

# **Project Report**

## **Employee Absenteeism**

***DEEKSHITHA R***

***8 MAY 2019***

# Contents

## 1. Introduction

1.1 Problem Statement . . . . .	3
1.2 Data . . . . .	3
1.3 Exploratory Data Analysis . . . . .	5

## 2. Methodology

2.1 Pre Processing . . . . .	6
2.1.1 Missing Value Analysis . . . . .	7
2.1.2 Outlier Analysis . . . . .	8
2.1.3 Feature Selection . . . . .	9
2.1.4 Feature Scaling . . . . .	10
2.1.5 Principal Component Analysis . . . . .	10
2.2 Modeling . . . . .	11
2.2.1 Decision Tree . . . . .	11
2.2.2 Random Forest . . . . .	11
2.2.3 Linear Regression . . . . .	11

## 3. Conclusion

3.1 Model Evaluation . . . . .	12
3.2 Model Selection . . . . .	12

## Appendix

Extra Figures . . . . .	17
-------------------------	----

## References

# Chapter 1

## Introduction

### 1.1 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 its 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?

### 1.2 Data

There are 21 variables in our data in which 20 are independent variables and 1 (Absenteeism time in hours) is dependent variable. Since our target variable is continuous in nature, this is a regression problem.

#### Variables Information:

1. Individual identification (ID)
2. Reason for absence (ICD) -

Absences attested by the **International Code of Diseases (ICD)** stratified into 21 categories (I to XXI) as follows:

- I. Certain infectious and parasitic diseases
- II. Neoplasms
- III. Diseases of the blood and blood-forming organs and certain disorders involving the immune mechanism
- IV. Endocrine, nutritional and metabolic diseases
- V. Mental and behavioral disorders
- VI. Diseases of the nervous system
- VII. Diseases of the eye and adnexa
- VIII. Diseases of the ear and mastoid process
- IX. Diseases of the circulatory system
- X. Diseases of the respiratory system
- XI. Diseases of the digestive system
- XII. Diseases of the skin and subcutaneous tissue
- XIII. Diseases of the musculoskeletal system and connective tissue

- XIV.** Diseases of the genitourinary system
- XV.** Pregnancy, childbirth and the puerperium
- XVI.** Certain conditions originating in the perinatal period
- XVII.** Congenital malformations, deformations and chromosomal abnormalities
- XVIII.** Symptoms, signs and abnormal clinical and laboratory findings, not elsewhere classified
- XIX.** Injury, poisoning and certain other consequences of external causes
- 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).

- 3.** Month of absence
- 4.** Day of the week (Monday (2), Tuesday (3), Wednesday (4), Thursday (5), Friday (6))
- 5.** Seasons (summer (1), autumn (2), winter (3), spring (4))
- 6.** Transportation expense
- 7.** Distance from Residence to Work (kilometers)
- 8.** Service time
- 9.** Age
- 10.** Work load Average/day
- 11.** Hit target
- 12.** Disciplinary failure (yes=1; no=0)
- 13.** Education (high school (1), graduate (2), postgraduate (3), master and doctor (4))
- 14.** Son (number of children)
- 15.** Social drinker (yes=1; no=0)
- 16.** Social smoker (yes=1; no=0)
- 17.** Pet (number of pet)
- 18.** Weight
- 19.** Height
- 20.** Body mass index
- 21.** Absenteeism time in hours (target)

### **1.3 Exploratory Data Analysis**

Exploratory Data Analysis (EDA) is an approach to analyzing data sets to summarize their main characteristics. In the given data set there are 21 variables and data types of all variables are either float64 or int64. There are 740 observations and 21 columns in our data set. Missing value is also present in our data.

**List of columns and their number of unique values -**

ID	36
Reason for absence	28
Month of absence	13
Day of the week	5
Seasons	4
Transportation expense	24
Distance from Residence to Work	25
Service time	18
Age	22
Work load Average/day	38
Hit target	13
Disciplinary failure	2
Education	4
Son	5
Social drinker	2
Social smoker	2
Pet	6
Weight	26
Height	14
Body mass index	17
Absenteeism time in hours	19

**From EDA we have concluded that there are 10 continuous variable and 11 categorical variable in nature.**

## **Chapter 2**

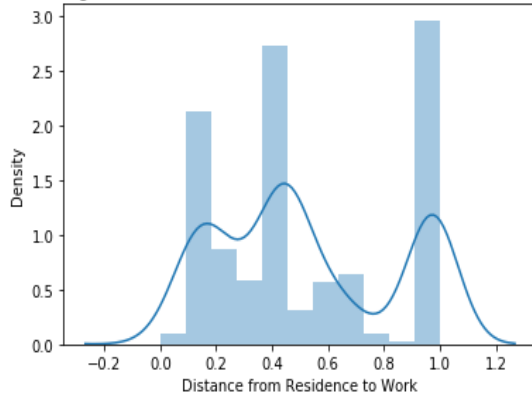
# Methodology

Before feeding the data to the model we need to clean the data and convert it to a proper format. It is the most crucial part of data science project we spend almost 80% of time in it.

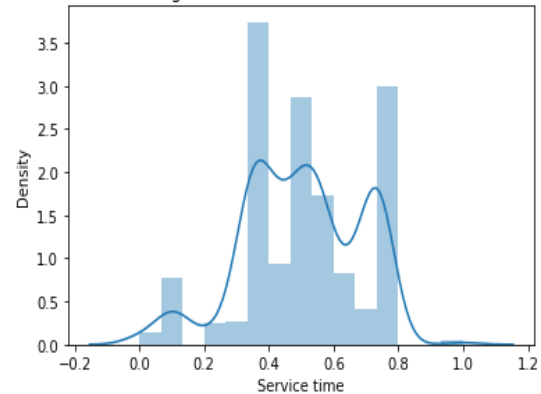
## 2.1 Pre Processing

Any predictive modeling requires that we look at the data before we start modeling. 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. To start this process we will first try and look at all the probability distributions of the variables. Most analysis like regression, require the data to be normally distributed. We can visualize that in a glance by looking at the probability distributions or probability density functions of the variable.

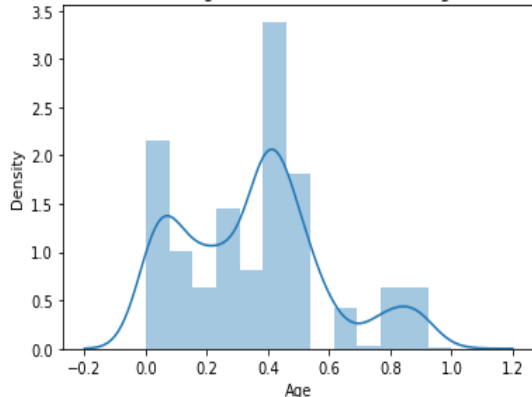
Checking Distribution for Variable Distance from Residence to Work



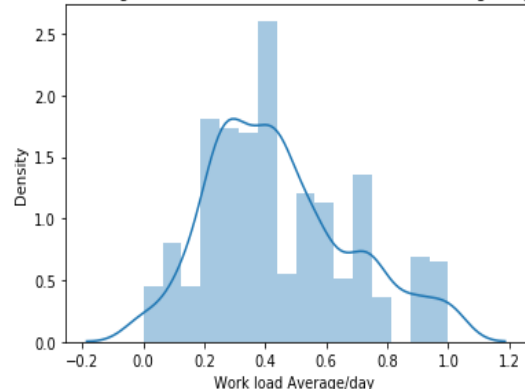
Checking Distribution for Variable Service time



