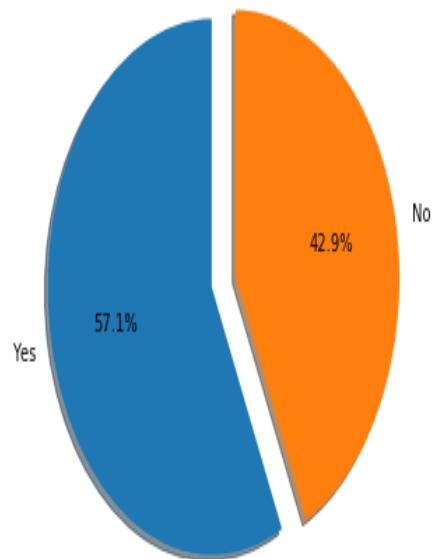
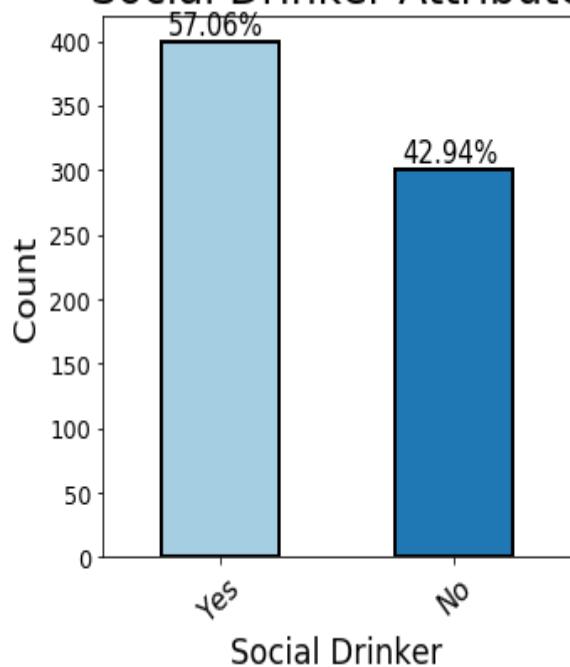
Categorical data are easier to interpret as compared to numerical variables. The best ways to analyze categorical variables are through bar graphs and pie charts. Bar graphs are mostly preferred to visualize the frequency of each category that falls into that variable, whereas pie charts are used to visualize the percentage of each category. Plots for analyzing categorical variables are shown below:-



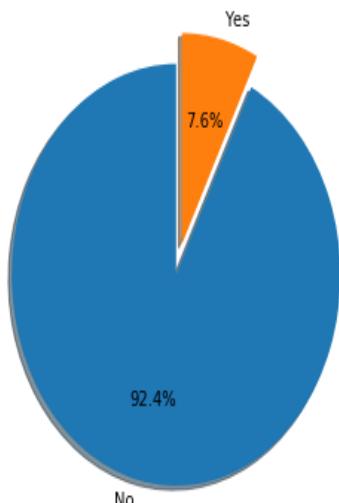
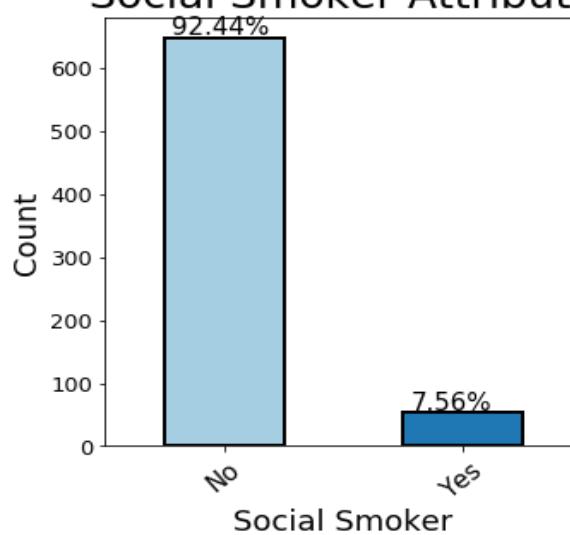


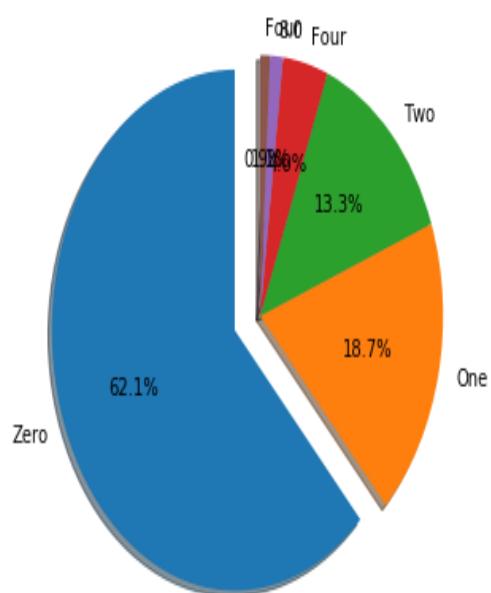
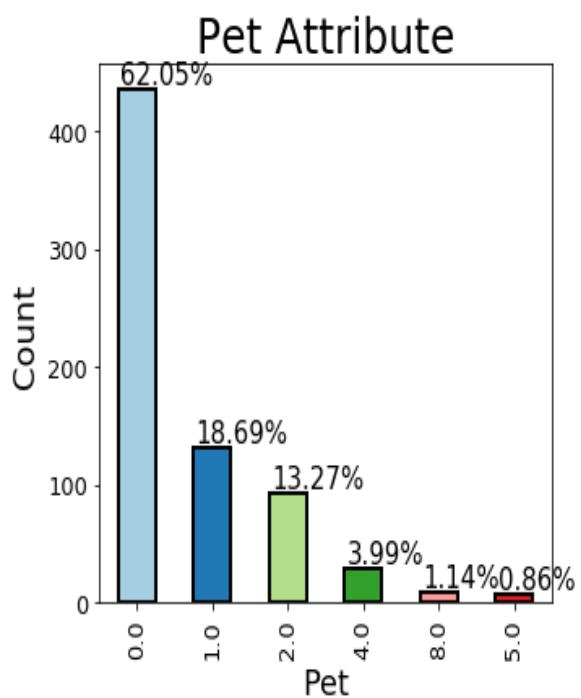
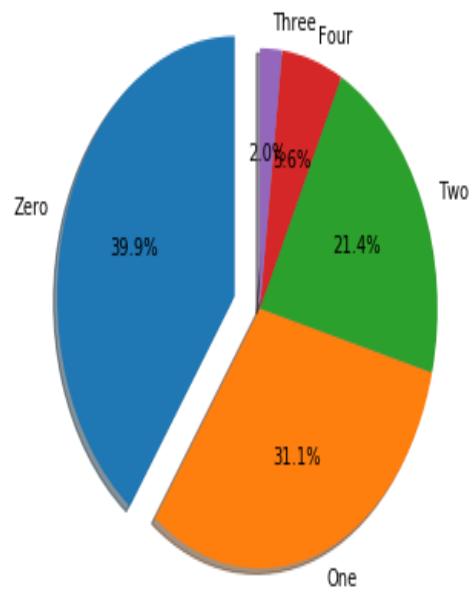
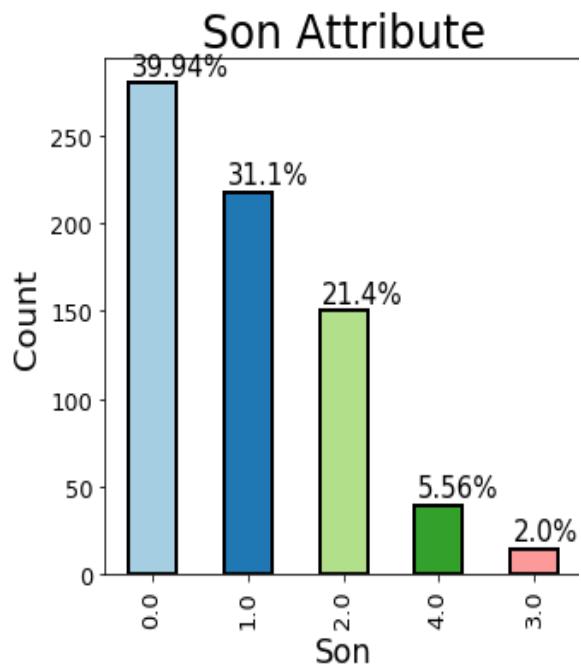


### Social Drinker Attribute

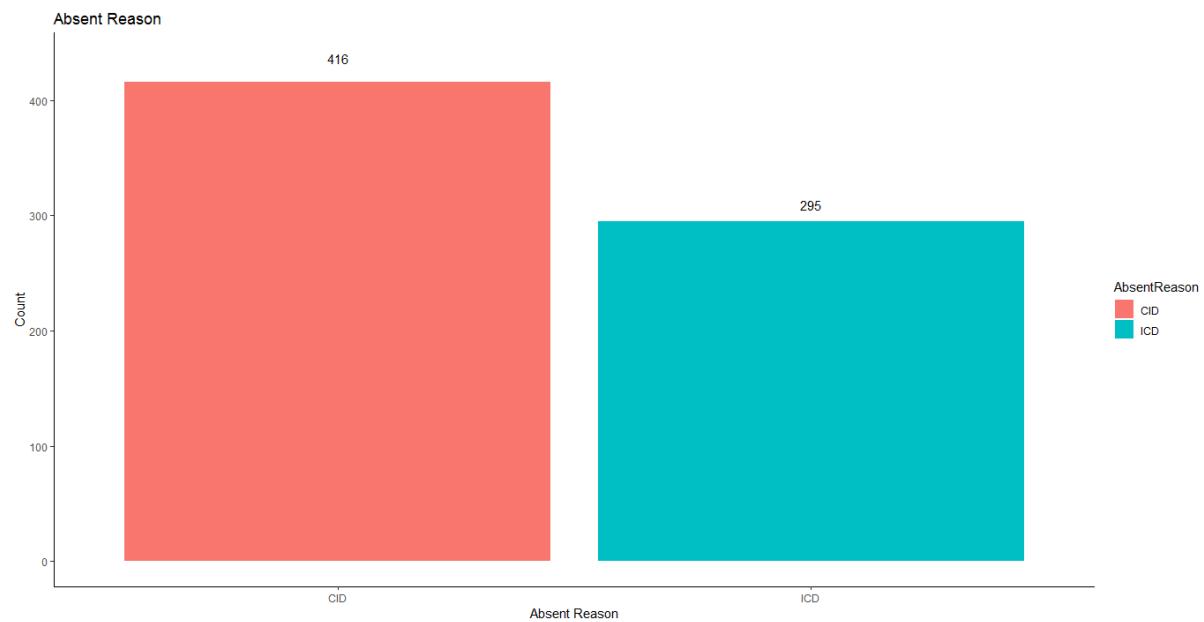
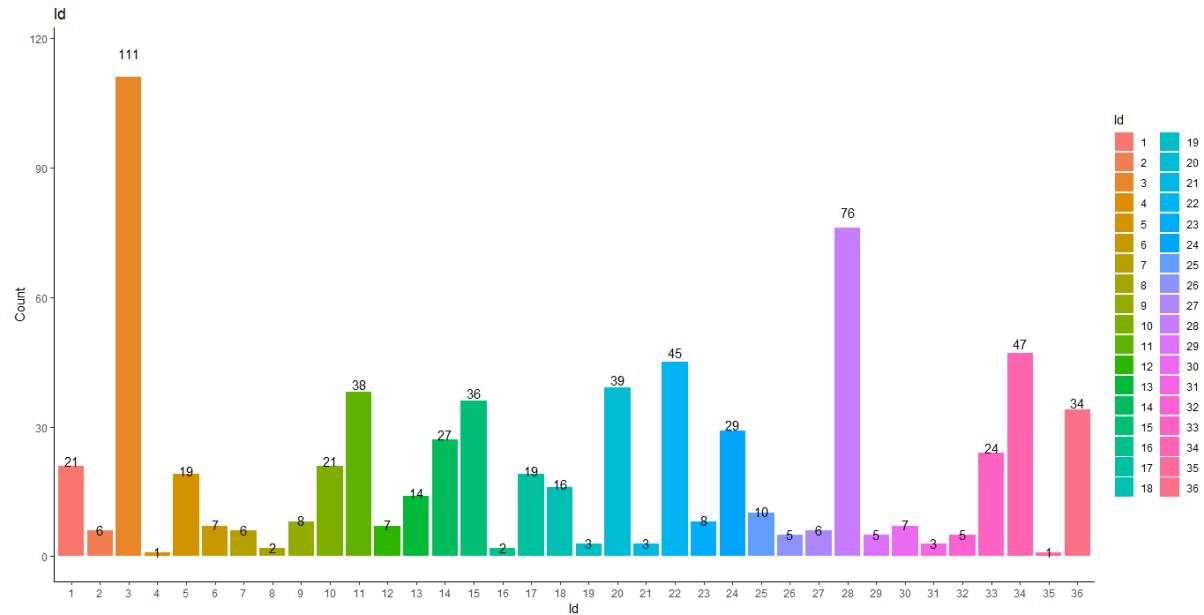


### Social Smoker Attribute

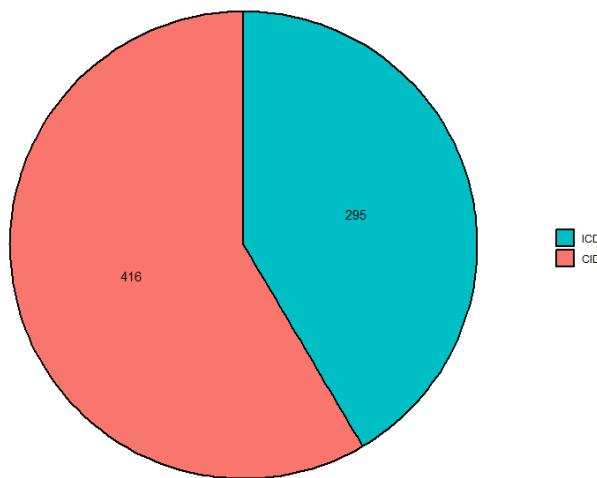




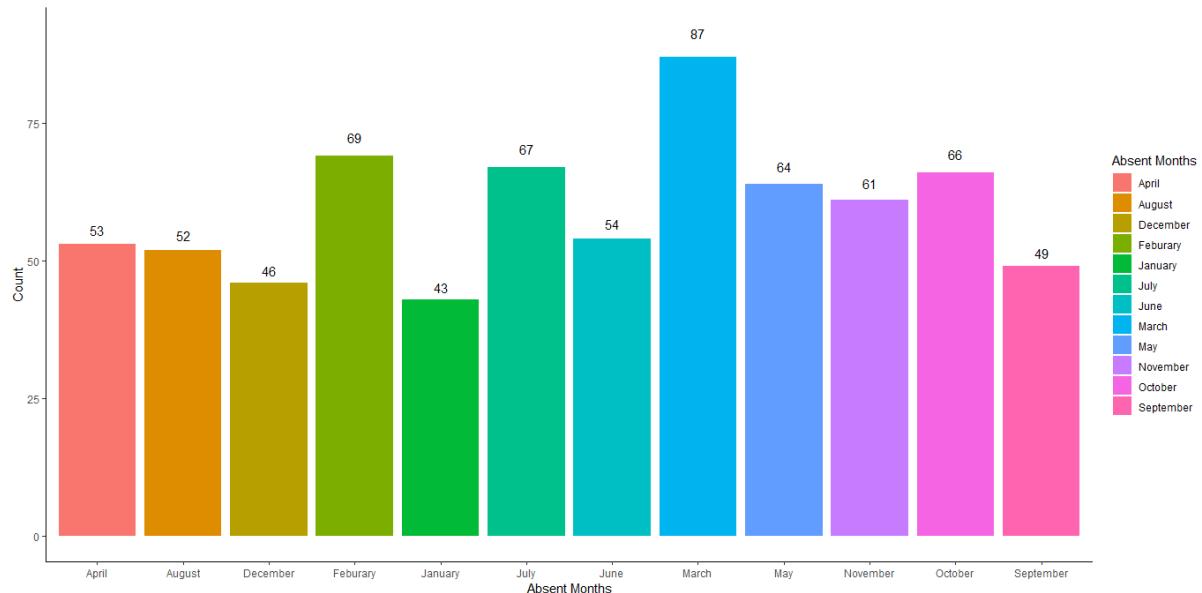
ggplot is very strong library written in R language and the graphs with better visualizations are shown below for categorical variables :-

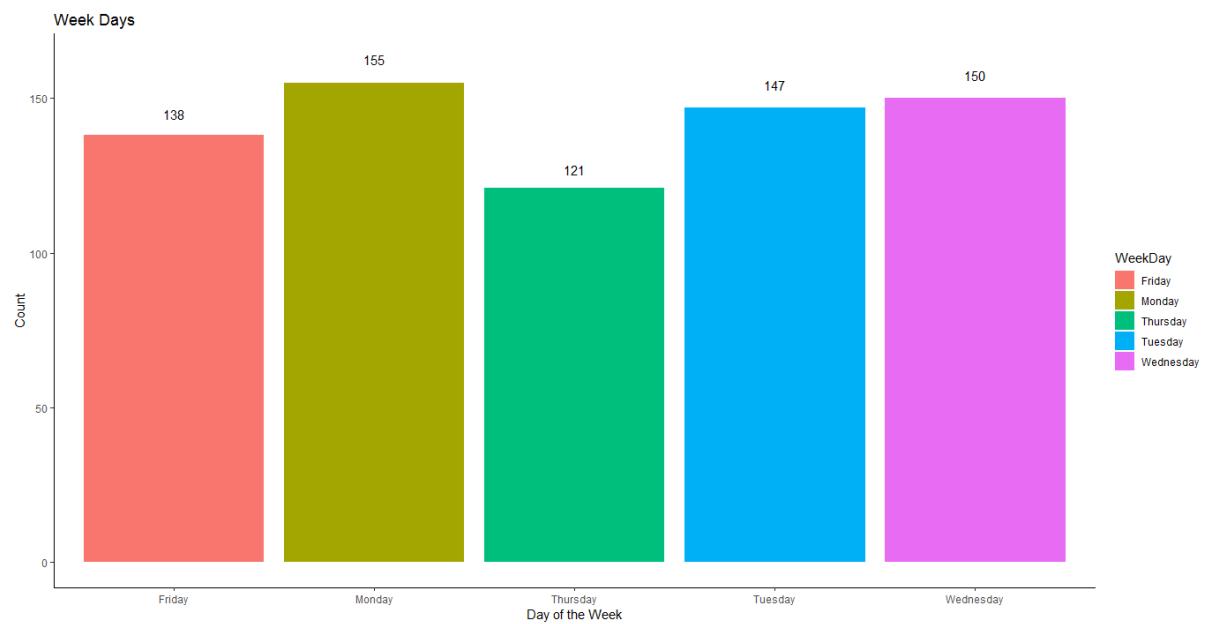


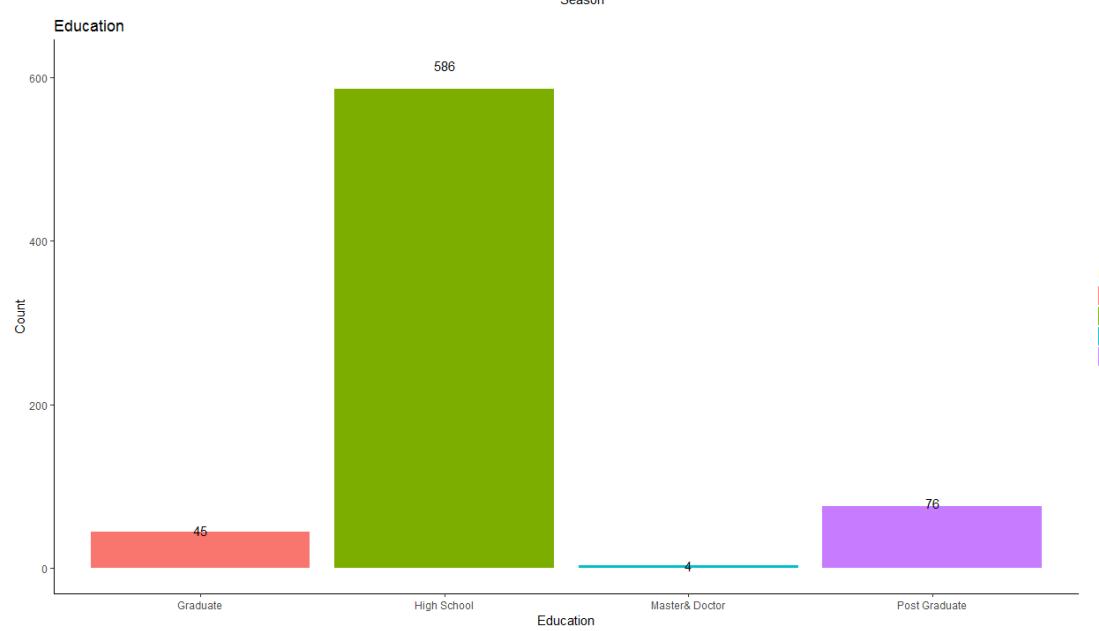
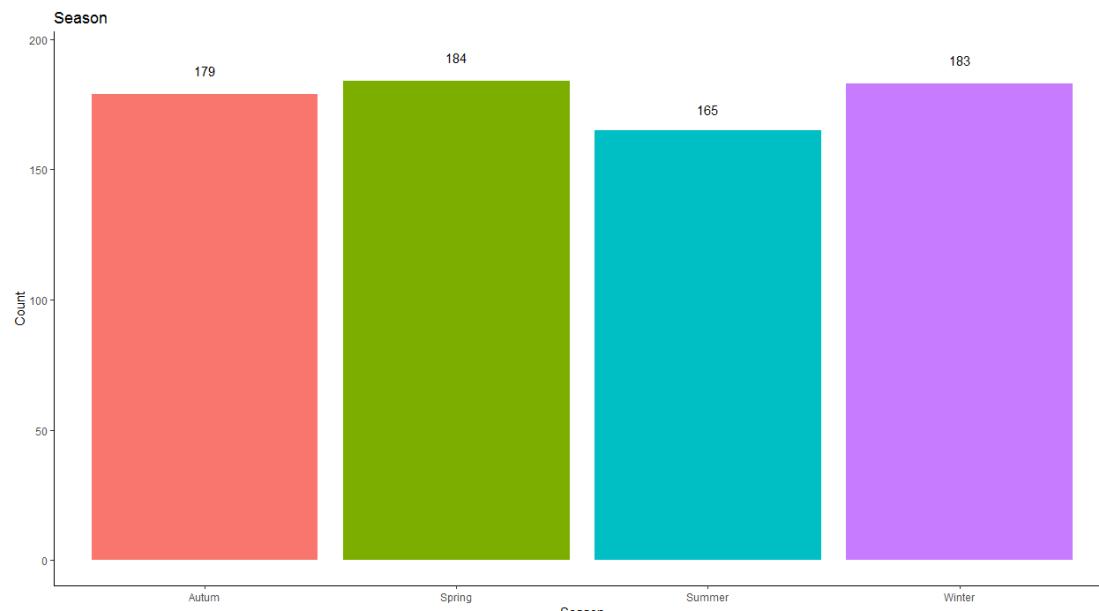
Absent Reason

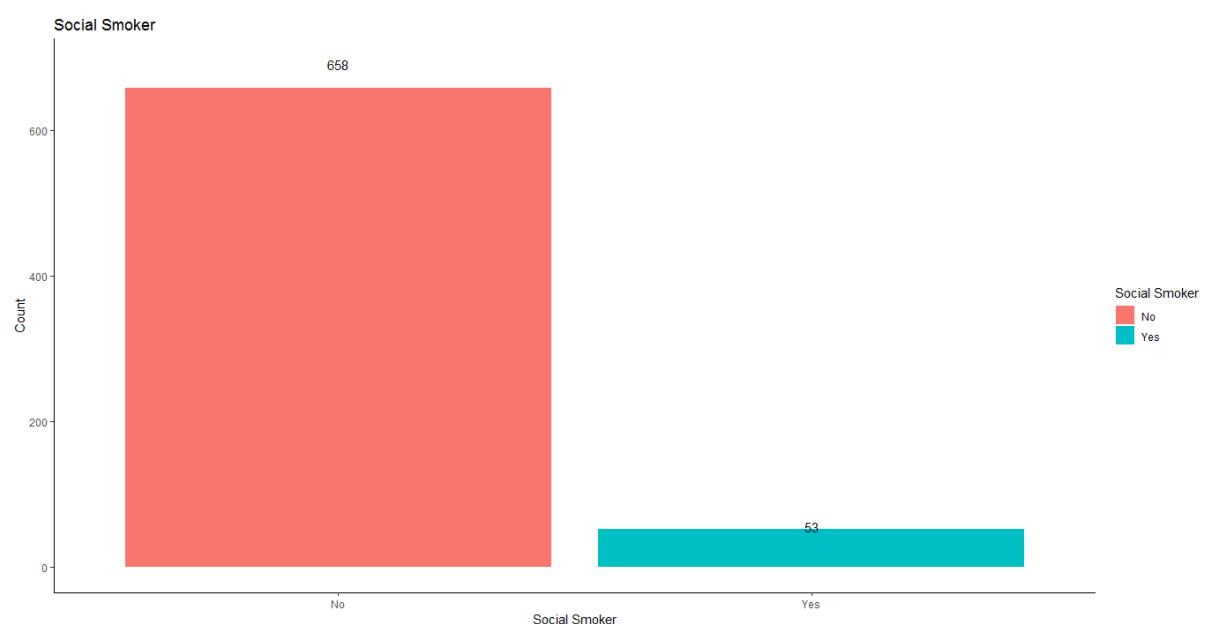
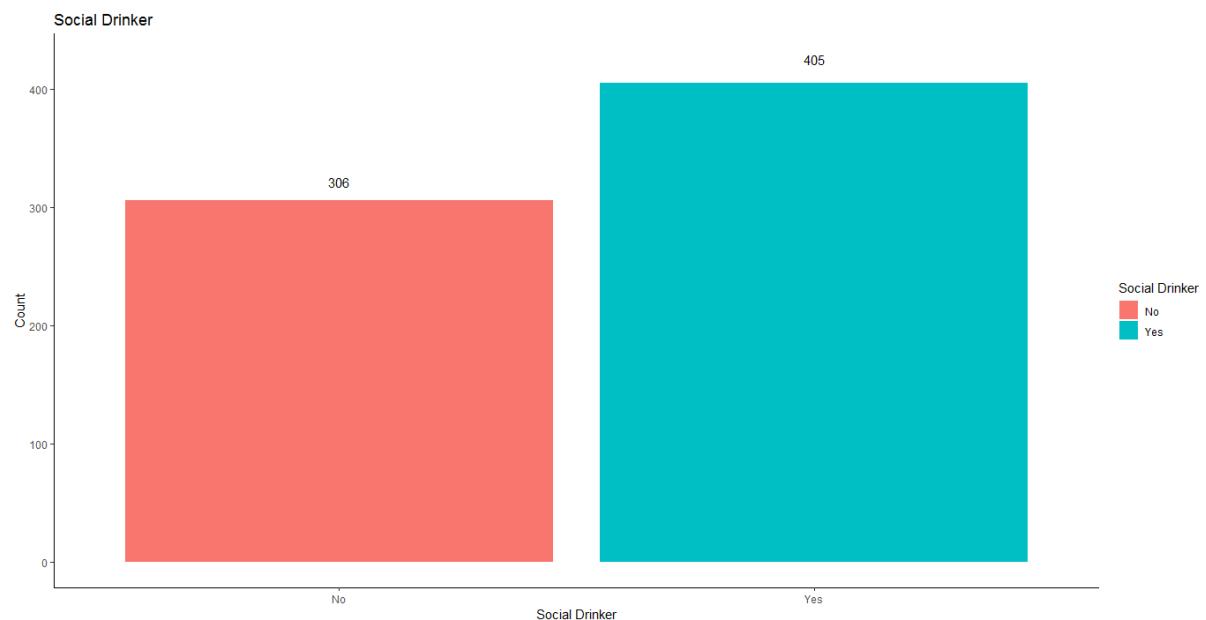


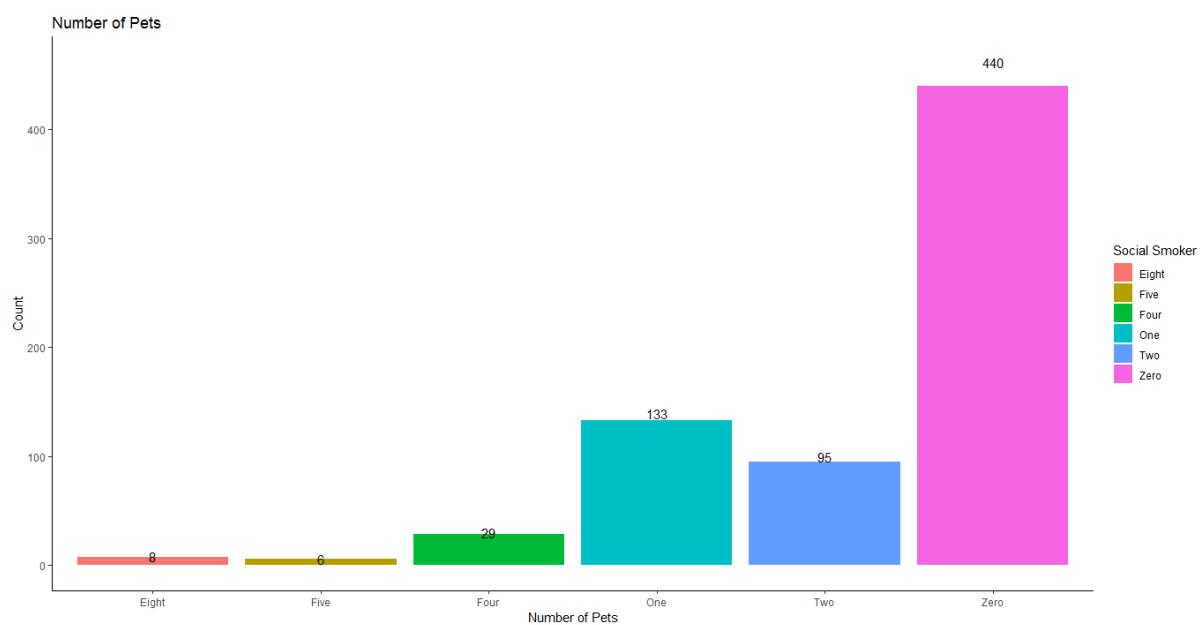
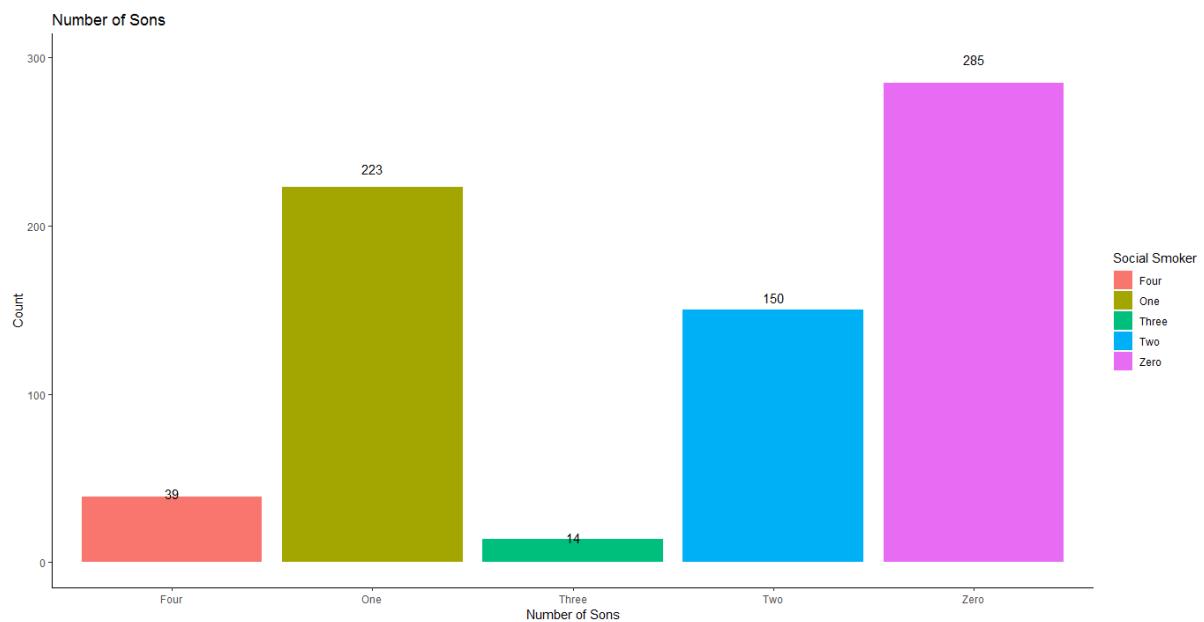
Absent Month











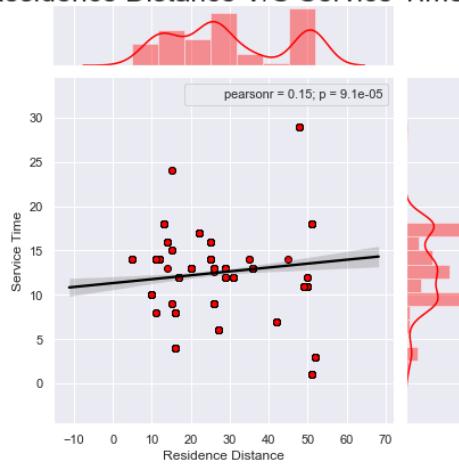
## Bivariate Analysis

Bivariate analysis is the simultaneous analysis of two variables. It explores the concept of relationship between two variables, whether there exists an association and the strength of this association, or whether there are differences between two variables and the significance of these differences.

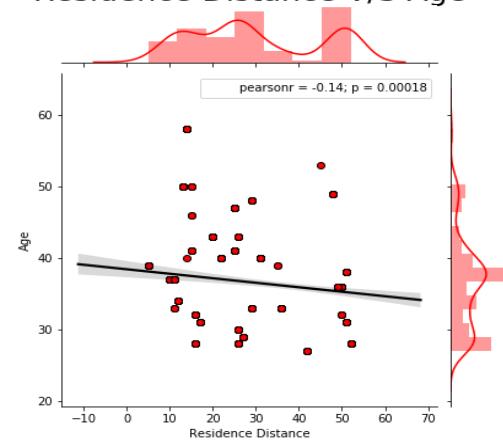
## Continuous Variables

Scatter plot for all the continuous variables are shown below. Correlation value is shown in the legend and it describes the relationship between variables against which graphs are plotted.

