

# **Employee Absenteeism**

*Akash Kumar Dubey*

*17 Sept 2018*

# Contents

|   |           |
|---|-----------|
| <b>1. Introduction</b>                    | <b>3</b>  |
| 1.1 Problem Description .....             | 3         |
| 1.2 Problem Statement .....               | 3         |
| 1.3 Data .....                            | 3         |
| 1.4 Performance Metric .....              | 5         |
| <b>2. Methodology</b>                     | <b>6</b>  |
| 2.1 Exploratory Data Analysis .....       | 6         |
| 2.1.1 Data Visualisation .....            | 6         |
| 2.1.1.1 Univariate Analysis .....         | 6         |
| 2.1.1.2 Bivariate Analysis .....          | 10        |
| 2.1.1.3 Multivariate Analysis .....       | 14        |
| 2.1.2 Data Preparation And Cleaning ..... | 16        |
| 2.1.2.1 Missing Value Analysis.....       | 16        |
| 2.1.2.2 Outlier Analysis .....            | 16        |
| 2.1.2.3 Feature Selection .....           | 17        |
| 2.2 Modeling                              |           |
| 2.2.1 KNN.....                            | 18        |
| 2.2.2 Ordinary Least Squares.....         | 18        |
| 2.2.3 Ridge Regression.....               | 18        |
| 2.2.4 Lasso Regression.....               | 19        |
| 2.2.5 Support Vector Regression .....     | 19        |
| 2.2.6 Decision Tree .....                 | 20        |
| 2.2.7 Gradient Boosted Decision Tree..... | 20        |
| 2.2.8 Random Forest .....                 | 21        |
| <b>3. Conclusion</b>                      | <b>23</b> |
| 3.1 Model Evaluation .....                | 23        |
| 3.1.1 Root Mean Square Value.....         | 23        |
| 3.2 Model Selection .....                 | 24        |
| 3.3 Answer to asked Questions.....        | 25        |
| <b>Appendix A - Extra Figures</b> .....   | <b>28</b> |
| <b>References</b> .....                   | <b>30</b> |

# Chapter 1

## Introduction

### 1.1 Problem Description

Employee Absenteeism is the absence of an employee from work. Its a major problem faced by almost all employers of today. Employees are absent from work and thus the work suffers. Absenteeism of employees from work leads to back logs, piling of work and thus work delay.

Absenteeism can be of two types :

- **Innocent absenteeism** - Is one in which the employee is absent from work due to genuine cause or reason. It may be due to his illness or personal family problem or any other real reason
- **Culpable Absenteeism** - is one in which a person is absent from work without any genuine reason or cause. He may be pretending to be ill or just wanted a holiday and stay at home.

### 1.2 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 much losses every month can we project in 2011 if same trend of absenteeism continues?

### 1.3 Data

The data is a Time-Series data but instead we will approach it as Regression Problem. Our task is to build a regression model which will predict the absenteeism in hours per employee based on the employee attributes and information in their work place and general information available to the company about them.

Table 1.1 : Employee Attributes Data (Columns : 1-8)

| ID | Reason for<br>Absence | Day of Month | Day of the<br>week | Seasons | Transportation<br>Expense | Distance From<br>Residence to<br>work | Service Time |
|----|-----------------------|--------------|--------------------|---------|---------------------------|---------------------------------------|--------------|
| 11 | 26                    | 7            | 3                  | 1       | 289                       | 36                                    | 13           |
| 36 | 0                     | 7            | 3                  | 1       | 118                       | 13                                    | 18           |
| 3  | 23                    | 7            | 4                  | 1       | 179                       | 51                                    | 18           |
| 7  | 7                     | 7            | 5                  | 1       | 279                       | 5                                     | 14           |
| 11 | 23                    | 7            | 5                  | 1       | 289                       | 36                                    | 13           |

Table 1.2 : Employee Attributes Data (Columns : 9-15)

| Age | Work load<br>Average/day | Hit target | Disciplinary<br>failure | Education | Son | Social Drinker |
|-----|--------------------------|------------|-------------------------|-----------|-----|----------------|
| 33  | 239554                   | 97         | 0                       | 1         | 2   | 1              |
| 50  | 239554                   | 97         | 1                       | 1         | 1   | 1              |
| 38  | 239554                   | 97         | 0                       | 1         | 0   | 1              |
| 39  | 239554                   | 97         | 0                       | 1         | 2   | 1              |
| 33  | 239554                   | 97         | 0                       | 1         | 2   | 1              |

Table 1.3 : Employee Attributes Data (Columns : 16-21)

| Social Smoker | Pet | Weight | Height | Body mass Index | Absenteeism time<br>in hours |
|---------------|-----|--------|--------|-----------------|------------------------------|
| 0             | 1   | 90     | 172    | 30              | 4                            |
| 0             | 0   | 98     | 178    | 31              | 0                            |
| 0             | 0   | 89     | 170    | 31              | 2                            |
| 1             | 0   | 68     | 168    | 24              | 4                            |
| 0             | 1   | 90     | 172    | 30              | 2                            |

As we can see in the table above, we have we have the following 20 variables, using which we have to correctly predict the ‘Absenteeism in hours’ for the employees.

Table 1.4 : Predictor Variables

| S.No | Predictor                       | S.No | Predictor              |
|------|---------------------------------|------|------------------------|
| 1    | ID                              | 11   | Social drinker         |
| 2    | Reason for absence              | 12   | Social smoker          |
| 3    | Distance from residence to work | 13   | Pet                    |
| 4    | Service time                    | 14   | Day of the week        |
| 5    | Age                             | 15   | Weight                 |
| 6    | Work load average/day           | 16   | Height                 |
| 7    | Hit target                      | 17   | Body mass index        |
| 8    | Disciplinary failure            | 18   | Transportation expense |
| 9    | Education                       | 19   | Month of absence       |
| 10   | Son                             | 20   | Seasons                |

## 1.4 Performance Metric

RMSE : Root Mean Square Error (RMSE) is the standard deviation of the residuals (prediction errors). Residuals are a measure of how far from the regression line data points are, RMSE is a measure of how spread out these residuals are. In other words, it tells you how concentrated the data is around the line of best fit.

Also, Since the errors are squared before they are averaged, the RMSE gives a relatively high weight to large errors. So, RMSE becomes more useful when large errors are particularly undesirable. So, Roost Mean Square value seems like a perfect choice for our problem at hand.

# Chapter 2

## Methodology

### 2.1 Exploratory Data Analysis

Exploratory Data Analysis (EDA) is the first step in our data analysis process. We do this by taking a broad look at patterns, trends, outliers, unexpected results and so on in our existing data, using visual and quantitative methods to get a sense of the story this tells. To start with this process, we will first have a look at univariate analysis like plotting Box plot and whiskers for individual features, Histogram plots, Bar plots and Kernel Density Estimation for the same for the same. Then we will proceed to Multivariate analysis like Bar and Histogram and Bar plot using group-by function and Pivot table for the features with respect to the target variable.