Checking Distribution for Variable Age



Checking Distribution for Variable Work load Average/day

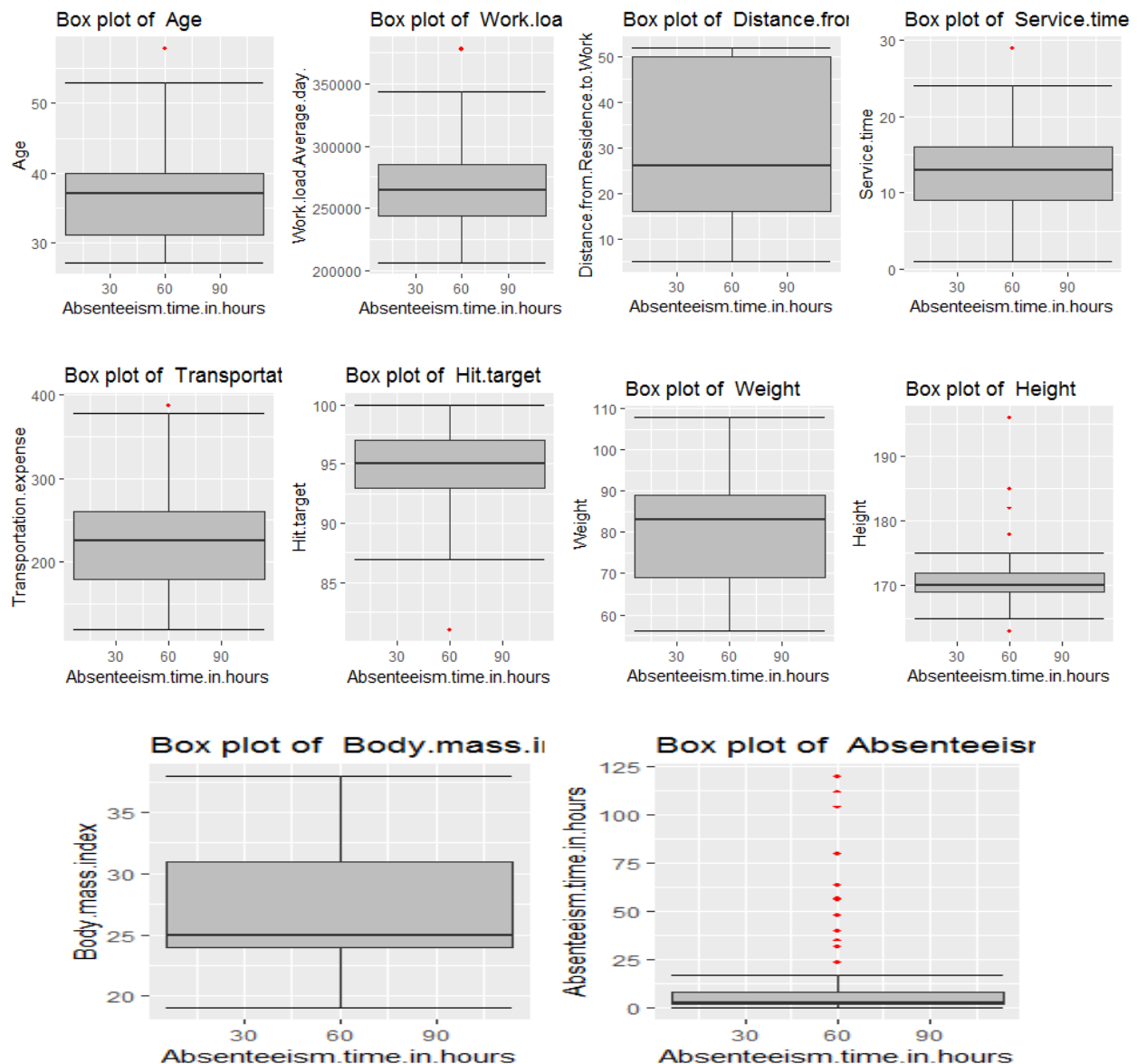




## 2.1.2 Outlier Analysis

We can clearly observe from these probability distributions that most of the variables are skewed. The skew in these distributions can be most likely explained by the presence of outliers and extreme values in the data. One of the other steps of pre-processing apart from checking for normality is the presence of outliers. In this case we use a classic approach of removing outliers. We visualize the outliers using boxplots.

In figure we have plotted the boxplots of the 11 predictor variables with respect to **Absenteeism time in hour**. A lot of useful inferences can be made from these plots. First as you can see, we have a lot of outliers and extreme values in each of the data set.

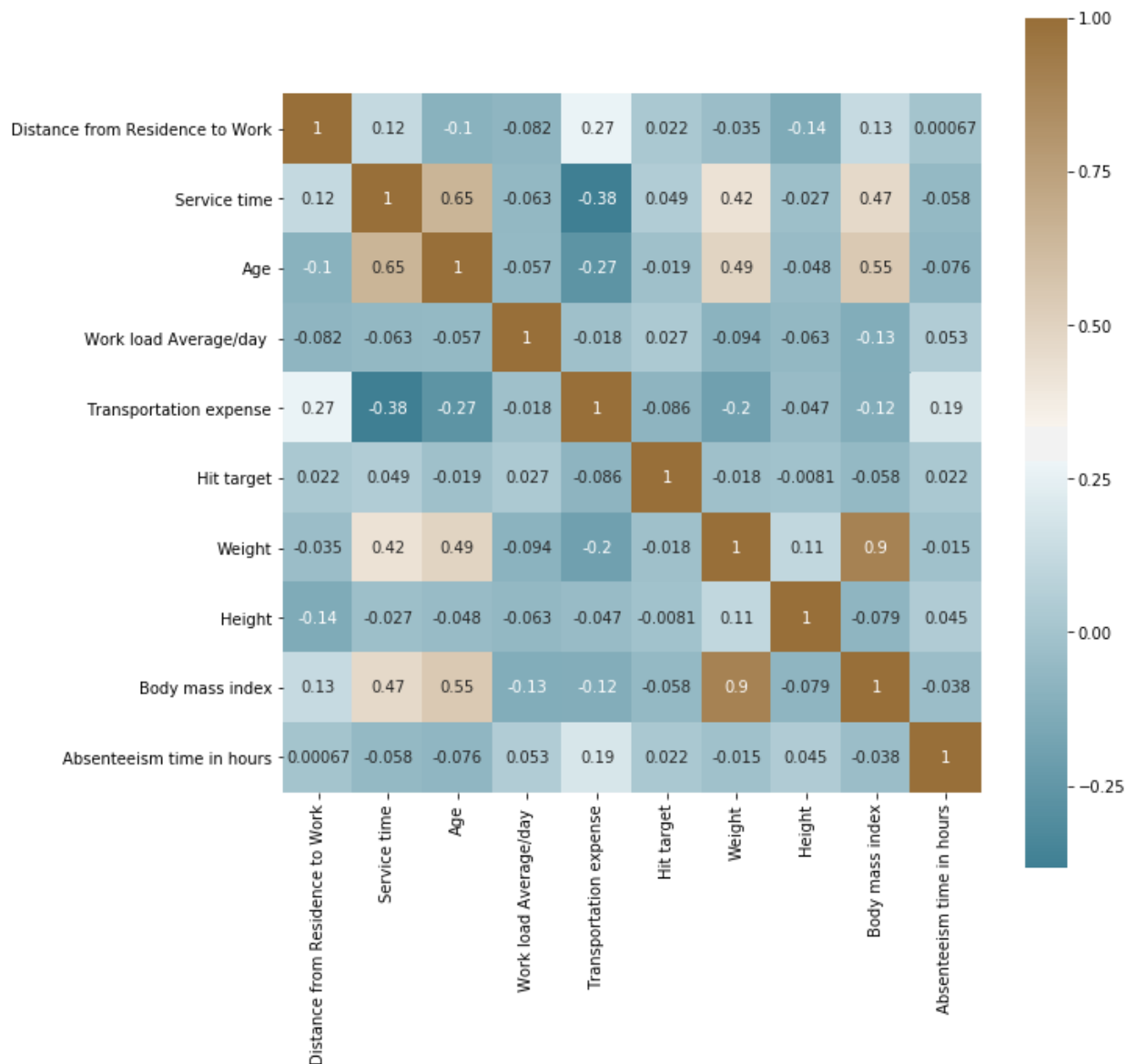


From the boxplot almost all the variables **except “Distance from residence to work”, “Weight” and “Body mass index”** consists of outliers. We have converted the outliers (data beyond minimum and maximum values) as NA i.e. missing values and fill them by **KNN** imputation method.