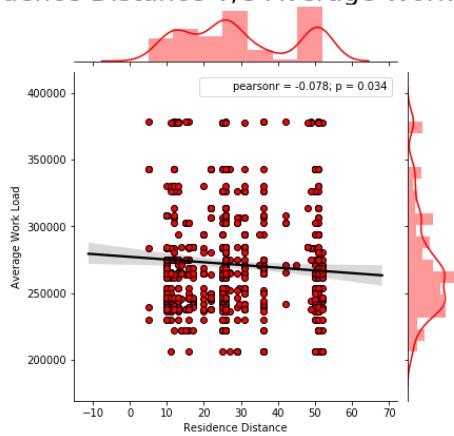
Residence Distance V/S Service Time



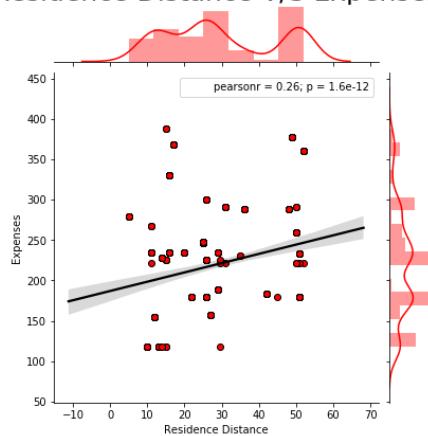
Residence Distance V/S Age



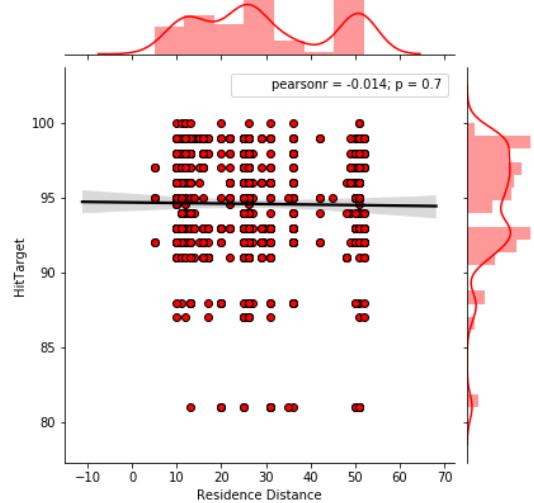
Residence Distance V/S Average Work Load



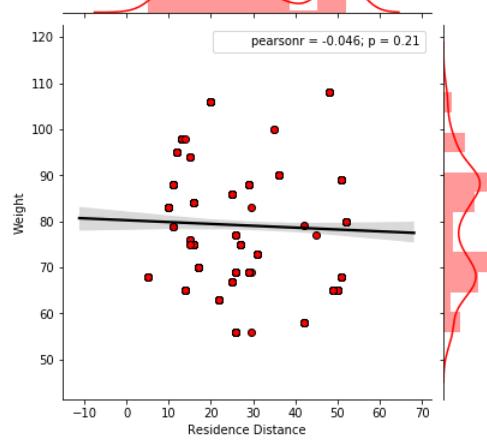
Residence Distance V/S Expenses



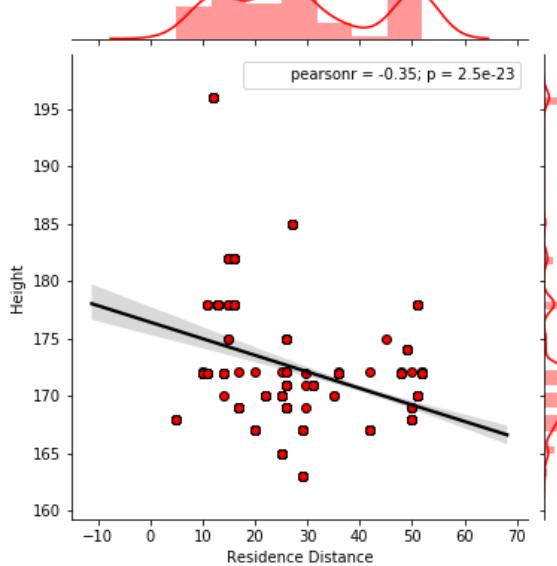
### Residence Distance V/S Hit Target



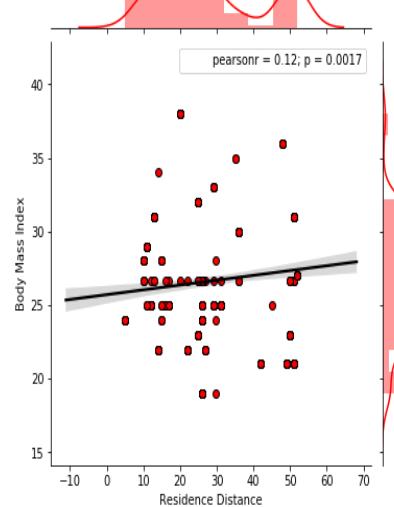
### Residence Distance v/s Weight



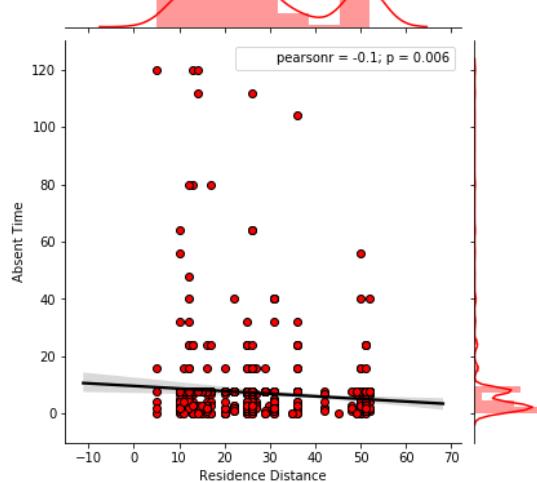
### Residence Distance V/S Height



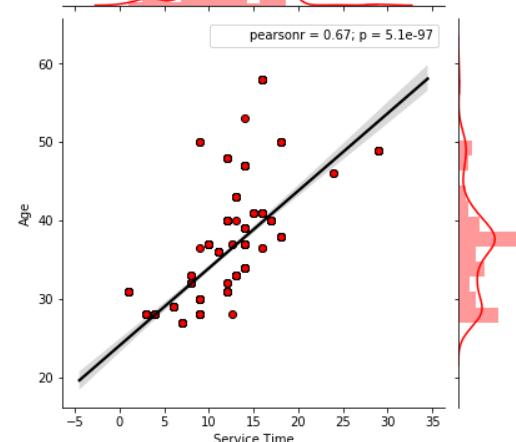
### Residence Distance v/s Body Mass Index



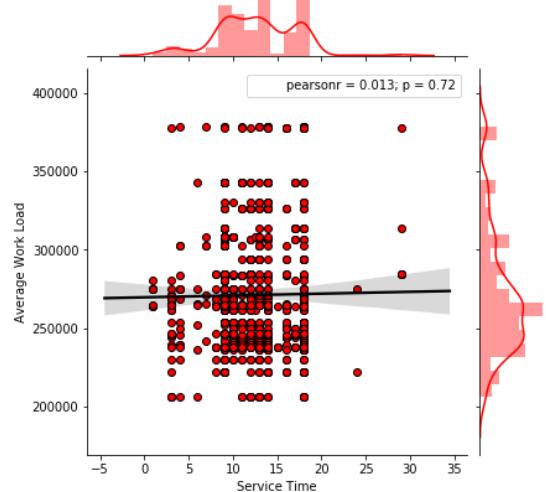
### Residence Distance v/s Absent Tim



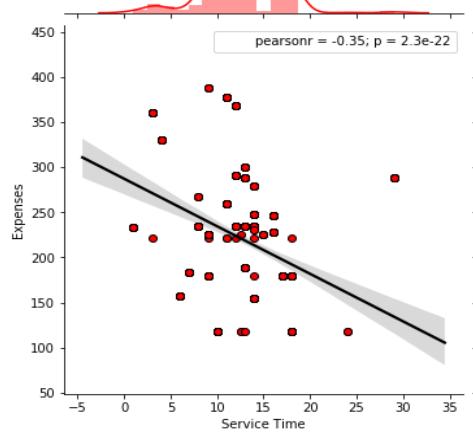
### Service Time v/s Age



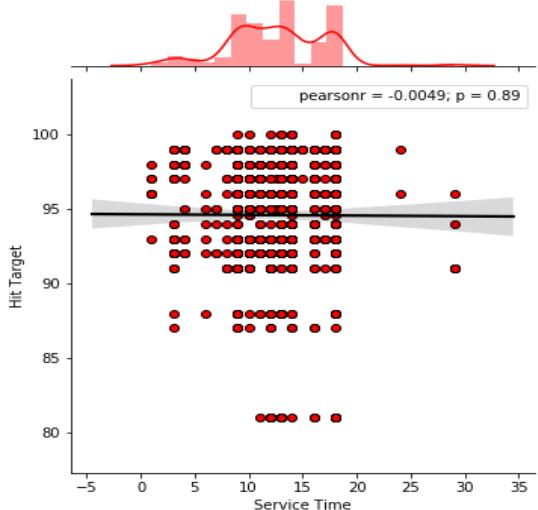
### Service Time v/s Average Work Load



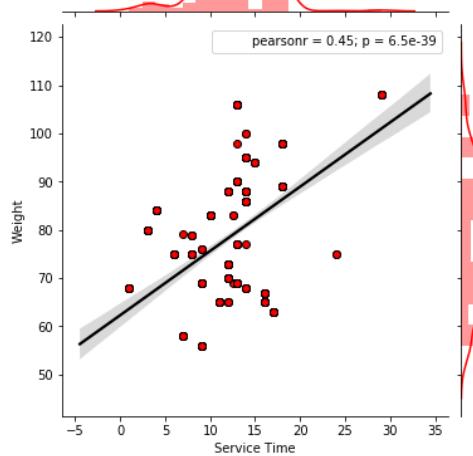
### Service Time v/s Expenses



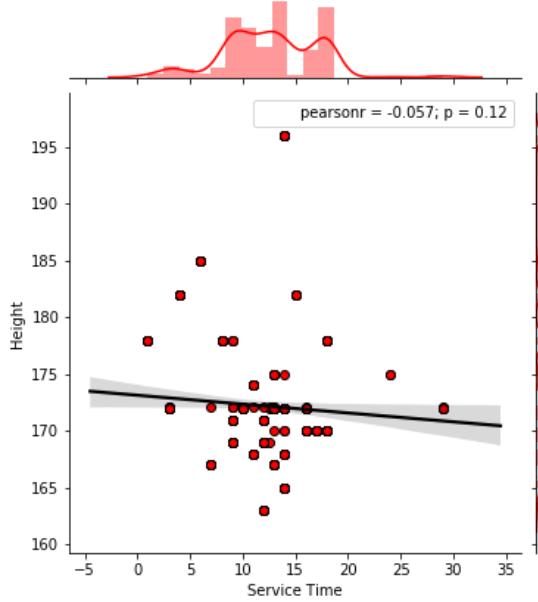
### Service Time v/s Hit Target



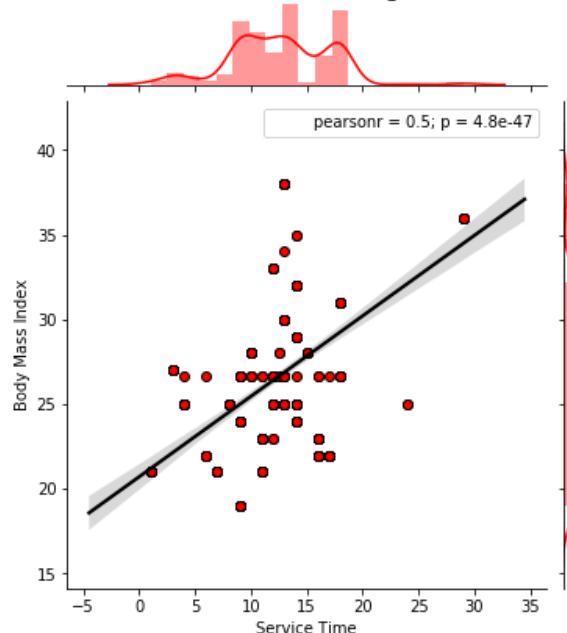
### Service Time v/s Weight



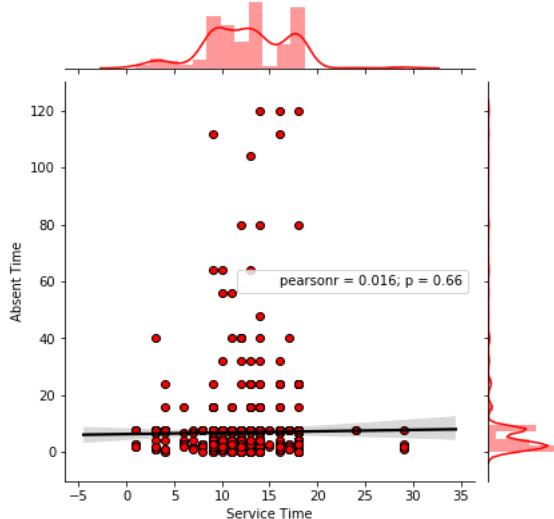
### Service Time v/s Height



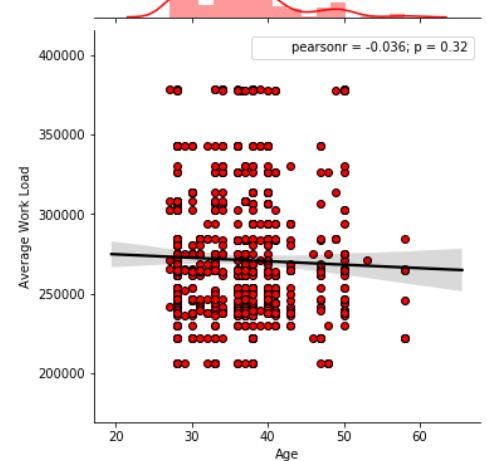
### Service Time v/s Body Mass Ind



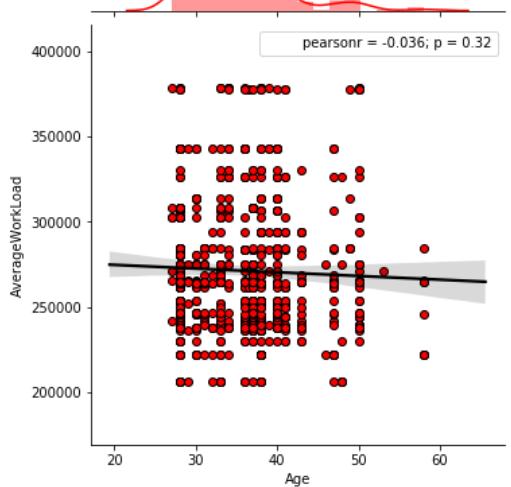
### Service Time v/s Absent Time



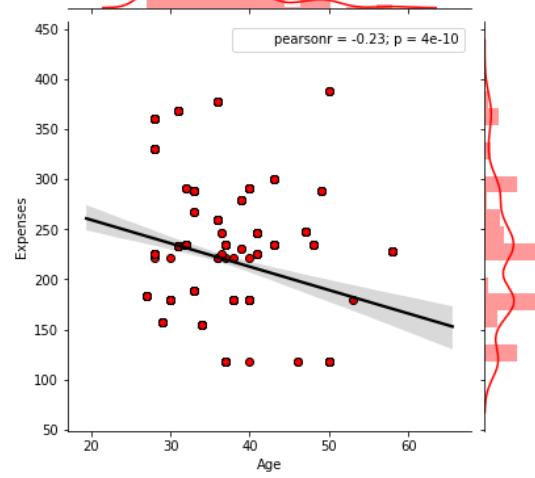
### Age v/s Average Work Load



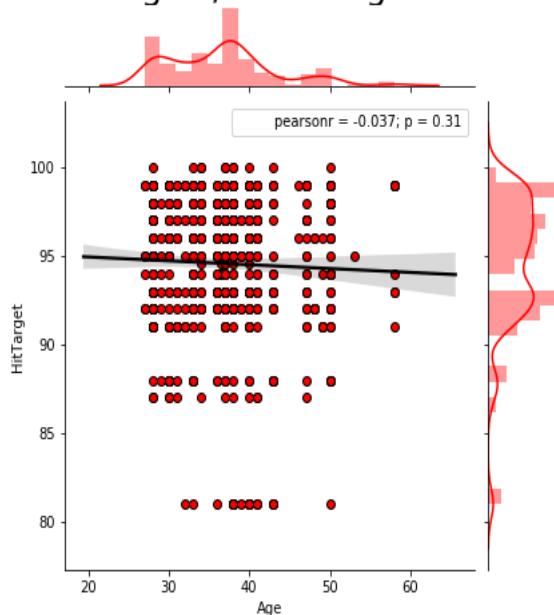
### Age v/s Average Work Load



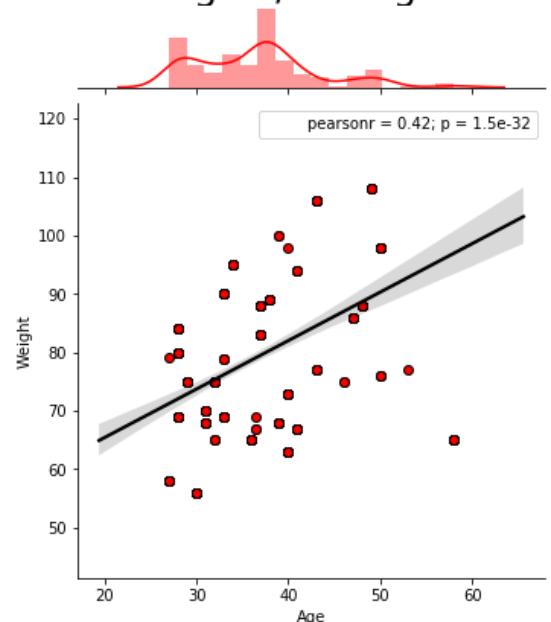
### Age v/s Expenses



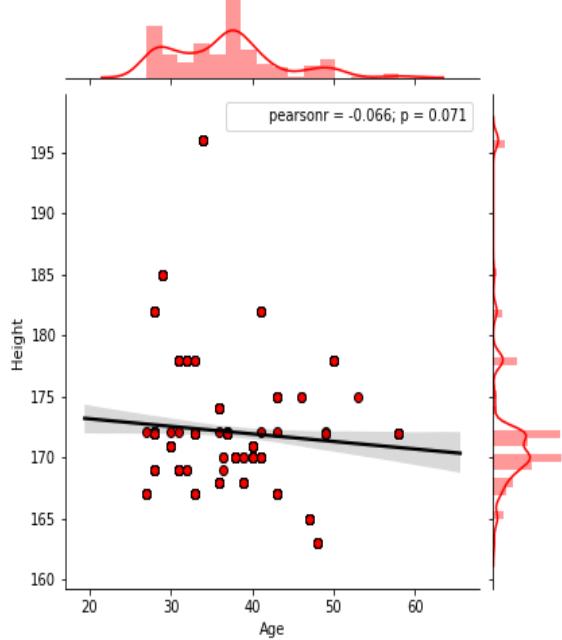
### Age v/s Hit Target



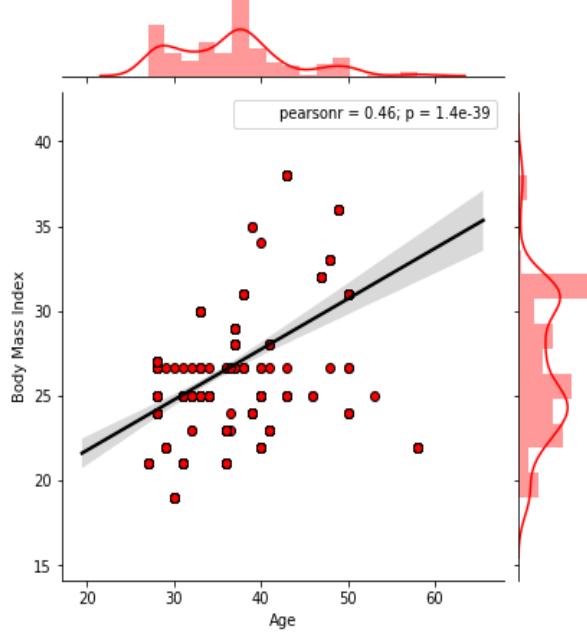
### Age v/s Weight



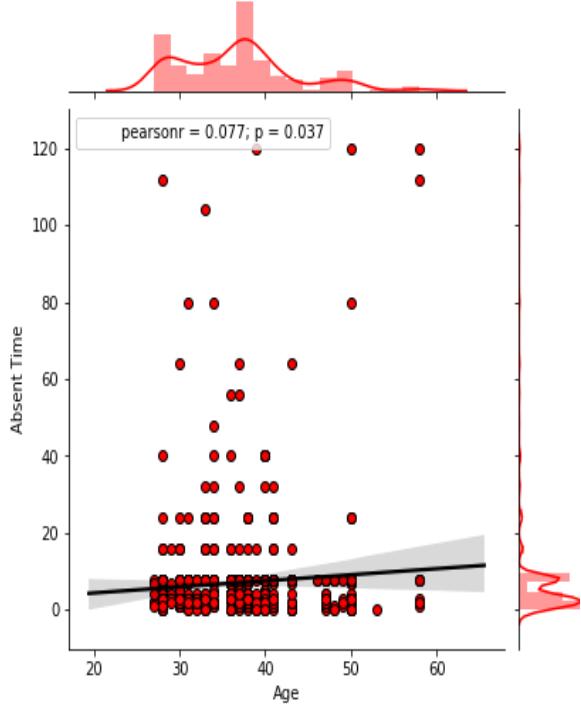
### Age v/s Height



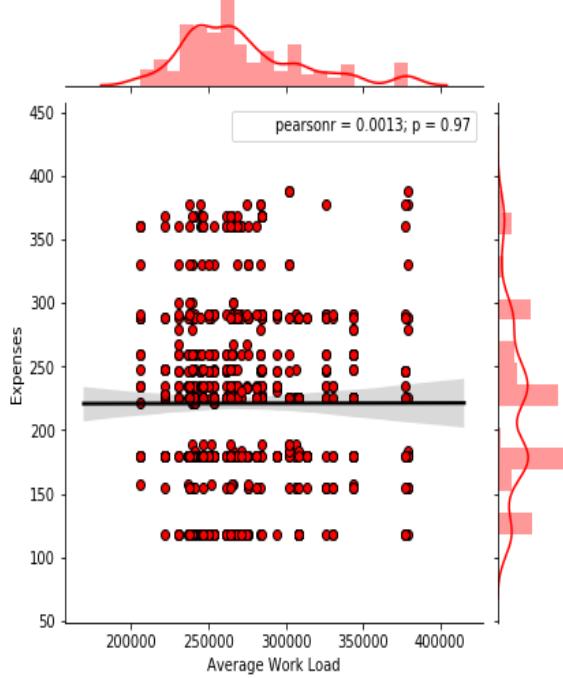
### Age v/s Body Mass Index



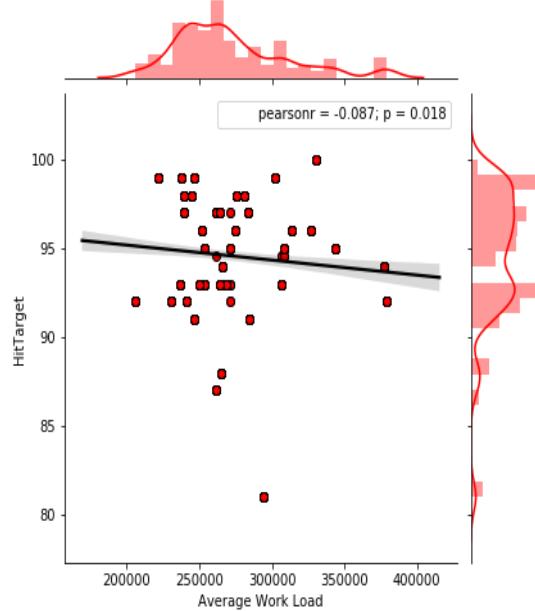
### Age v/s Absent Time



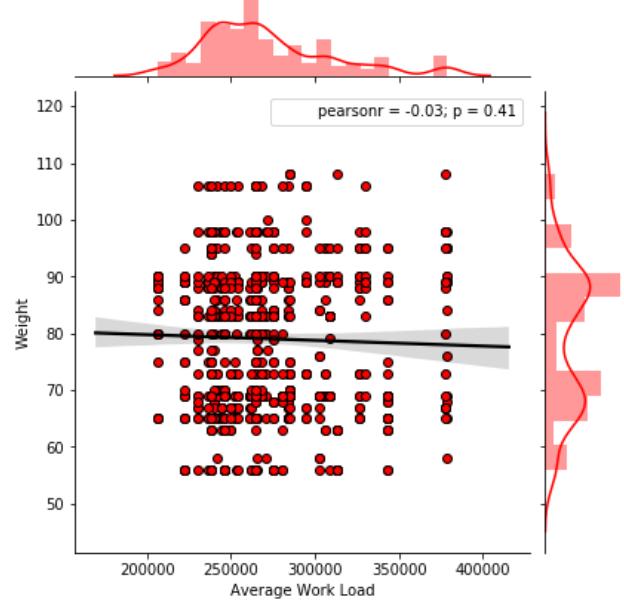
### Average Work Load v/s Expenses



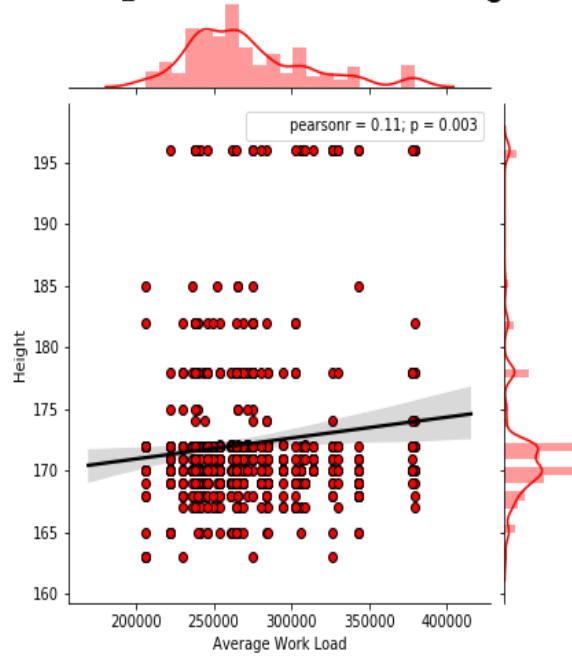
### Average Work Load v/s Hit Target



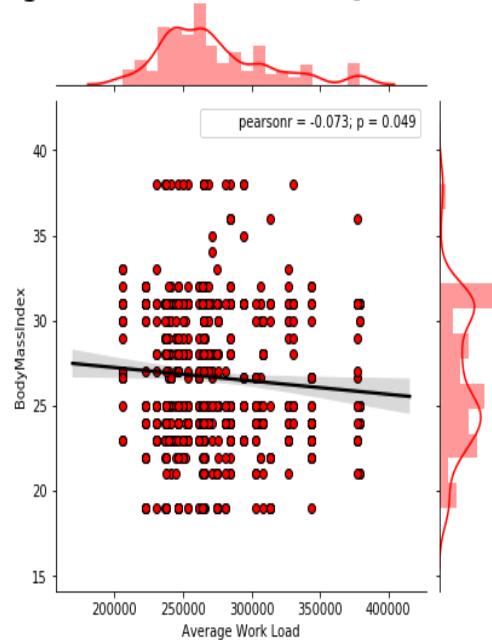
### Average Work Load v/s Weight



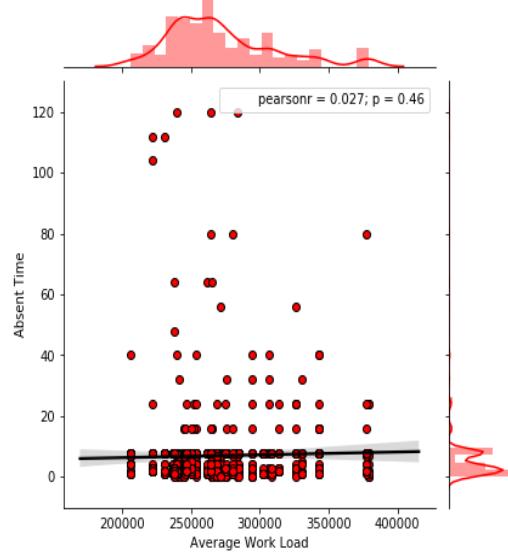
### Average Work Load v/s Height



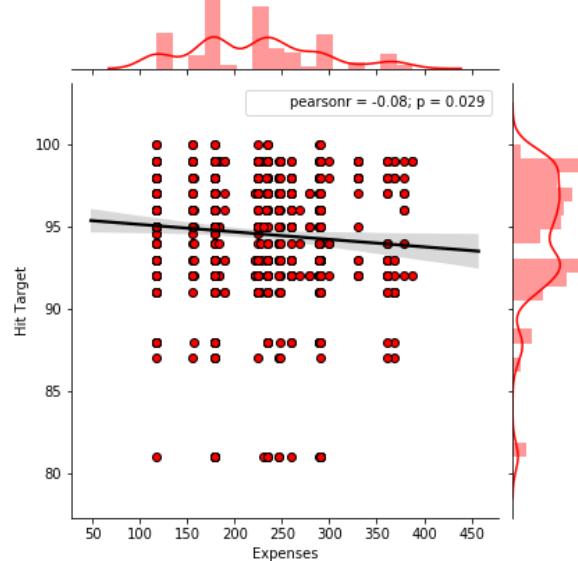
### Average Work Load v/s Body Mass Index