#### 2.1.1 Data Visualisation

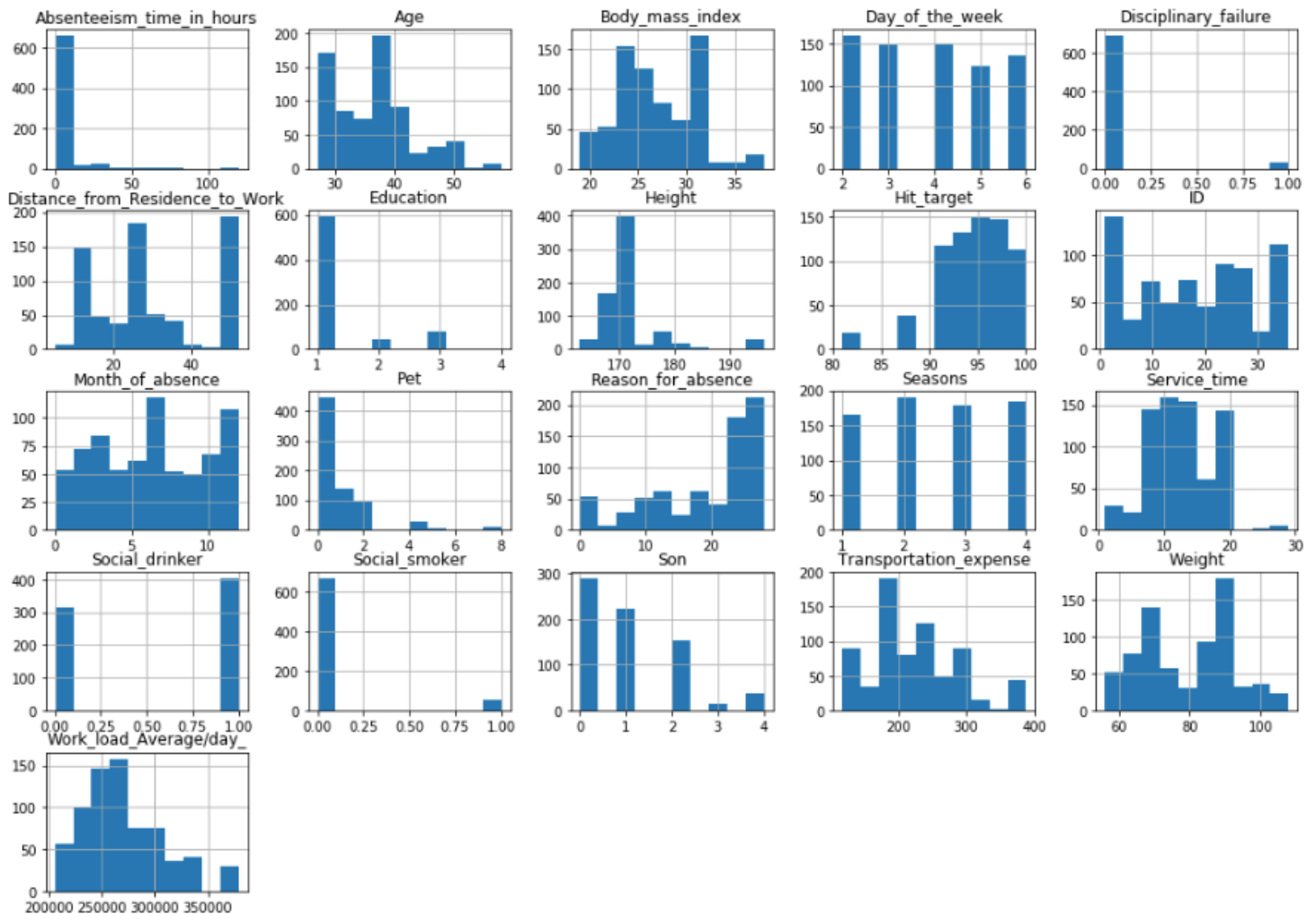
Data visualisation helps us to get better insights of the data. By visualising data, we can identify areas that need attention or improvement and also clarifies which factors influence customer behaviour and how the resources are used by the customers.

##### 2.1.1.1 Univariate Analysis

Univariate analysis is the simplest form of data analysis where the data being analysed contains only one variable. Since it's a single variable it doesn't deal with causes or relationships. The main purpose of univariate analysis is to describe the data and find patterns that exist within it.

So, Lets have a look at histogram plot, to identify the characteristic of the features and the data.

Figure 2.1 : Histogram plot for distribution of features in the data

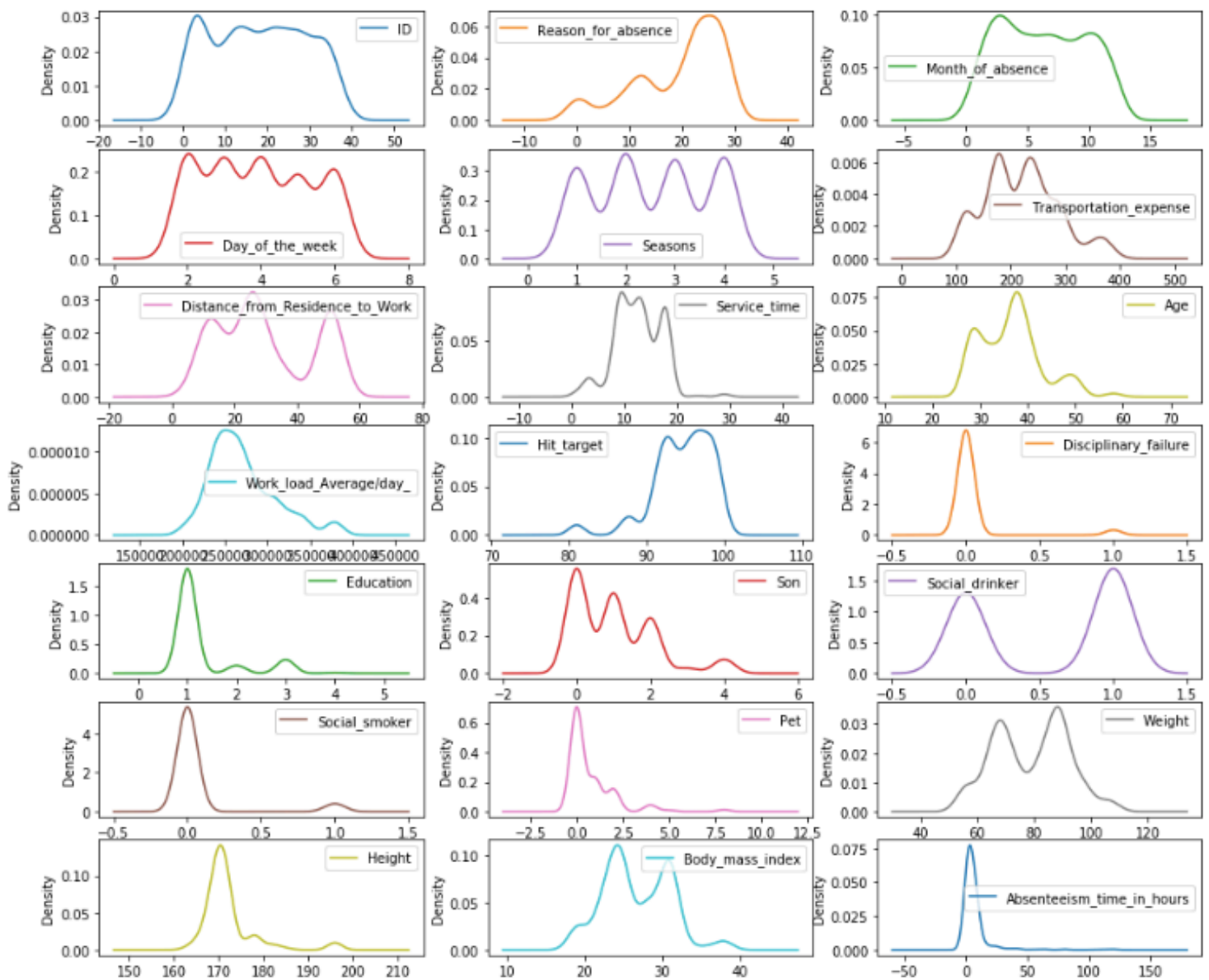


Histograms are constructed by binning the data and counting the number of observations in each bin. The objective of plotting Histogram plot is usually to visualise the shape of the distribution. The number of bins needs to be large enough to reveal interesting features and small enough not to be too noisy.

From the above histogram plot, we can clearly observe that none of the features in our data are actually skewed. Although feature like 'work load average/day' seems like it is right skewed a little. Also if observed properly, It is worth noting the following points :

1. Majority of the employees working in the company have age below 40 years.
2. A very large portion of the population have only passed 'High School'.
3. More than half of the employees in the company are 'social drinker'.
4. Only a very few portion of the employees in the company are 'social smoker'.