### 2.1.3 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. Selecting subset of relevant columns for the model construction is known as Feature Selection. We cannot use all the features because some features may be carrying the same information or irrelevant information which can increase overhead. To reduce overhead we adopt feature selection technique to extract meaningful features out of data. This in turn helps us to avoid the problem of multi collinearity. In this project we have selected **Correlation Analysis** for numerical variable and **ANOVA** (Analysis of variance) for categorical variable.



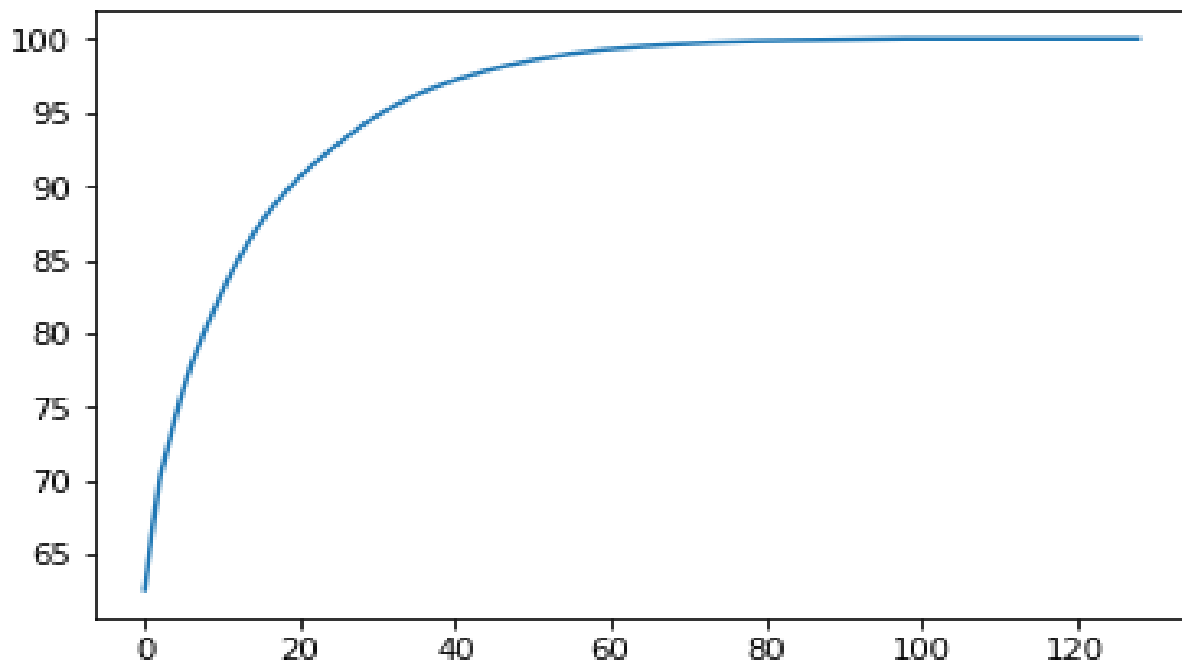
From correlation analysis we have found that **Weight** and **Body mass index** has high correlation ( $>0.7$ ), so we have excluded the **Weight** column.

### 2.2.4 Feature Scaling

**Feature scaling** is a method used to standardize the range of independent variables or features of data. In data processing, it is also known as data normalization and is generally performed during the data preprocessing step. Since the range of values of raw data varies widely, in some machine learning algorithms, objective functions will not work properly without normalization. For example, the majority of classifiers calculate the distance between two points by the Euclidean distance. If one of the features has a broad range of values, the distance will be governed by this particular feature. Therefore, the range of all features should be normalized so that each feature contributes approximately proportionately to the final distance. Since our data is not uniformly distributed we will use **Normalization** as Feature Scaling Method.

### 2.2.5 Principal Component Analysis

Principal component analysis is a method of extracting important variables (in form of components) from a large set of variables available in a data set. It extracts low dimensional set of features from a high dimensional data set with a motive to capture as much information as possible. With fewer variables, visualization also becomes much more meaningful. PCA is more useful when dealing with 3 or higher dimensional data. After creating dummy variable of categorical variables the shape of our data became 107 columns and 714 observations, this high number of columns leads to bad accuracy.



We have applied PCA algorithm on our data and from the above graph we have concluded that 45 variables out of 107 explains more than 95% of data. So we have selected only those 45 variables to feed our models.

## 2.2 Modeling

After a thorough preprocessing we will be using some regression models on our processed data to predict the target variable. Following are the models which we have built –

### 2.2.1 Decision Tree

A decision tree is a decision support tool that uses a tree-like graph or model of decisions and their possible consequences, including chance event outcomes, resource costs, and utility. Each branch connects nodes with “and” and multiple branches are connected by “or”. It can be used for classification and regression. It is a supervised machine learning algorithm. Accept continuous and categorical variables as independent variables. Extremely easy to understand by the business users. Split of decision tree is seen in the below tree. The RMSE value and  $R^2$  value for our project in R and Python are –

Decision Tree	R	PYTHON
RMSE Train	0.410	0.573
RMSE Test	0.363	0.568
$R^2$ Test	0.98	0.96

### 2.2.2 Random Forest

Random Forest is an ensemble technique that consists of many decision trees. The idea behind Random Forest is to build n number of trees to have more accuracy in dataset. It is called random forest as we are building n no. of trees randomly. In other words, to build the decision trees it selects randomly n no of variables and n no of observations to build each decision tree. It means to build each decision tree on random forest we are not going to use the same data. The RMSE value and  $R^2$  value for our project in R and Python are –

Random Forest	R	PYTHON
RMSE Train	0.369	0.033
RMSE Test	0.614	0.029
$R^2$ Test	0.96	0.99

### 2.2.3 Liner Regression

Linear Regression is one of the statistical methods of prediction. It is applicable only on continuous data. To build any model we have some assumptions to put on data and model. Here are the assumptions to the linear regression model.

Linear Regression	R	PYTHON
RMSE Train	0.001	4.34e-15
RMSE Test	0.001	4.14e-15
$R^2$ Test	0.99	1

# Chapter 3

## Conclusion

In this chapter we are going to evaluate our models, select the best model for our dataset and try to get answers of the asked questions.

### **3.1 Model Evaluation**

In the previous chapter we have seen the **Root Mean Square Error (RMSE)** and **R-Squared Value** of different models. **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. Whereas **R-squared** is a relative measure of fit, **RMSE** is an absolute measure of fit. As the square root of a variance, **RMSE** can be interpreted as the standard deviation of the unexplained variance, and has the useful property of being in the same units as the response variable. Lower values of **RMSE** and higher value of **R-Squared Value** indicate better fit.