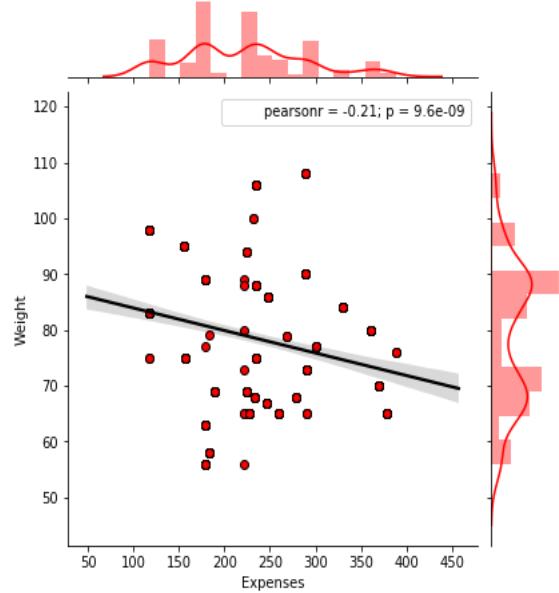
### Average Work Load v/s Absent Time



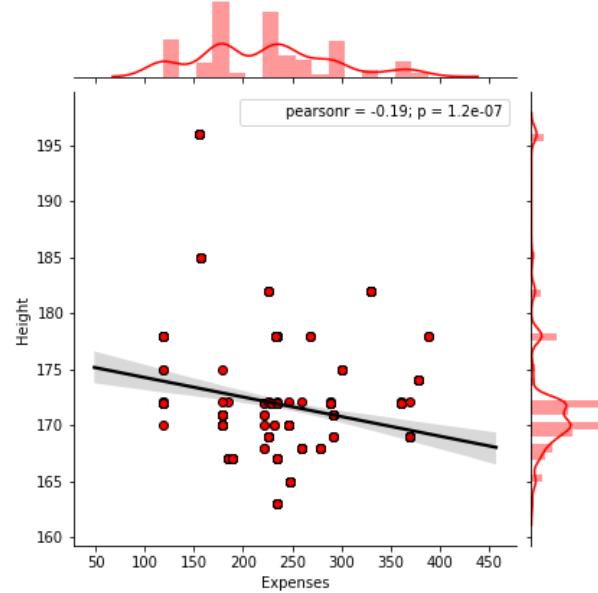
### Expenses v/s Hit Target



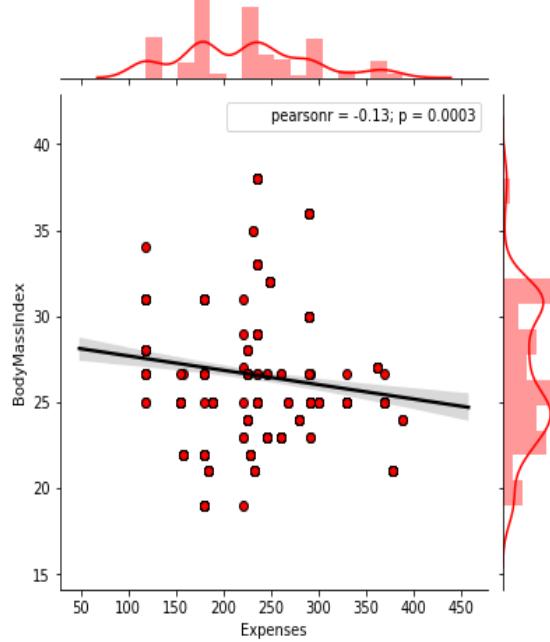
### Expenses v/s Weight



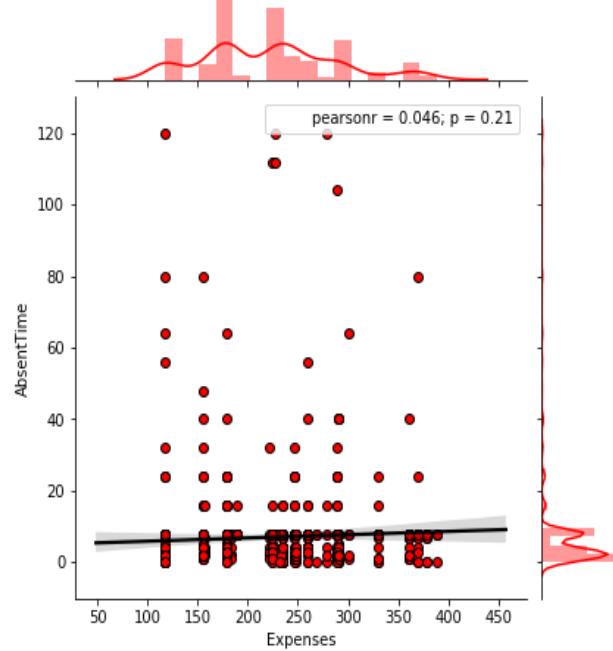
### Expenses v/s Height



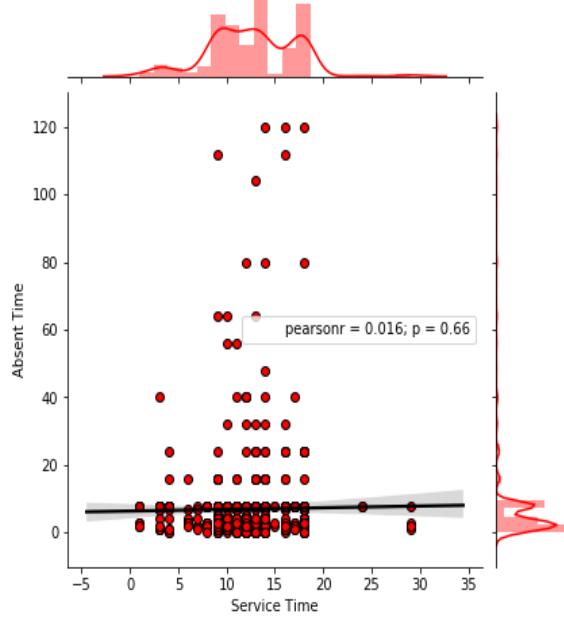
### Expenses v/s Body Mass Index



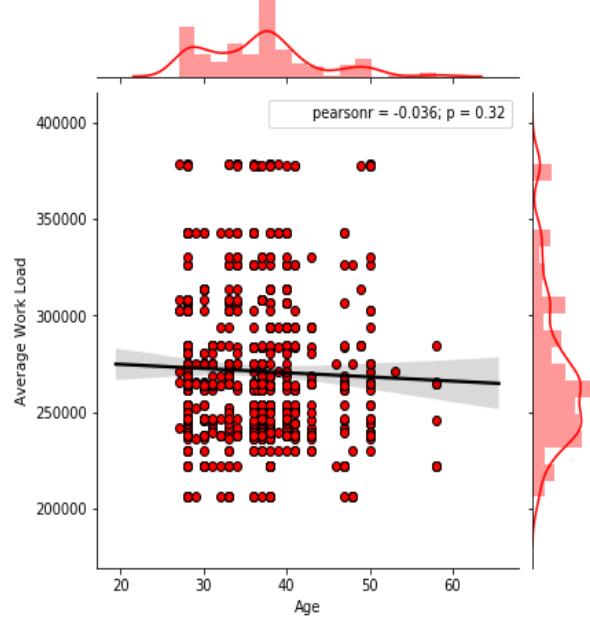
### Expenses v/s Absent Time



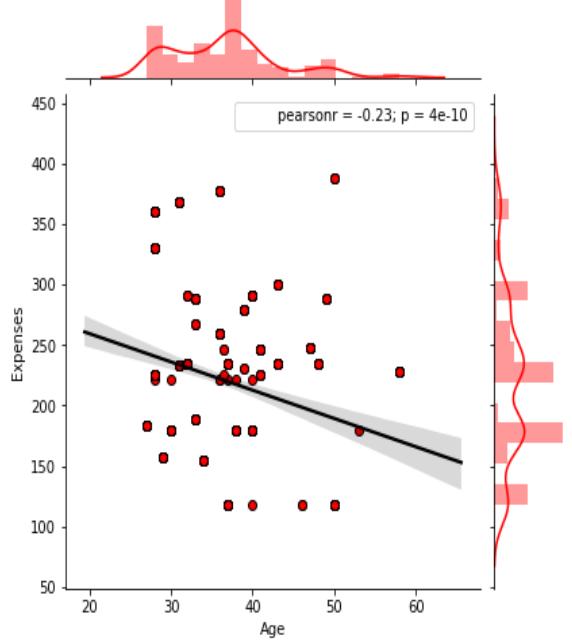
### Service Time v/s Absent Time



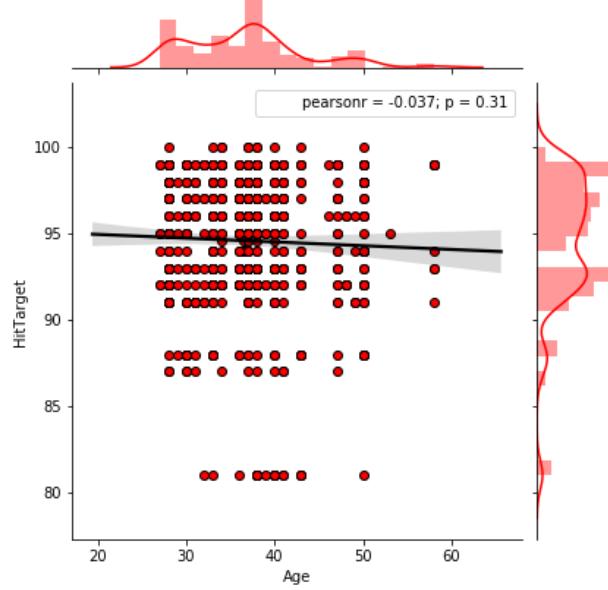
### Age v/s Average Work Load



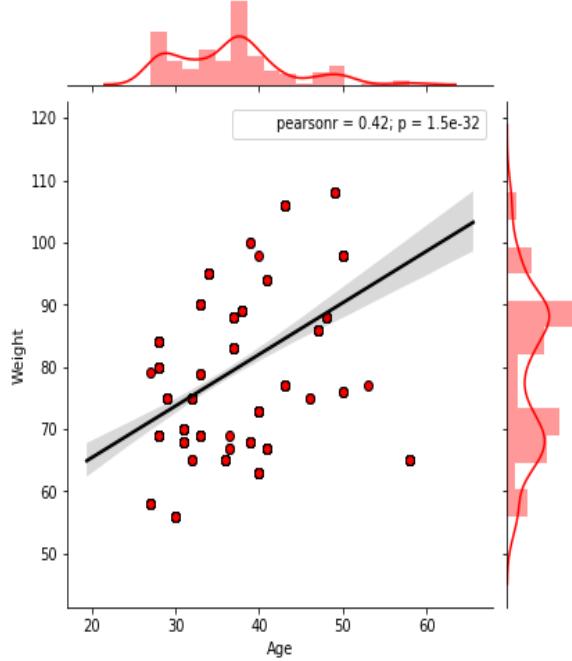
### Age v/s Expenses



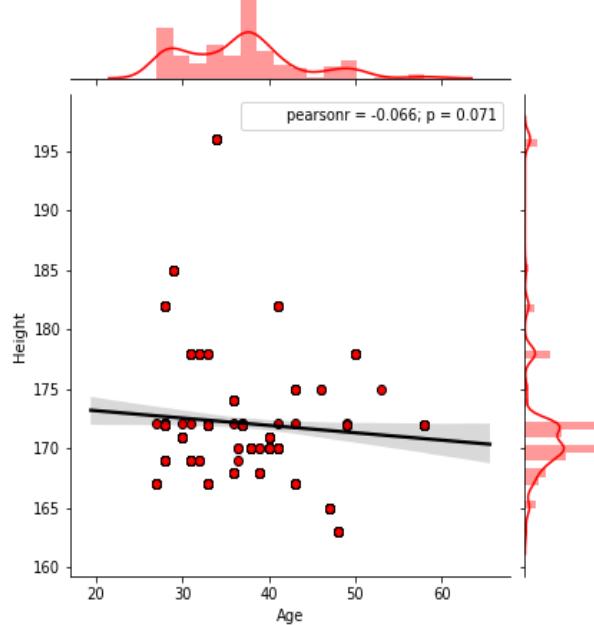
### Age v/s Hit Target



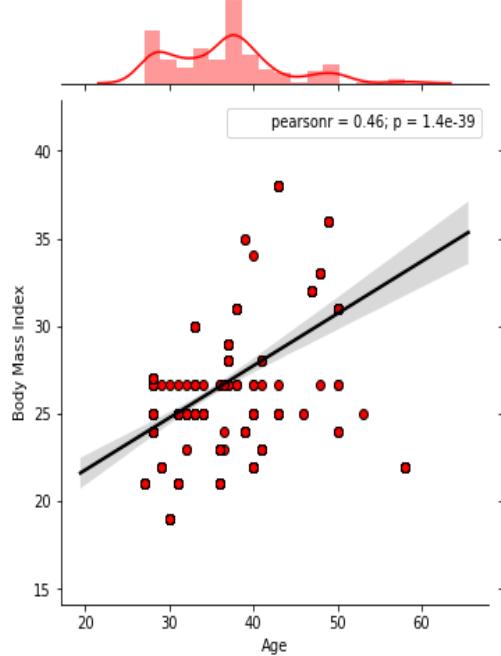
### Age v/s Weight



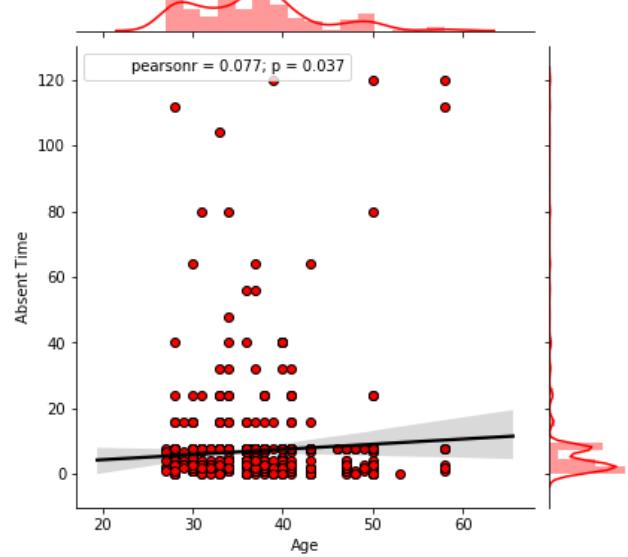
### Age v/s Height



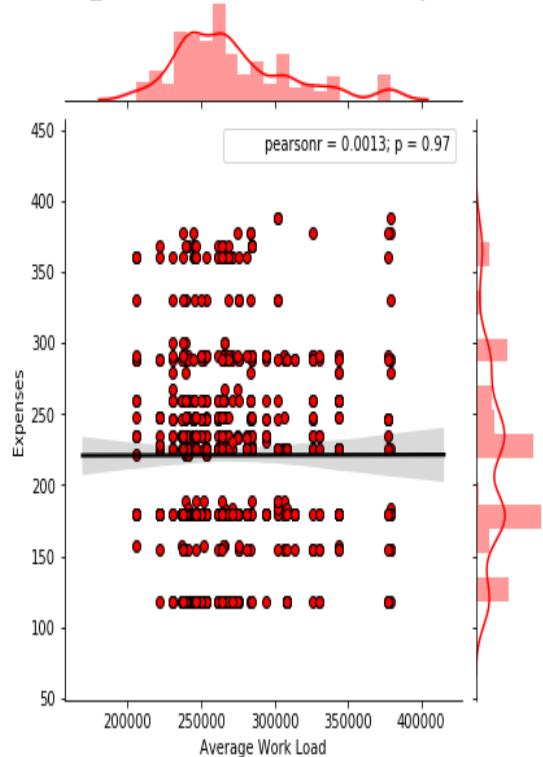
### Age v/s Body Mass Index



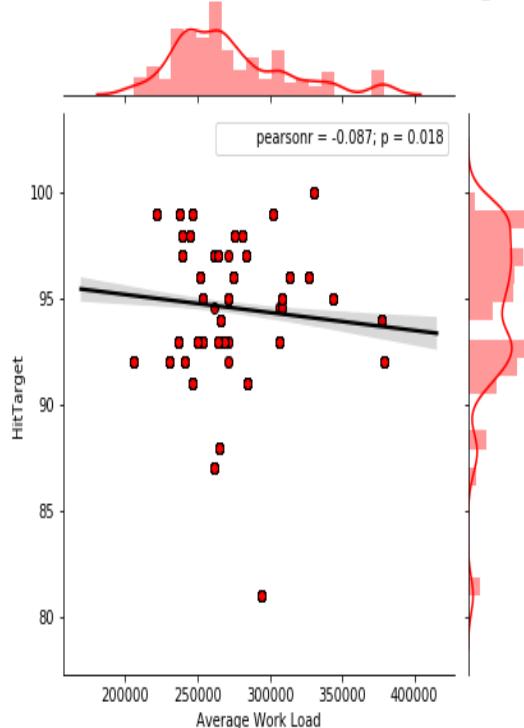
### Age v/s Absent Time



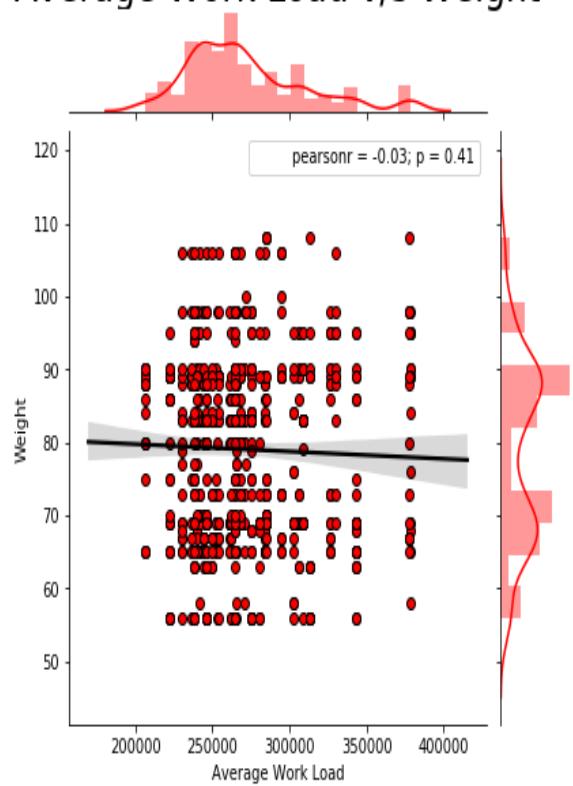
### Average Work Load v/s Expenses



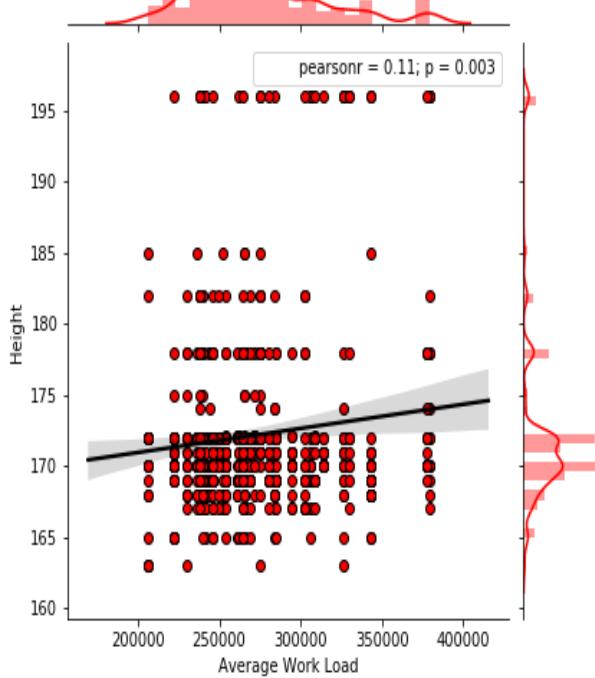
### Average Work Load v/s Hit Target



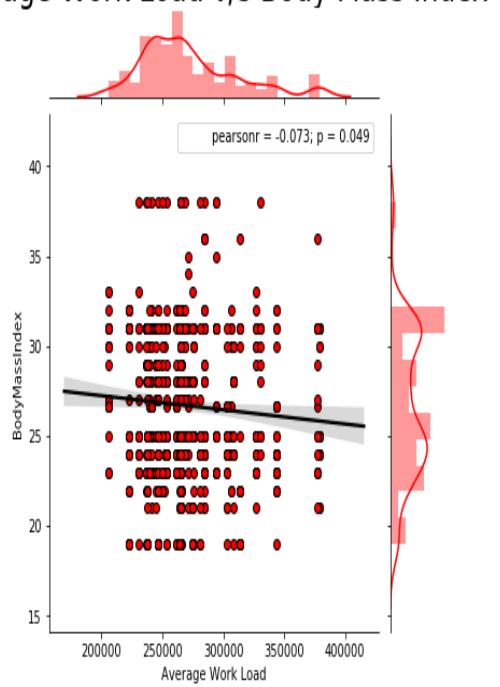
### Average Work Load v/s Weight



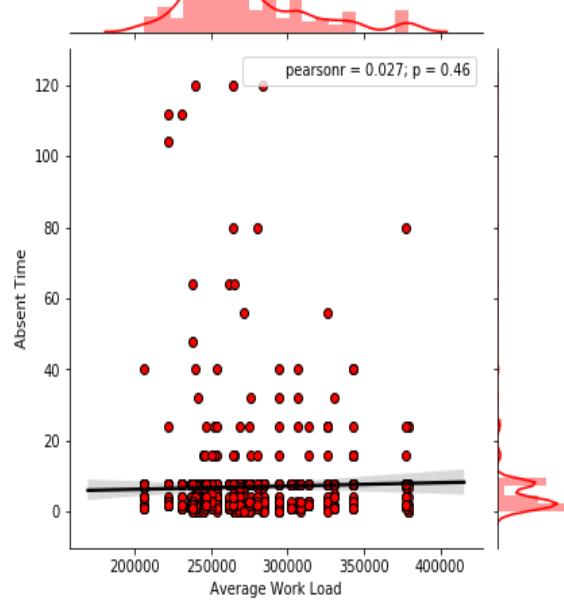
### Average Work Load v/s Height



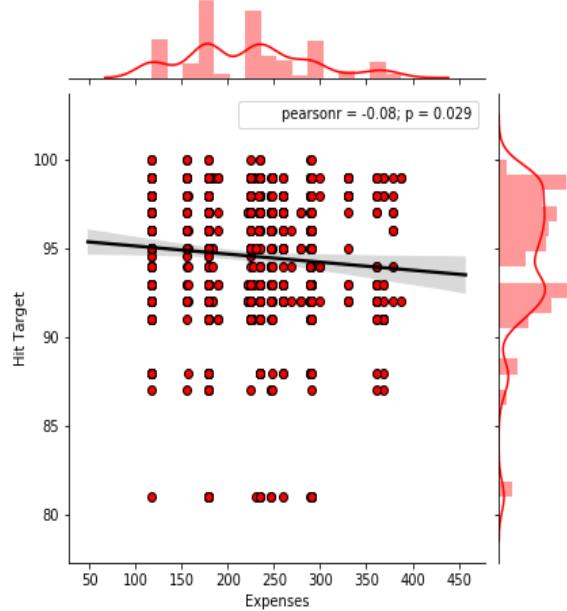
### Average Work Load v/s Body Mass Index