Figure 2.2 : KDE plot for distribution of features in the data

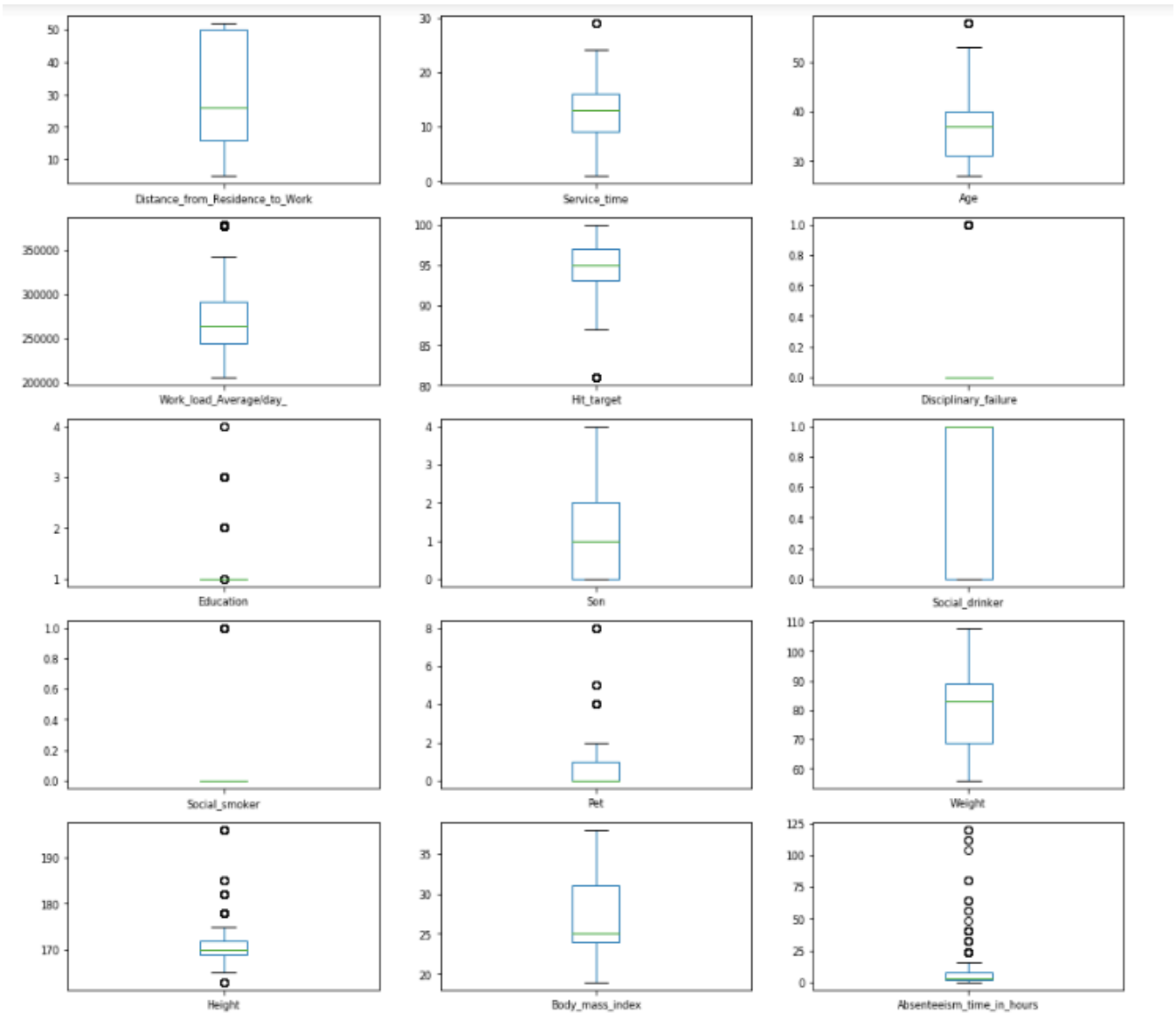


A Density Plot visualises the distribution of data over a continuous interval or time period. Density plots can be thought of as plots of smoothed histograms. An advantage Density Plots have over Histograms is that they're better at determining the distribution shape because they're not affected by the number of bins used.

So, Looking at the above density plot , we can observe that none of the features follow Gaussian distribution. Few of the features like 'Disciplinary failure', 'Social smoker', 'Work load average/day' seems to follow gaussian distribution at first sight but they either have long tail at the left or right or they are either jagged at the end.



Figure 2.3 : Box and Whiskers plot of features in the data



From the above Box and whisker plots , we can observe that not all the features contains outliers. Continuous features like ‘Weight’ , ‘Distance from residence to work’ does not contain any outliers at all. Few features like ‘work load average/day’, ‘Hit target’ and ‘Height’ have a very few outliers.

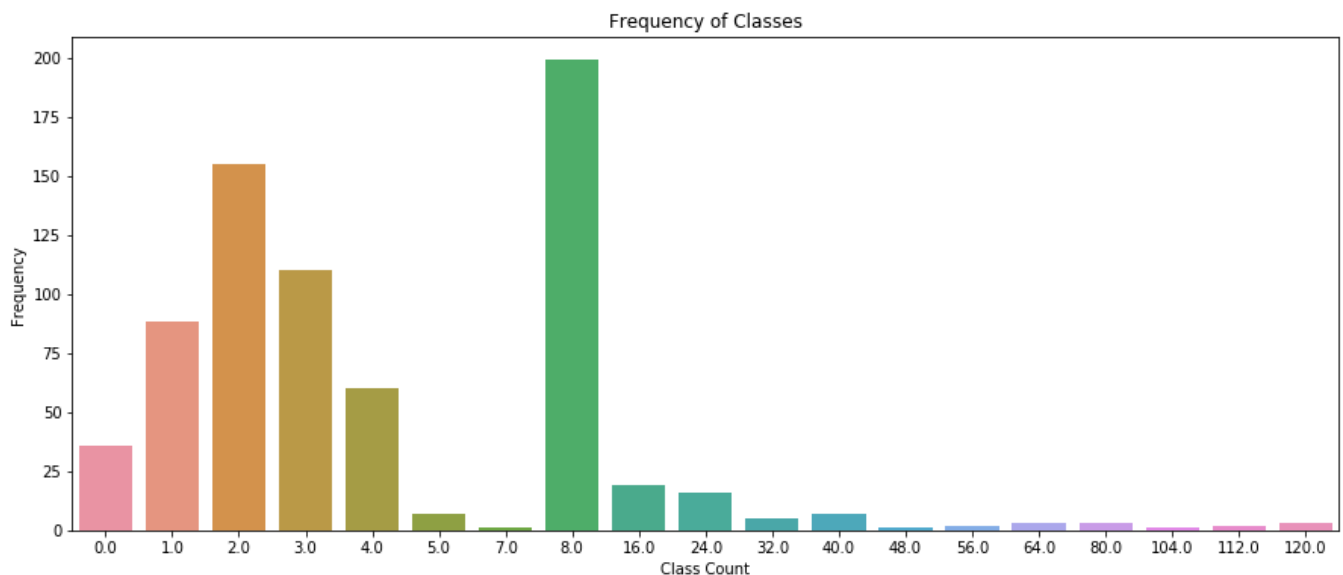
It is also evident from the above plot that none of the features are symmetric to the median and it can easily be interpreted that none of the features follow symmetric distribution. Also, it can also be observed that Median of the feature ‘Body mass index’ is very close to 25th percentile value which means median of this feature is almost equal to 25th percentile.

### 2.1.1.2 Bivariate Analysis

Bivariate analysis refers to the analysis of bivariate data. It is one of the simplest forms of statistical analysis, used to find out if there is a relationship between two sets of values. It usually involves one predictor variable and one target variable.

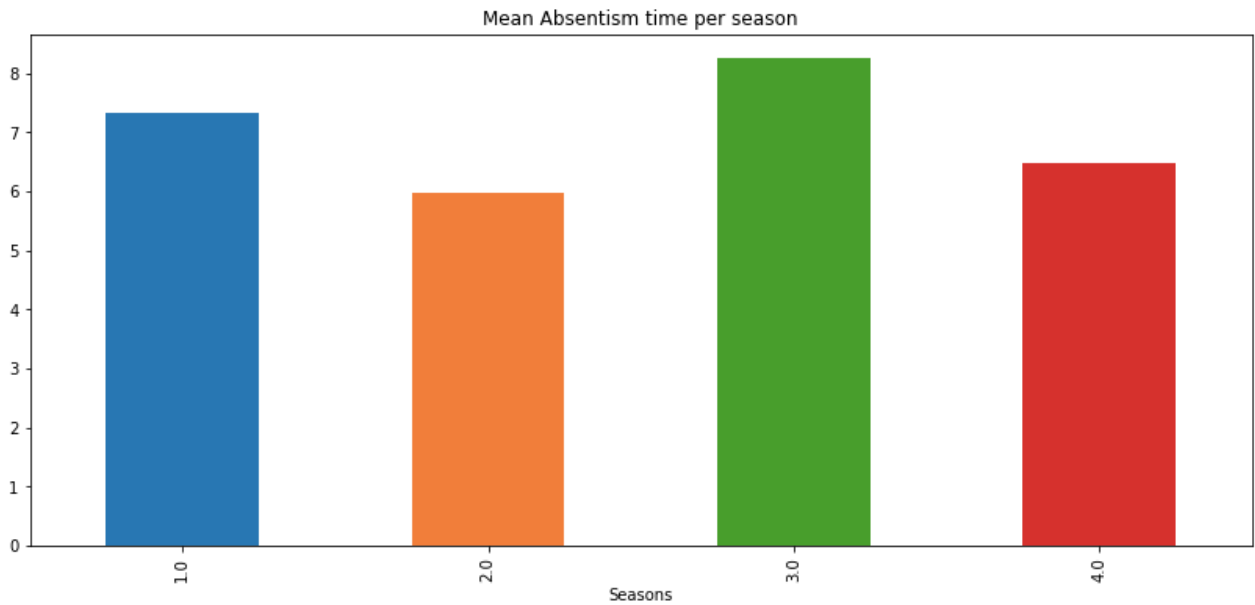
So, Lets have a look at the Histogram and Bar Plots to understand the Employee behaviour better.