### **3.2 Model Selection**

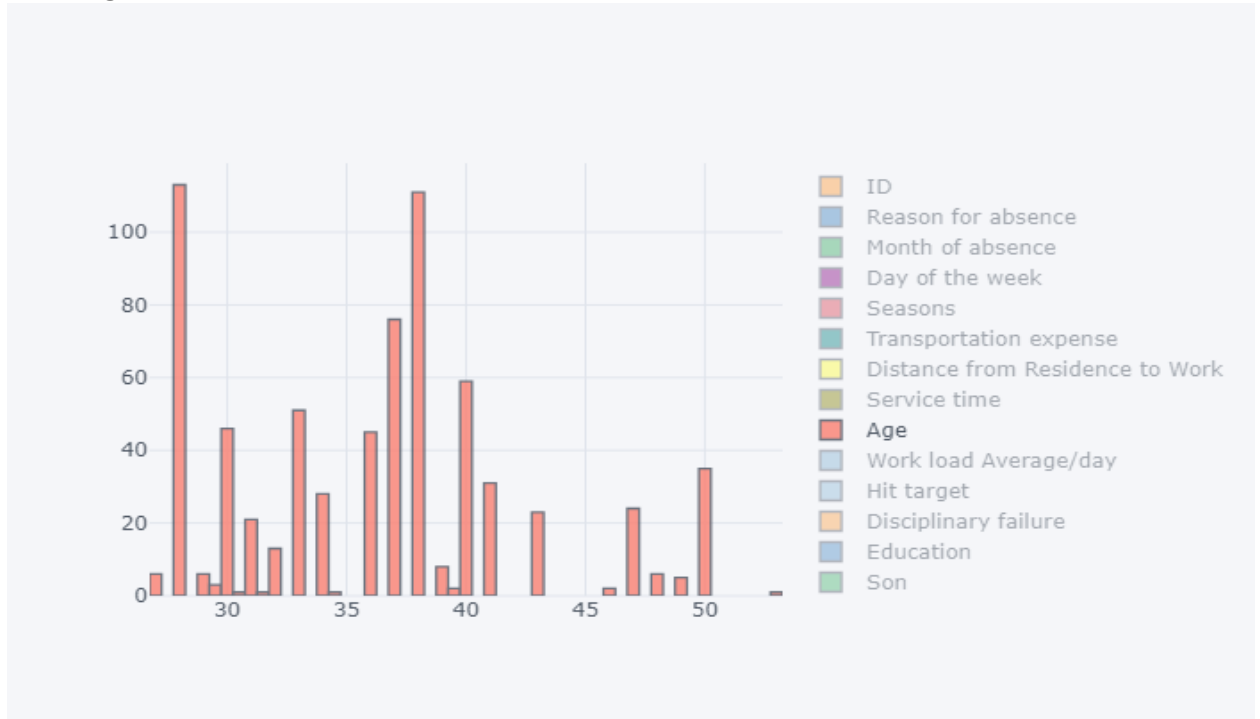
From the observation of all **RMSE Value** and **R-Squared Value** we have concluded that **Linear Regression Model** has minimum value of RMSE and it's **R-Squared Value** is also maximum (i.e. 1). The RMSE value of testing data and Training does not differ a lot this implies that it is not the case of overfitting.

# Appendix

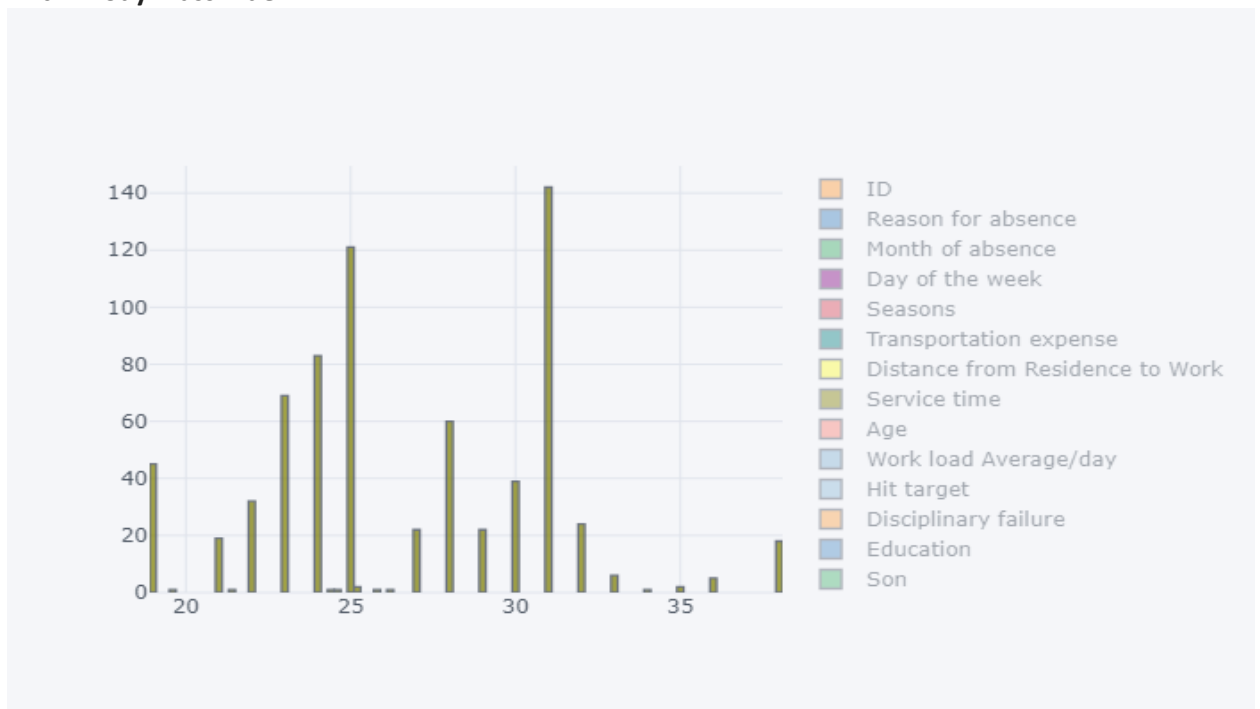
## Extra Figures

Relationship of our target variable (Absentee time in hour) with other variables.

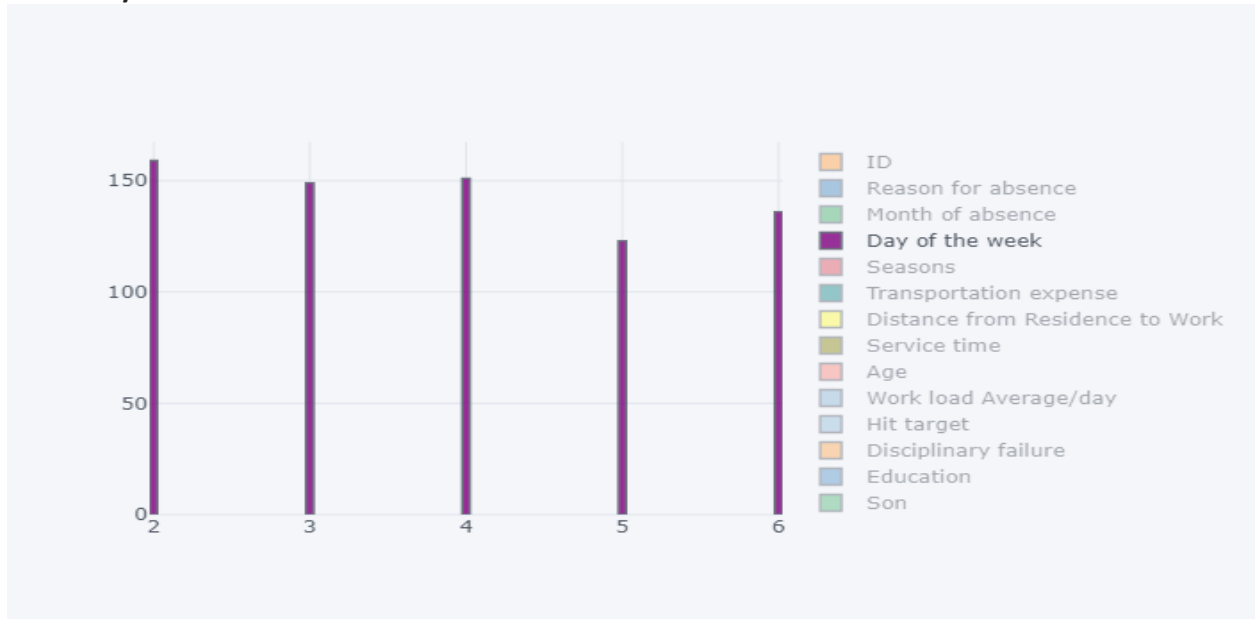
### 1. With “Age”