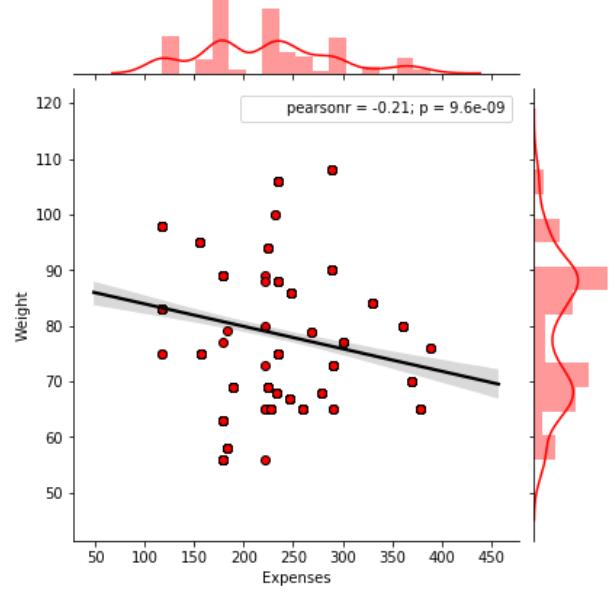
### Average Work Load v/s Absent Time



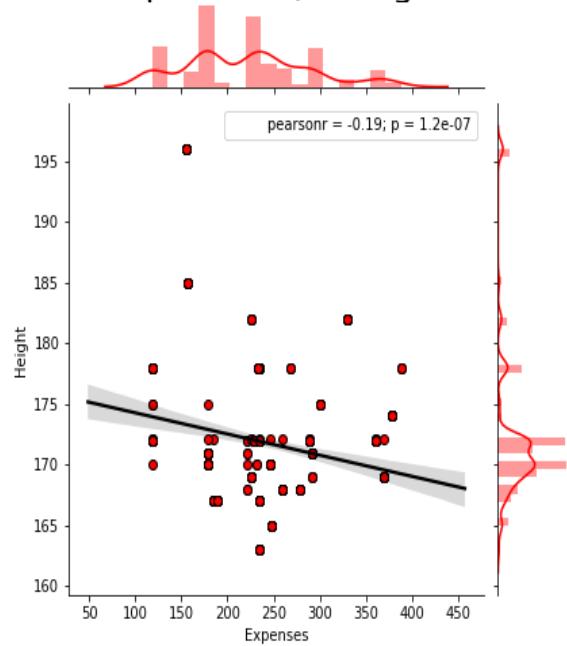
### Expenses v/s Hit Target



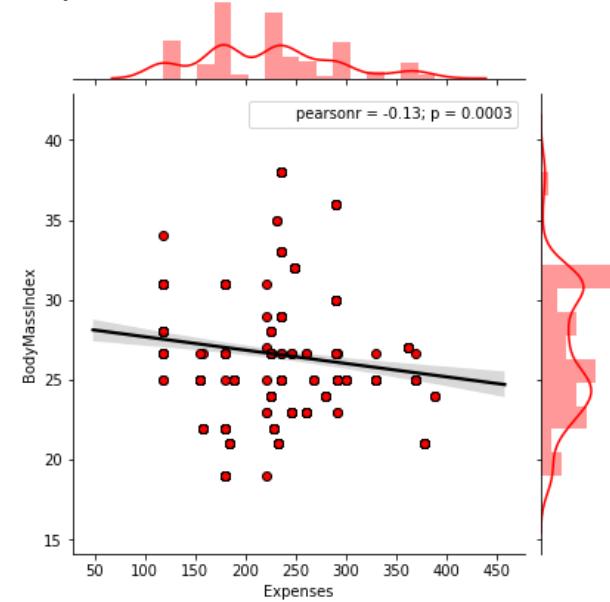
### Expenses v/s Weight



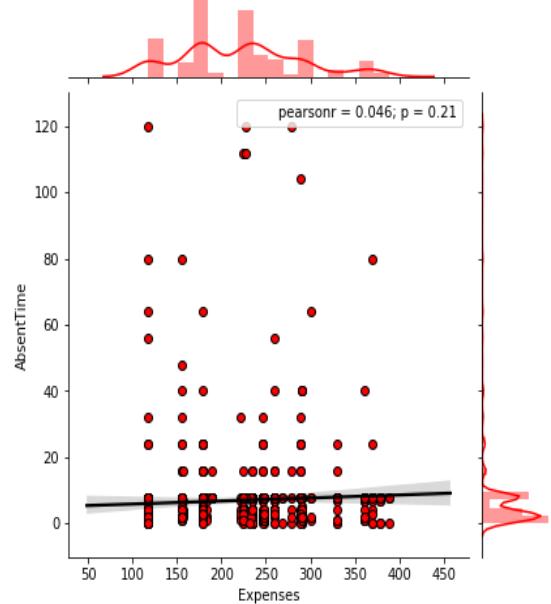
### Expenses v/s Height



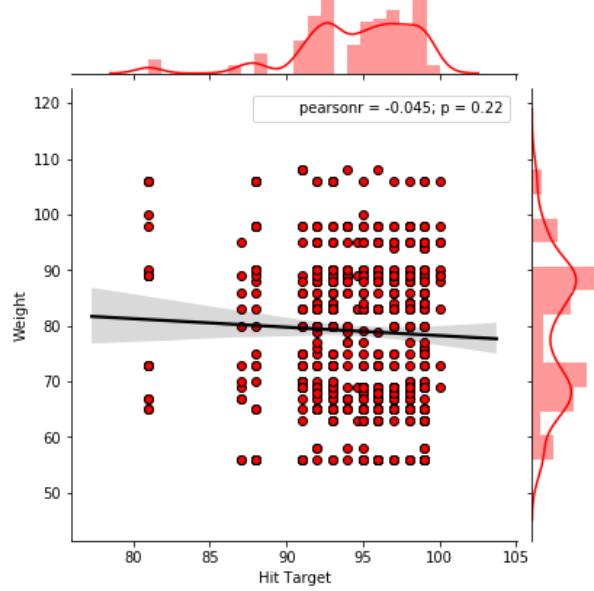
### Expenses v/s Body Mass Index



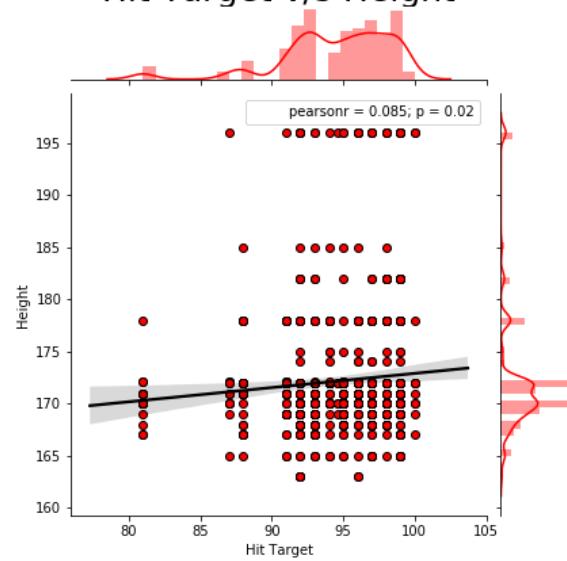
### Expenses v/s Absent Time



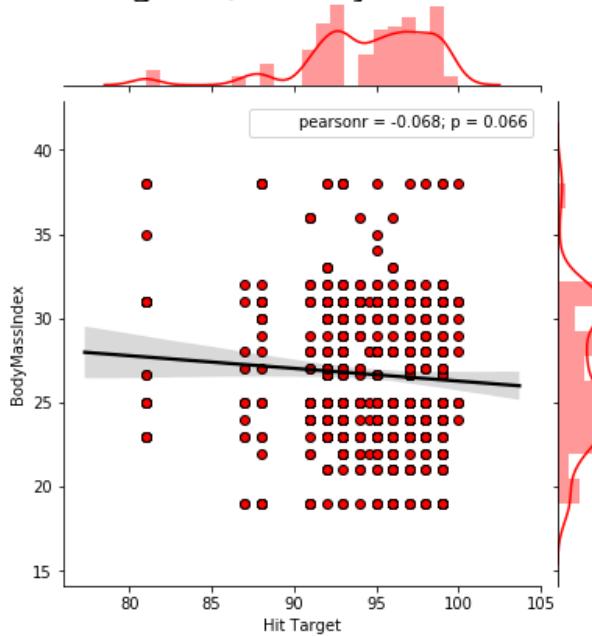
### Hit Target v/s Weight



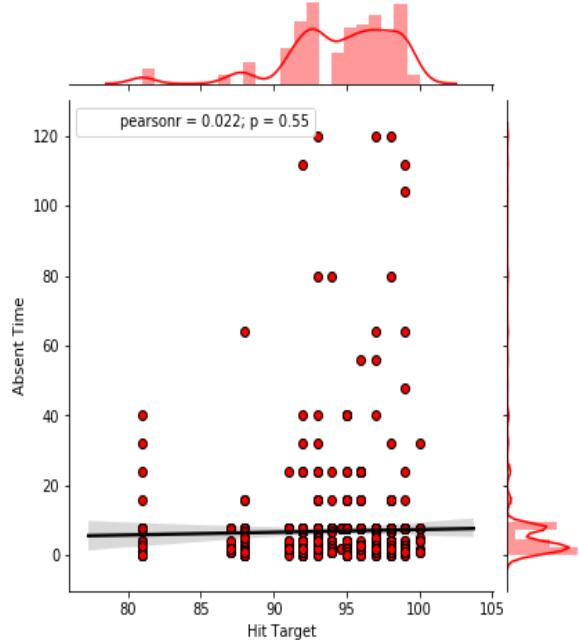
### Hit Target v/s Height



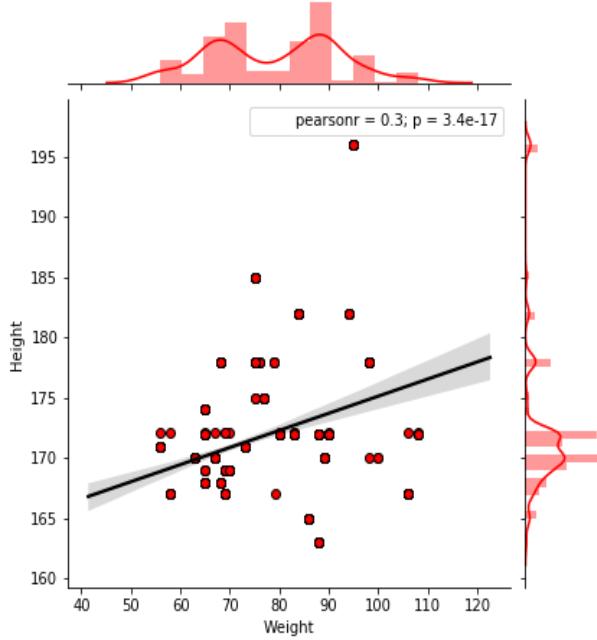
### Hit Target v/s Body Mass Index



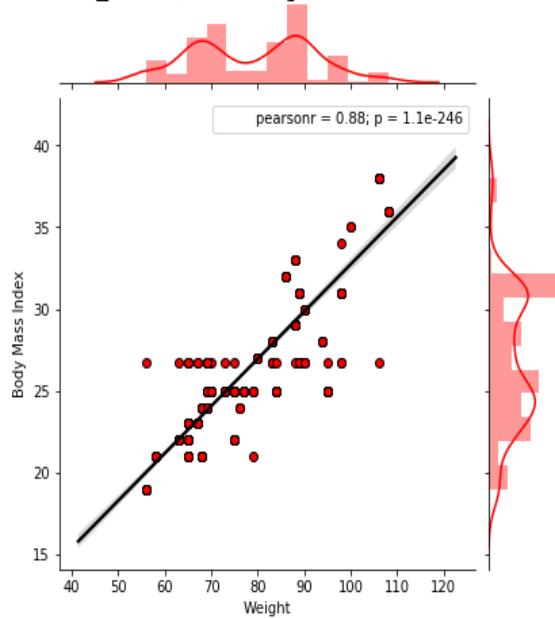
### Hit Target v/s Absent Time



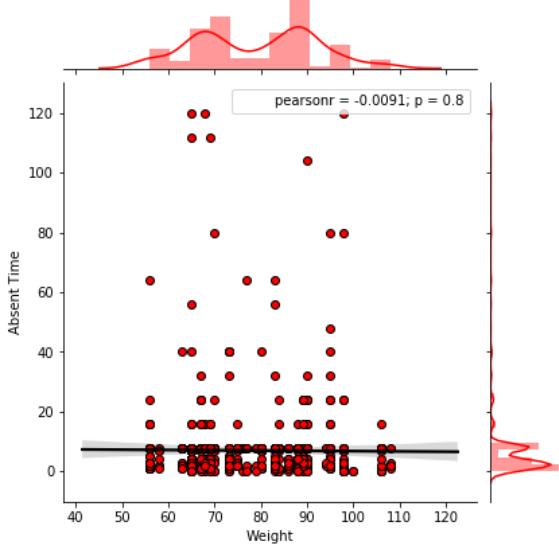
### Weight v/s Height



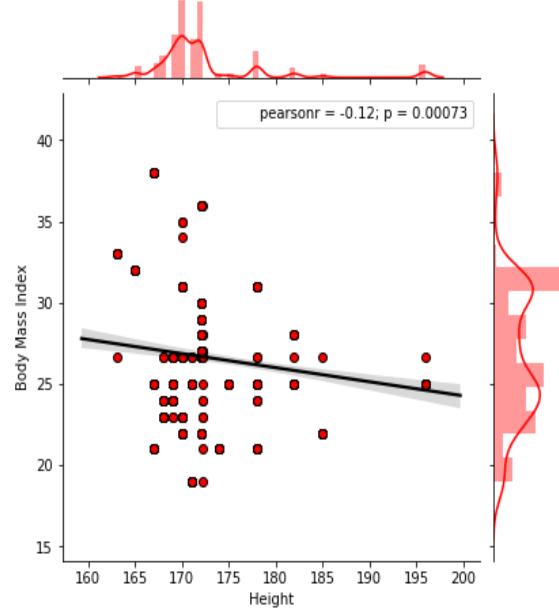
### Weight v/s Body Mass Index



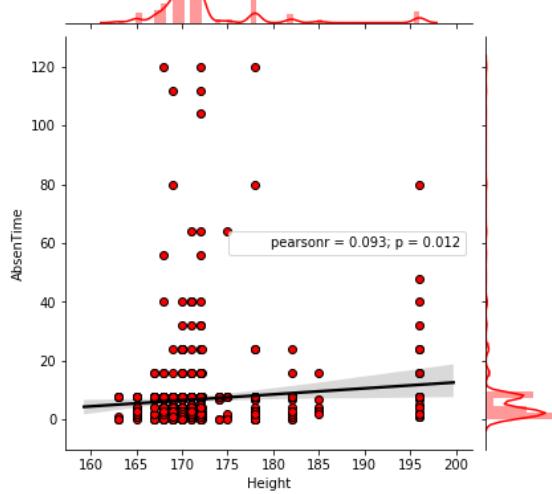
### Weight v/s Absent Time



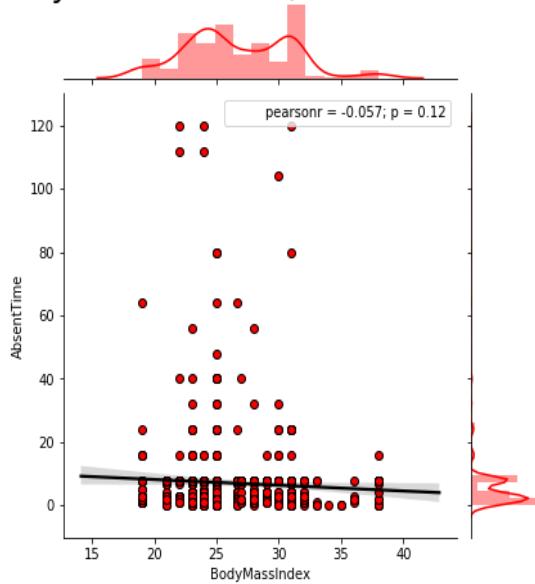
### Height v/s Body Mass Index



### Height v/s Absent Time

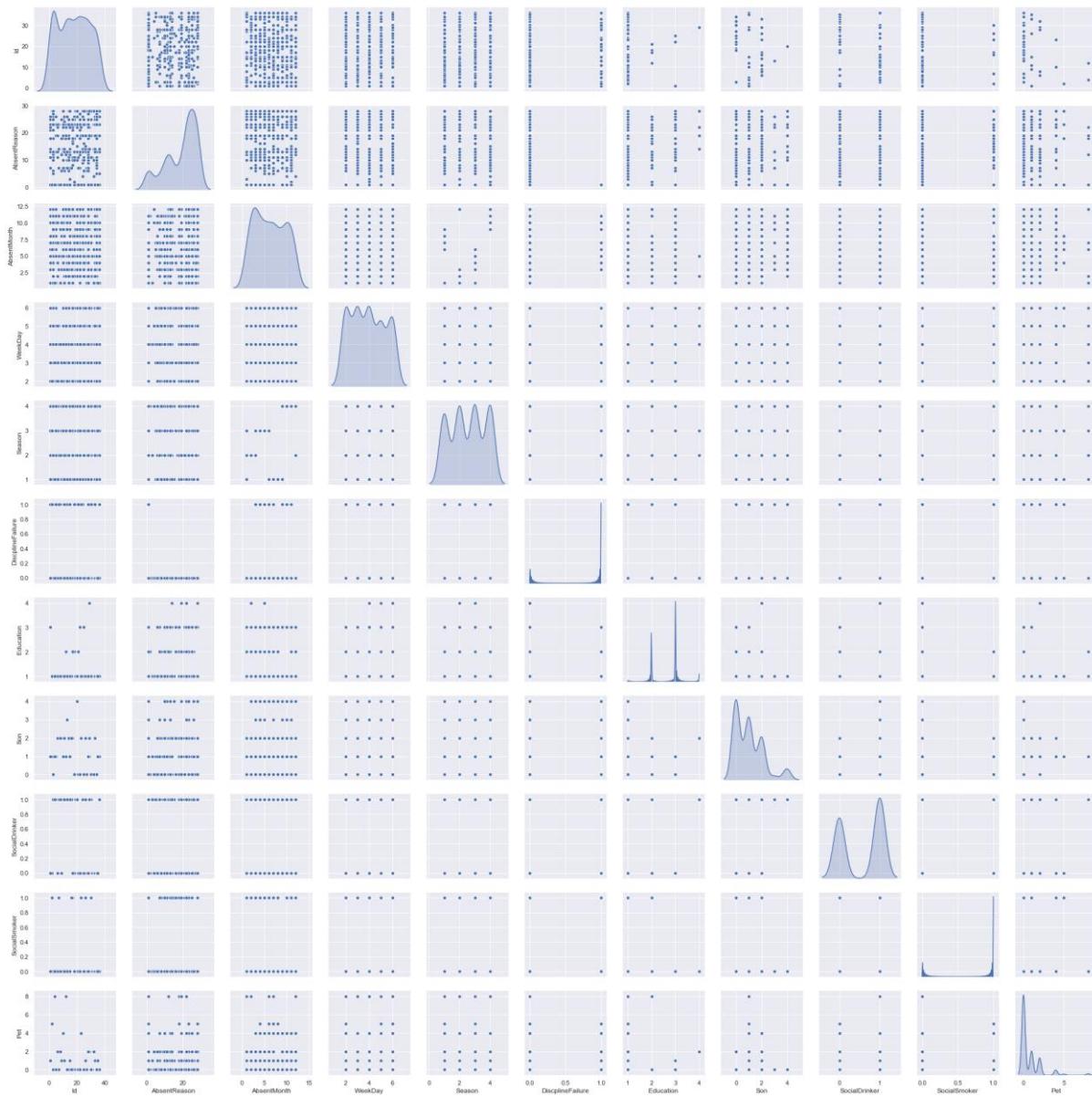


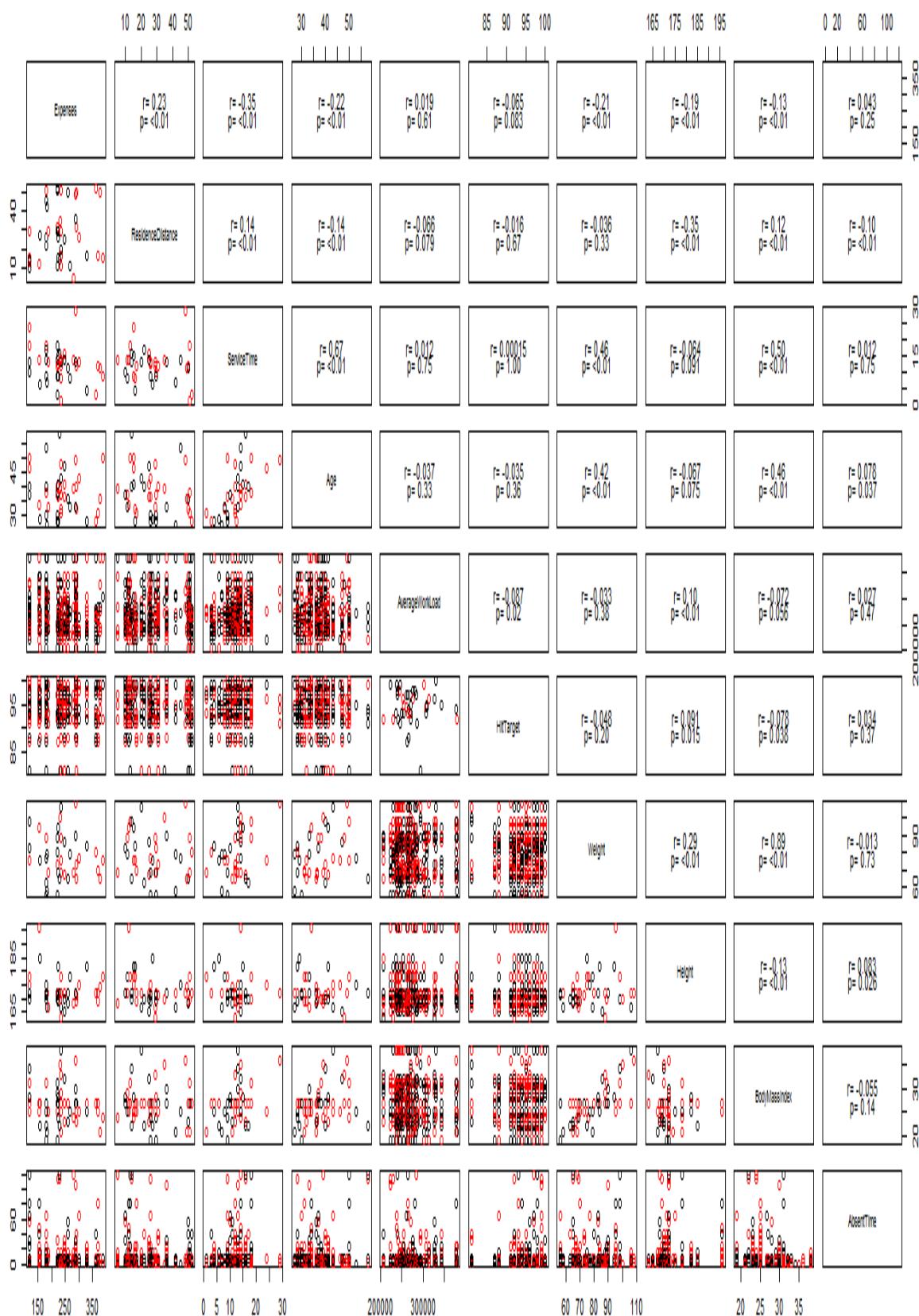
### Body Mass Index v/s Absent Time



## Pair Plot

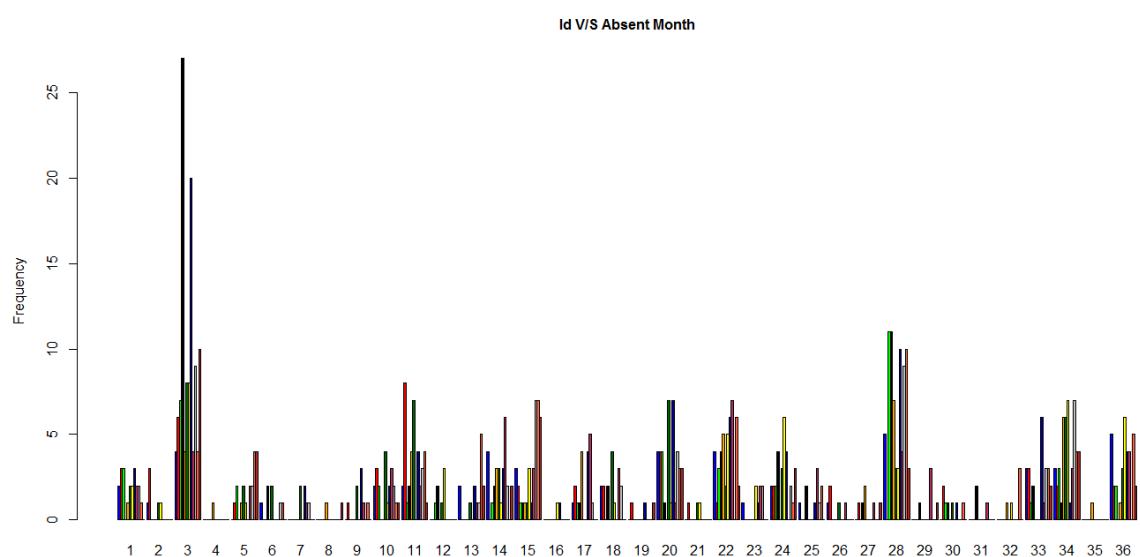
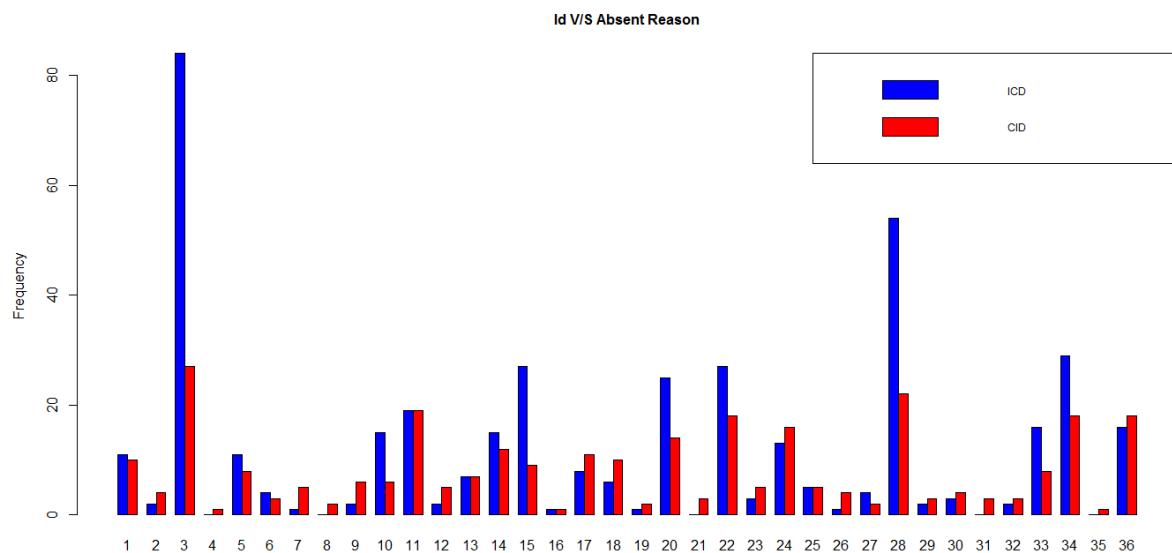
Pair plot is used to understand the best set of features to explain a relationship between two variables or to form the most separated clusters. It also helps to form some simple classification models by drawing some simple lines or make linear separation in our dataset. Pair plots for all our numerical variables are shown below:-



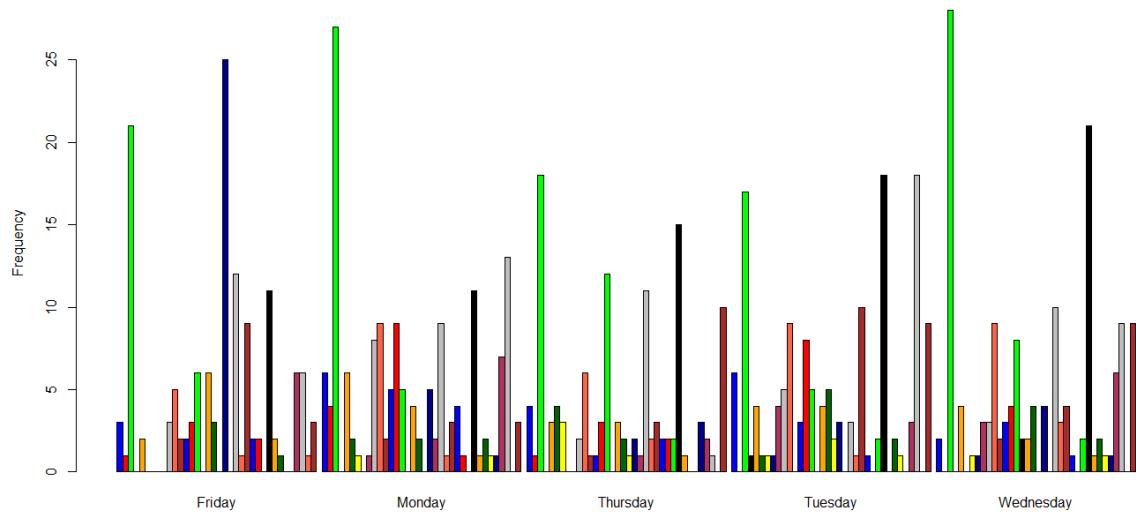


## Bivariate Analysis of Categorical Variables

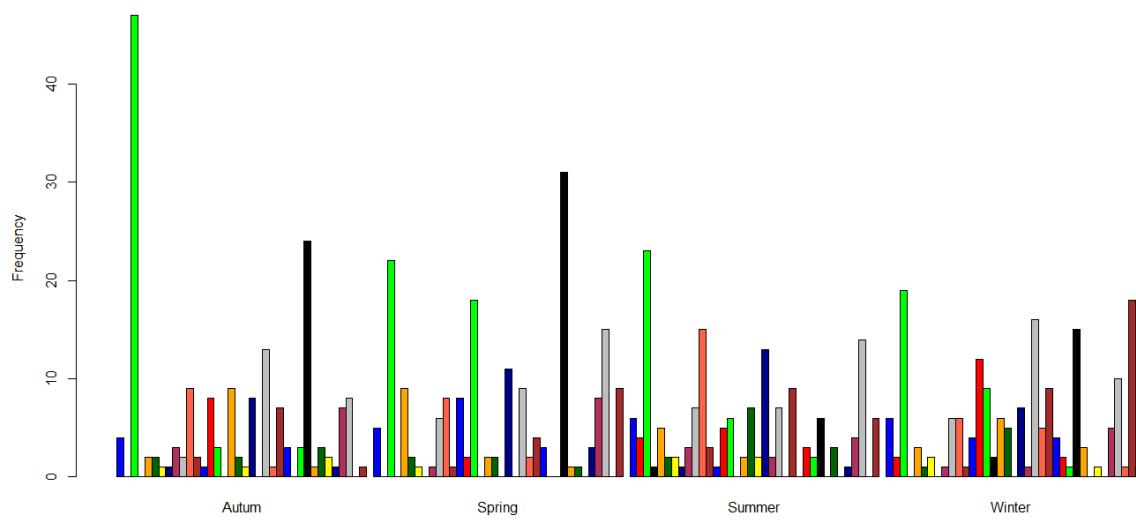
Bivariate analyses of categorical variables are shown below in the form of visualizations:-



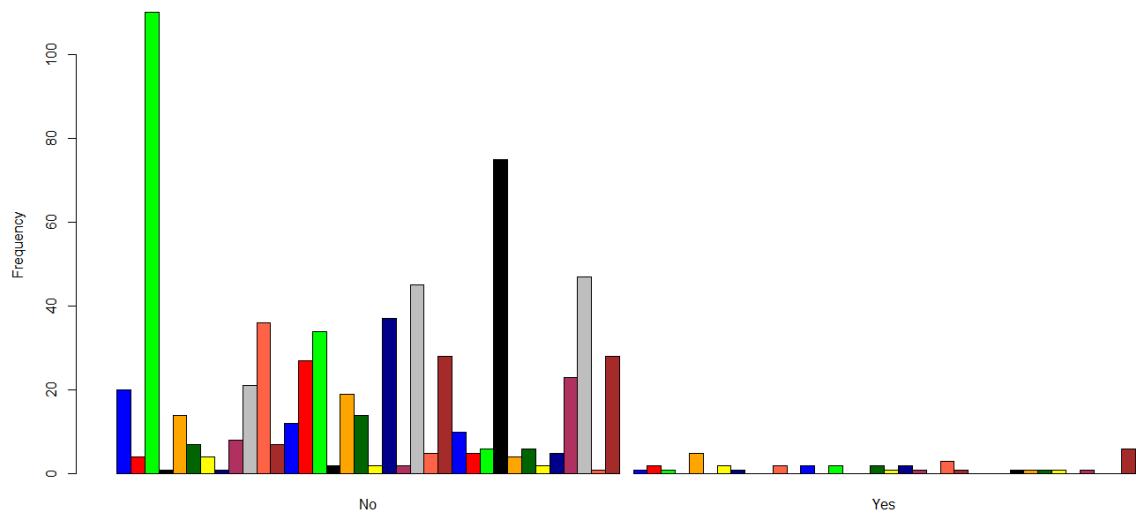
Id V/S Week Day



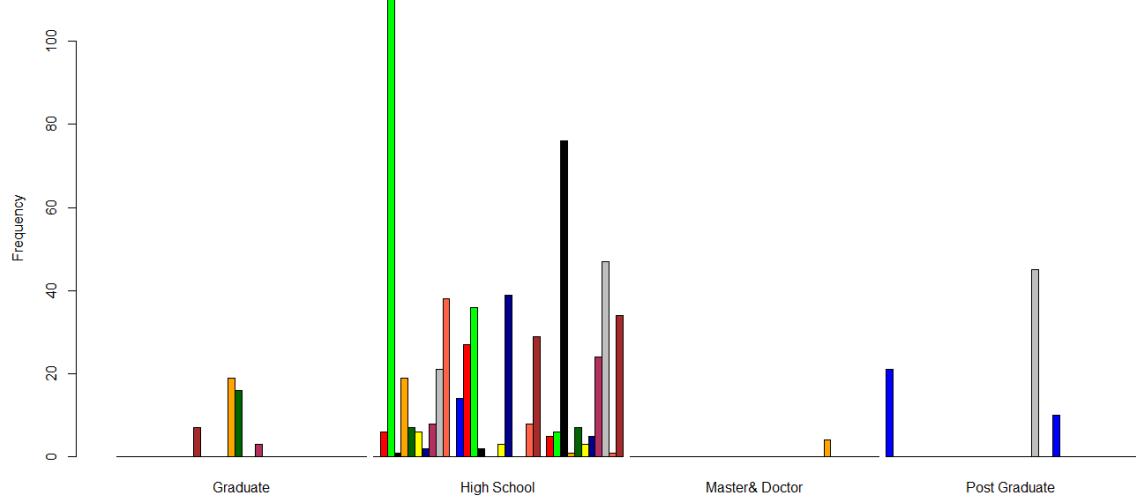
Id V/S Season

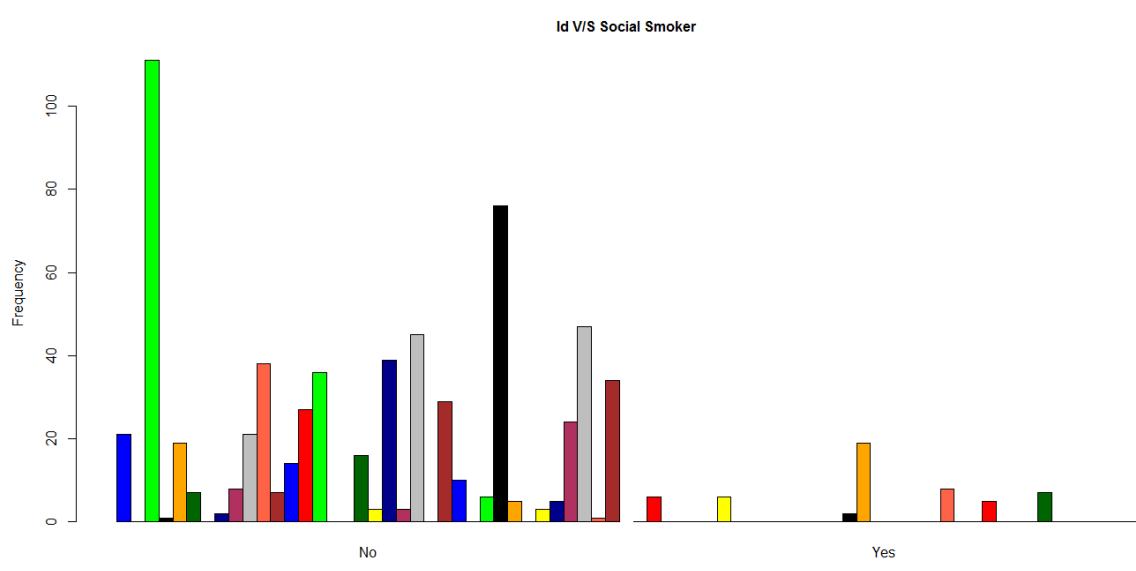
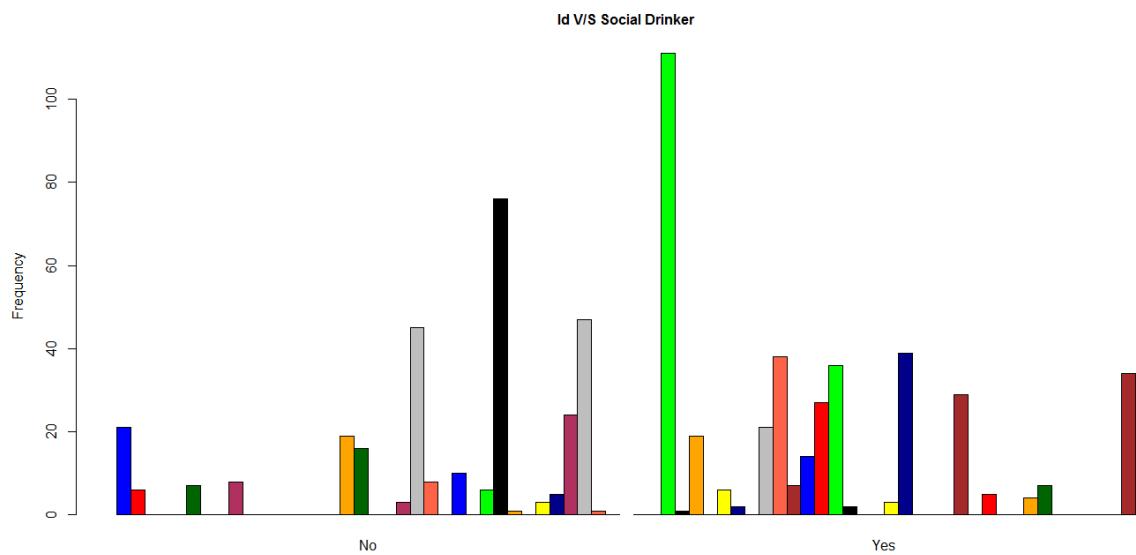


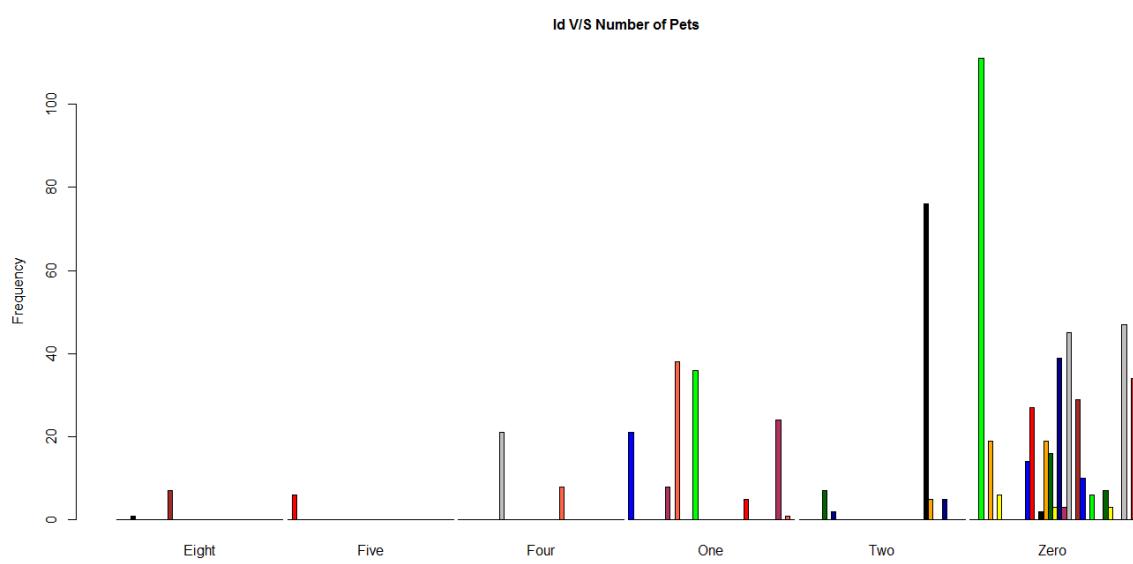
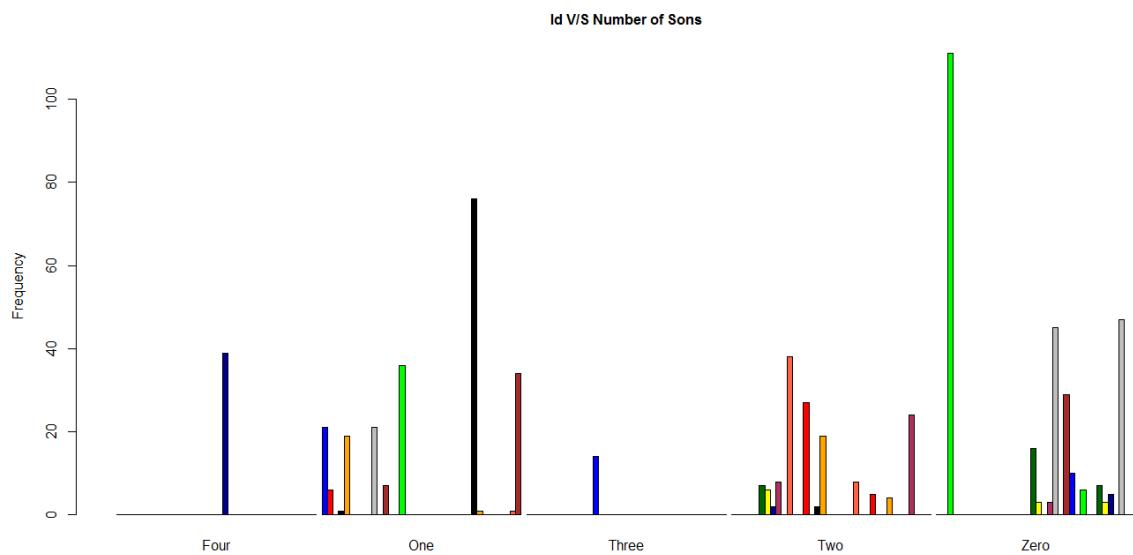
**Id V/S Discipline Failure**

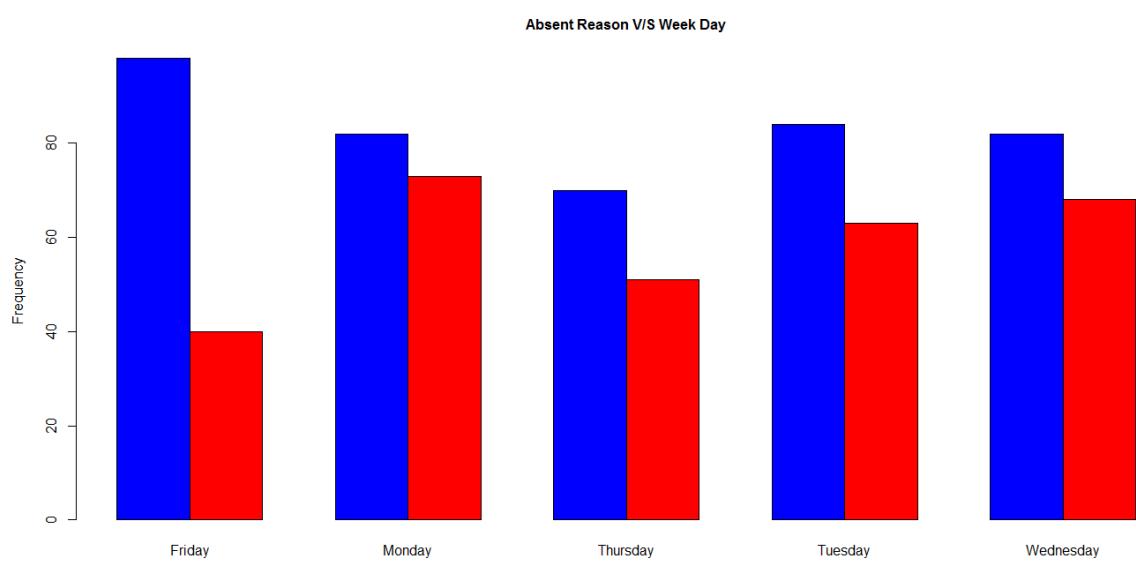
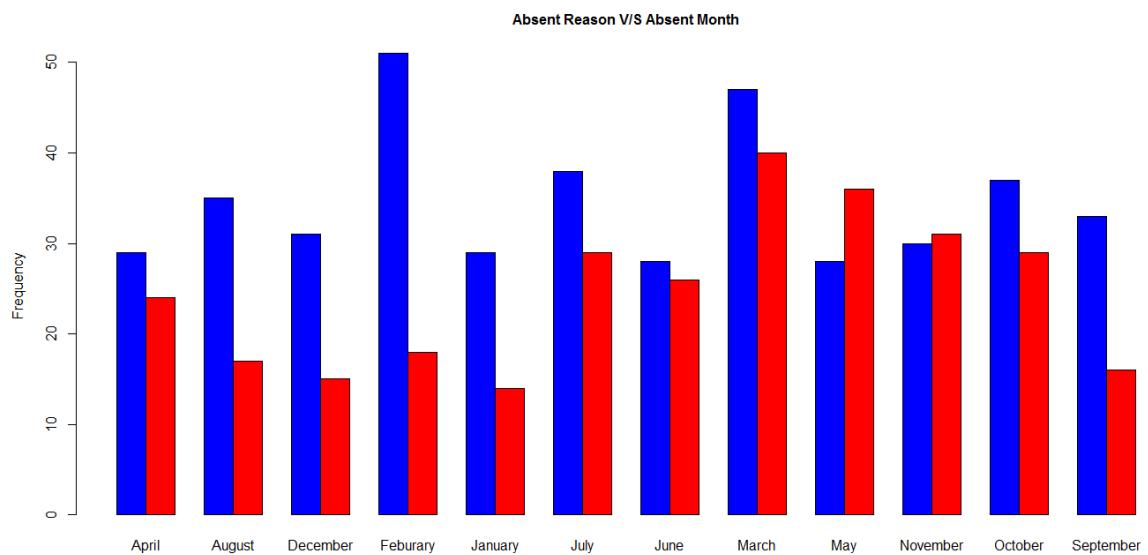


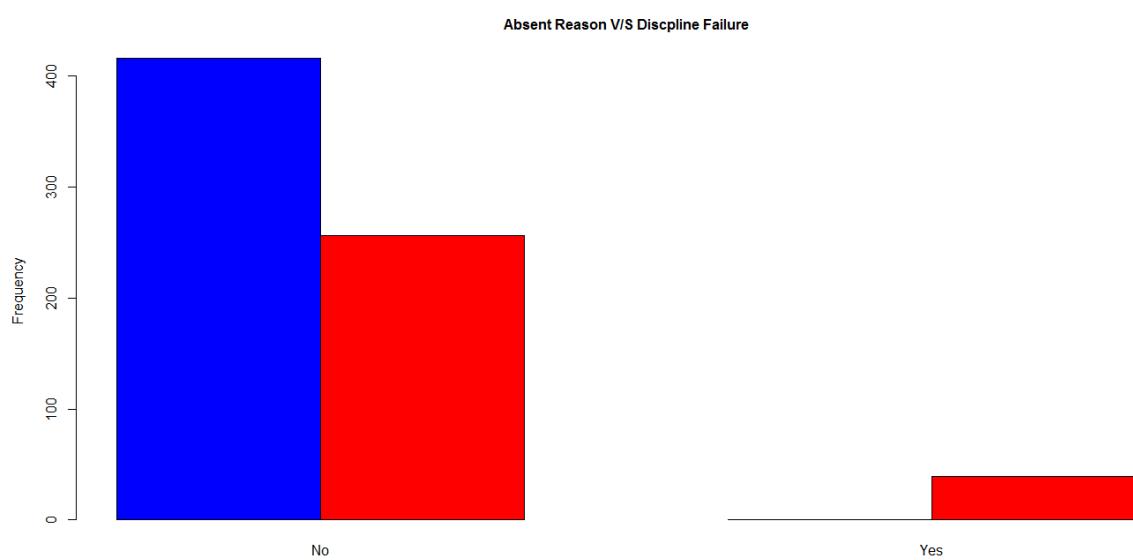
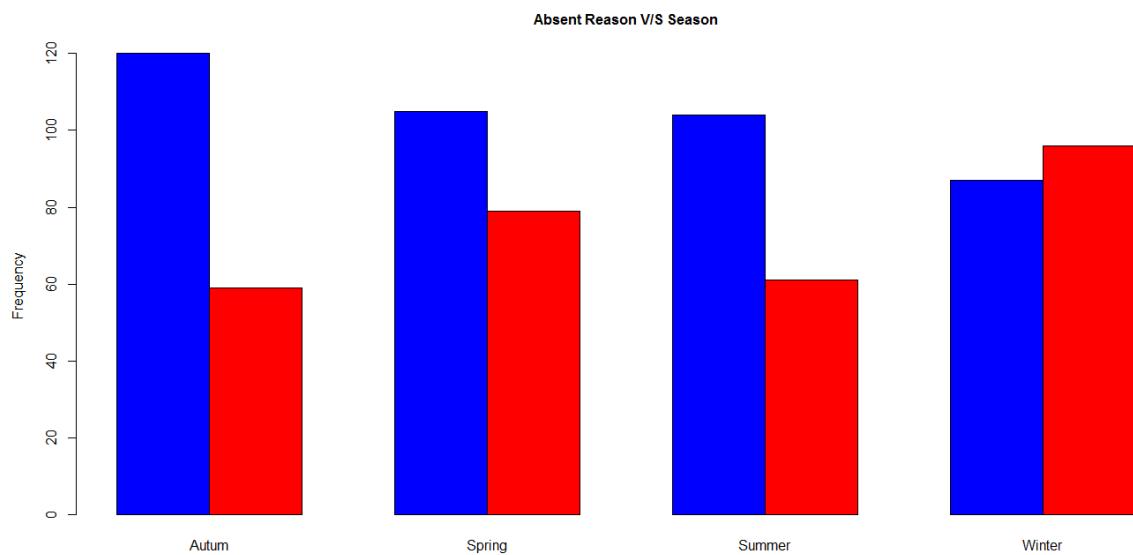
**Id V/S Education**

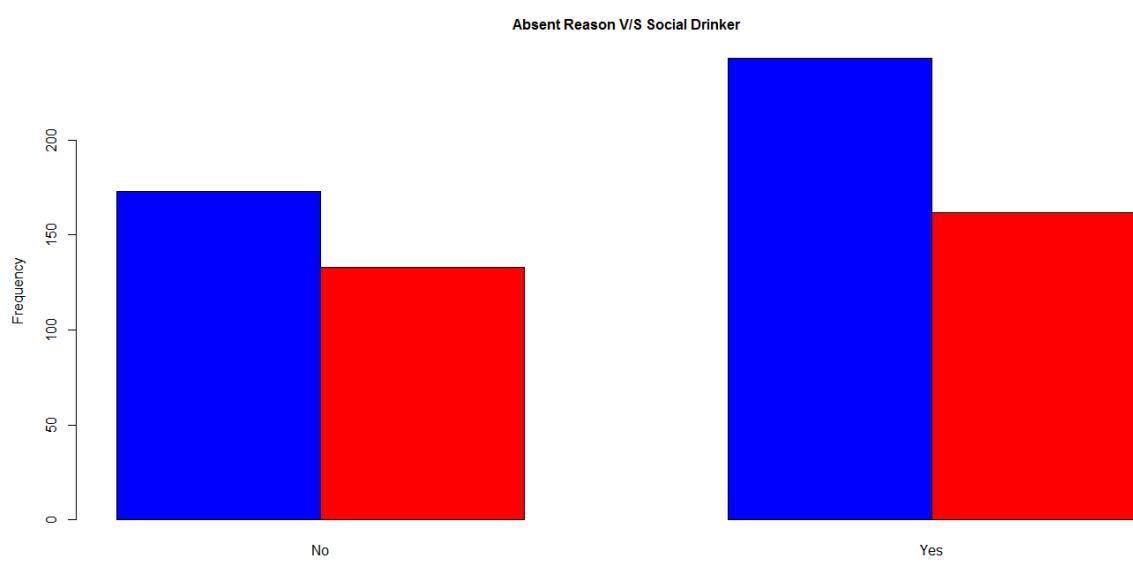
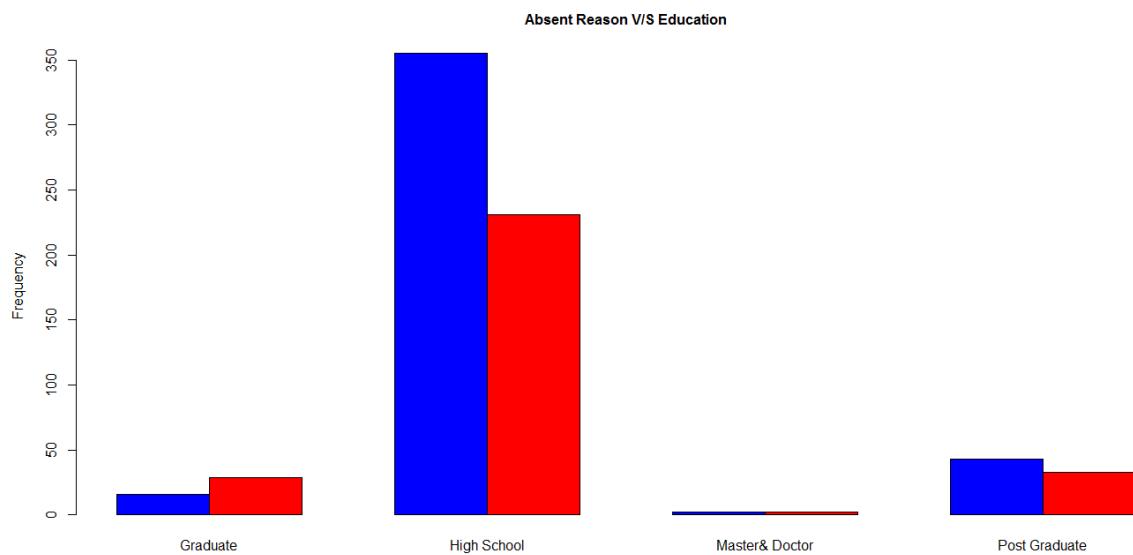




