# Data Mining Assignment

GROUP: P1DM6

# Business Context & Problem Statement

**Context**

Important features of Amazing Zone's business model are

- Various products sold online
- Retailer takes responsibility of the shipment of the product
- Products delivered within promised timeline increase happiness among the customer
- Customers are more inclined to purchase the products that has "Retailer fulfilled" tag associated with them

**Problem Statement**

The company Amazing Zone is facing issues with deliveries due to frequent absenteeism of the delivery personnel.

Due to heavy workload these delivery personnel are not able to perform their optimum limits, which results into their absenteeism from the work, which in turn increases the load on other delivery personnels.

# Opportunity Identifications

In order to improve productivity, we need to take a close look at absenteeism record of delivery personnel's, we must
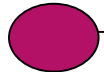
- **Identify the factors causing the absenteeism at work**
- **Suggest a model that can help to determine the absenteeism hours for an employee**
- **Suggest ways by which absenteeism can be reduced**

# Data Understanding & Preparations

Data Acquisition                                    Data Pre-processing and EDA                                    Model Creation
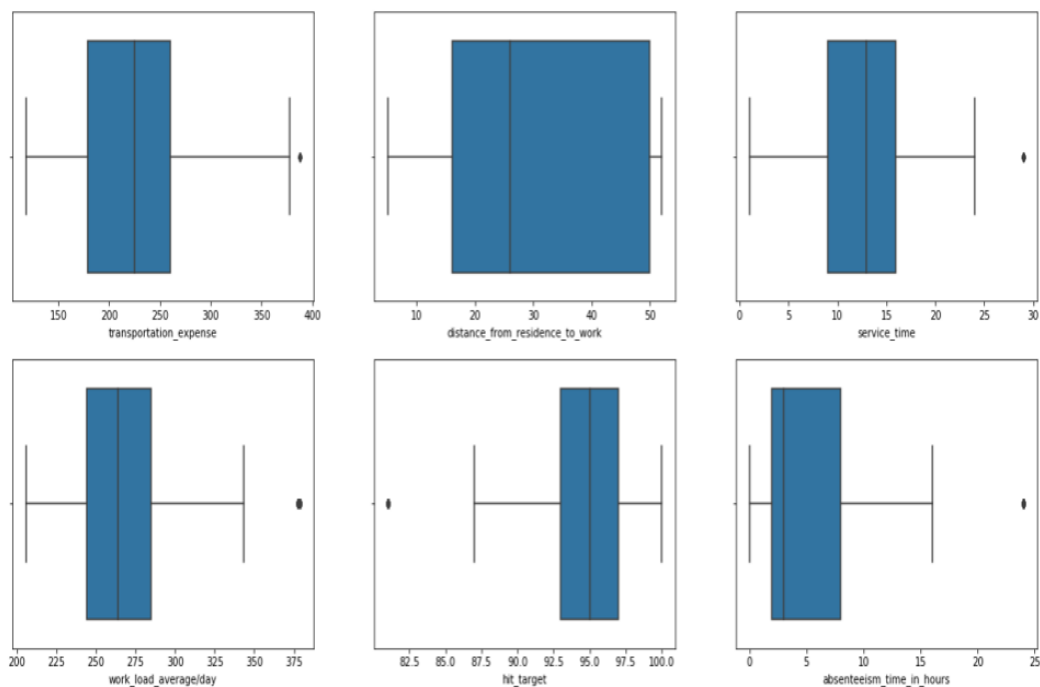
- **Data Acquisition** – Reading the data from CSV file and creating a dataframe
- **Data Pre-Processing and EDA**
    - Check for null values
    - Identify and remove outliers
    - Encoding categorical features
    - Scaling
    - EDA with bar graphs, heatmap and pairplots
- **Model Creation**
    - Splitting the data into training and test data
    - Linear Regression
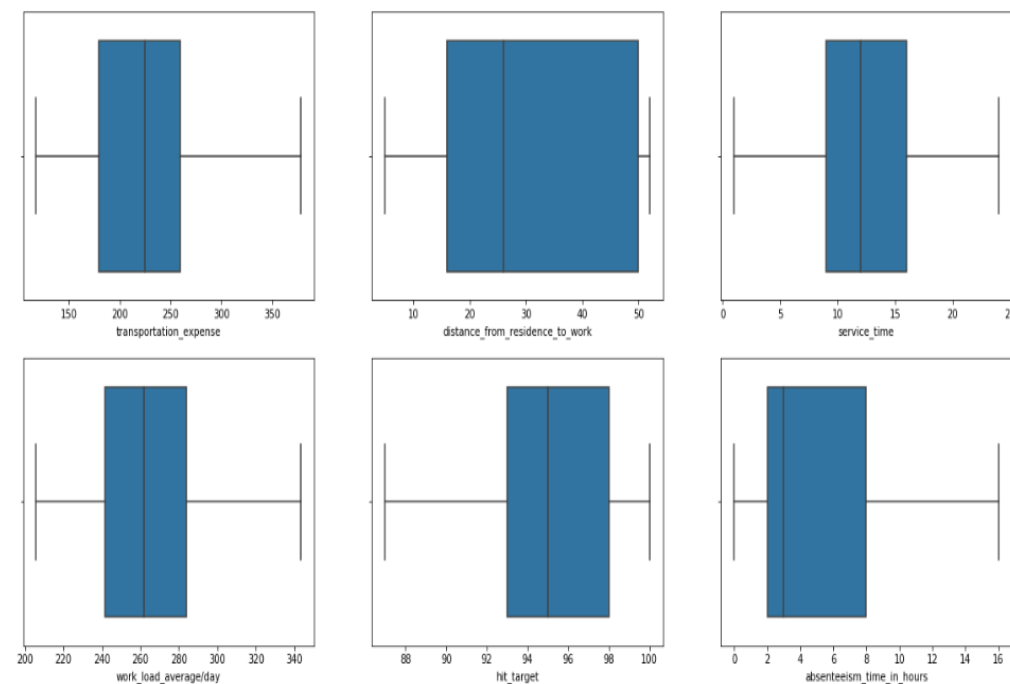    - Lasso
    - Model evaluation RMSE and R2