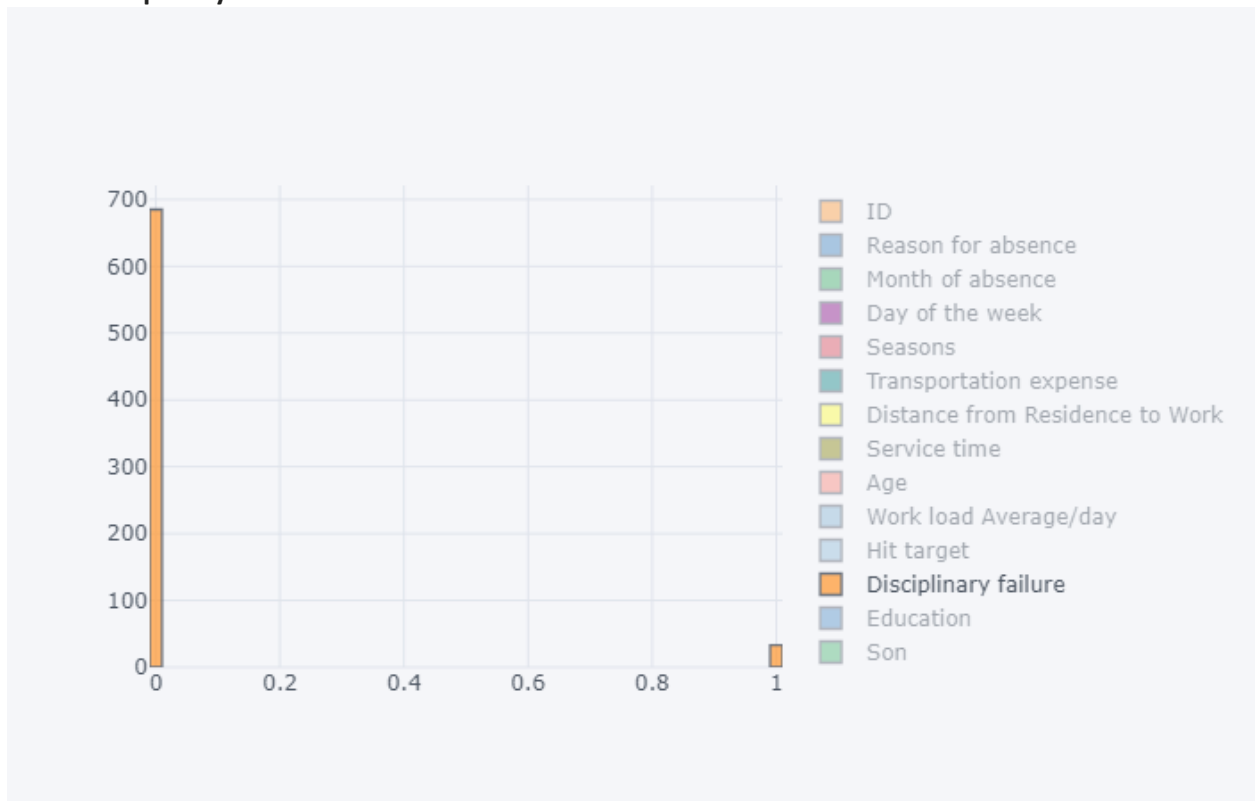
### 2. With “Body mass index”



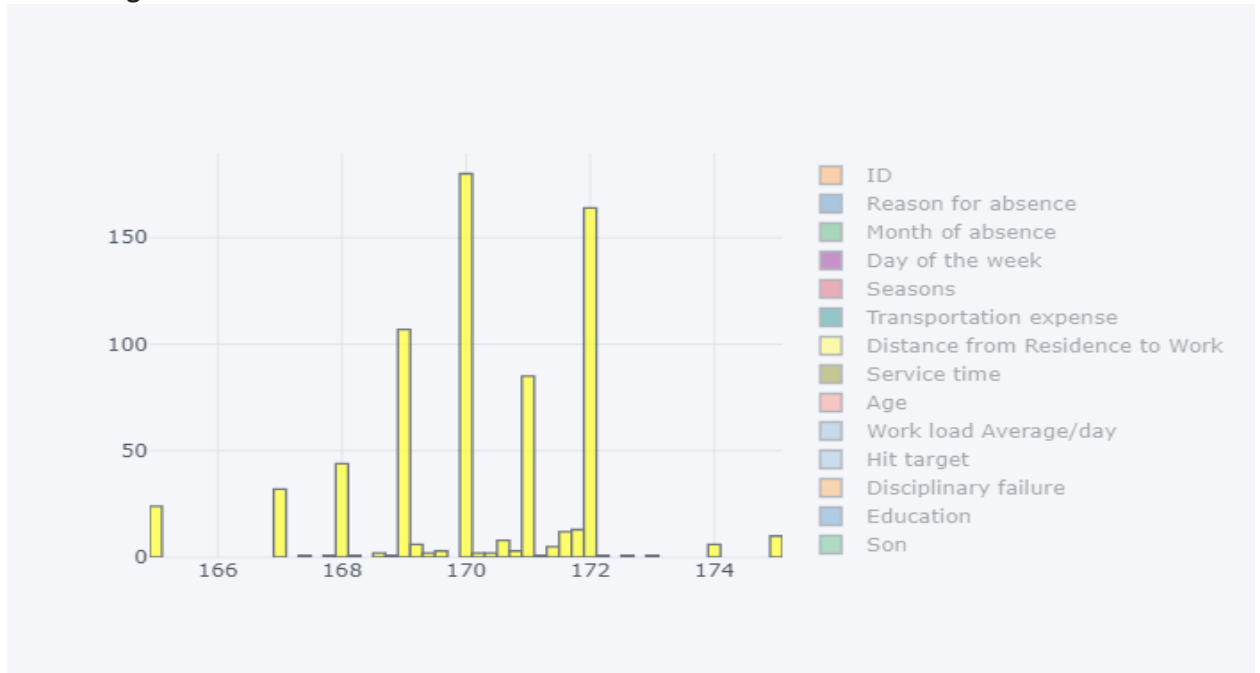
### 3. With "Day of week"



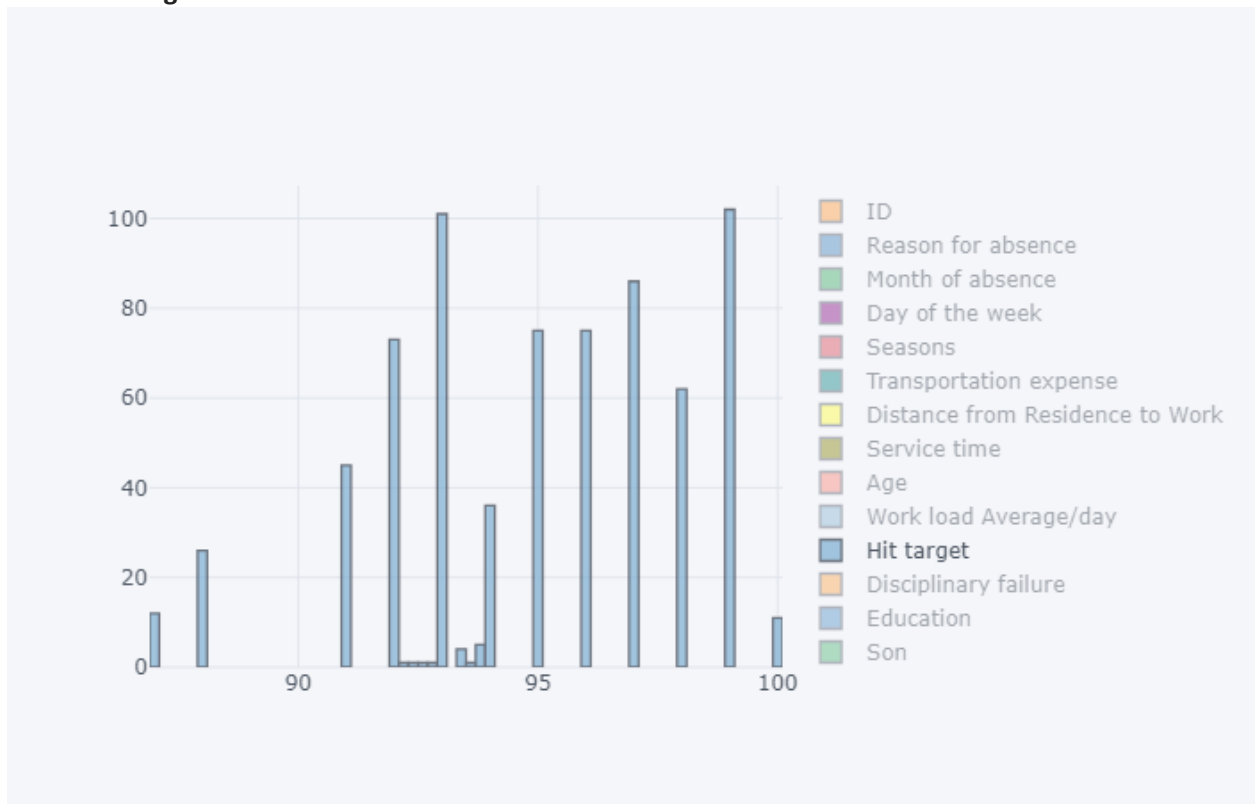
### 4. With "Disciplinary failure"