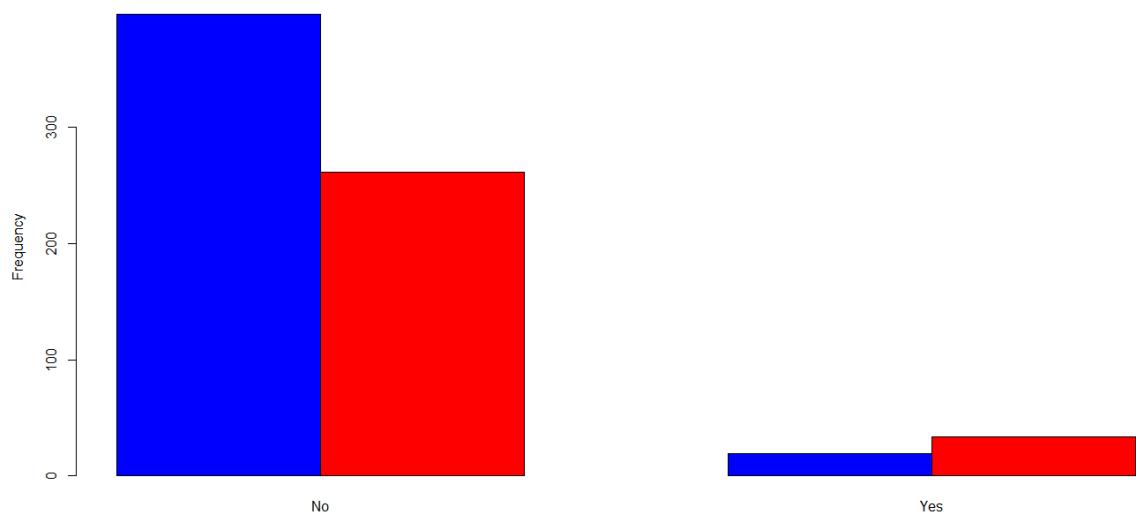




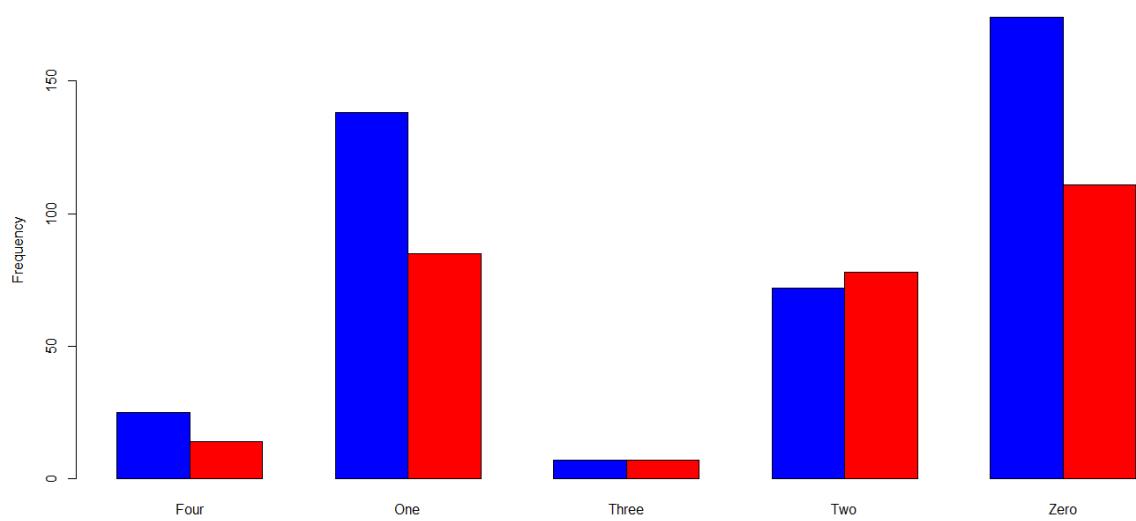


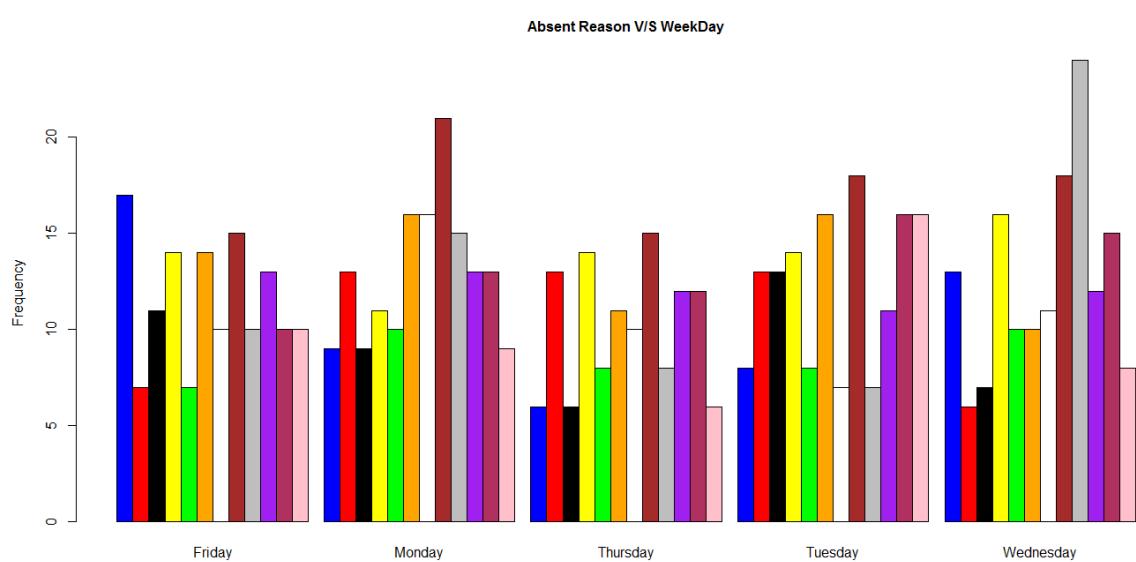
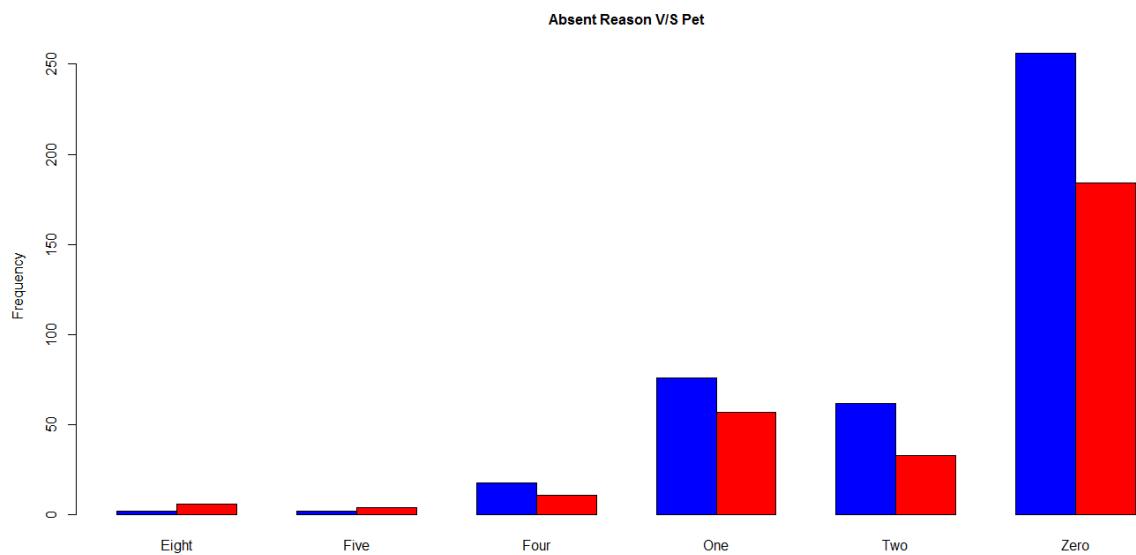


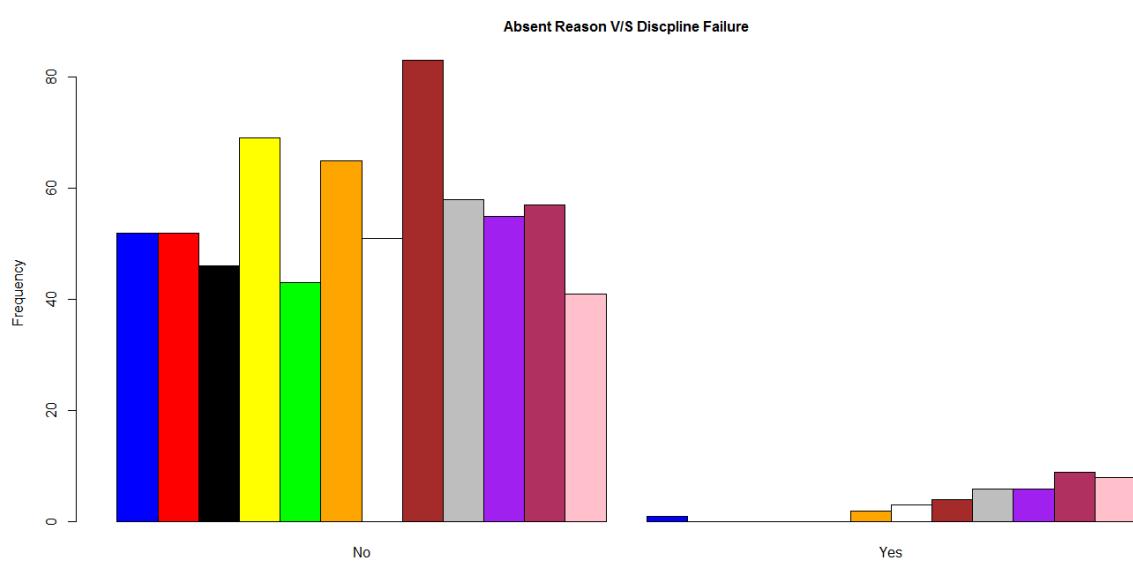
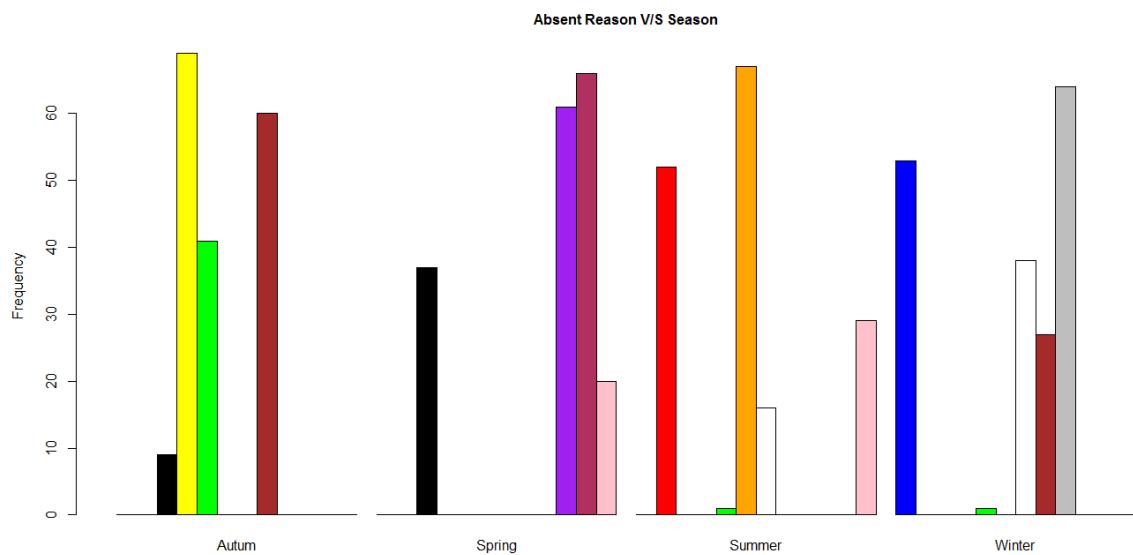
Absent Reason V/S Social Smoker

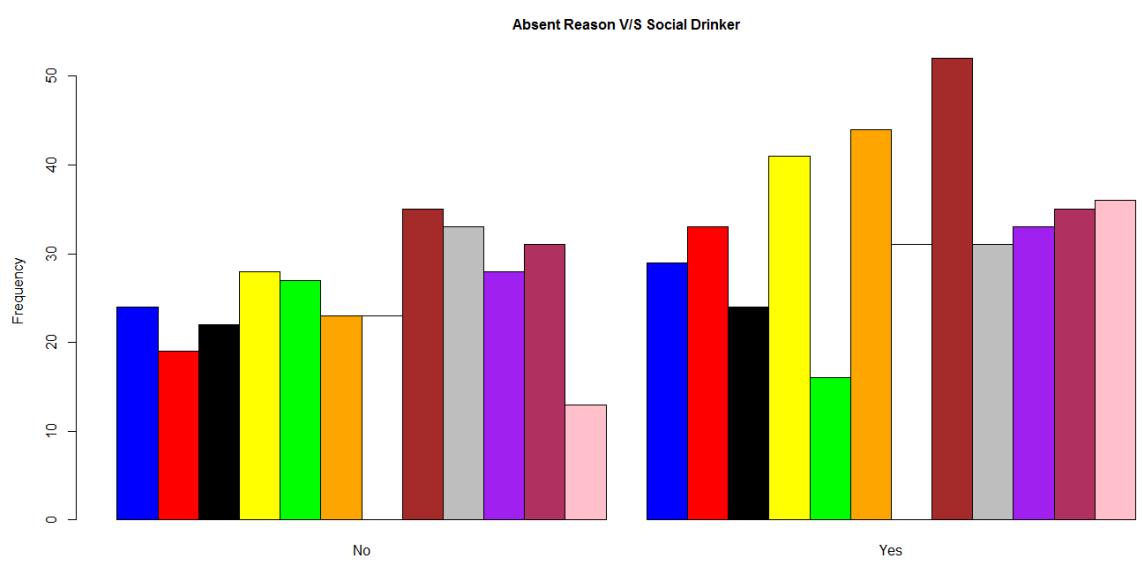
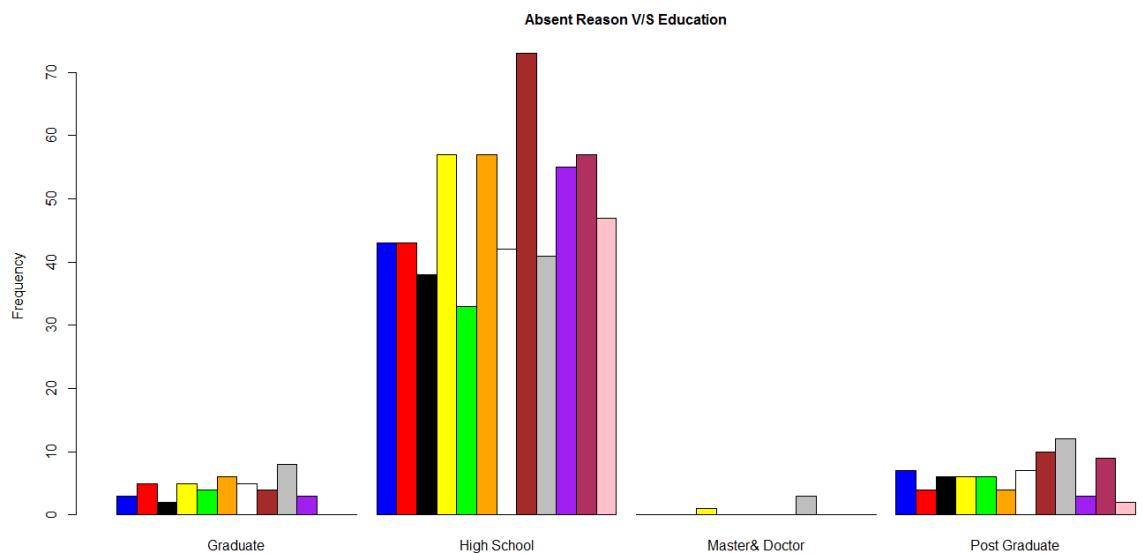


Absent Reason V/S Son

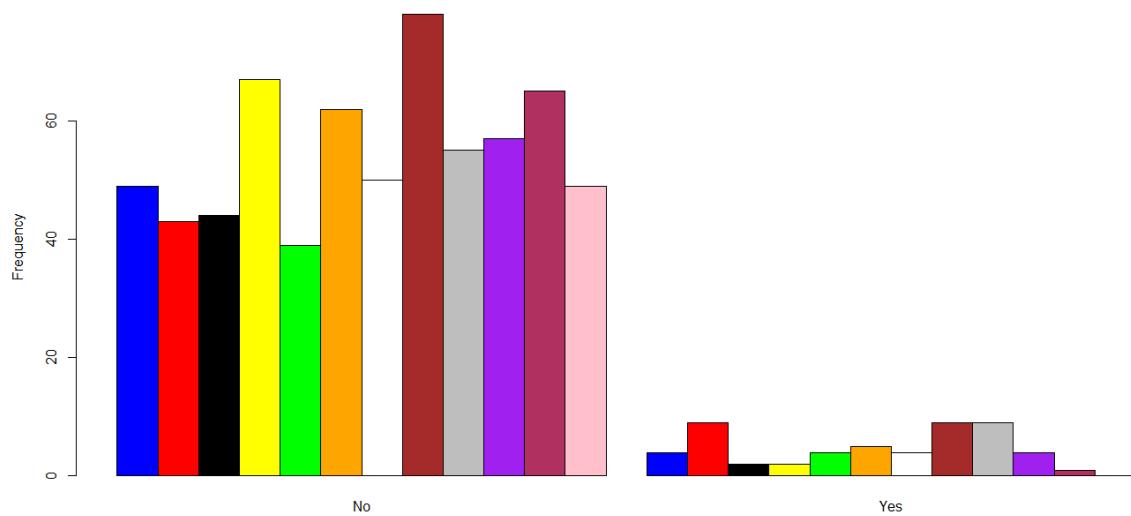




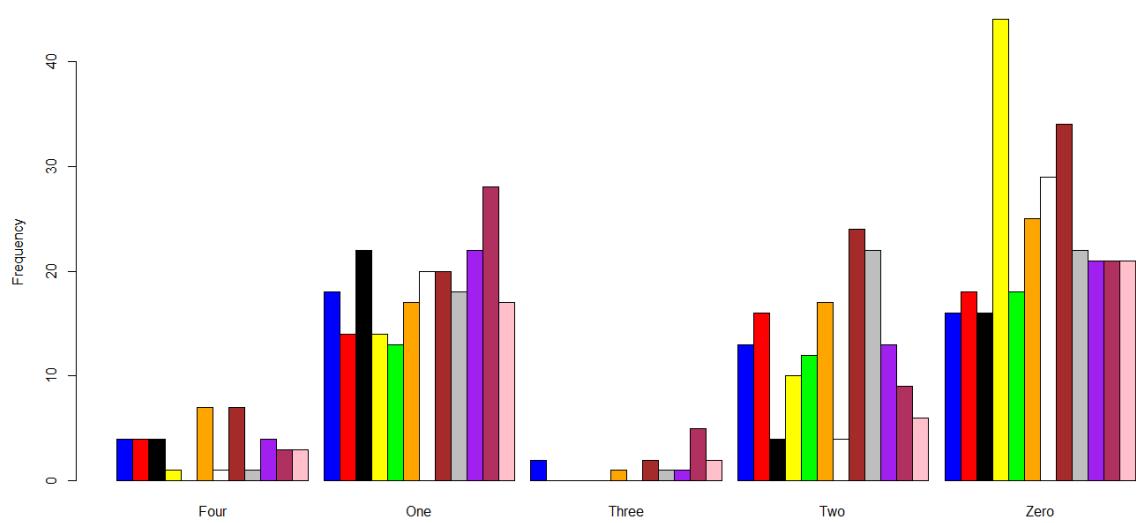




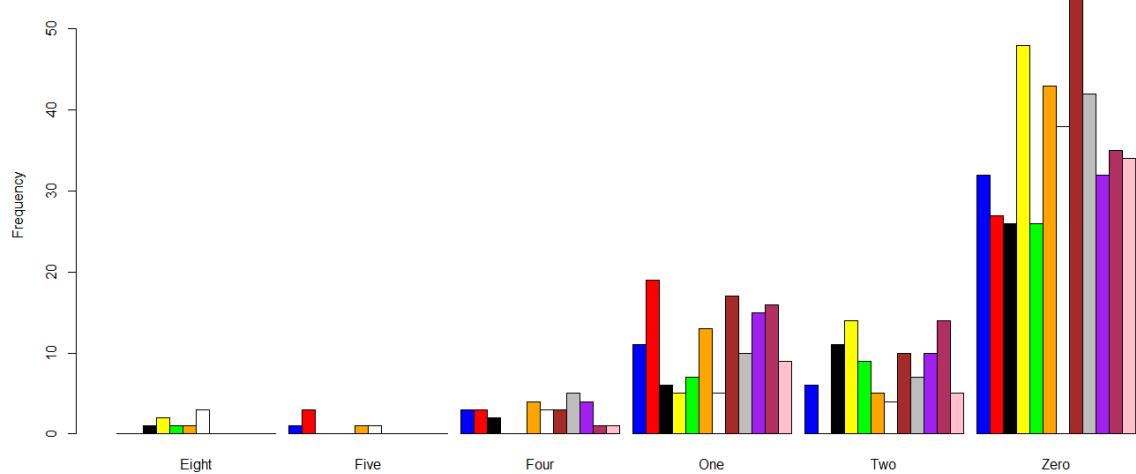
Absent Reason V/S Social Smoker



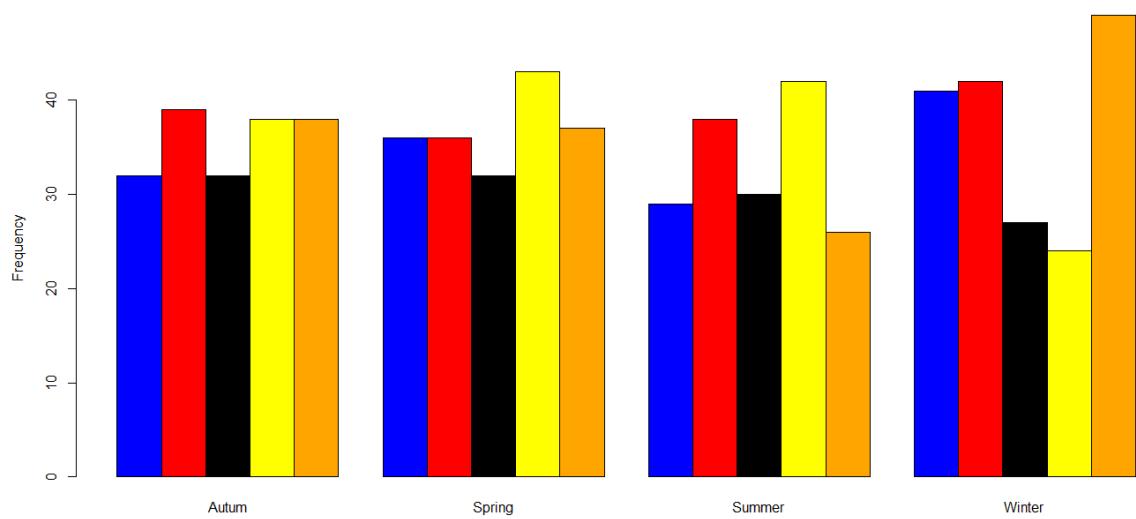
Absent Reason V/S Son



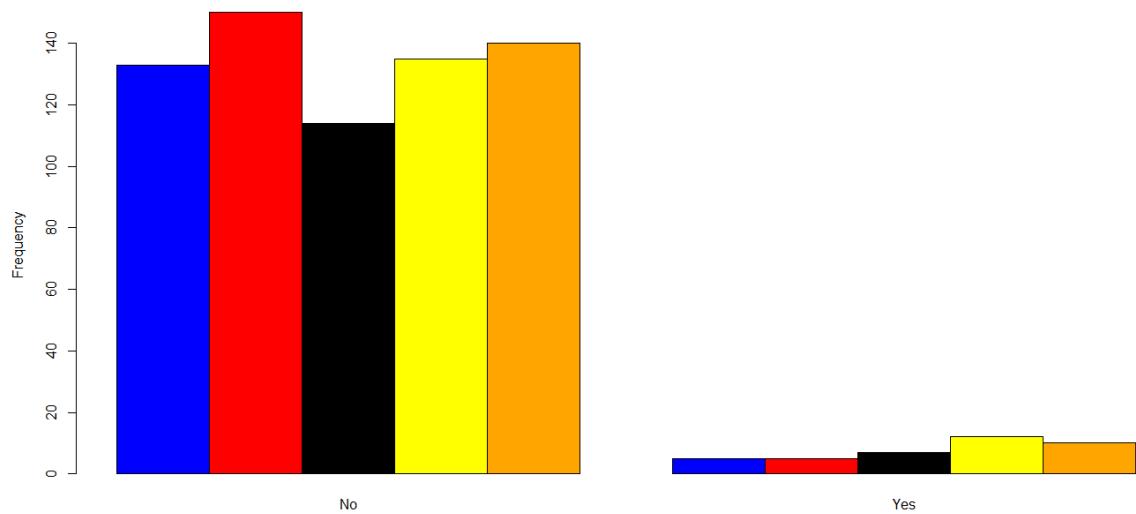
Absent Reason V/S Pet



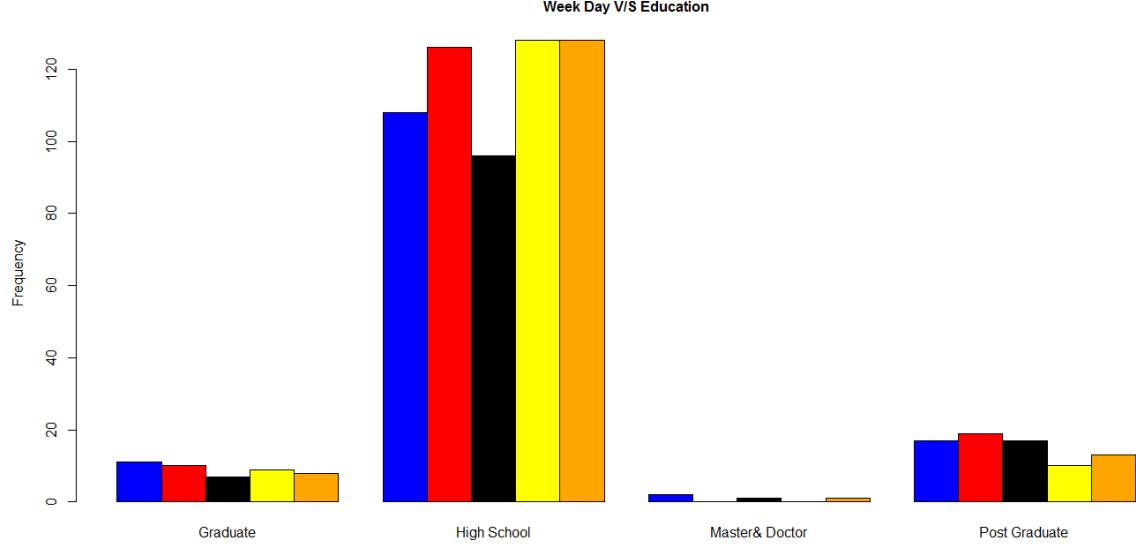
Week Day V/S Season

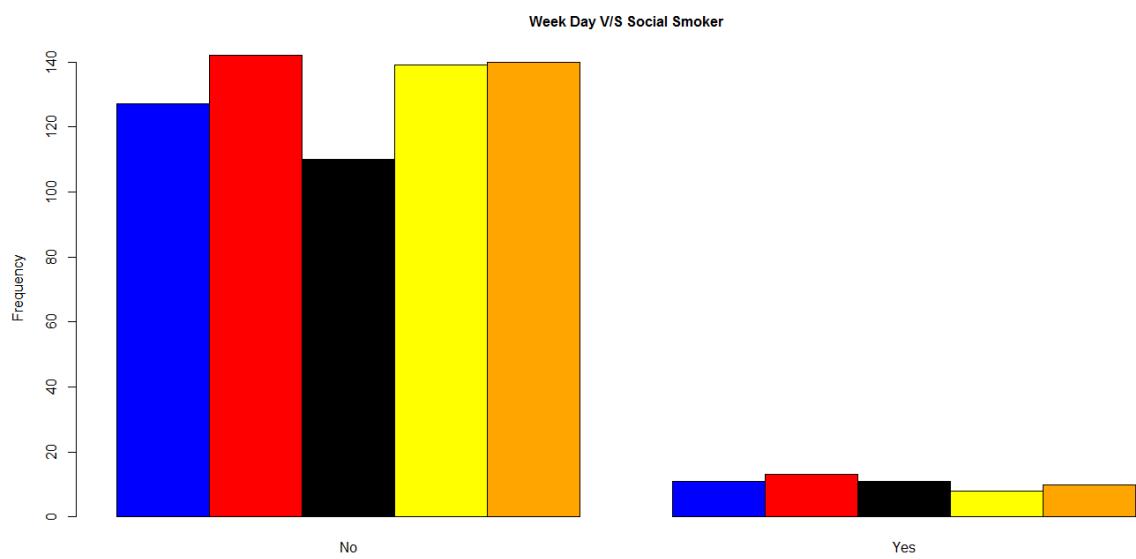
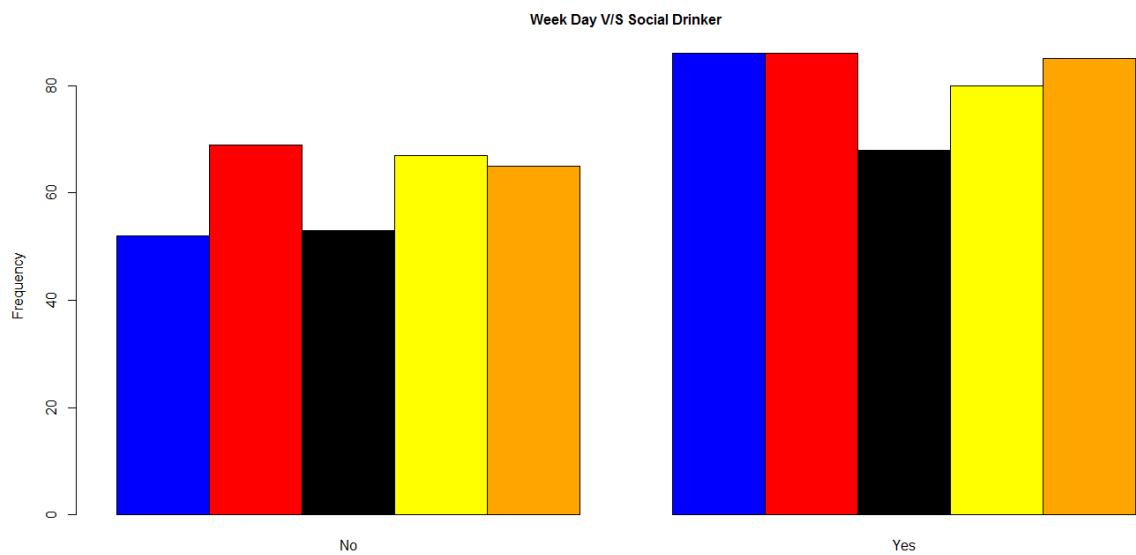


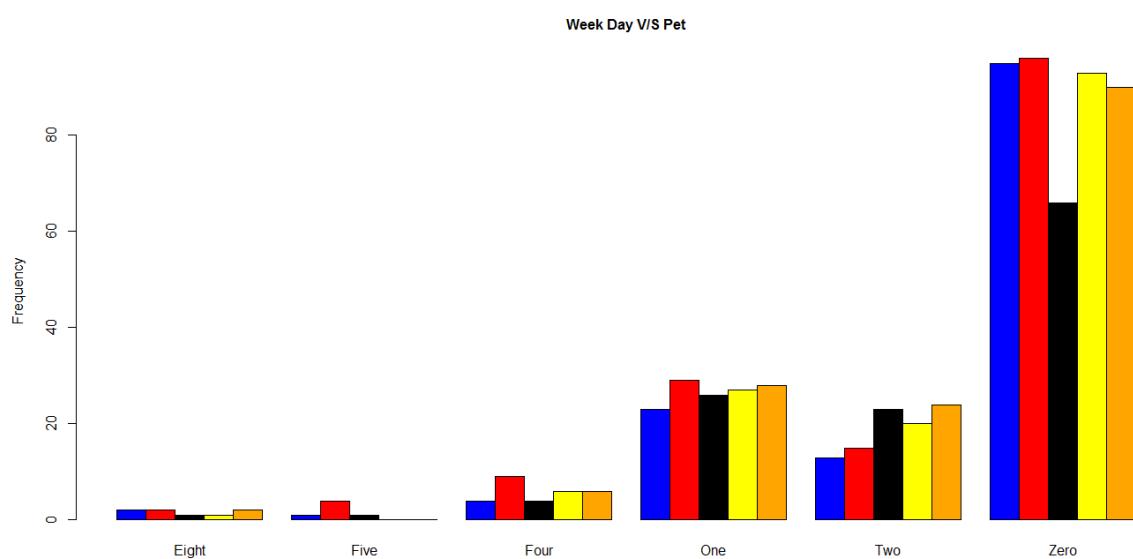
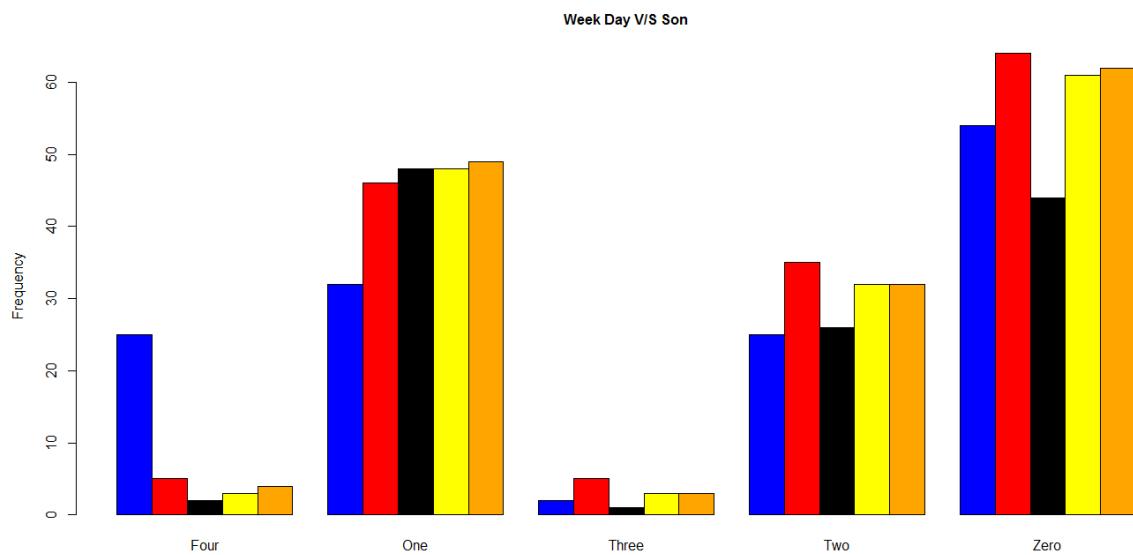
Week Day V/S Discipline Failure