Figure 2.4 : Bar plot for Distribution of Target class



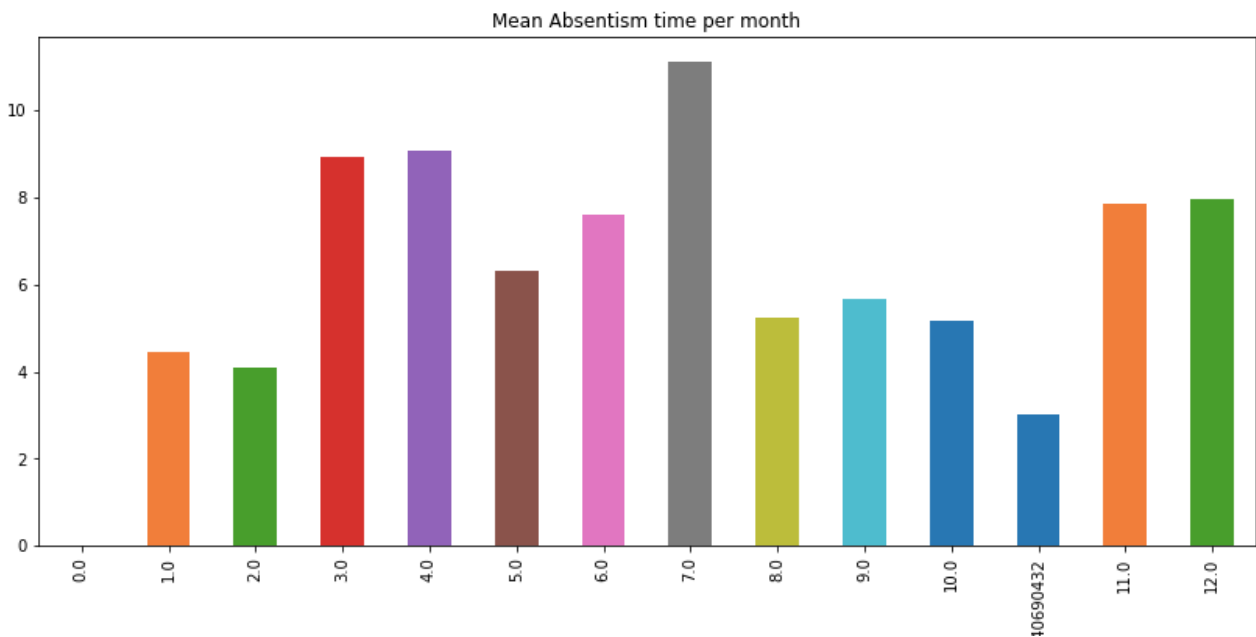
So, we can clearly observe from the above histogram plot that, most of the data points in our data are labelled as '8', followed by '2', '3', '1', '4', '0' etc. Also, It is worth noticing that any label after '8' is a multiple of 8 and have the least occurrence within the data.

Figure 2.5 Bar Plot for Mean absenteeism time per season



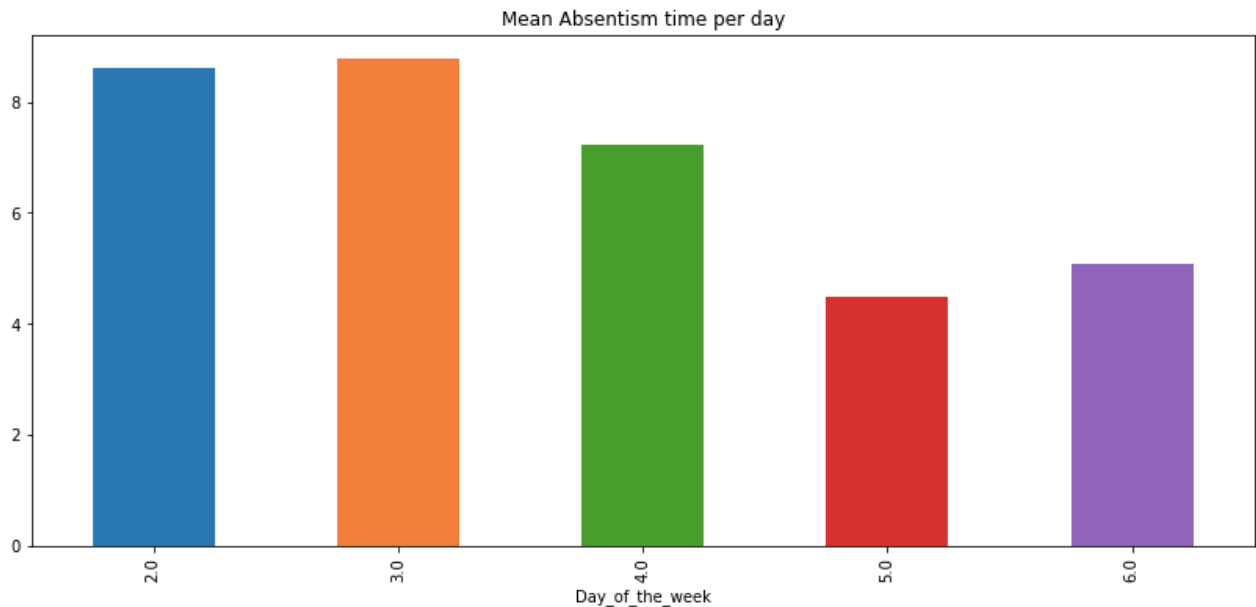
From the above Bar plot , it can be observed that the ‘*Absenteeism rate*’ is maximum in *Season 3 : Winter* followed by *Season 1 : Summer* , *Season 4 : Spring*, *Season 2 : Autumn*.

Figure 2.6 Bar Plot for Mean absenteeism time per month



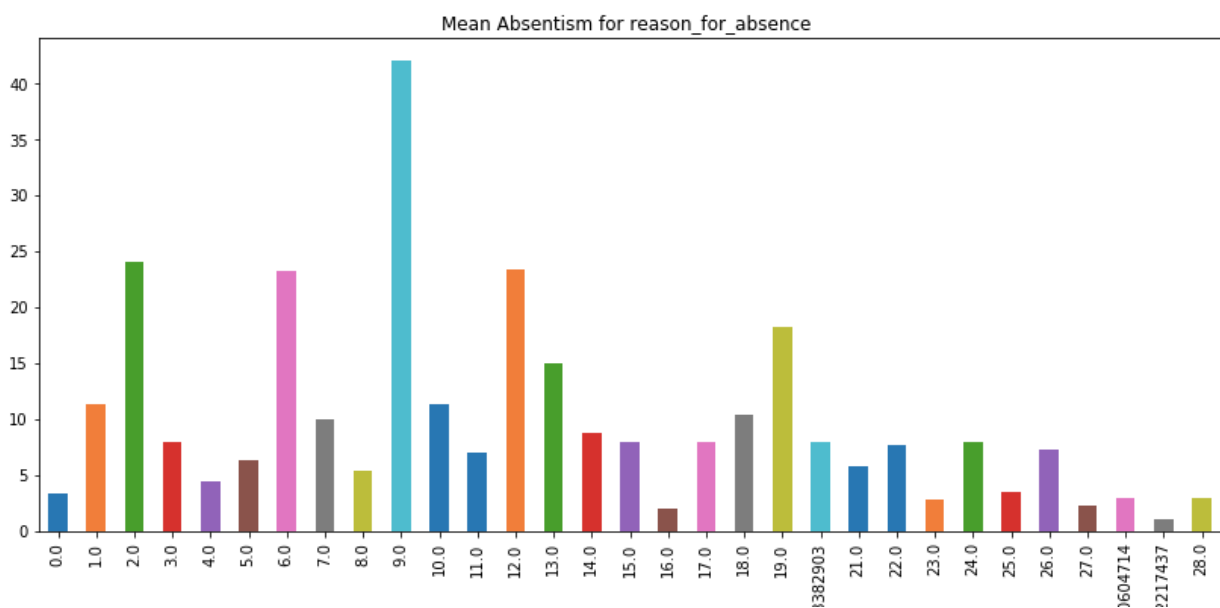
From the above Bar plot , we can clearly observe that the ‘*Absenteeism rate*’ is maximum in *Month 7 : July* followed by *Month 4 : April* , *Month 3 : March*, *Month 12 : December*, *Month 11 : November* , *Month 6 : June* , *Month 5 : May* etc.