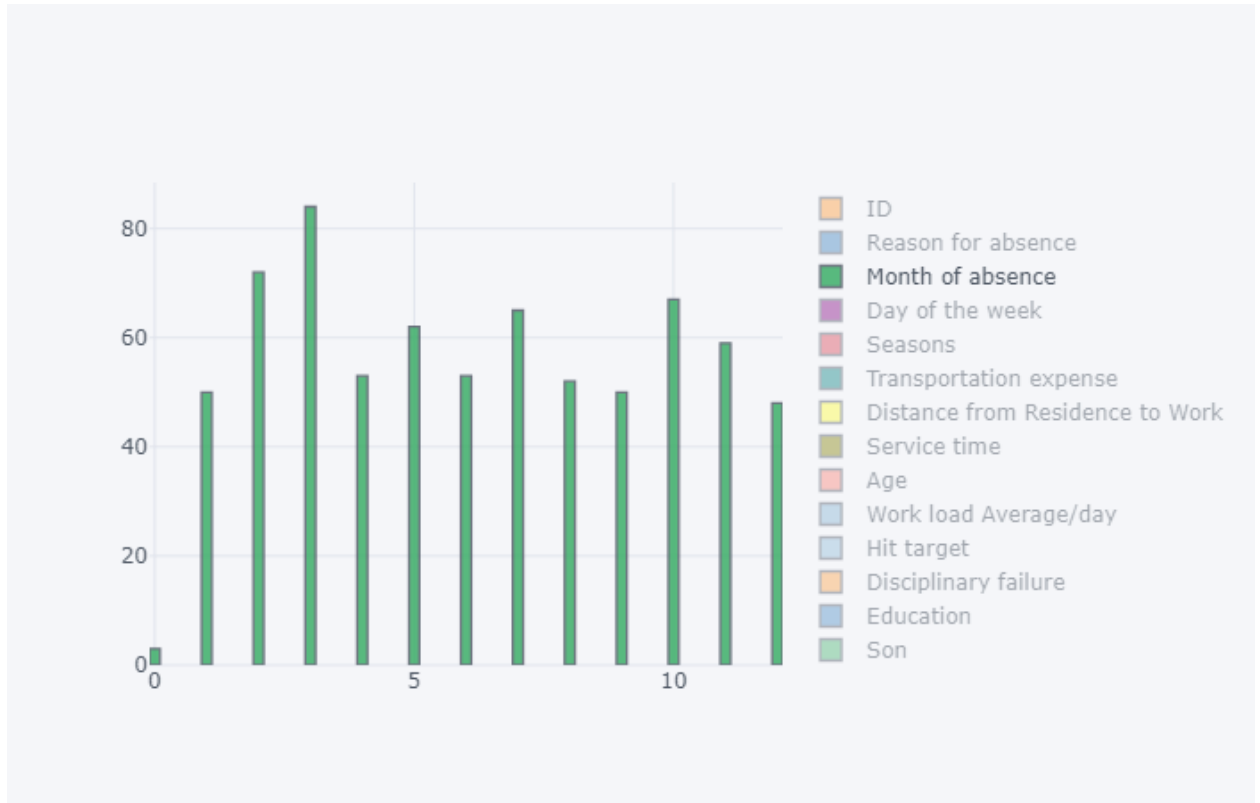
## 5. With "Height"



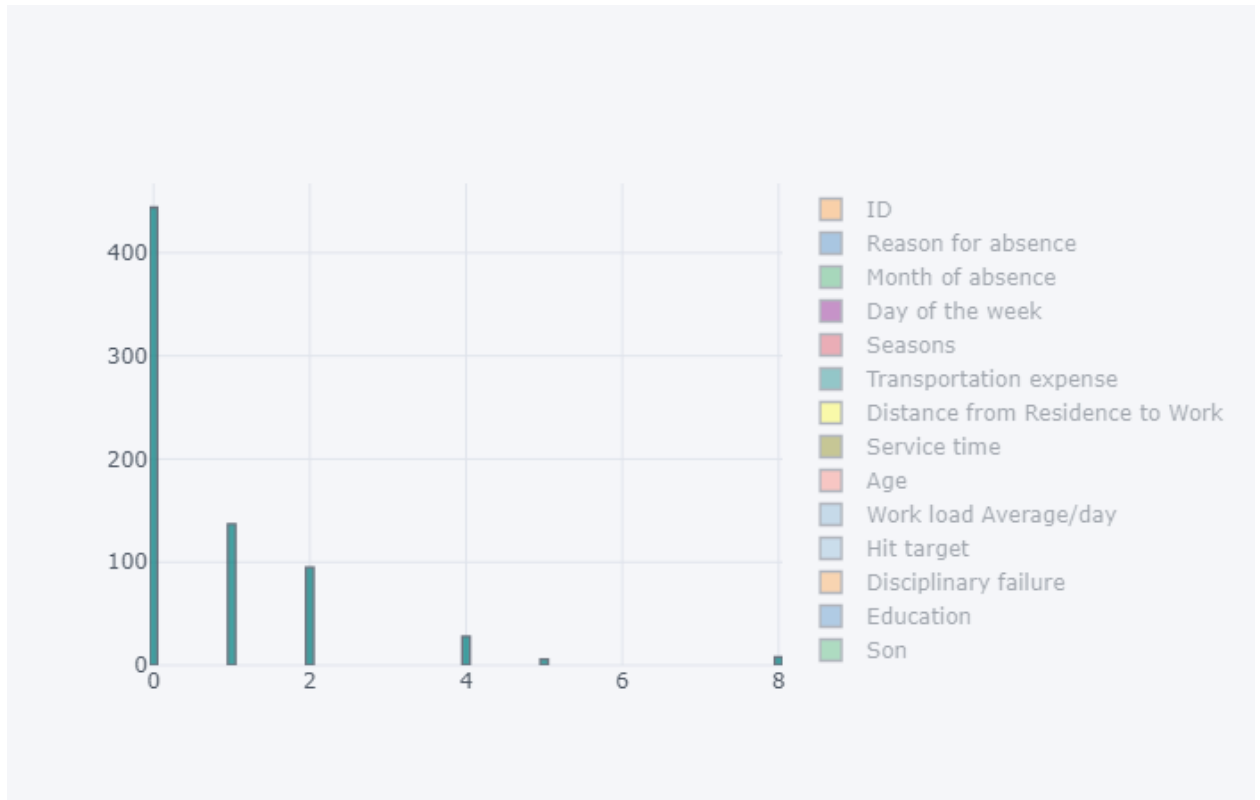
## 6. With "Hit target"