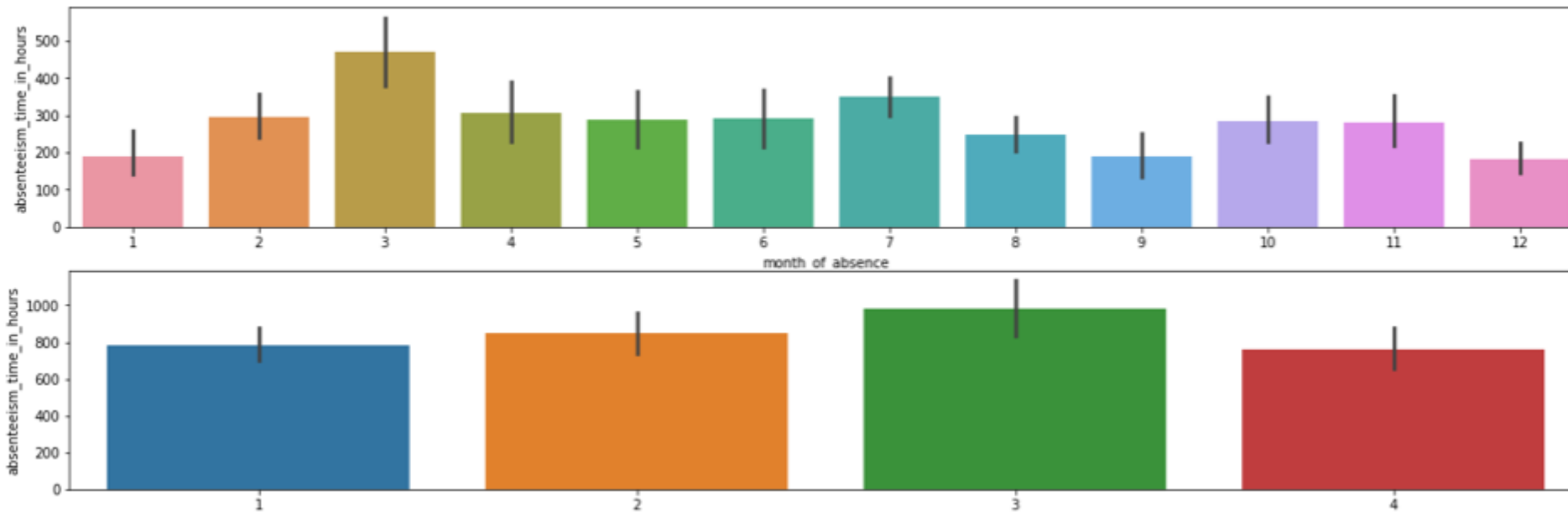# Data Understanding & Preparations