Figure 2.7 Bar Plot for Mean absenteeism time per day



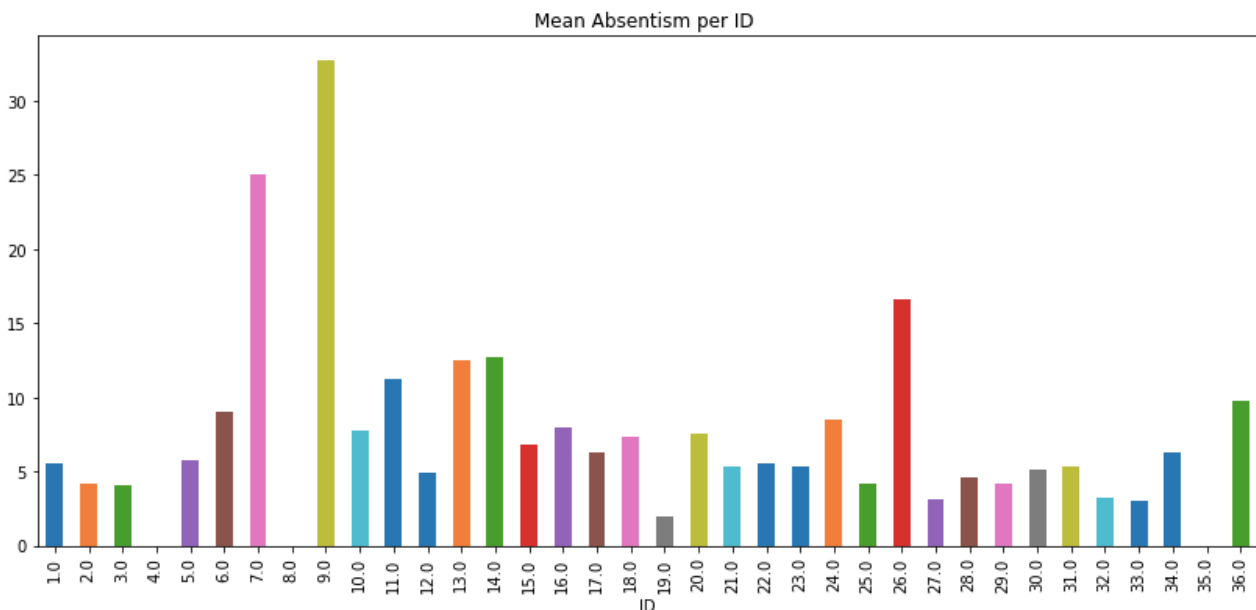
From the above Bar plot, it can clearly be observed that the 'Absenteeism rate' is maximum on the third day of the week i.e *Day 3 : Tuesday* followed by *Day 2 : Monday*, *Day 4 : Wednesday*. Also, the 'absenteeism rate' is lowest on *Day 6 : Saturday* followed by *Day 5 : Friday*.

Figure 2.8 Bar Plot for Mean absenteeism for reason\_of\_absence



From the above plot, we can observe that ‘9 : *Diseases of the circulatory system*’ is the most frequent reason for the absence of the employees. The second most frequent reason given by the employees for their absence is ‘2 : *Neoplasms*’ followed by ‘6 : *Diseases of the nervous system*’, ‘12: *Diseases of the skin and subcutaneous tissue*’, ‘19 : *Injury, poisoning and certain other consequences of external causes*’ etc.

Figure 2.9 Bar Plot for Mean absenteeism time per ID



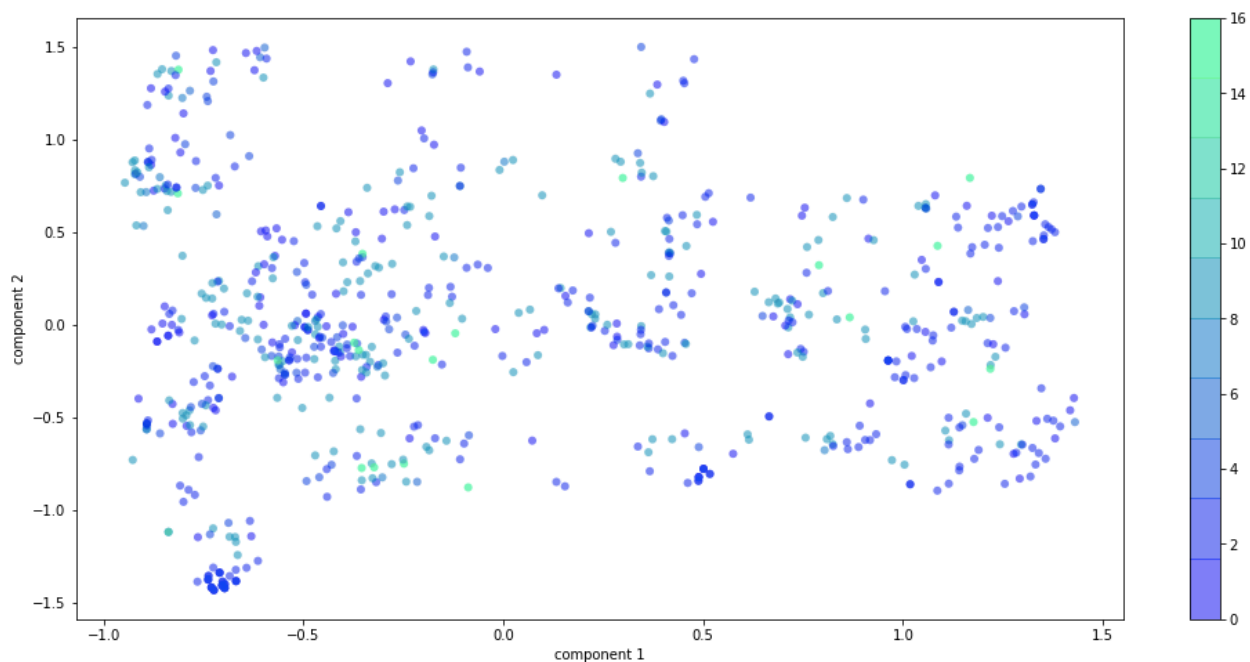
From the above plot, It can be observed that the absence rate is maximum for employee with *ID* : 9, followed by employees with *ID* : 7,26,14,13, 36, 11 and 6. Also, it can be observed that Employee with employee *ID* : 4,8 and 35 never absents and are very much regular to work.

### 2.1.1.3 Multivariate Analysis

Multivariate analysis is the analysis of more than one variable in a dataset. Multivariate analysis becomes important when we have large dimensional data to visualise and it becomes very difficult to visualise every predictor variable individually. It also helps us to identify the dominant patterns and clusters in the data.