## 7. With "Month of absent"



## 8. With "Pet"

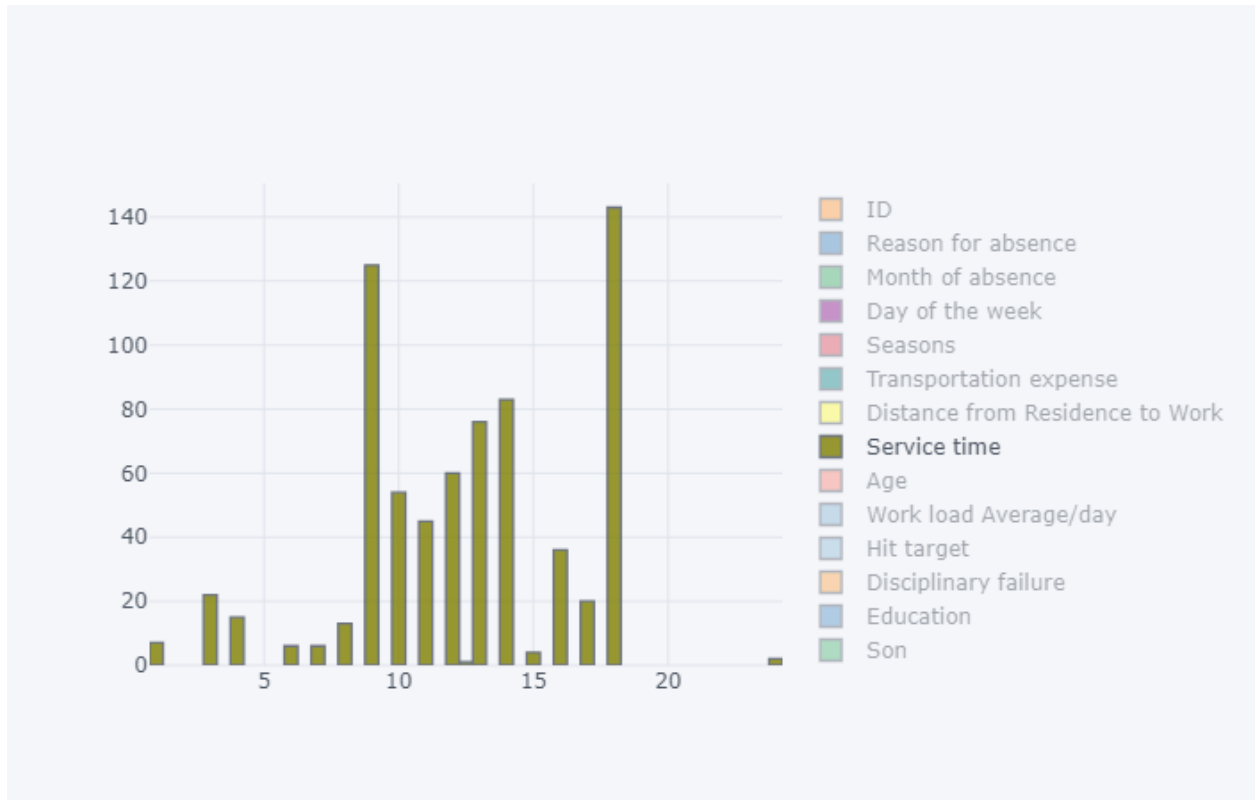




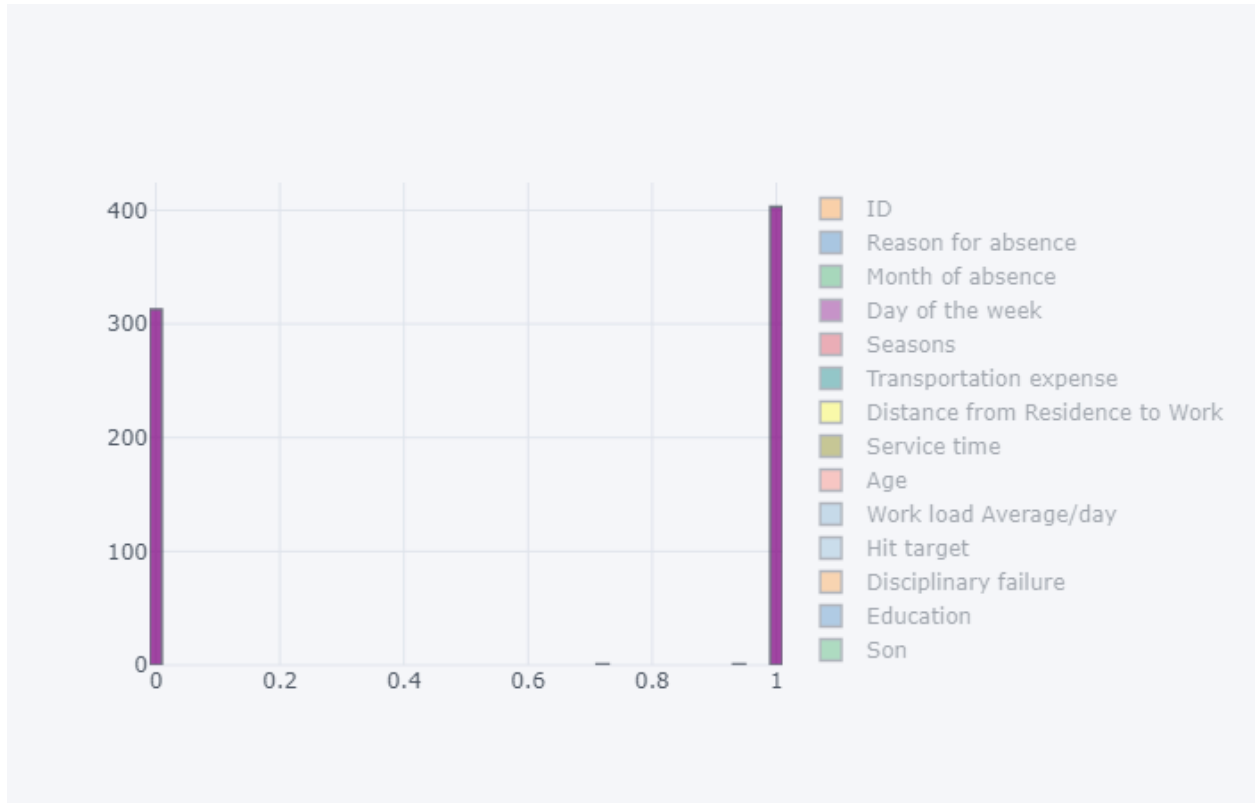
## 9. With "Seasons"



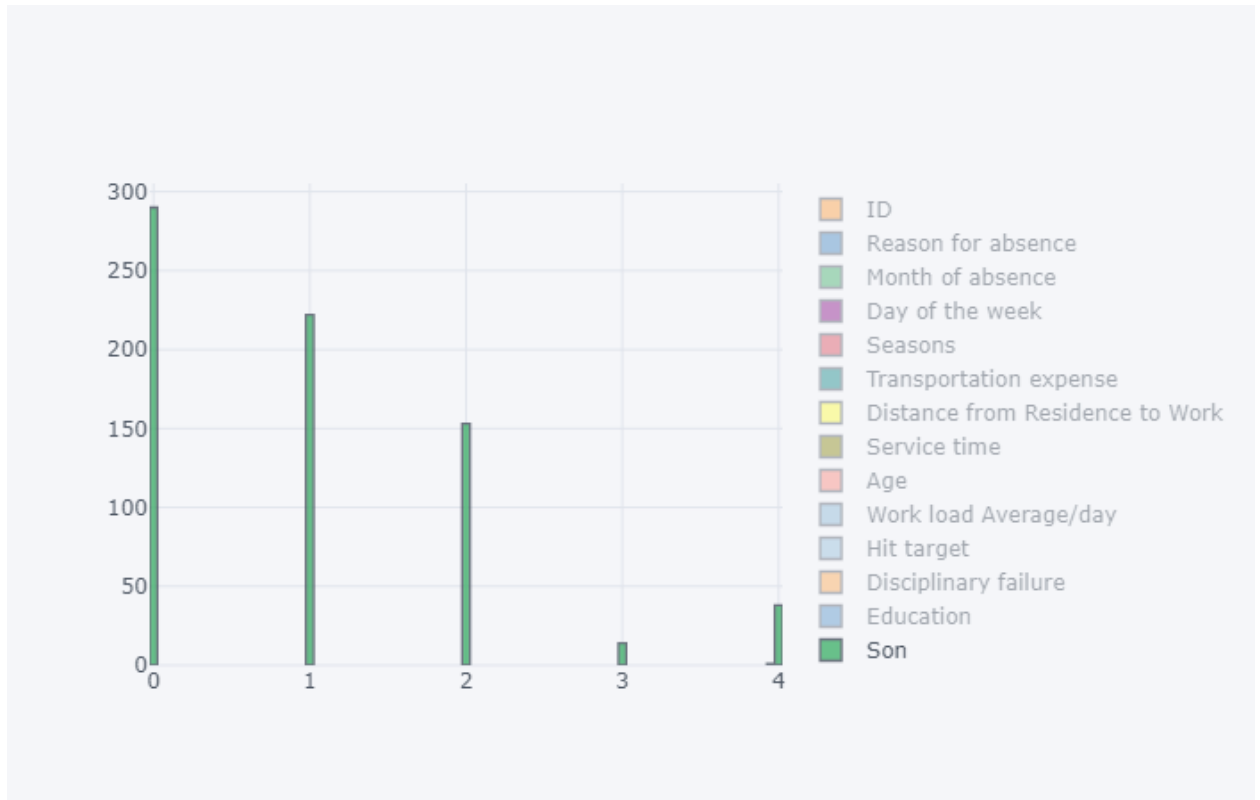
## 10. With "Service time"