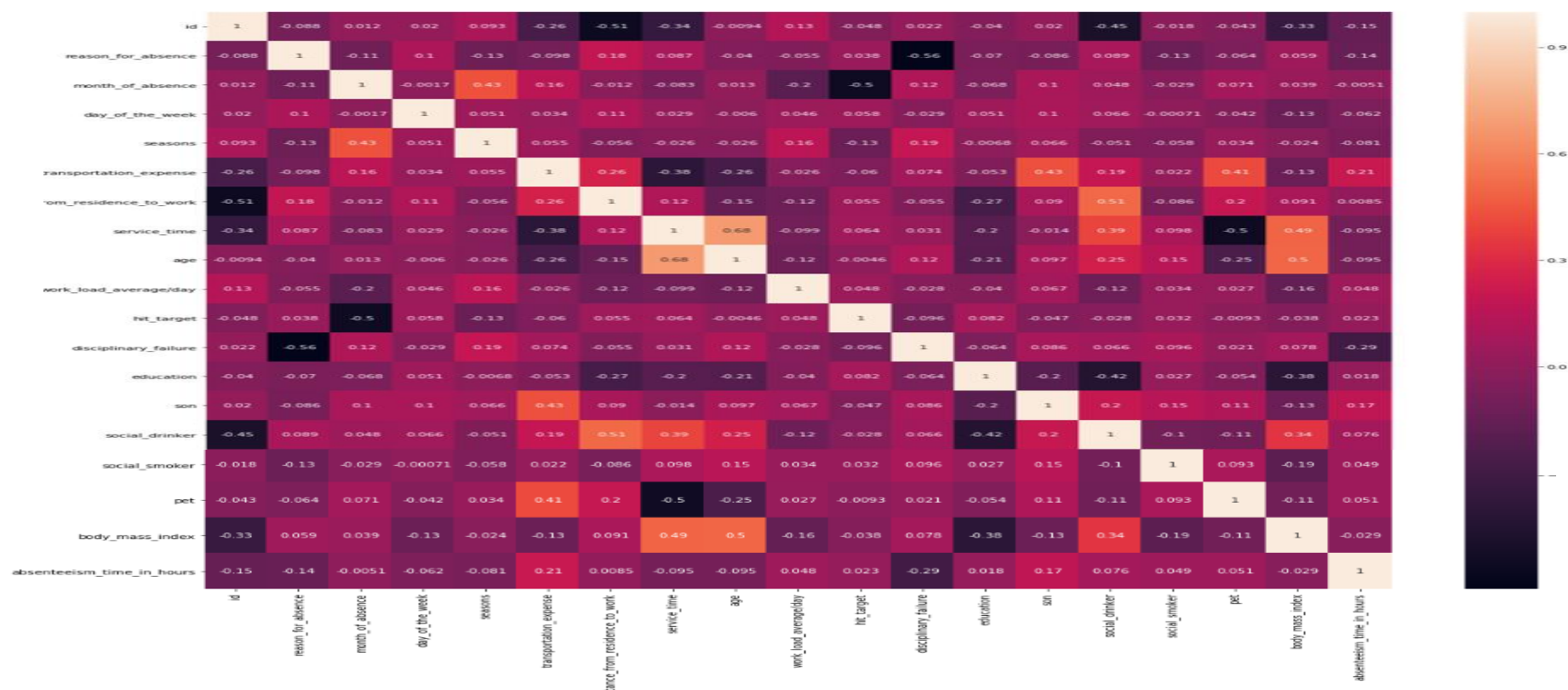
**Outlier Analysis**