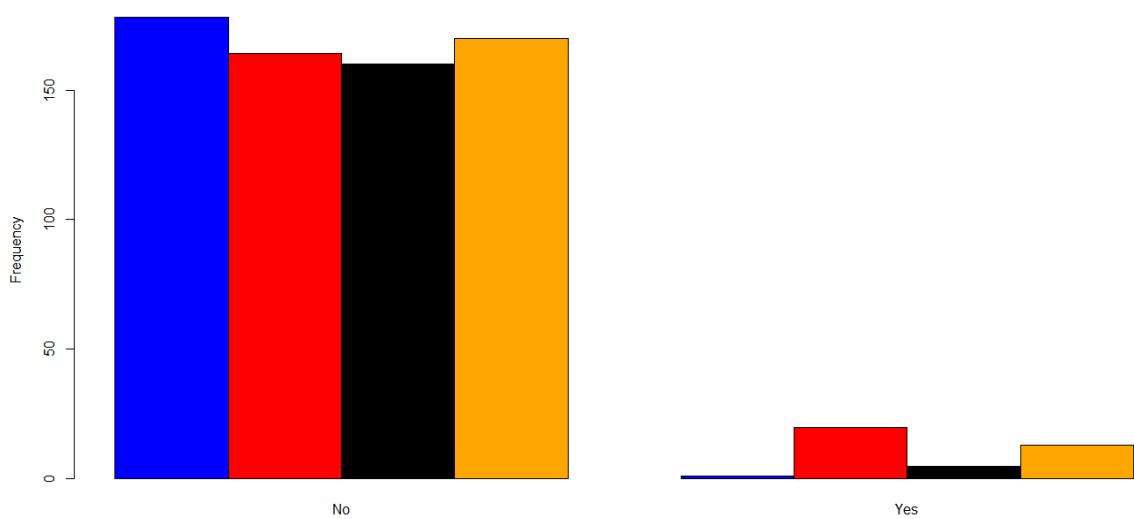
Week Day V/S Education



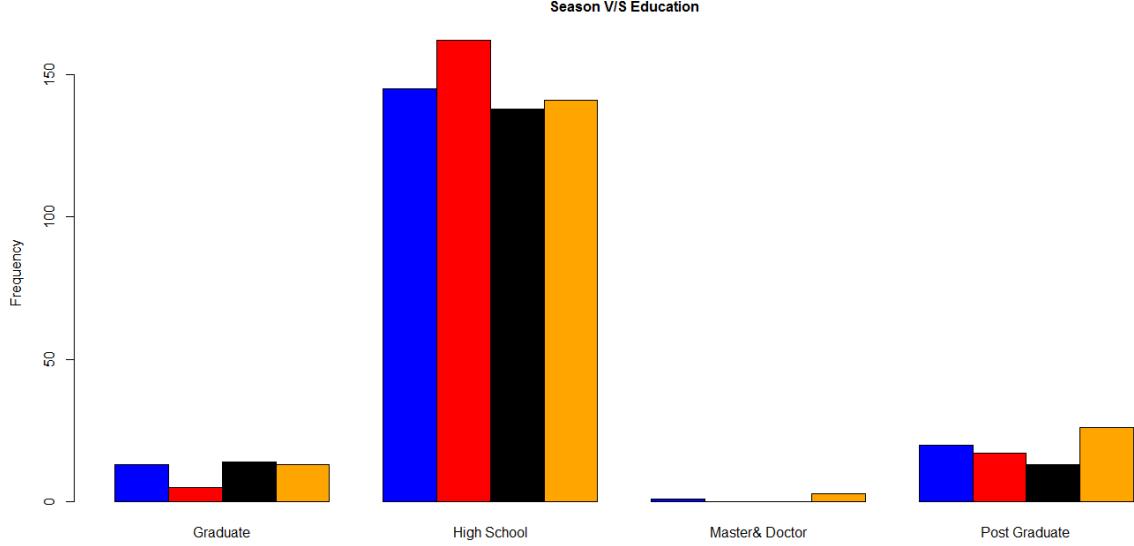




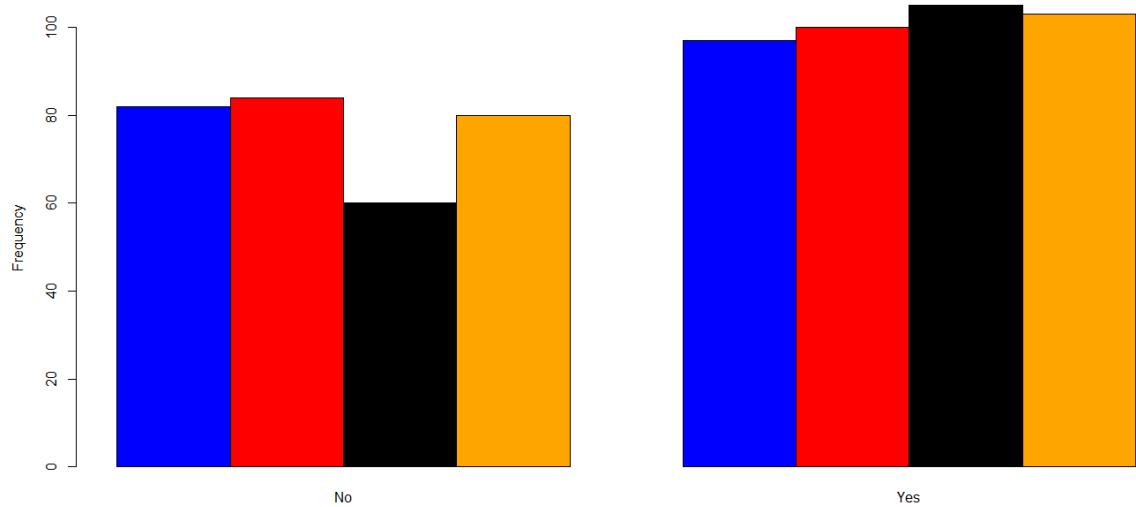
**Season V/S Discipline Failure**



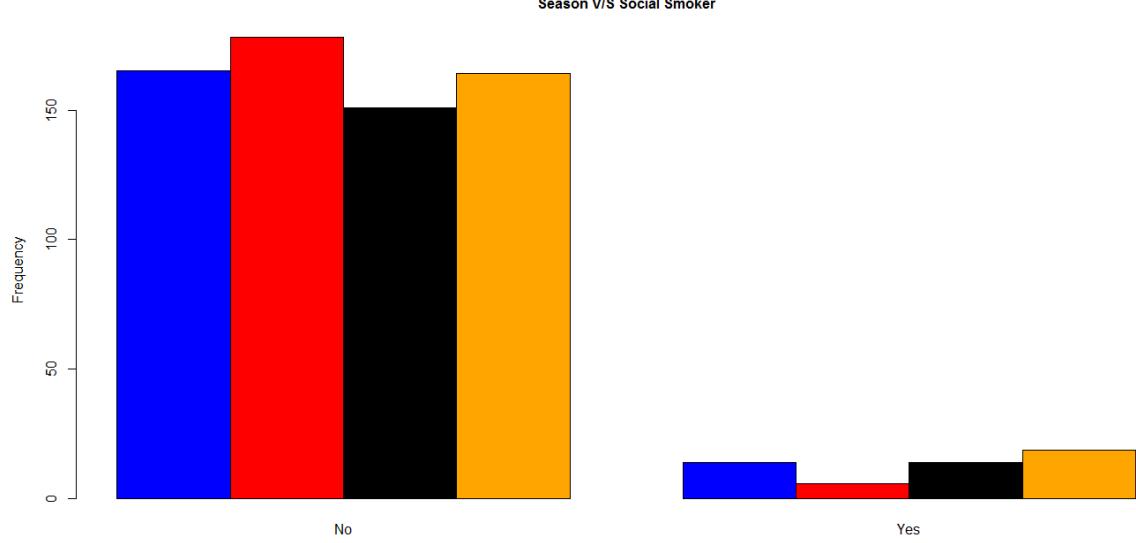
**Season V/S Education**

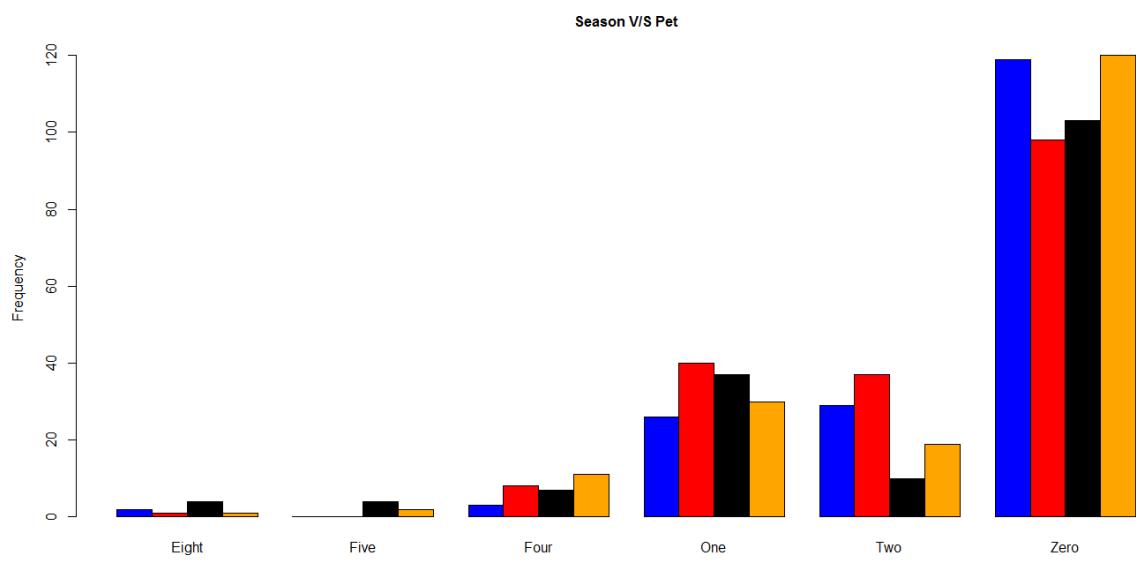
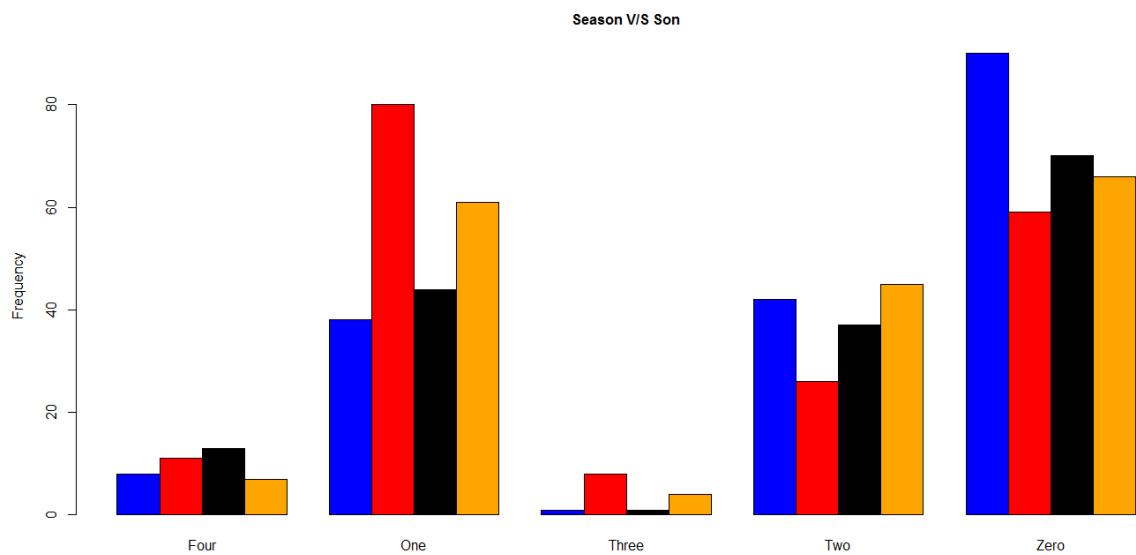