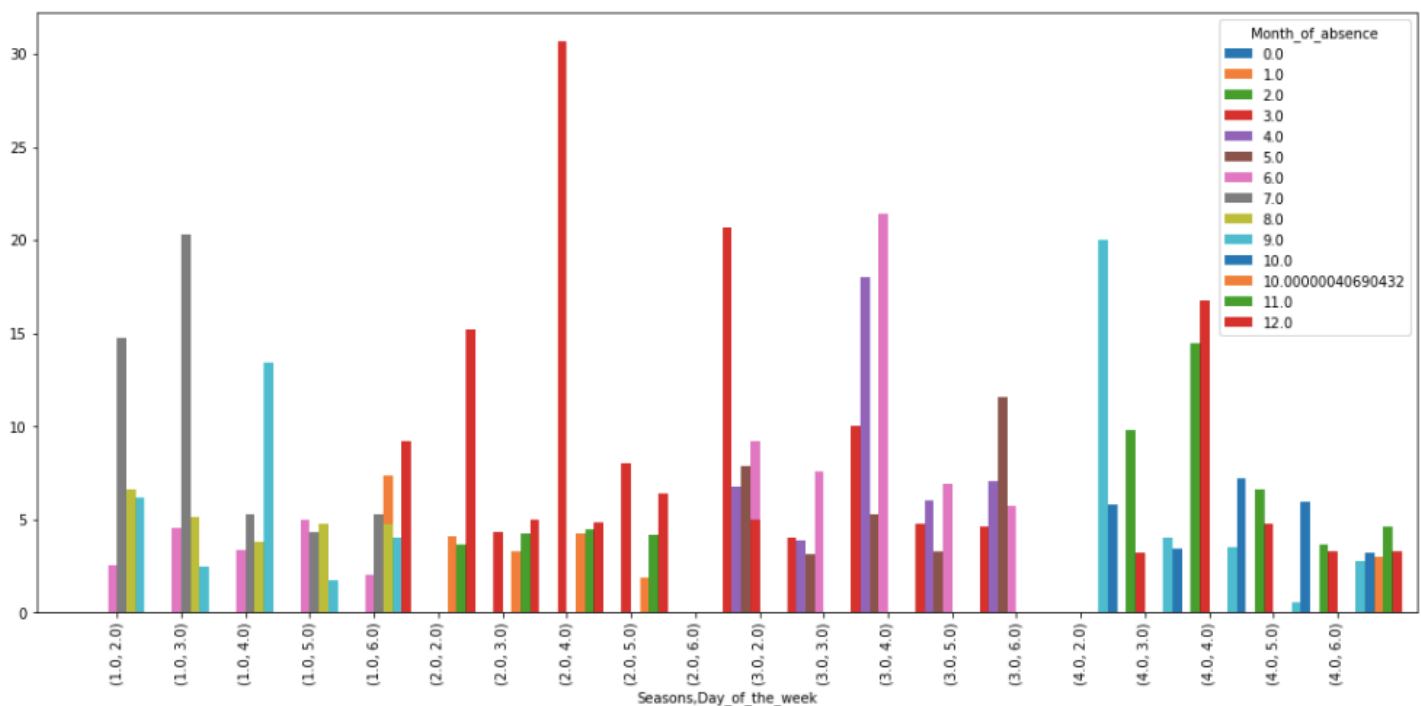
So, we will first look at one of the most widely used algorithm to visualise high dimensional data based on Eigen Value and Eigen vectors which is called as *Principal Component Analysis*. In PCA visualisation, we have plotted the scatter plot along two dimensions.

Figure 2.10 PCA Visualisation of the data



In the above plot, 100% of the variance in the data is explained by the principal component 1 and 0% of the variance in the data is explained by the principal component 2. Also, it is hard to interpret anything from the above plot except the fact that the data is spread unevenly across the space. The data does not seem to follow any pattern or any kind of linear or polynomial relationship.

Figure 2.11 Season,Month of absence,Day of week visualisation of the data



The above plotted bar graph is plotted using Pivot table. The graph shows relation between 'Seasons', 'Day of the week', 'Month of absence' with 'Absenteeism in hours' in a single plot. We can draw many insights from the above plot, like : In 2nd Season, 4th day of the week and in the month of December, the absenteeism rate was the highest. So, the above plot narrows down the visualisation to Absenteeism rate for even days of a particular season for a particular month.

## **2.1.2 Data Preparation and Cleaning**

### **2.1.2.1 Missing Value Analysis**

One of the most common problems I have faced in Data Cleaning/Exploratory Data Analysis is handling the missing values. Firstly, there is no good way to deal with missing data. But still missing value analysis helps address several concerns caused by incomplete data. If cases with missing values are systematically different from cases without missing values, the results can be misleading. Also, missing data may reduce the precision of calculated statistics because there is less information than originally planned. Another concern is that the assumptions behind many statistical procedures are based on complete cases, and missing values can complicate the theory required.

So, In our data, there are plenty of missing values available in different variables. So, after computing the percentage of missing data that is available to us in the dataset, it accounts to around 12% of the data. It is also important to note that, the missing value has been calculated after removing the missing values within the target variable. Also, as we have very less data available to us, we impute the missing values in other columns using KNN imputation, because that fits the best after trying various other imputation techniques like Mean, Median and Random value imputation.

### **2.1.2.2 Outliers Analysis**

In statistics, an outlier is an observation point that is distant from other observations. In layman terms, we can say that an outlier is something which is separated/different from the crowd. Also, Outlier analysis is very important because they affect the mean and median which in turn affects the error (absolute and mean) in any data set. When we plot the error we might get big deviations if outliers are in the data set.

In Box plots analysis of individual features, we can clearly observe from these box-plots that, not every feature contains outliers and many of them even have very few outliers. Also, given the constraint that, we have only 640 data-points and after removing the outliers, the data gets decreased by almost 25%. So, dropping the outliers is probably not the best idea.

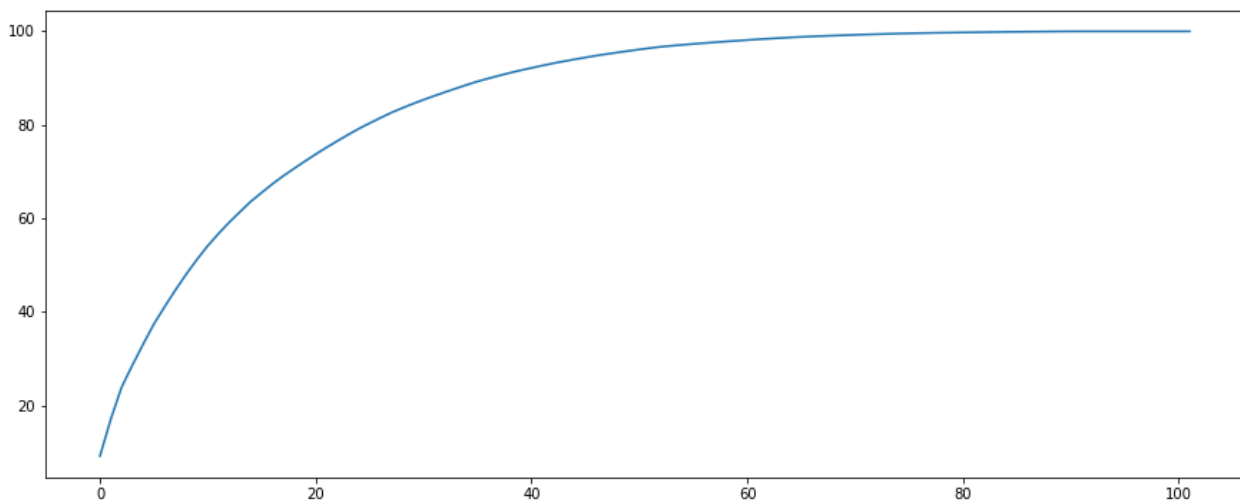
Instead we will try to visualise and find out the outliers using box plots and will fill them with NA, that means we have created 'missing values' in place of outliers within the data. Now, we can treat these outliers like missing values and impute them using standard imputation techniques. In our case, we use KNN imputation to impute these missing values.



### 2.1.2.3 Feature Selection

Before performing any type of modeling we need to assess the importance of each predictor variable in our analysis. There is a possibility that many variables in our analysis are not important at all to the problem of class prediction. There are several methods of doing that. Here, we do this in two steps : Firstly, we find and remove the correlated features and then we use a more advanced technique for dimensionality reduction called PCA .While doing this, we first plot a cumulative distribution function plot to observe how much percentage of variance is explained by how many variables (Principle Components). The CDF plot for the same is plotted below :

Figure 2.12 CDF Plot for ‘Variance Explained’ Vs ‘Principle components’



It is very clear from the above CDF plot for ‘*Variance Explained*’ Vs ‘*Principle components*’, that almost 95%+ variance is explained by just 45 variables (Principle components). We, can imagine how powerful PCA is, It just shrank down our feature space to just 45 from a total of 107 features. So, we will keep only 45 principle components in the data and will perform modeling on it.

## 2.2 Modeling

We always start our model building from the most simplest to more complex. Therefore we start with KNN Regressor.