# EDA

**Data Visualization with Bar Graph**

# EDA

**Heatmap**

# EDA

**Pair plots**

# Data Preparation

**Encoding Categorical Data**

One Hot encoding has been performed on categorical features of the data to make it viable for analysis.
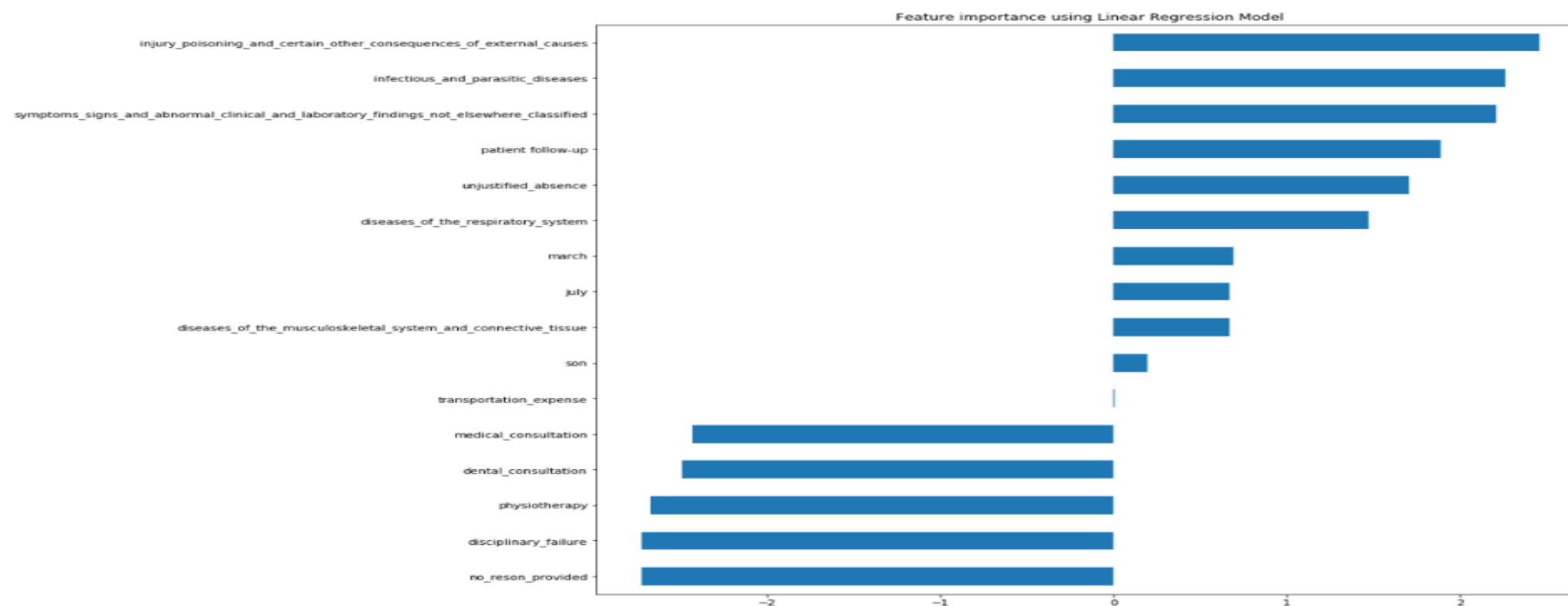
Encoding performed on below columns
- Seasons
- Reason for absence
- Month of absence
- Day of week

# Feature selection and model creation

- Calculating correlation coefficient for all the features
- Selecting features with correlation coefficient greater than +/- 0.1
- Splitting the data into dependent and independent features
- Implemented Linear regression model to predict absenteeism
- Implemented LASSO (Least Absolute Shrinkage Selector Operator) Model
- Calculated "Root Mean Square Error" to evaluate the accuracy of the model
- Calculated R2 Score to evaluate strength of correlation between dependent and independent variables
- Displayed linear regression coefficient & intercept
- Plotted feature importance with linear regression model

# Feature selection and model creation

# Suggestions for Improvement

- ▶ Provide periodic medical consultation to the delivery persons which will reduce the occurrences of ailments

- ▶ Provide physiotherapy sessions which will not only reduce stress but increase productivity

- ▶ Provide transportation or reimburse transportation expense

- ▶ Provide periodic dental consultation