Season V/S Social Drinker

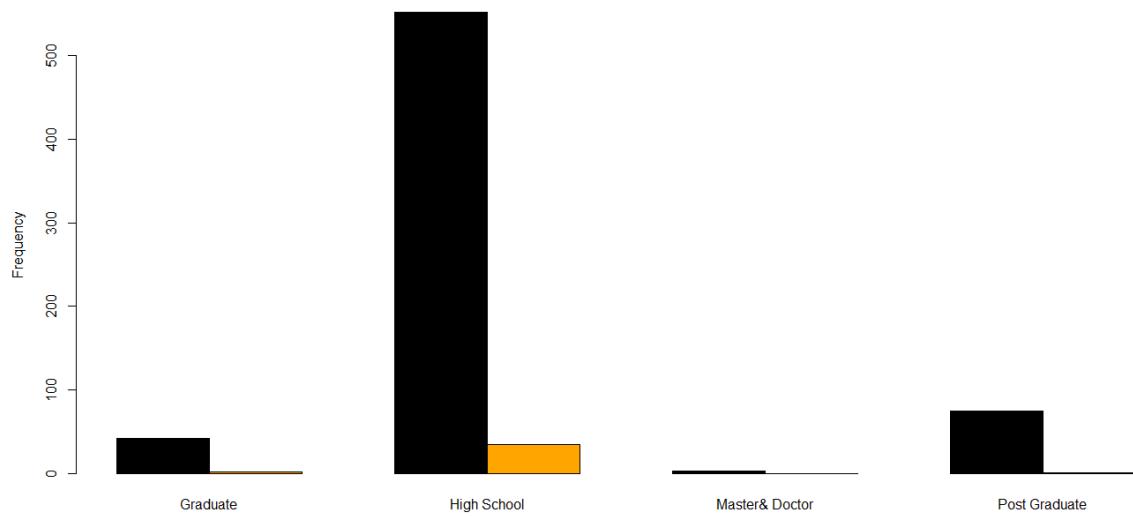


Season V/S Social Smoker

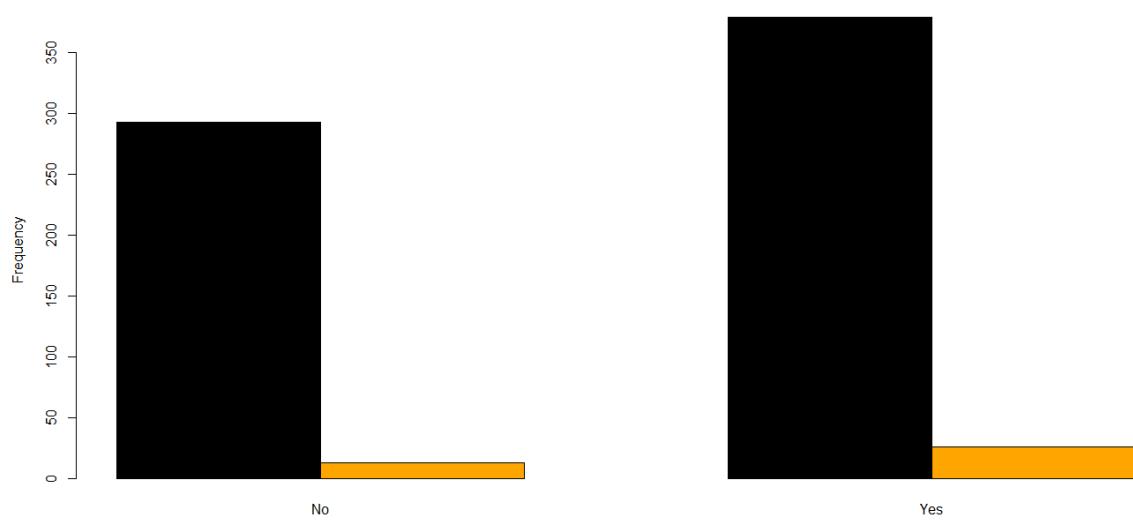




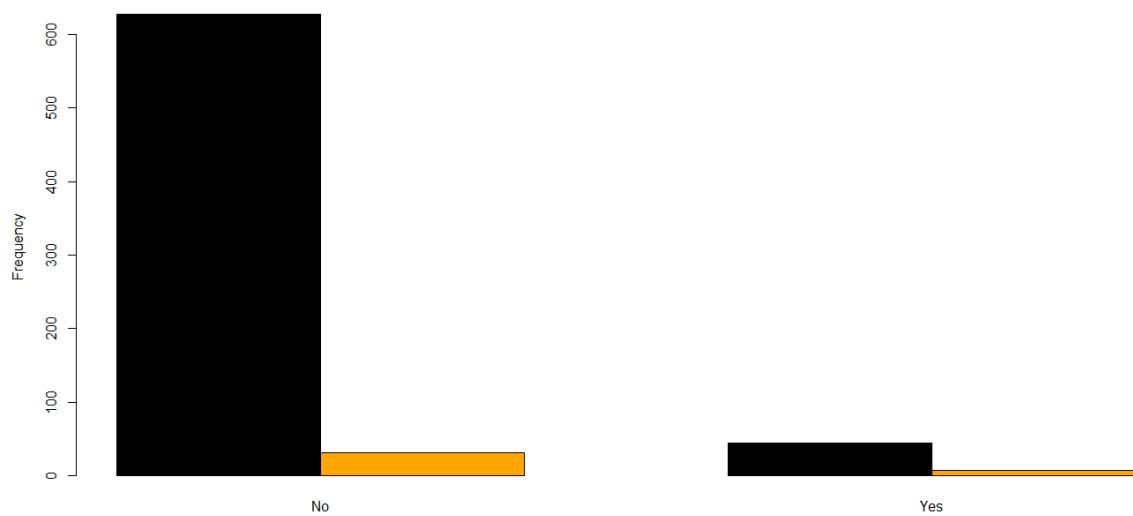
Discipline Failure V/S Education



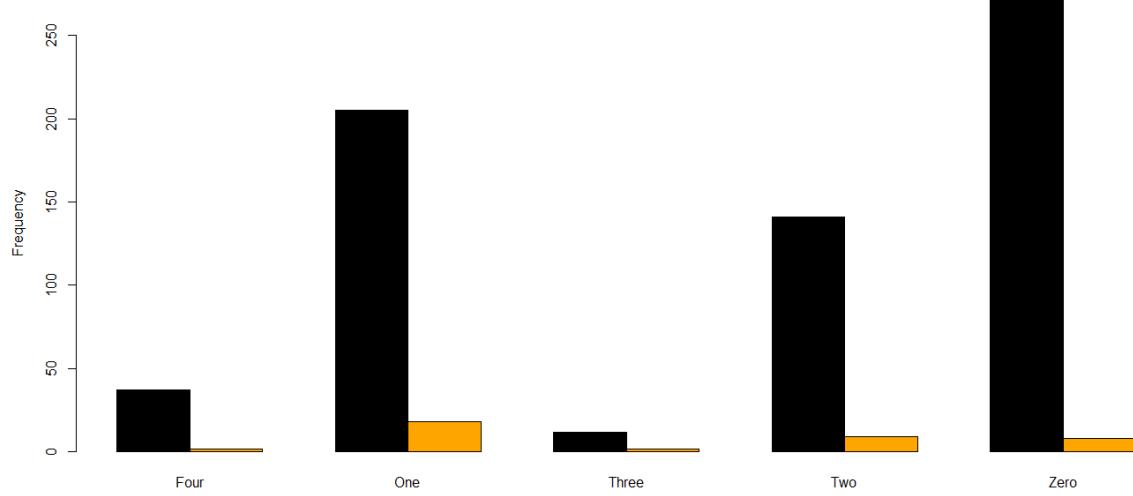
Discipline Failure V/S Social Drinker



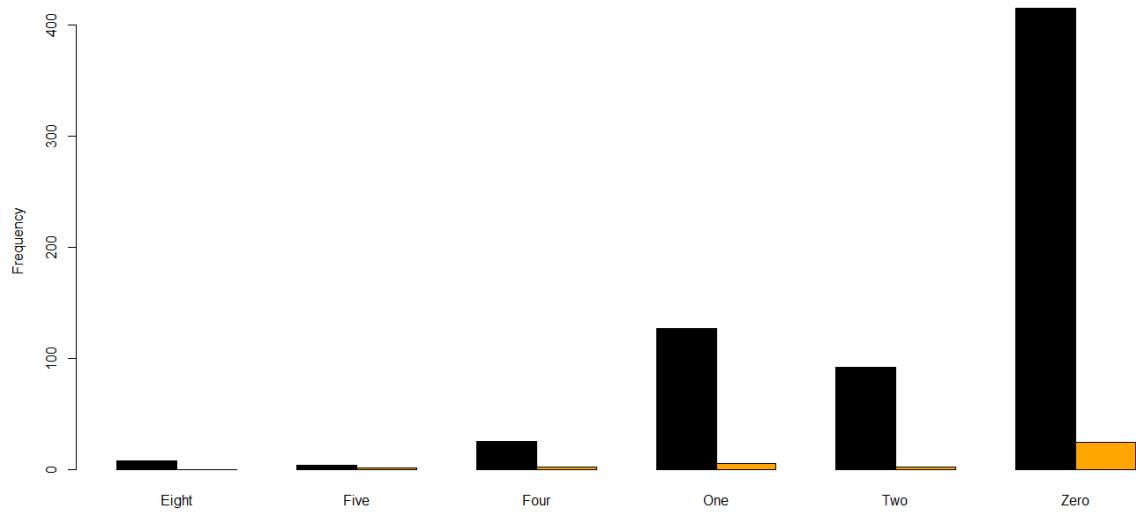
Discipline Failure V/S Social Smoker



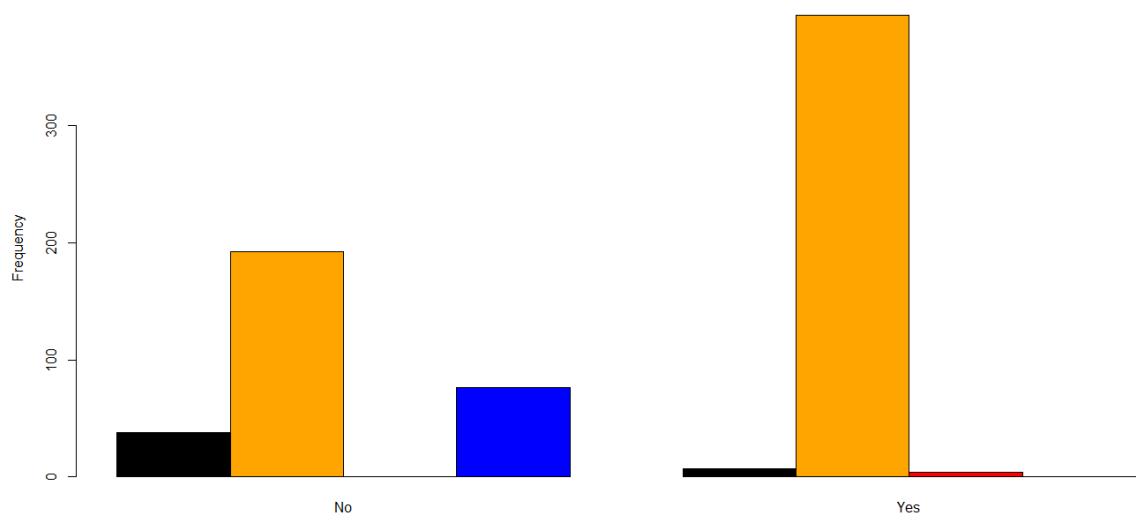
Discipline Failure V/S Son



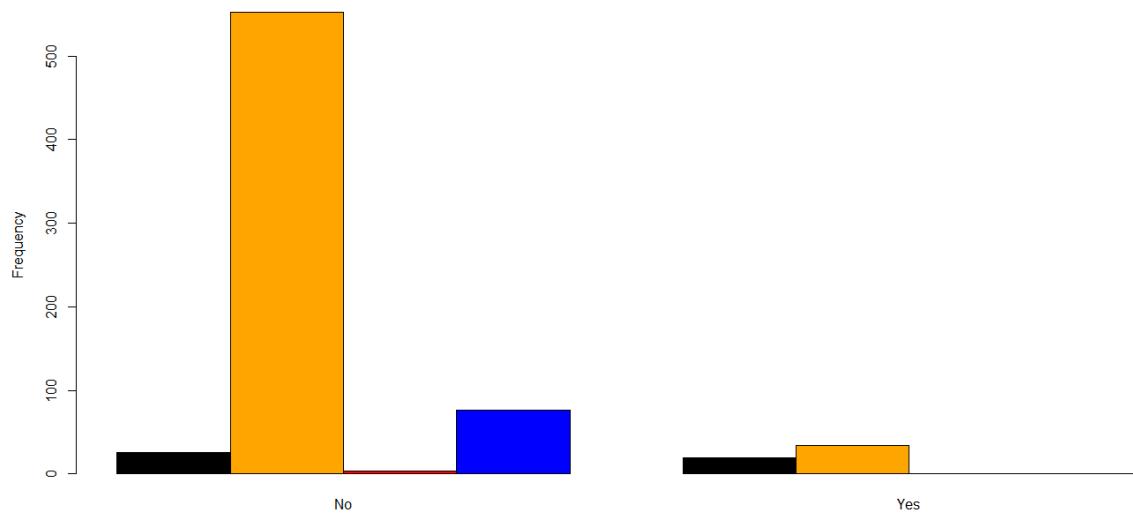
Discipline Failure V/S Pet



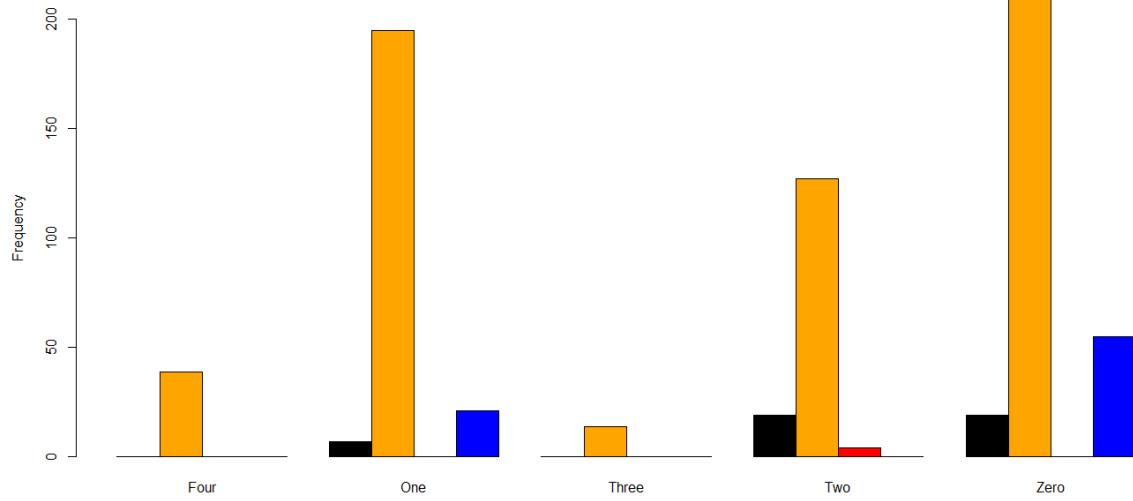
Education V/S Social Drinker



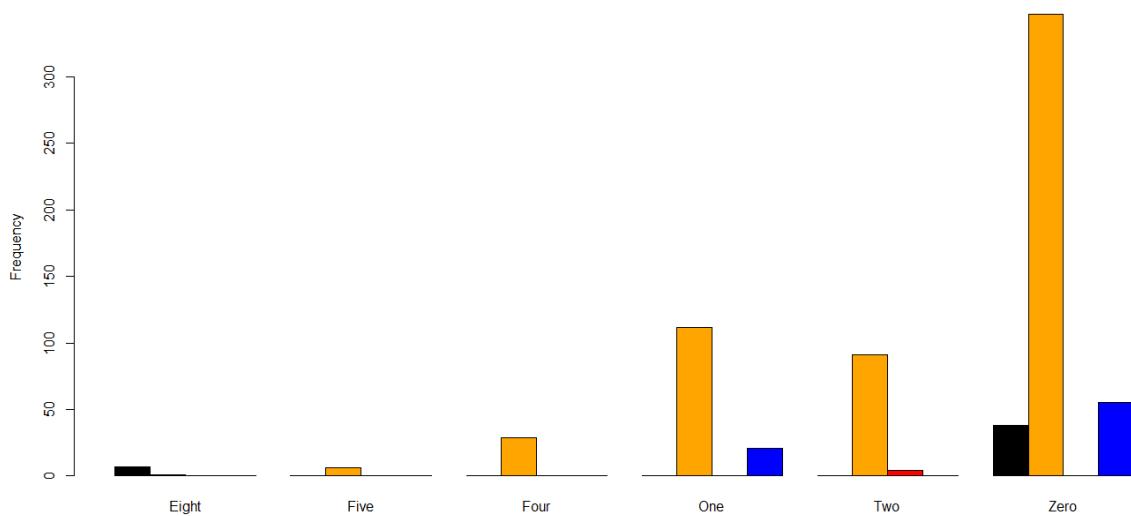
Education V/S Social Smoker



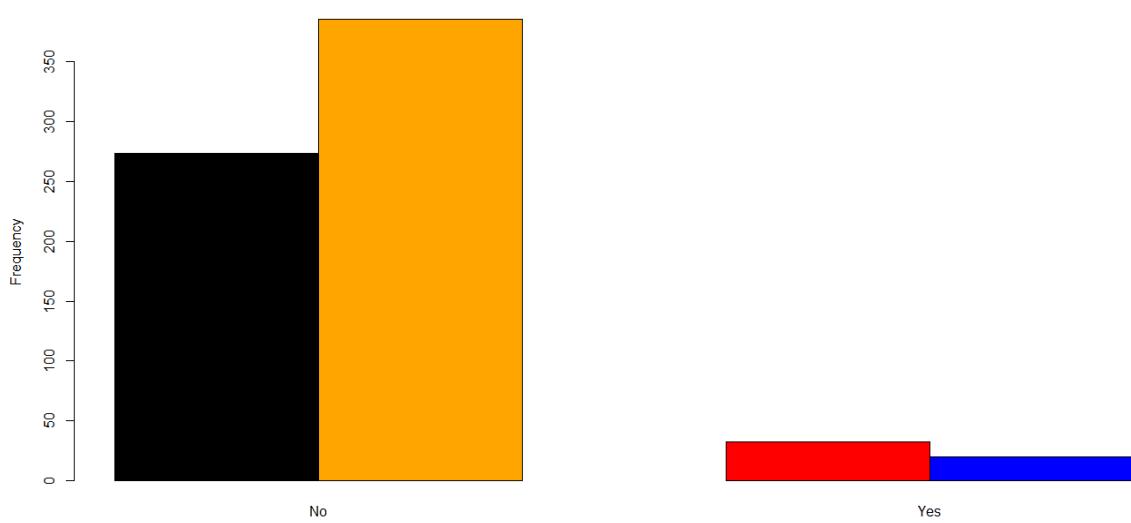
Education V/S Son



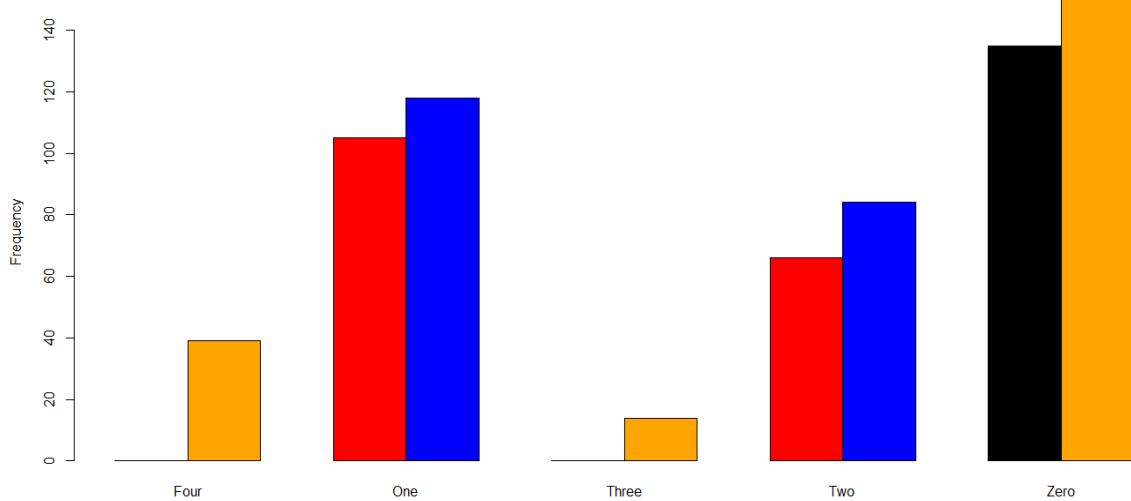
**Education V/S Pet**



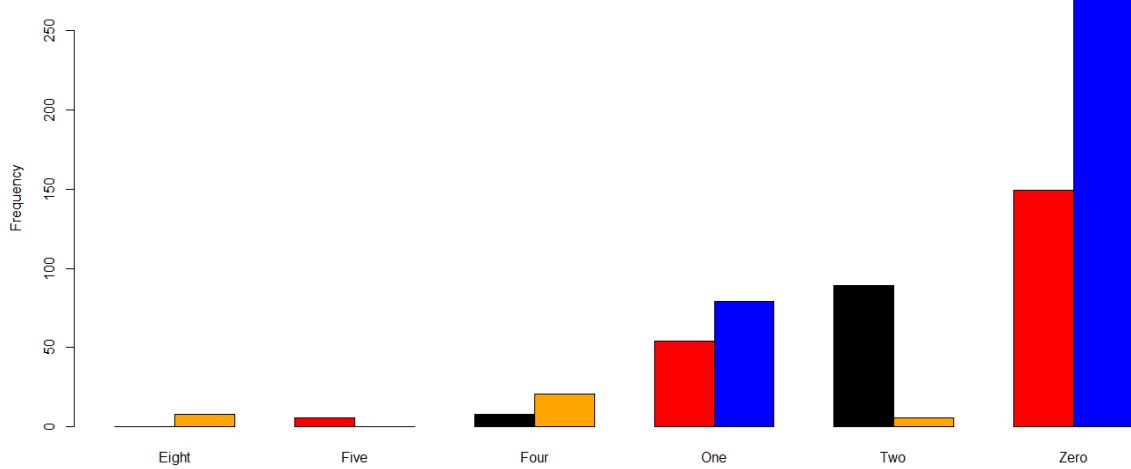
**Social Drinker V/S Socail Smoker**



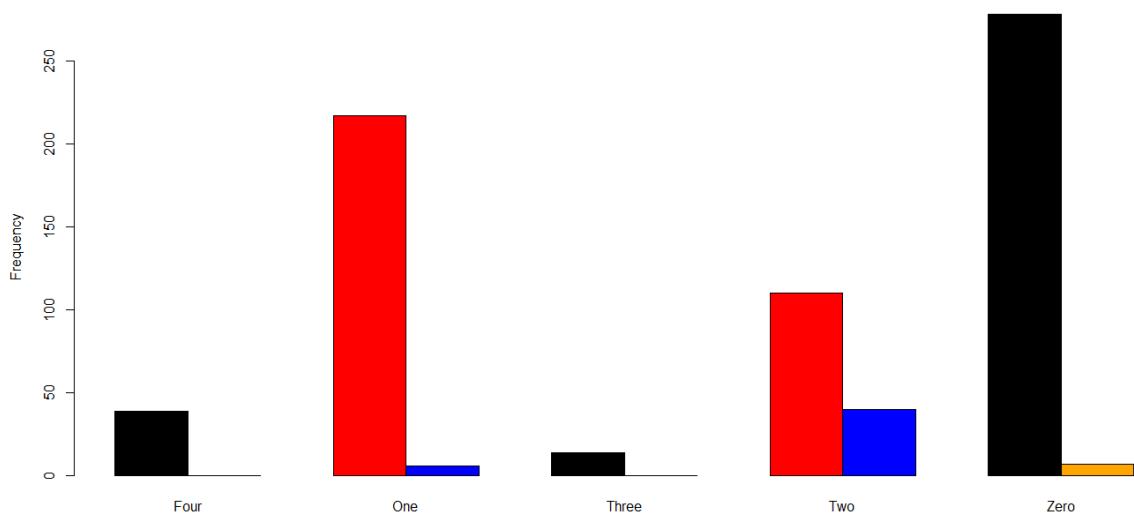
Social Drinker V/S Son



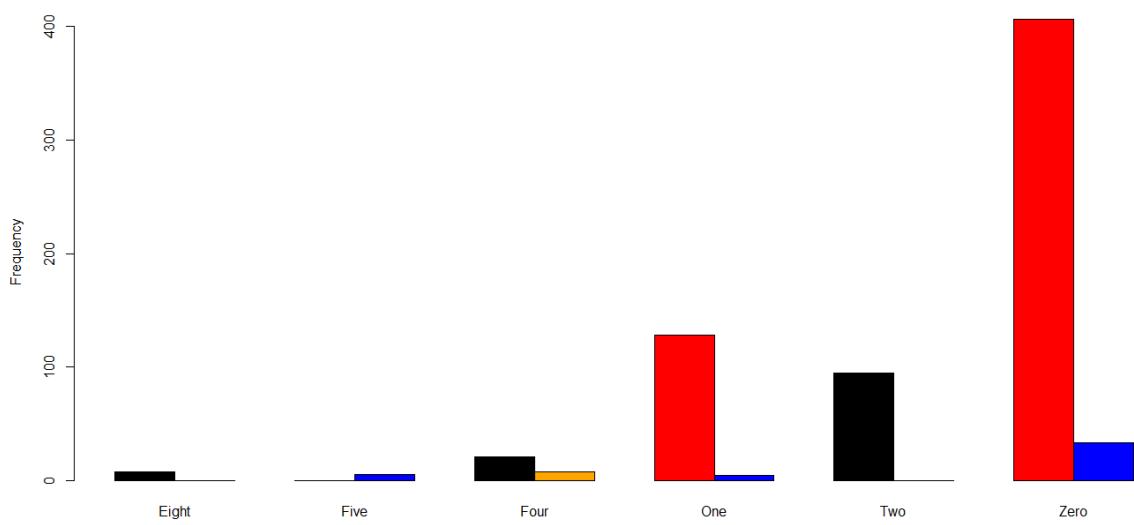
Social Drinker V/S Pet

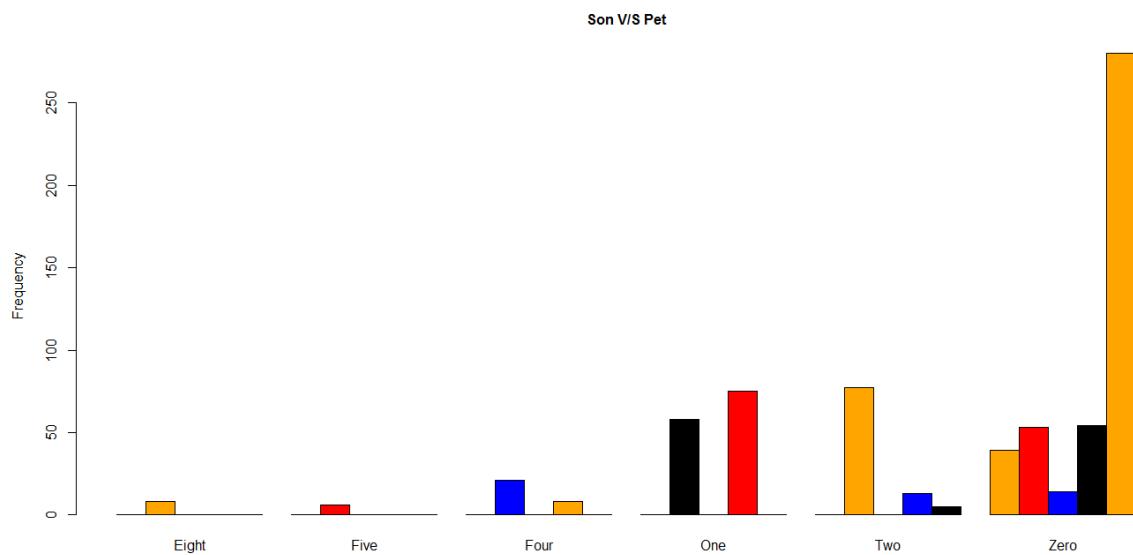


Social Smoker V/S Son



Social Smoker V/S Pet





## Missing Values

In statistics, missing data, or missing values, occur when no data value is stored for the variable in an observation. Missing data are a common occurrence and can have a significant effect on the conclusions that can be drawn from the data.

Missing data can occur because of nonresponse: no information is provided for one or more items or for a whole unit ("subject"). Some items are more likely to generate a nonresponse than others

In our dataset, there is no missing value present. A table below will throw light on our dataset:-

```
In [69]: #GETting sum of all missing values
AbsentData.isnull().sum()

Out[69]: Id          0
AbsentReason      3
AbsentMonth       1
WeekDay           0
Season             0
Expenses           7
ResidentDistance   3
ServiceTime        3
Age                3
AverageWorkLoad    10
HitTarget          6
DisciplineFailure  6
Education          10
Son                6
SocialDrinker      3
SocialSmoker        4
Pet                 2
Weight              1
Height              14
BodyMassIndex       31
AbsentTime          22
dtype: int64
```

	Variables	MissingPercentage
0	BodyMassIndex	4.189189
1	AbsentTime	2.972973
2	Height	1.891892
3	AverageWorkLoad	1.351351
4	Education	1.351351
5	Expenses	0.945946
6	Son	0.810811
7	DisciplineFailure	0.810811
8	HitTarget	0.810811
9	SocialSmoker	0.540541
10	Age	0.405405
11	AbsentReason	0.405405
12	ServiceTime	0.405405
13	ResidentDistance	0.405405
14	SocialDrinker	0.405405
15	Pet	0.270270
16	Weight	0.135135
17	AbsentMonth	0.135135
18	Season	0.000000
19	WeekDay	0.000000
20	Id	0.000000

All the missing values must have to me imputed before training the models. There are different methods available to impute the missing values and it depends on the data scientist to choose the appropriate method. After imputing the missing values all NA in the dataset will be removed. All count of missing values must be zero after imputing as shown below:-

```
In [72]: #Checking Missing Values exists after filling or not.
AbsentData.isnull().any()
#ALL Missing Values are filled.
```

```
Out[72]: Id      False
AbsentReason  False
AbsentMonth   False
WeekDay       False
Season        False
Expenses      False
ResidentDistance  False
ServiceTime   False
Age          False
AverageWorkLoad  False
HitTarget    False
DisciplineFailure  False
Education    False
Son          False
SocialDrinker False
SocialSmoker  False
Pet          False
Weight       False
Height       False
BodyMassIndex False
AbsentTime   False
dtype: bool
```

## Outliers

An Outlier is a rare chance of occurrence within a given data set. In Statistics and Data Science, an Outlier is an observation point that is distant from other observations. An Outlier may be due to variability in the measurement or it may indicate experimental error.

Outliers, being the most extreme observations, may include the sample maximum or sample minimum, or both, depending on whether they are extremely high or low. However, the sample maximum and minimum are not always outliers because they may not be unusually far from other observations.

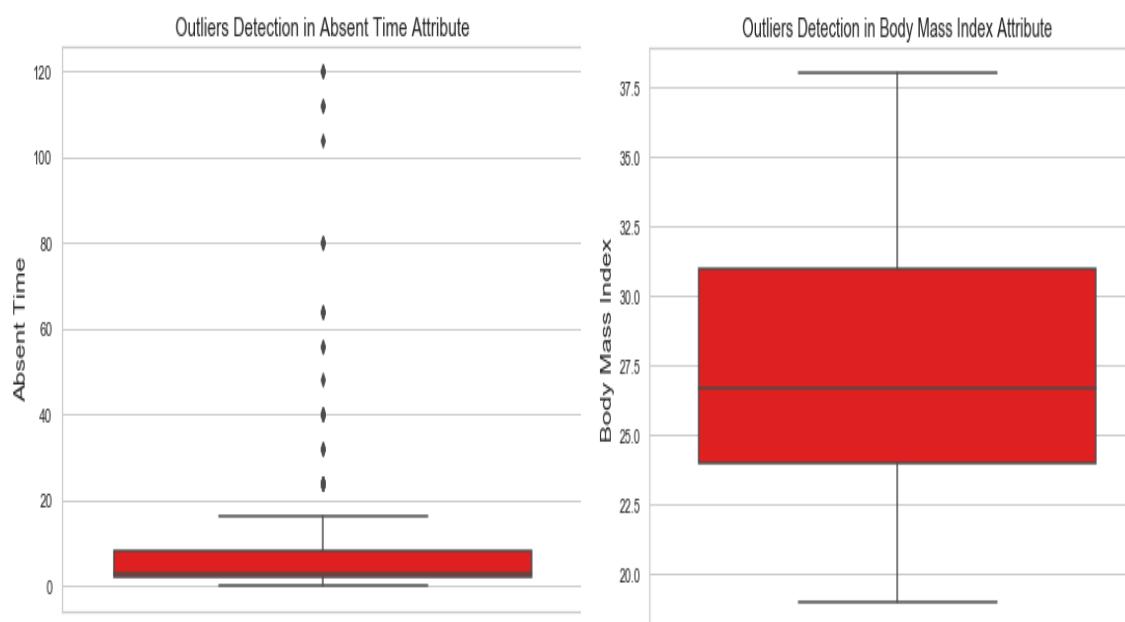
While outliers are attributed to a rare chance and may not necessarily be fully explainable, Outliers in data can distort predictions and affect the accuracy, if you don't detect and handle them.

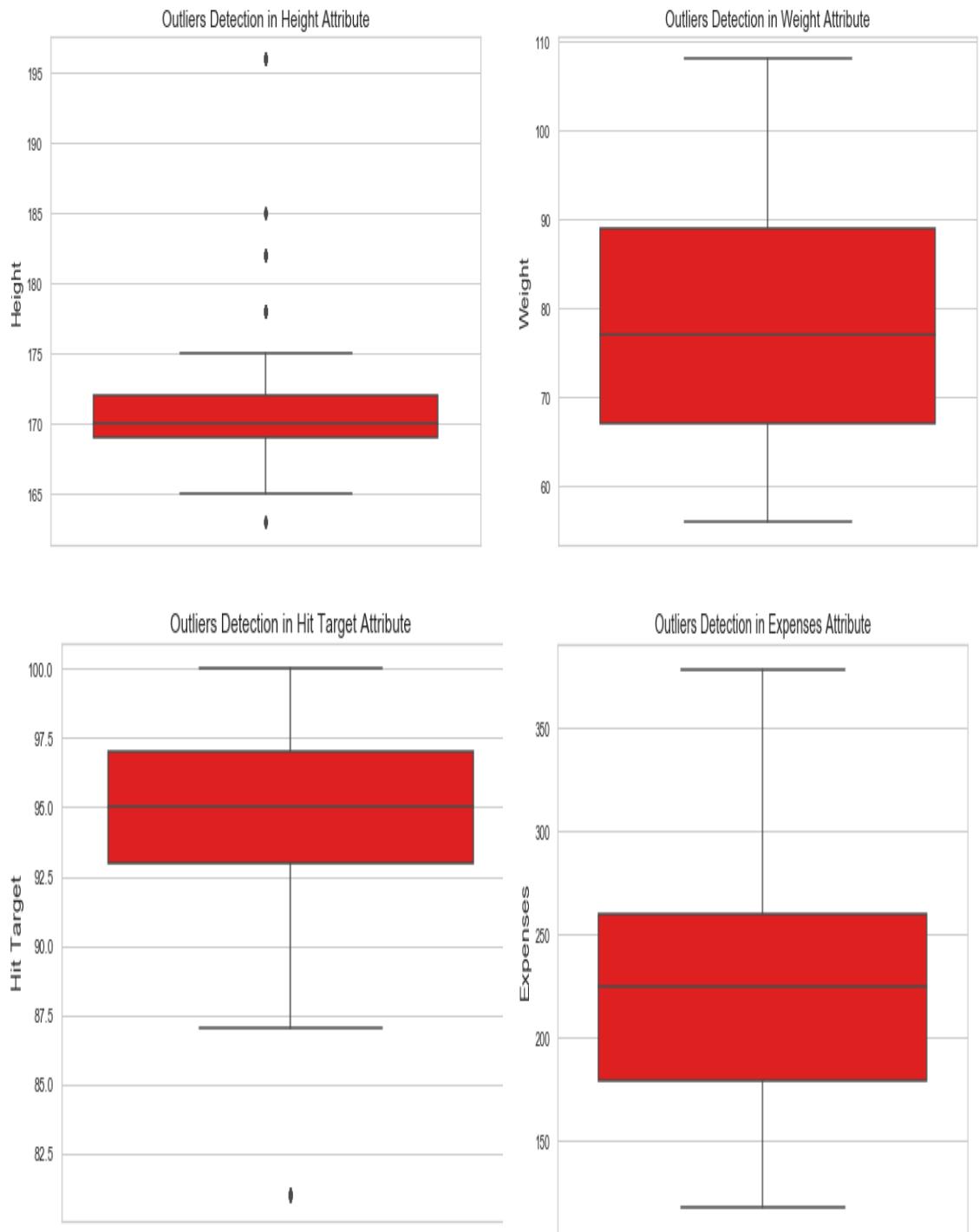
The contentious decision to consider or discard an outlier needs to be taken at the time of building the model. Outliers can drastically bias/change the fit estimates and predictions.

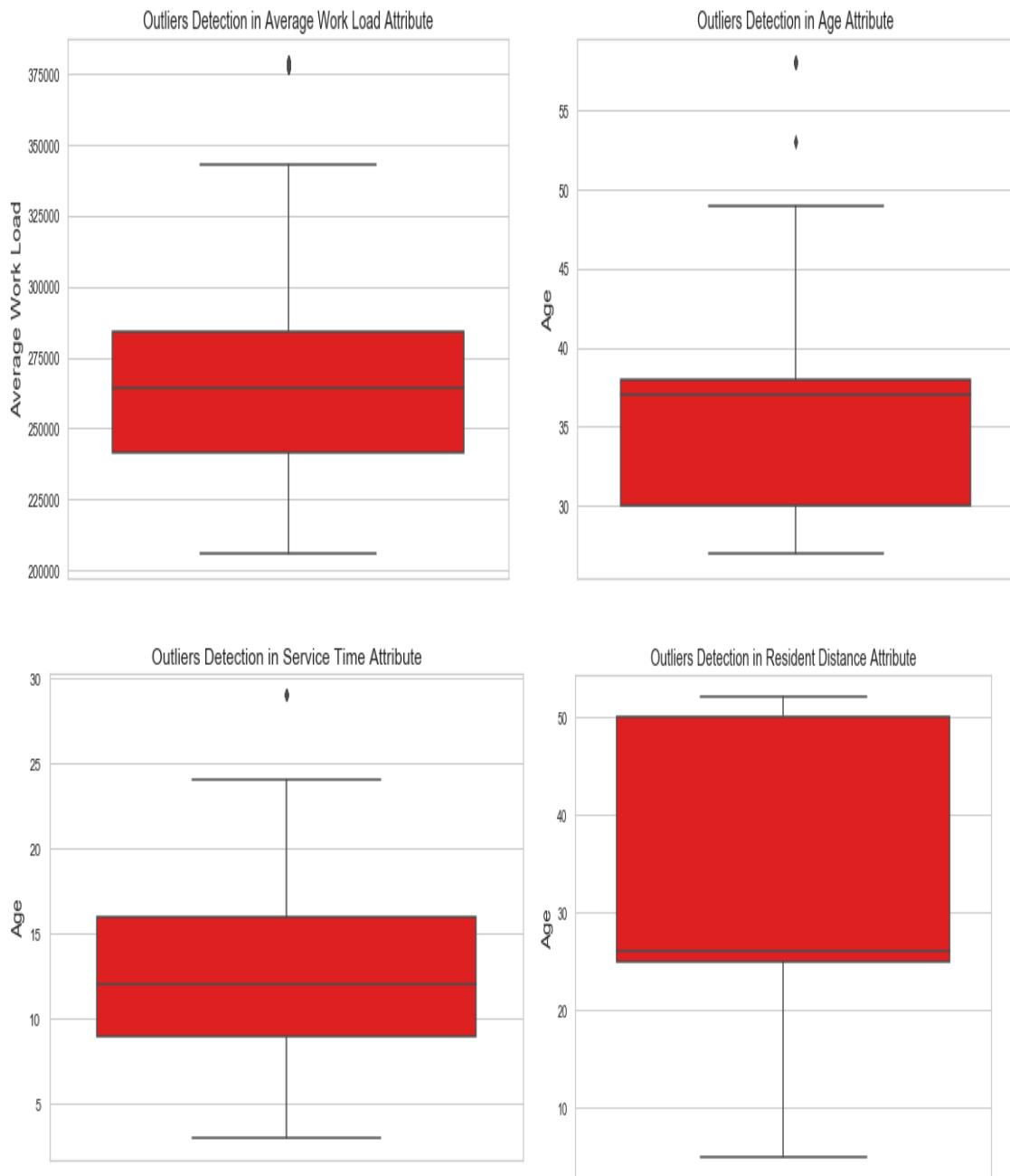
## Detecting and Removing Outliers

Mostly outliers are present in the continuous variables and box plot method is best and easy way to detect and remove outliers. Moreover, our dataset contains categorical variables that are already encoded so we will perform outlier detections only on continuous variables.

Box plot of all the variables are shown below:-







## Feature Selection

Machine learning works on a simple rule – if we put garbage in, we will only get garbage to come out.

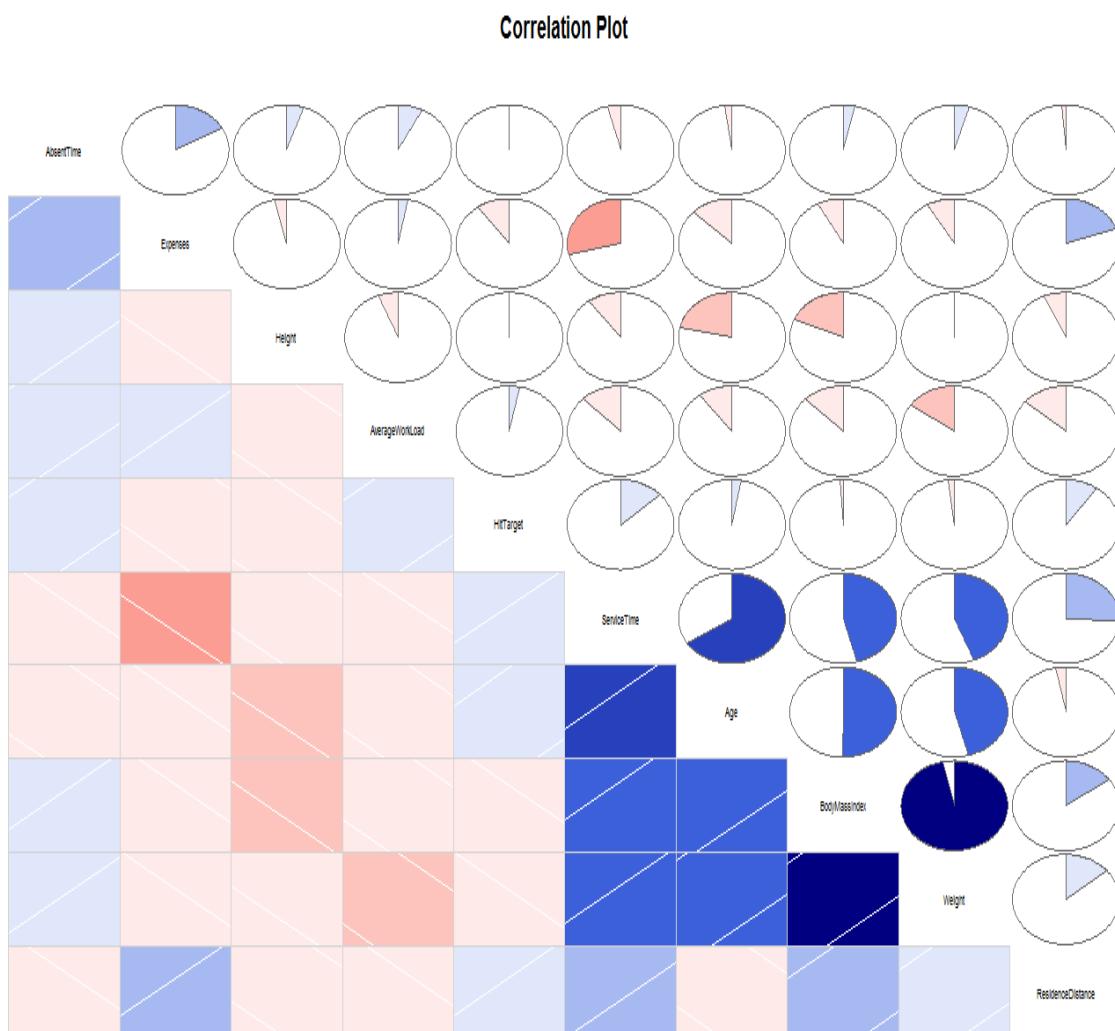
This becomes even more important when the number of features is very large. We need not use every feature at our disposal for creating an algorithm. We can assist our algorithm by feeding in only those features that are really important. Feature subsets giving better results than complete set of feature for the same algorithm or – “Sometimes, less is better!”.

We should consider the selection of feature for model based on below criteria:-

1. The relationship between two independent variable should be less and
2. The relationship between Independent and target variables should be high.

Below figure shows a graphical display of a correlation matrix, called a correlogram. The cells of the matrix are coloured to show the correlation value.

	AbsentTime	BodyMassIndex	Height	Age	Weight	HitTarget	Expenses	AverageWorkLoad	ServiceTime	ResidentDistance
AbsentTime	1.0	-0.056	0.083	0.076	-0.013	0.035	0.042	0.028	0.012	-0.11
BodyMassIndex	-0.056	1.0	-0.13	0.46	0.89	-0.08	-0.13	-0.068	0.51	0.12
Height	0.083	-0.13	1.0	-0.068	0.29	0.091	-0.19	0.11	-0.084	-0.36
Age	0.076	0.46	-0.068	1.0	0.42	-0.035	-0.22	-0.037	0.67	-0.14
Weight	-0.013	0.89	0.29	0.42	1.0	-0.051	-0.21	-0.024	0.46	-0.038
HitTarget	0.035	-0.08	0.091	-0.035	-0.051	1.0	-0.067	-0.089	-7.9e-05	-0.02
Expenses	0.042	-0.13	-0.19	-0.22	-0.21	-0.087	1.0	0.015	-0.38	0.24
AverageWorkLoad	0.028	-0.068	0.11	-0.037	-0.024	-0.089	0.015	1.0	0.0051	-0.064
ServiceTime	0.012	0.51	-0.064	0.67	0.46	-7.9e-05	-0.38	0.0051	1.0	0.15
ResidentDistance	-0.11	0.12	-0.38	-0.14	-0.038	-0.02	0.24	-0.064	0.15	1.0



## Feature Scaling

Feature scaling is done to reduce unwanted variation either within or between variables and to bring all of the variables into proportion with one another. I will use Normalization process to perform feature scaling. Formula for Normalization is given below:-

$$x_{new} = \frac{x - x_{min}}{x_{max} - x_{min}}$$

Attribute after normalization are given below:-

Expenses	ResidentDistance	ServiceTime	Age	AverageWorkLoad	...
0.657692	0.659574	0.476190	0.30	0.312568	...
0.234615	0.978723	0.714286	0.55	0.312568	...
0.619231	0.000000	0.523810	0.60	0.312568	...
0.657692	0.659574	0.476190	0.30	0.312568	...
0.234615	0.978723	0.714286	0.55	0.312568	...

## Modeling

### Model Selection

In the case of this dataset we have to predict the count of total bike rented on basis of environmental and seasonal condition. The target variable here is a continuous variable and for a continuous variable we can use various Regression models. Trained model having less error rate and more accuracy will be our final model. Different machine learning methods which will be used to train our final model are mentioned below:-

1. Decision Tree Regression Model
2. Random Forest Model
3. Liner Regression Model

Final model will be with the higher accuracy which we will able to decide at the end of the modelling process.

### Decision Tree Regression Model

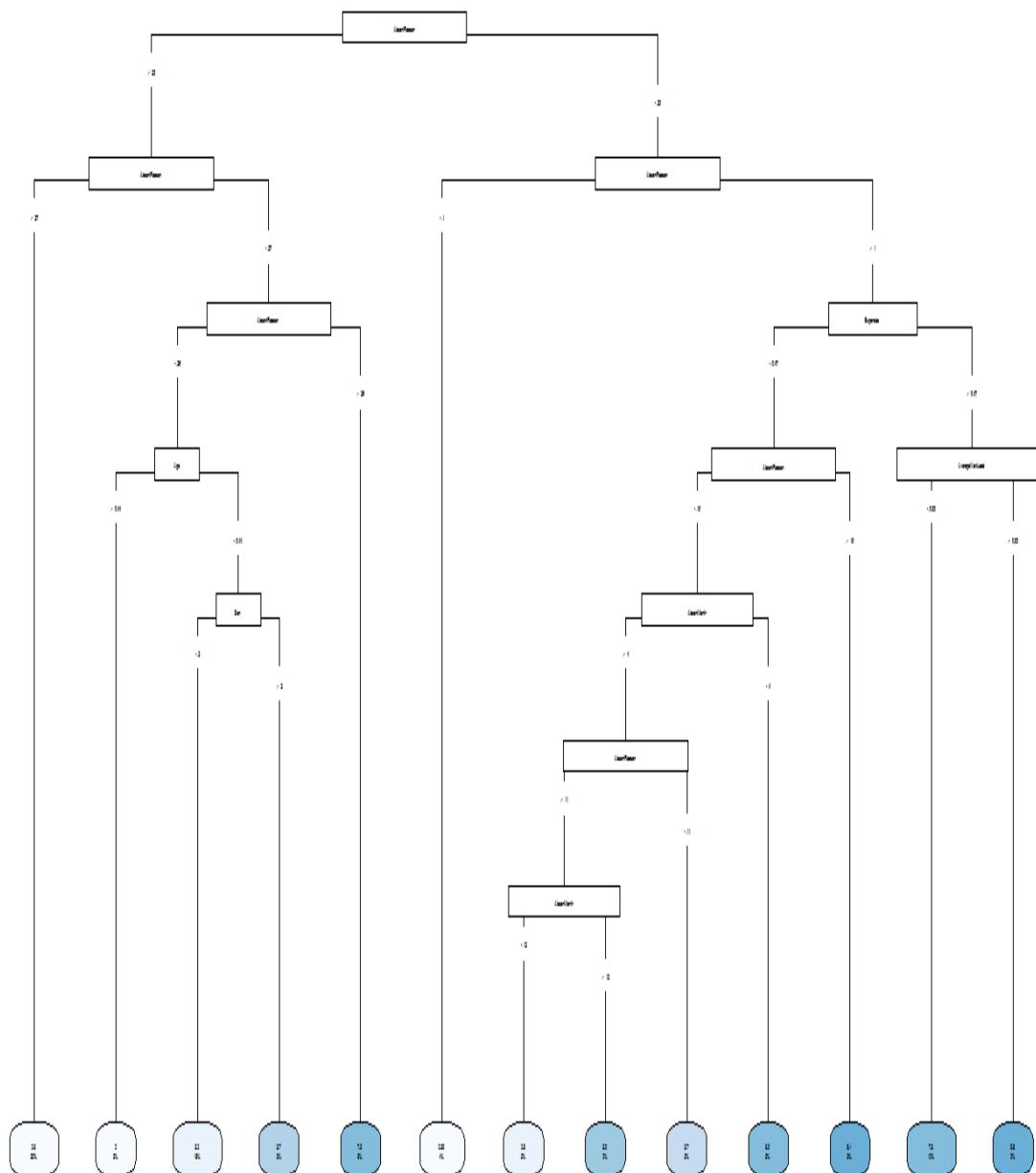
Decision tree builds regression models in the form of a tree structure. It breaks down a dataset into smaller and smaller subsets while at the same time an

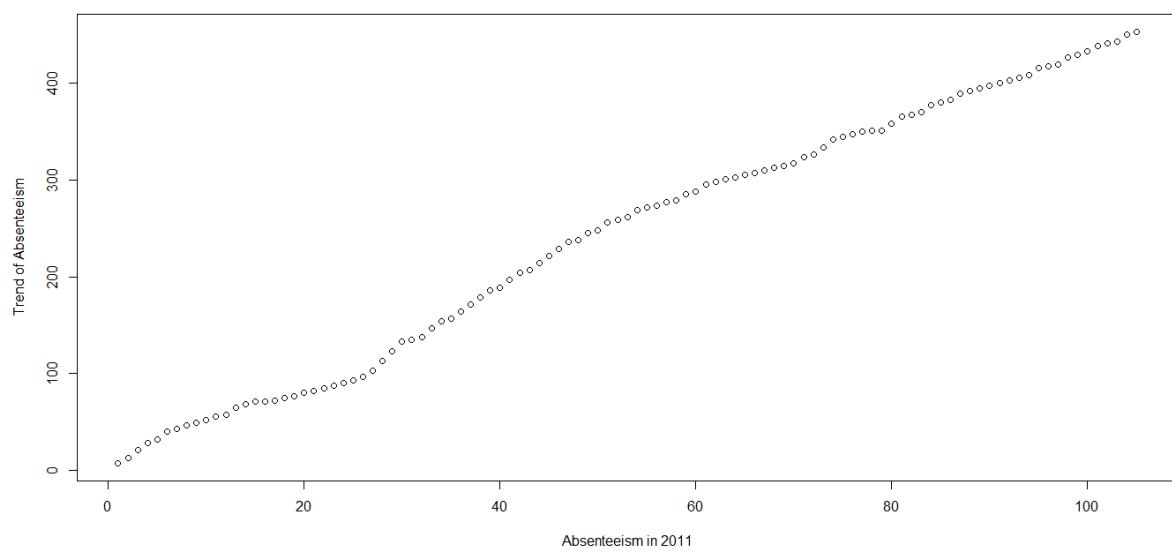
associated decision tree is incrementally developed. The final result is a tree with decision nodes and leaf nodes.

```
#Training Decision tree
DecisionTree = DecisionTreeRegressor(max_depth=10,random_state=0)
DecisionTree.fit(X_Train,Y_Train)

DecisionTreeRegressor(criterion='mse', max_depth=10, max_features=None,
max_leaf_nodes=None, min_impurity_decrease=0.0,
min_impurity_split=None, min_samples_leaf=1,
min_samples_split=2, min_weight_fraction_leaf=0.0,
presort=False, random_state=0, splitter='best')
```

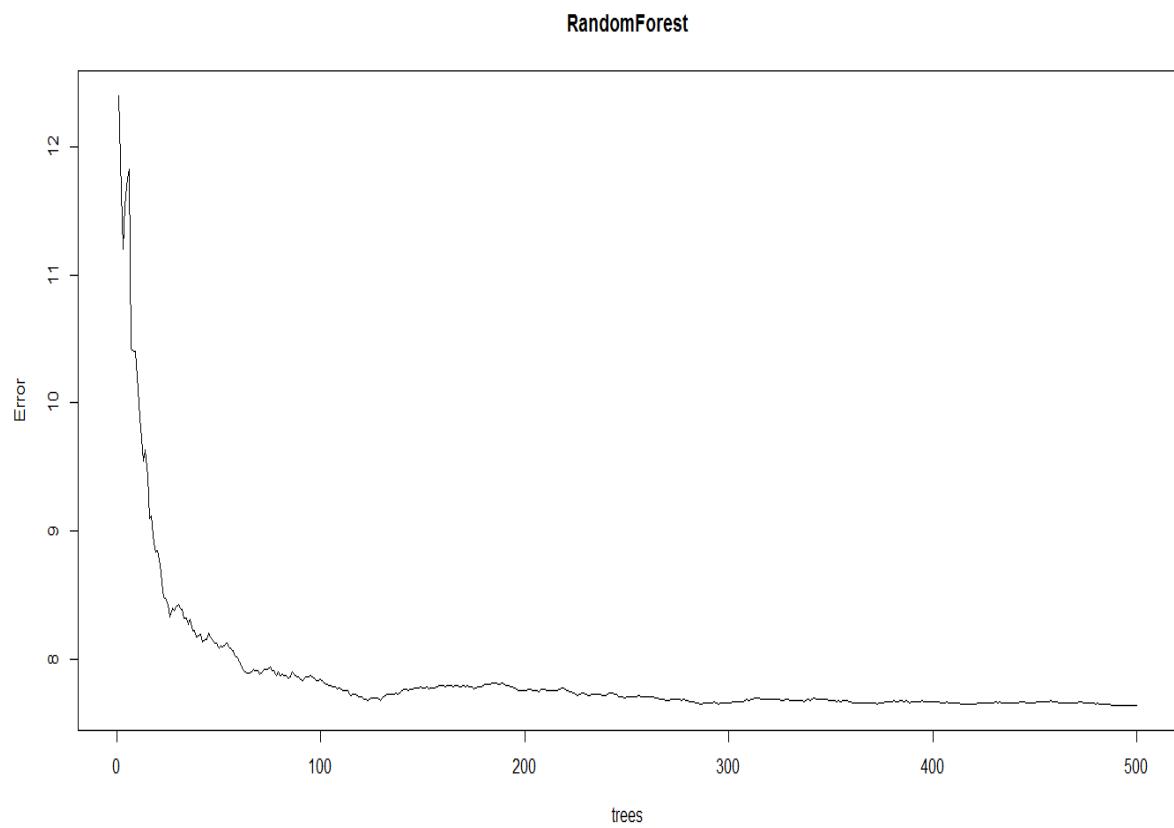
Graphical Visualization of a trained modal is shown below:



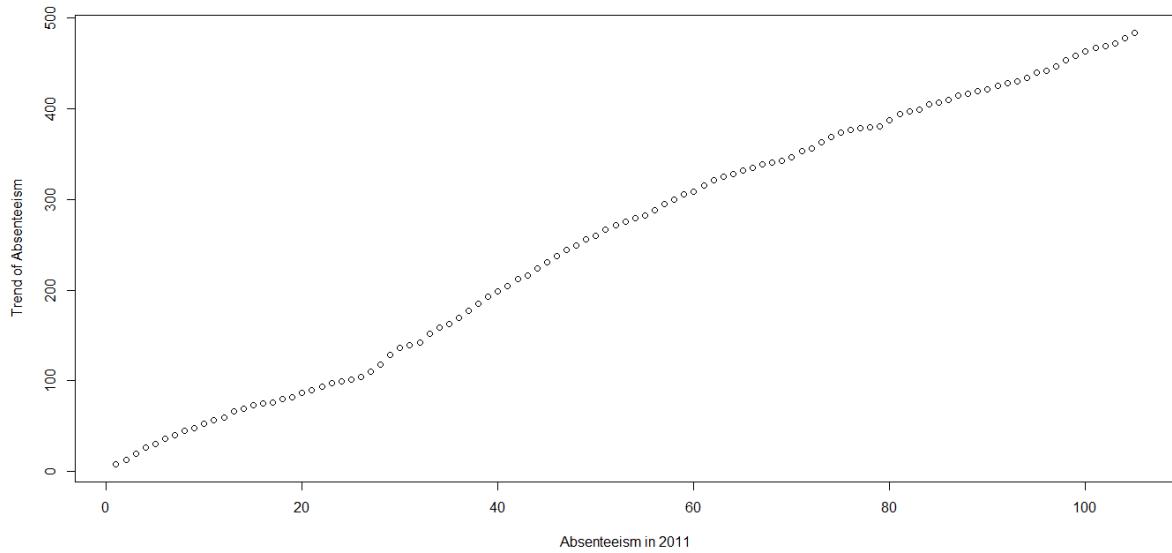


### **Random Forest Model**

Random Forest is a flexible, easy to use machine learning algorithm that produces, even without hyper-parameter tuning, a great result most of the time. It is also one of the most used algorithms, because its simplicity and the fact that it can be used for both classification and regression tasks.