### 2.2.1 KNN Regression

KNN regression is one of the simplest algorithm in the whole of Machine learning. It gives a weighted average of the regression function in a local space (k nearest points to a given point). So, we first try to implement and fit KNN regression to our Data and got following results after tuning the hyper-parameter k :

```
Train Data
n_neighbours : 30 ----KNN rmse: 2.2706908685134333
Test Data
n_neighbours : 30 ----KNN rmse: 2.756796390255527
```

So, as we can see from the above results that, our model gives a RMSE (Root Mean Square Value) of 2.27 for the train data and a RMSE value of 2.75 for the test data. Looking at the train and test error, We can say that our model seems like overfitting a little. But still we can consider it as a pretty good score given the shortage of data to us.

### 2.2.2 Ordinary Least Squares

Now we will try to implement Multiple Linear Regression algorithm using Ordinary Least Squares, the simplest of all. Ordinary least squares (OLS) minimises the squared distances between the observed and the predicted dependent variable. So, we get the following results after implementing the model :

```
Train Data
Ordinary Least Squares rmse: 2.551463189617839
Test Data
Ordinary Least Squares rmse: 2.7433491590595773
```

So, as we can see from the above results that, our model gives a RMSE (Root Mean Square Value) of 2.55 for the train data and a RMSE value of 2.74 for the test data. Looking at the train and test error, We can say that our model seems like overfitting a little but it still it overfits less then KNN. Also, we can observe that Ordinary Least Squares perform a little but better then the KNN Regression model.

### 2.2.3 Ridge Regression

Ridge Regression essentially is an instance of Linear Regression with regularisation. Ridge regression is that it enforces the  $\beta$  coefficients to be lower, but it does not enforce them to be zero. That is, it will not get rid of irrelevant features but rather minimise their impact on the trained model. So, after we implement Ridge Regression on our data, we get the following results :

```
Train Data
Ridge rmse: 2.6021156562086447
Test Data
Ridge rmse: 2.669065175428975
```

So, we can see from the above results that, our model gives a RMSE (Root Mean Square Value) of 2.60 for the train data and a RMSE value of 2.66 for the test data. Looking at the train and test error, we can say that the model doesn't fit at all and that might be because of the regularisation term involve in the cost function. Also, we can claim that Ridge Regression performs better then all the algorithm implemented before.

### 2.2.4 Lasso Regression

Least absolute shrinkage and selection operator, abbreviated as LASSO, is an Linear Regression technique which also performs regularisation on variables in consideration. It sets the coefficients to zero thus reducing the errors completely. That is, It will get rid of irrelevant features completely. So, after we implement Lasso Regression on our data, we get the following results :

```
Train Data
Lasso rmse: 3.4620116434744213
Test Data
Lasso rmse: 3.1430230704614934
```

So, we can see from the above results that, our model gives a RMSE (Root Mean Square Value) of 3.46 for the train data and a RMSE value of 3.14 for the test data. Looking at the train and test error, we can say that the model seems to overfit because the difference in the train and test error is noticeable. Also, we see that, Lasso Regression performs worse than the other models that we have seen.

### **2.2.5 Support Vector Regression**

Support Vector Machine can be applied not only to classification problems but also to the case of regression. Still it contains all the main features that characterise maximum margin algorithm: a non-linear function is learned by linear learning machine mapping into high dimensional kernel induced feature space. So, after we implement Support Vector Regression on our data, we get the following results :

```
Train Data
Support Vector Regression rmse: 2.6206055428869166
Test Data
Support Vector Regression rmse: 2.6912727332945185
```

So, we can see from the above results that, our model gives a RMSE (Root Mean Square Value) of 2.62 for the train data and a RMSE value of 2.69 for the test data. Looking at the train and test error, we can say that the model does not overfit at all because the difference in the train and test error is very low. Also, we see that, Support Vector Regression performs better than the other models that we have seen except Ridge Regression.

### **2.2.6 Decision Tree Regression**

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. So, after we implement Decision Tree Regression on our data, we get the following results :

```
Train Data
Decision Tree rmse: 3.31035213395237
Test Data
Decision Tree rmse: 3.272149876154806
```

So, we can see from the above results that, our model gives a RMSE (Root Mean Square Value) of 3.31 for the train data and a RMSE value of 3.27 for the test data. Although, it seems like the model doesn't overfit on the train data, Decision Tree Regression does not give impressive results and our other linear algorithms gives much lower error.

### 2.2.7 Gradient Boosting Decision Tree Regression

Gradient boosting is a machine learning technique for regression and classification problems, which produces a prediction model in the form of an ensemble of weak prediction models, typically decision trees. It builds the model in a stage-wise fashion like other boosting methods do, and it generalises them by allowing optimisation of an arbitrary differentiable loss function. So, after we implement Gradient Boosting Decision Trees on our data, we get the following results :

```
Train Data
GBDT rmse: 3.31035213395237
Test Data
GBDT rmse: 3.272149876154806
```

So, we can see from the above results that, our model gives a RMSE (Root Mean Square Value) of 3.31 for the train data and a RMSE value of 3.27 for the test data. Although, it seems like the model doesn't overfit on the train data, Gradient Boosting Decision Tree Regression does not give impressive results and our other linear algorithms gives much lower error. Also, It can be noticed that, the result of Decision Trees and GBDT are exactly same.

## 2.2.8 Random Forest Regression

Random Forest Regression or Regression Trees are known to be very unstable, in other words, a small change in our data may drastically change your model. The Random Forest uses this instability as an advantage through bagging resulting in a very stable model. So, after we implement Random Forest Regression or Regression Trees on our data, we get the following results :

```
Train Data
Random Forest rmse: 1.1209862159378214
Test Data
Random Forest rmse: 2.7611642280442057
```

So, we can see from the above results that, our model gives a RMSE (Root Mean Square Value) of 1.12 for the train data and a RMSE value of 2.76 for the test data. Looking at the train and test error, We can say that our model seems like overfitting a lot. But still we can consider it gives a pretty good score on the test data. Also, it outperforms many algorithms that we have seen earlier.

# Chapter 3

## Conclusion

### 3.1 Model Evaluation

Now that we have a few models for predicting the target variable, we need to decide which one to choose. There are several criteria that exist for evaluating and comparing models. We can compare the models using any of the following criteria:

1. Predictive Performance
2. Interpretability
3. Computational Efficiency

In our case of Employee Absenteeism, the latter two, *Interpretability* and *Computation Efficiency*, do not hold much significance. Therefore we will use *Predictive performance* as the criteria to compare and evaluate models.

Predictive performance can be measured by comparing Predictions of the models with real values of the target variables, and calculating some average error measure.

#### 3.1.1 Root Mean Square Value

Root Mean Square Error (RMSE) is the standard deviation of the residuals (prediction errors). Residuals are a measure of how far from the regression line data points are, RMSE is a measure of how spread out these residuals are. In other words, it tells you how concentrated the data is around the line of best fit.

Also, Since the errors are squared before they are averaged, the RMSE gives a relatively high weight to large errors. So, RMSE becomes more useful when large errors are particularly undesirable. So, Root Mean Square value seems like a perfect choice for our problem at hand.

## 3.2 Model Selection

We saw that both models Random Forest along with SVR and Ridge Regression perform comparatively on RMSE (Root Mean Square Error) , Although Random Forest gives the best results on the test data but it is unstable and it overfits on the train data. So, It will be a wise decision to use either SVR or Ridge Regression for deployment. Model comparison table is given below.

|   | Model                     | rmse |
|---|---------------------------|------|
| 0 | KNN Regression            | 2.75 |
| 1 | Ordinary Least Square     | 2.74 |
| 2 | Ridge Regression          | 2.66 |
| 3 | Lasso Regression          | 3.14 |
| 4 | Support Vector Regression | 2.65 |
| 5 | Decision Trees            | 3.19 |
| 6 | GBDT                      | 3.19 |
| 7 | Random Forest             | 2.58 |

So, It is obvious from above model performance comparison table, both Random Forest, SVR and Ridge Regression perform comparatively on RMSE (Root Mean Square Error) and can be used for deployment.