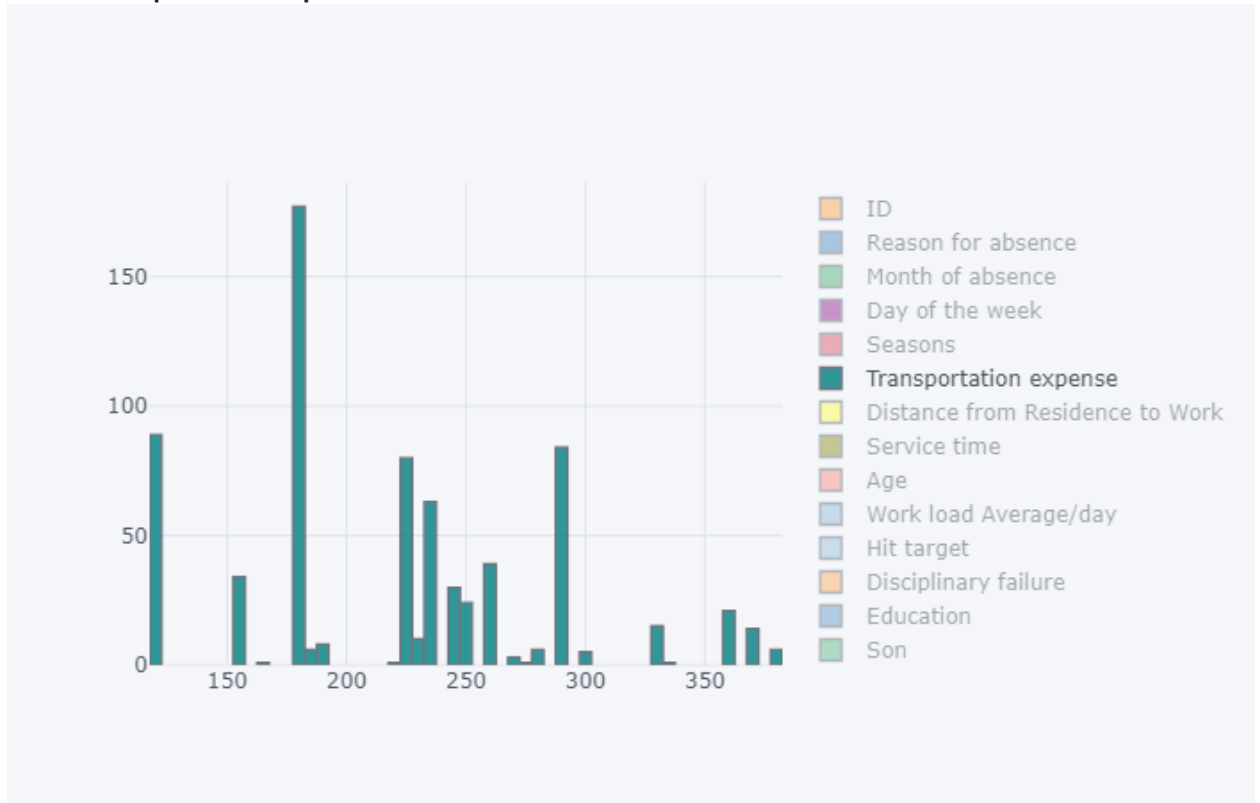
### 11. With "Social drinker"



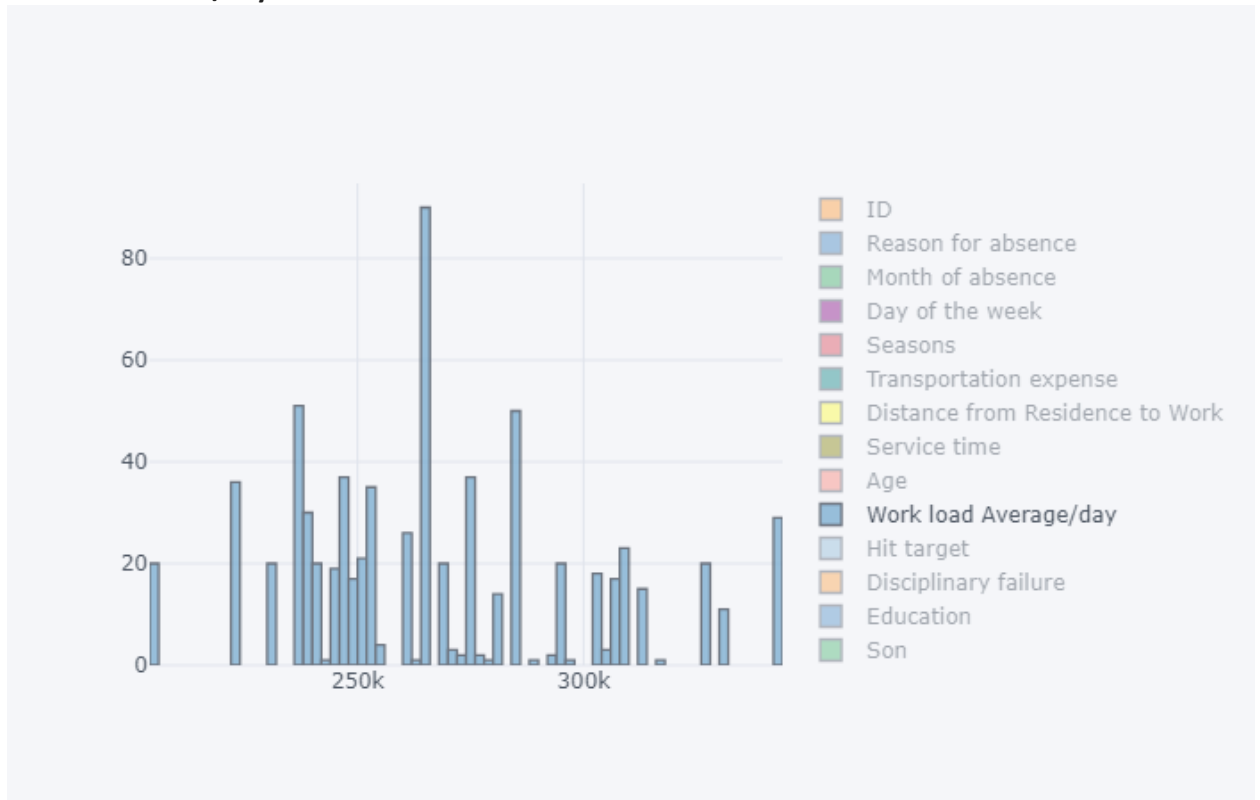
### 12. With "Son"



### 13. With "Transportation Expense"



### 14. With "Work load/day"



## References

1. For R-
  - <http://www.r-bloggers.com/>
  - <http://www.statmethods.net/>
  - <http://rfunction.com/>
2. For Data Cleaning and Model Development -  
<https://edwisor.com/career-data-scientist>
3. For Visualization –  
<https://www.udemy.com/python-for-data-science-and-machine-learning-bootcamp/>