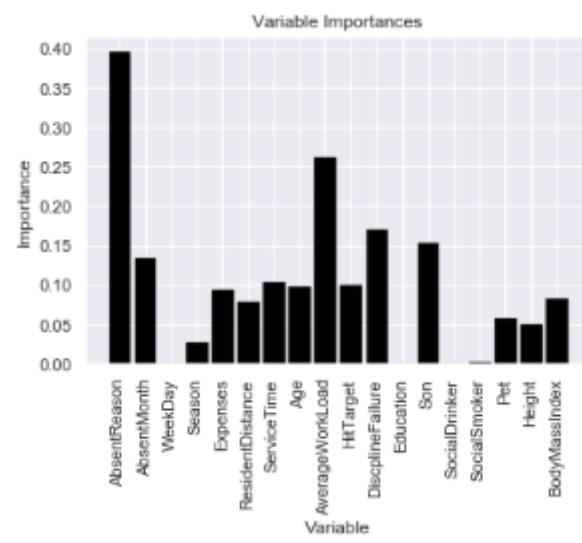


Visualization and Accuracy of the trained model is shown below:-

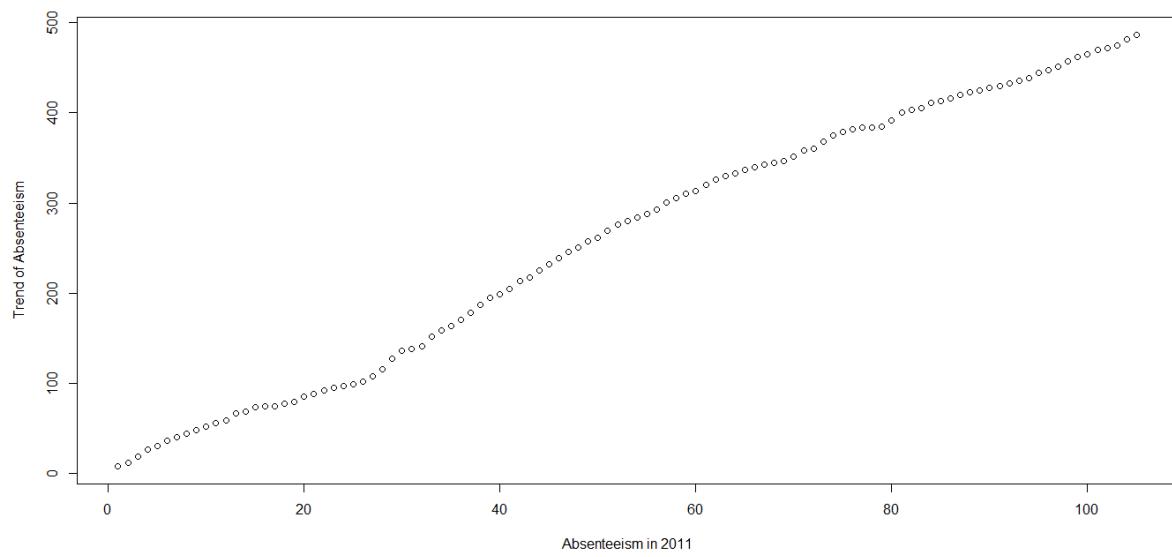


Over repetitions of iteration and training the model with all features were not increasing accuracy. So in order to increase the accuracy of the model we need to check importance of all features in the model then remove the features with least importance. Below graphs shows the importance of all the features present in the model.

Importance	
AbsentReason	0.289164
AverageWorkLoad	0.100921
DisciplineFailure	0.076935
HitTarget	0.071110
AbsentMonth	0.066839
WeekDay	0.065074
BodyMassIndex	0.057203
Expenses	0.050828
Age	0.040098
Height	0.036056
ResidentDistance	0.033223
Season	0.029051
ServiceTime	0.028654
Son	0.023746
Education	0.009844
Pet	0.009831
SocialDrinker	0.009393
Social Smoker	0.002031

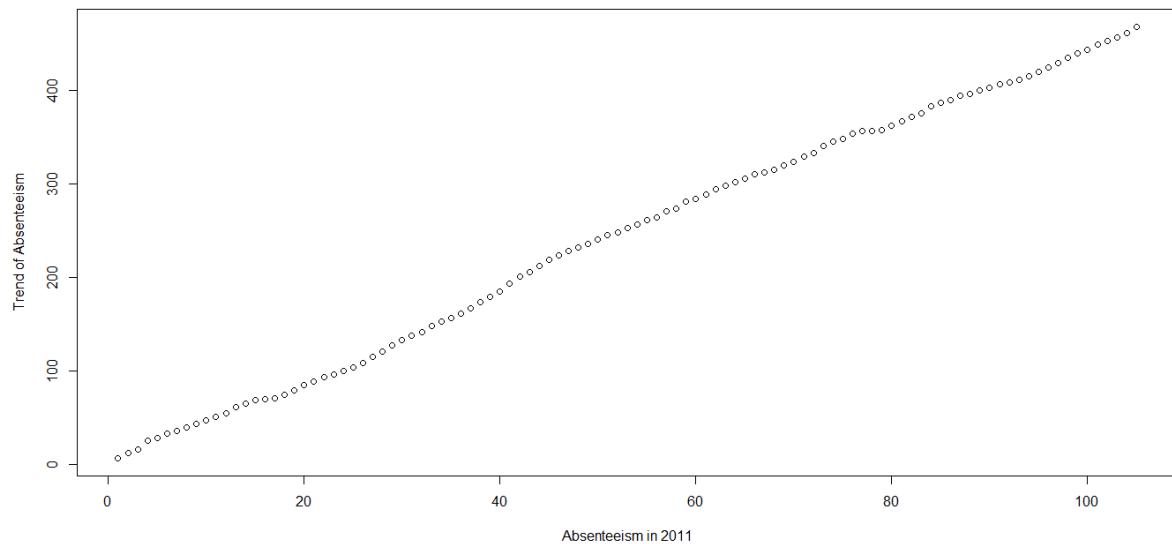


From the above graph, the next task is to remove the feature with less importance. So the graph of trained model after removing those features is shown below:-



### Linear Regression

Multiple Linear regression is the most common form of linear regression analysis. As a predictive analysis, the multiple linear regressions is used to explain the relationship between one continuous dependent variable and two or more independent variables. The independent variables can be continuous or categorical. The trained model is shown below:-

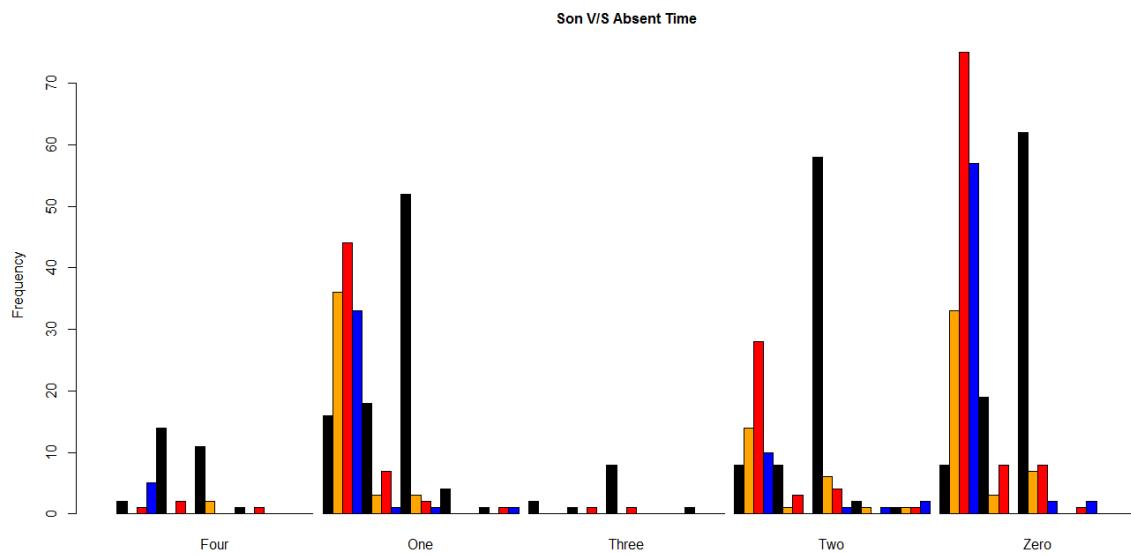


----- VIFs of the remained variables -----

Variables	VIF
Id	7.887307
AbsentReason	1.460496
AbsentMonth	1.523459
WeekDay	1.111049
Season	1.519623
Expenses	3.806174
ServiceTime	5.380615
Age	3.837641
AverageWorkLoad	1.232699
DisciplineFailure	1.438735
Education	4.858774
Son	2.123952
SocialDrinker	8.340348
SocialSmoker	1.979498
Pet	3.546069
Height	1.662905
BodyMassIndex	4.054443

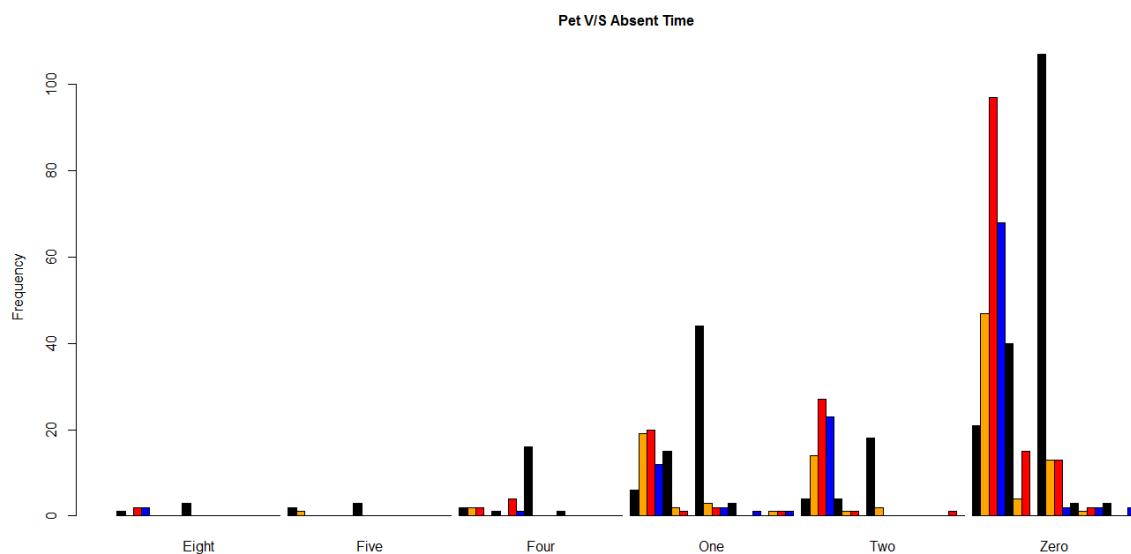
## What changes company should bring to reduce the number of absenteeism?

- Absent Time v/s Number of Sons



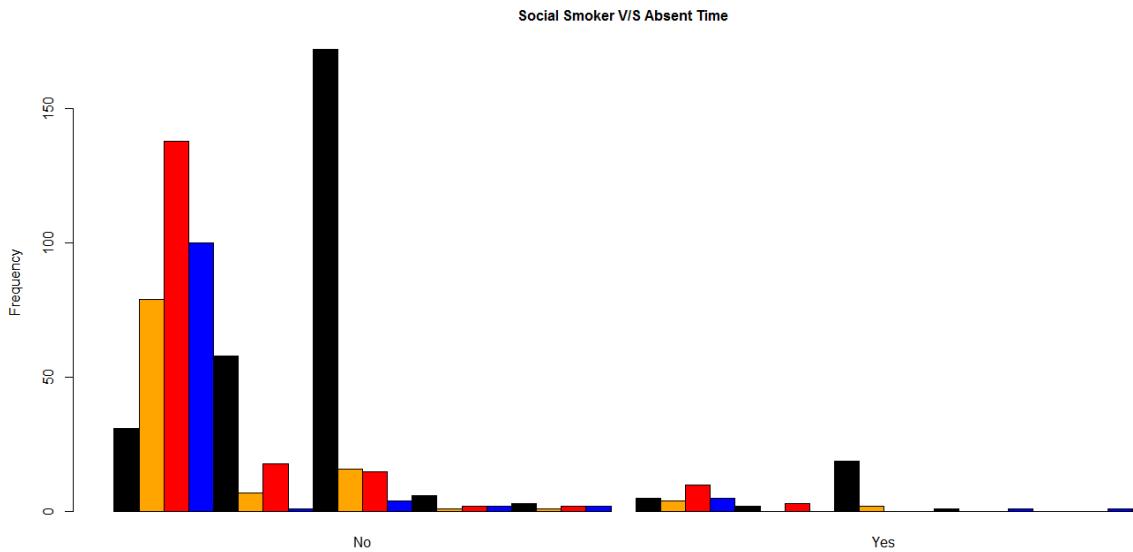
Employee with zero son are maximum time absent from their time. They may need some break from professional life so arranging some fun activities with the help of HR department can resolve this problem. Besides this employee with two and one child are also absent. It may be because their children are very small in age. Providing crèche facility can resolve this problem.

- Absent Time v./s Number of pets



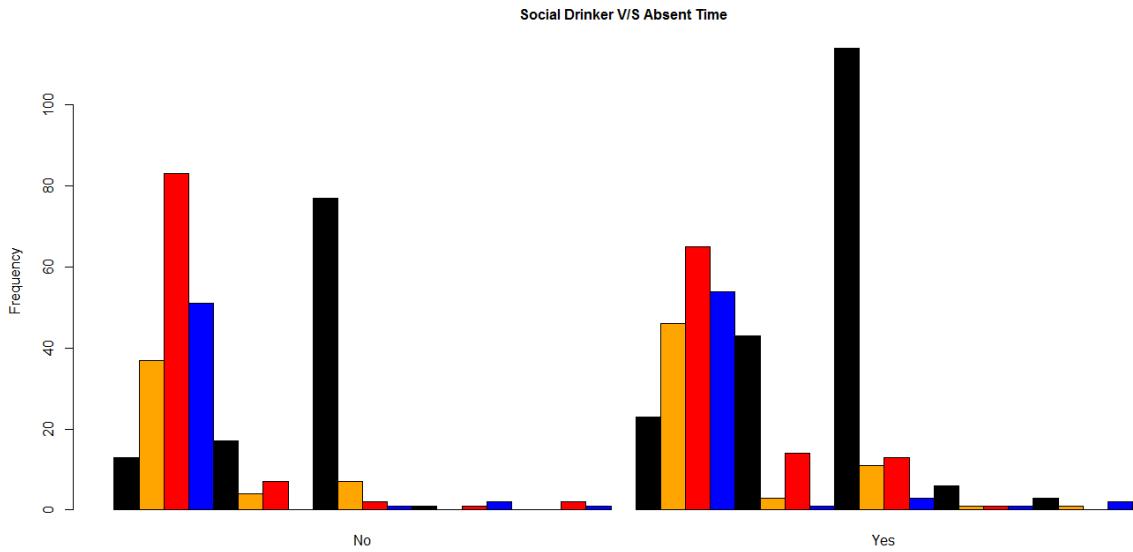
Employee with zero pets are most of the time absent.

- Social Smoker v/s Absent Time



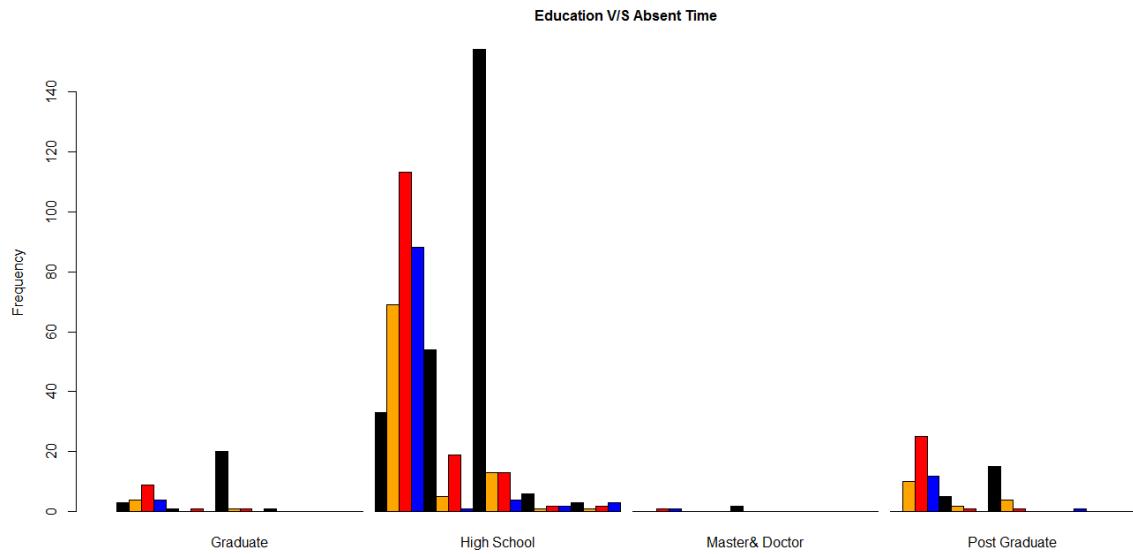
People who are not social smoker are mostly absent from their job.

- Social Drinker v/s Absent Time



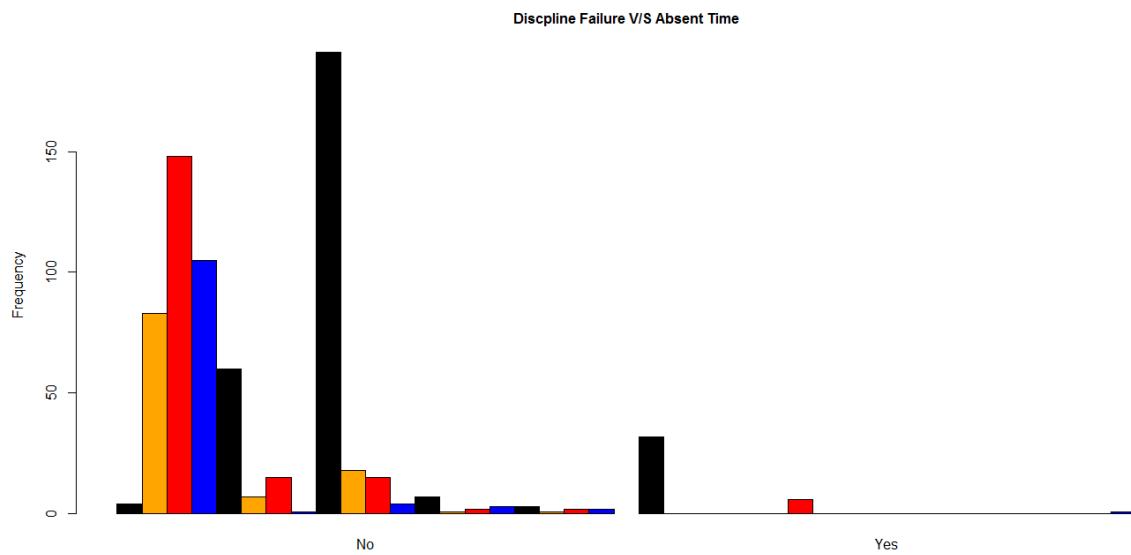
People who are not social smoker are mostly absent from their job.

- Absent Time v/s Education



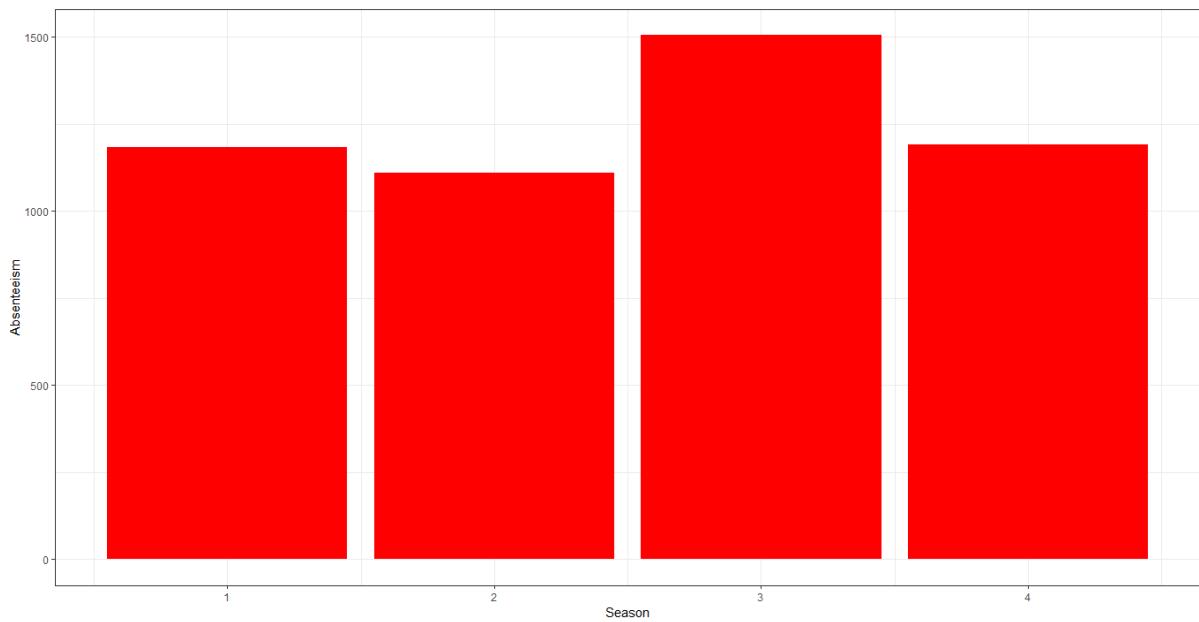
Employees who are less educated are mostly absent. Hiring more educated employees can resolve this issue.

- Absent Time v/s Discipline Failure



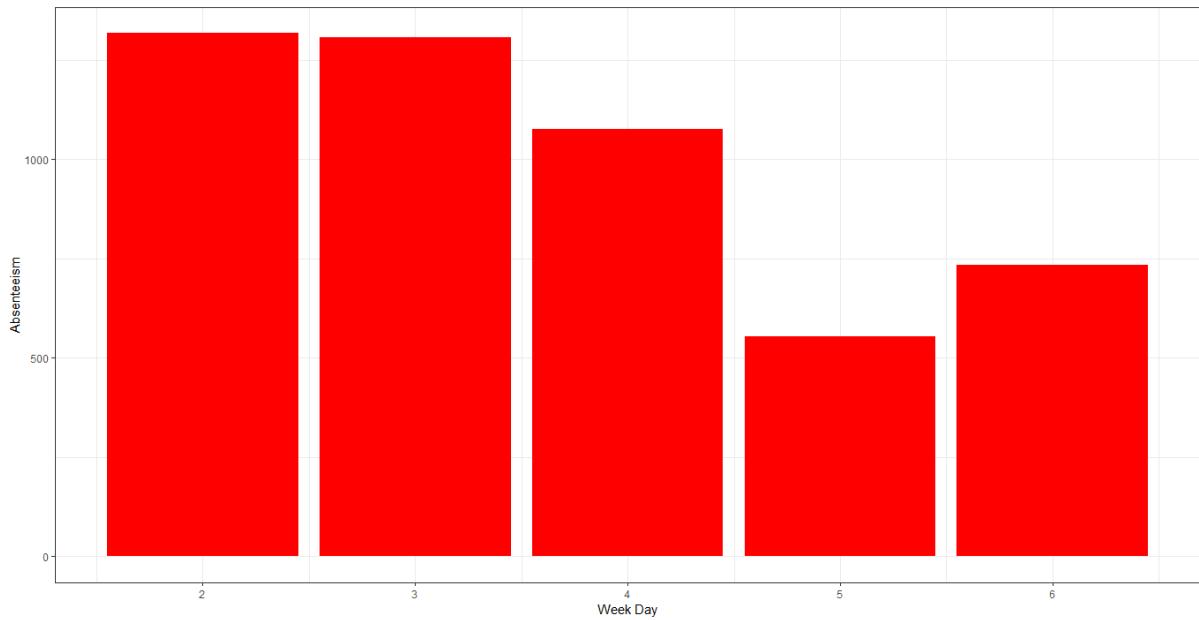
May be company is more strict about its policies. This is leading to more absent time. Giving some relaxation to their employee's can reduce employee absent time.

- Absent Time v/s Season



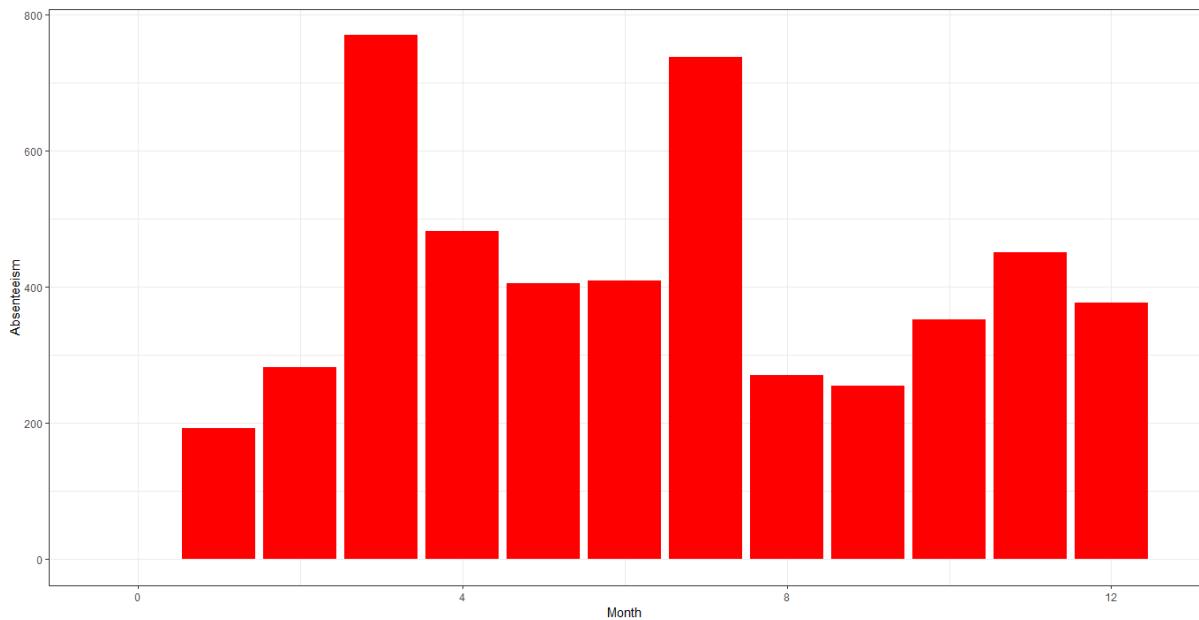
Most of the employees are absent in Winter Season. May be poor heating facility in office may leading to this. Providing blower can resolve this issue. Besides this providing them hot snacks with tea regular will pushes employee to join office daily.

- Absent Time v/s Week day



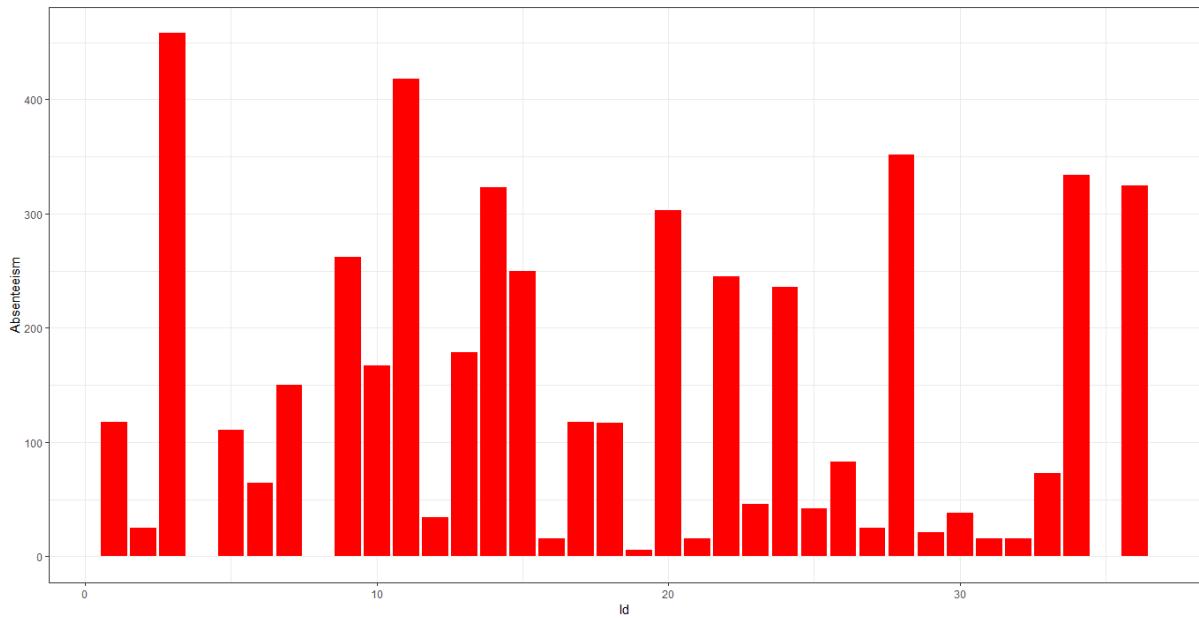
Employees are mostly absent on Monday and Tuesday. Creating meeting on these days will pushes employee to join office on these two days. Issuing warning can also help.

- Absent Time v/s Absent Month



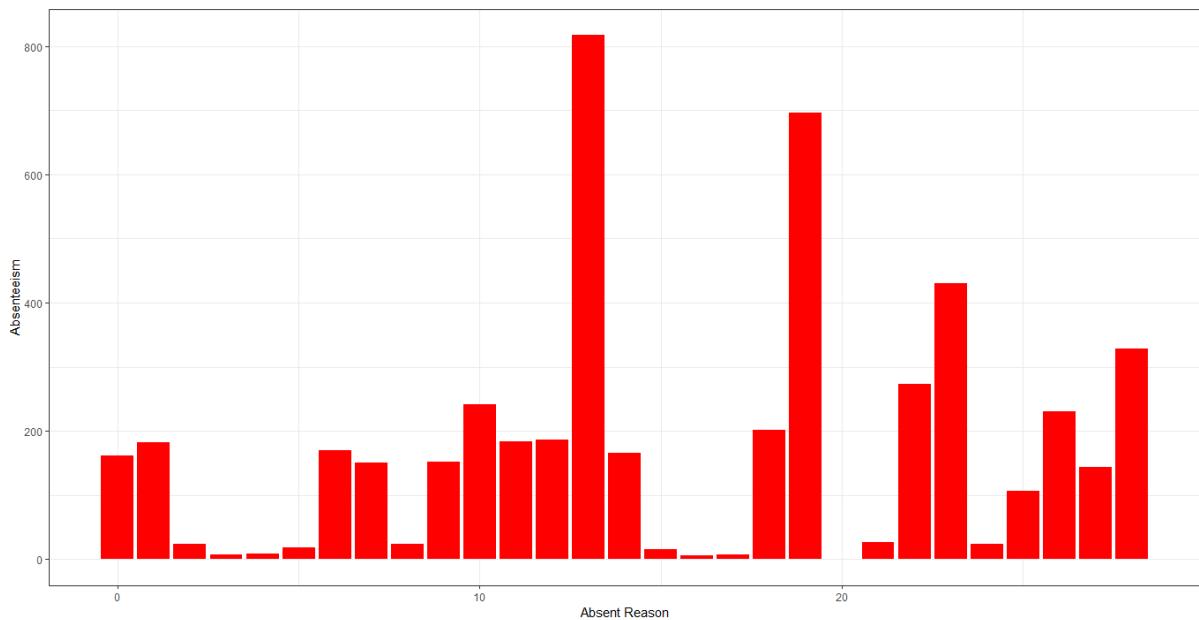
Employees are mostly absent in the month of March and July. Keeping Appraisal season around this month will resolve this issue.

- Absent Time v/s Id



Maximum Absent Rate is with 3,11,20,22,28,34,36 a. giving them warning to be regular can help the firm. Company must interact with the employees of these Ids and motivate them to be regular. Arranging of some motivational talk can increase their moral and helps the firm to deal to reduce absenteeism time of its employee.

- Absent Time v/s Absent Reason



Employees having maximum absent time if they are suffering from disease number 23,26,27,28,13,10,1. Proper medical check up and treatment can help them to reduce employee absent time. Company can provide regular medical check up to their employees. Providing employee's with ESI facilities can also help to deal with this type of problem.

### Graph of Prediction for future Business Decisions