## 3.3 Answer to the asked Questions

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

Looking at the Exploratory data analysis of the features, we observe and make following conclusion :

1. The rate of Absenteeism is maximum in *Season 3 : Winter* followed by *Season 1 : Summer* , *Season 4 : Spring*, *Season 2 : Autumn*.
2. Also, We can say that the '*Absenteeism rate*' is maximum in *Month 7 : July* followed by *Month 4 : April* , *Month 3 : March*, *Month 12 : December*, *Month 11 : November* , *Month 6 : June* , *Month 5 : May* etc.
3. Looking at the Bar plot of '*Absenteeism rate*' Vs '*Day of the week*', it can clearly be observed that the '*Absenteeism rate*' is maximum on the third day of the week i.e *Day 3 : Tuesday* followed by *Day 2 : Monday*, *Day 4 : Wednesday*. Also, the '*absenteeism rate*' is lowest on *Day 6 : Saturday* followed by *Day 5 : Friday*.
4. From the Bar plot of '*Absenteeism rate*' Vs '*Reason of absence*' we can observe that '*9 : Diseases of the circulatory system*' is the most frequent reason for the absence of the employees. The second most frequent reason given by the employees for their absence is '*2 : Neoplasms*' followed by '*6 : Diseases of the nervous system*', '*12: Diseases of the skin and subcutaneous tissue*', '*19 : Injury, poisoning and certain other consequences of external causes*' etc.
5. Looking at the Bar plot of '*Absenteeism rate*' Vs '*ID*' It can be observed that the absence rate is maximum for employee with *ID : 9*, followed by employees with *ID : 7,26,14,13, 36, 11* and *6*. Also, it can be observed that Employee with employee *ID : 4,8* and *35* never absents and are very much regular to work.

So, Now that we have understood the behaviour of Employee attributes against their Mean Absenteeism rate in the company, we can introduce following changes to reduce the number of Absenteeism :

1. Firstly, we can start by Increasing the employee morale, engagement, and commitment to the organisation.

2. We can set a certain threshold for minimum number of absence for employees during the workdays, and employees not meeting the criteria can be questioned.
3. As Absenteeism rate is maximum in Season of 'Winter' and in month of July, April and March respectively, We can issue special notices regarding the Absenteeism scenario around the company.
4. A Health care facility can be introduced in the company, so that the employees can have regular Medical check-ups to keep them fit and Working. Also, It would help with the company's reputation to have taken the responsibility of their Employees.
5. Also, we can introduce person to person phone calls if off sick and return to work interviews. This way, Employee would feel responsible for their action towards the company goal and achievements.
6. We can also we can come up with other ideas like : An incentive or conversion scheme for unused sick days.
7. Also, Strict action could be taken towards the employee with high absence rate in the workplace without any valid reason for the absence and Employees with no absence or a minimum absence can be rewarded with perks.
8. Lastly , we can apply performance policies to act at the root of the problem. In some cases, absence rate might be reduced by clear specification of employees' responsibilities and targets.

### **3.3.2 How much losses every month can we project in 2011 if same trend of absenteeism continues?**

Employee absence, whether caused by sickness or pressures outside of the workplace, can cost employers a large amount of money if not properly managed. Lowering absence levels across a business not only leads to a reduction in money being lost by the business, but also a happier, productive and content workforce.

To calculate Loss per month, We introduce the following formula :

$$\text{Loss} = \frac{\text{Work load average/day} * \text{Absenteeism time in hours}}{\text{Service Time}}$$

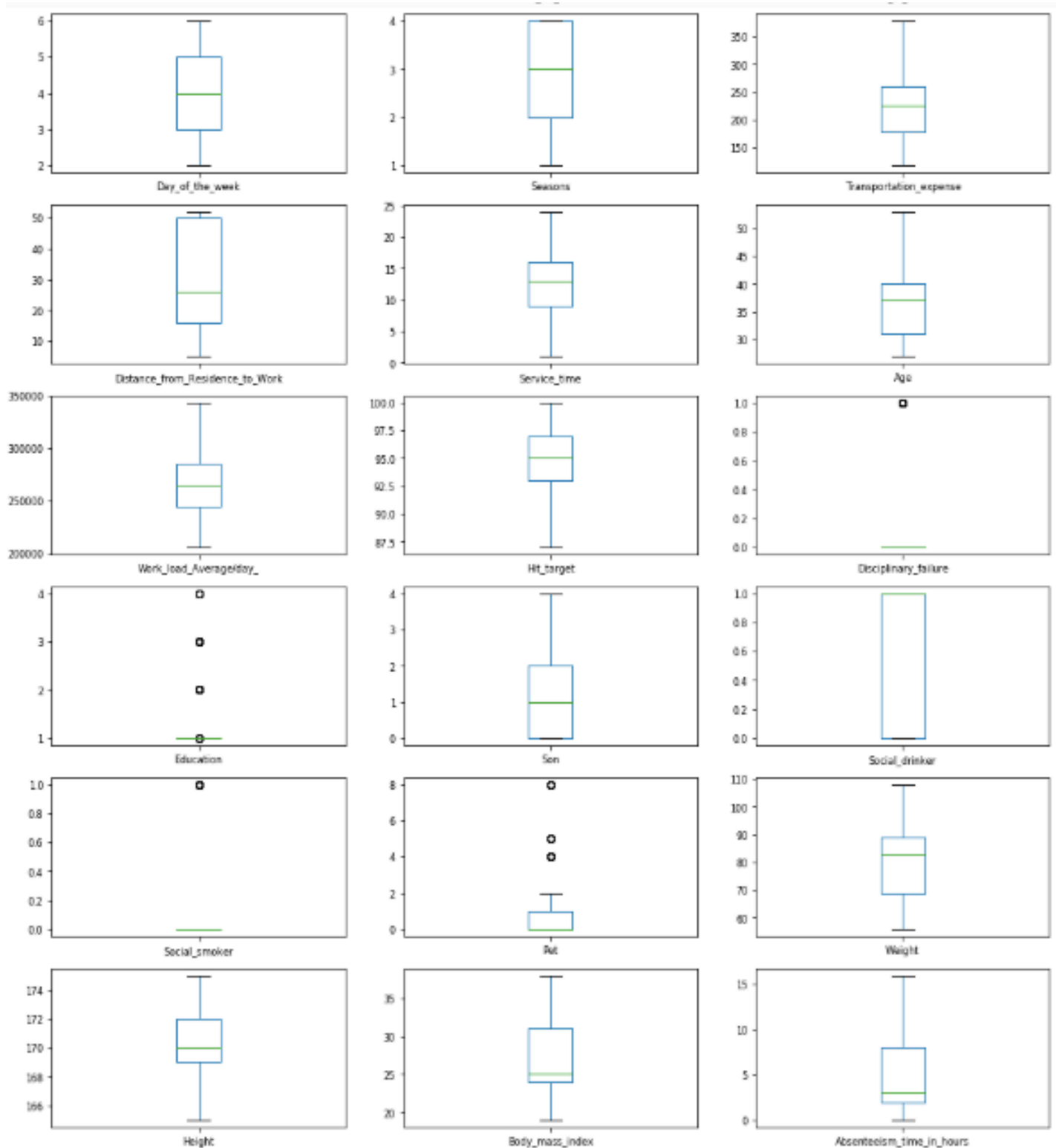
So, below chart represents loss per month and more likely, the same trend could follow in 2011 :

| <b>Month</b> | <b>Loss Per Month</b> |
|--------------|-----------------------|
| No Absent    | 0                     |
| January      | 4856265               |
| February     | 8003207               |
| March        | 10174728              |
| April        | 6350698               |
| May          | 6242207               |
| June         | 10254968              |
| July         | 11650434              |
| August       | 6400651               |
| September    | 4327941               |
| October      | 7227177               |
| November     | 6337747               |
| December     | 7692128               |

So, looking at the above table, we can say that the company incurred maximum loss in the month of July, followed by June and March etc. So, the same trend of loss can be expected in the year 2011 given the attributes of the employees.

## **Appendix A - Extra Figures**

*Fig A.1 Box Plots of Customer attributes after removal of outliers*



# References

Andrew Ng's Machine Learning course on Coursera (or, for more rigor, Stanford CS229).

An Introduction to Statistical Learning by Gareth James et al. Excellent reference for essential machine learning concepts, available free online.

All of Statistics: A Concise Course in Statistical Inference, by Larry Wasserman. Introductory text on statistics

MIT 18.05, Introduction to Probability and Statistics, taught by Jeremy Orloff and Jonathan Bloom. Provides intuition for probabilistic reasoning & statistical inference, which is invaluable for understanding how machines think, plan, and make decisions.













