

THANH PHONG, HO (PARKER)  
INSTITUTE OF DATA - 26 OCTOBER  
2021

# HYPOTHESIS TESTING PROJECT

---

## BUSINESS PROBLEM:

---

**The problem is that the work environment of today is more competitive, businesses want to:**

---

Reduce the cost of employee's absenteeism

---

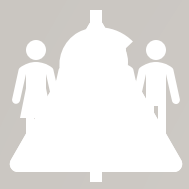
Achieve business goals.

---

---

Provide better healthcare for employees

---



---

## BUSINESS CONTEXT:

To understand whether employees with certain characteristics are expected to be away from work at some point in time or not.

We want to know how many working hours any employee could be away from work based on information like:

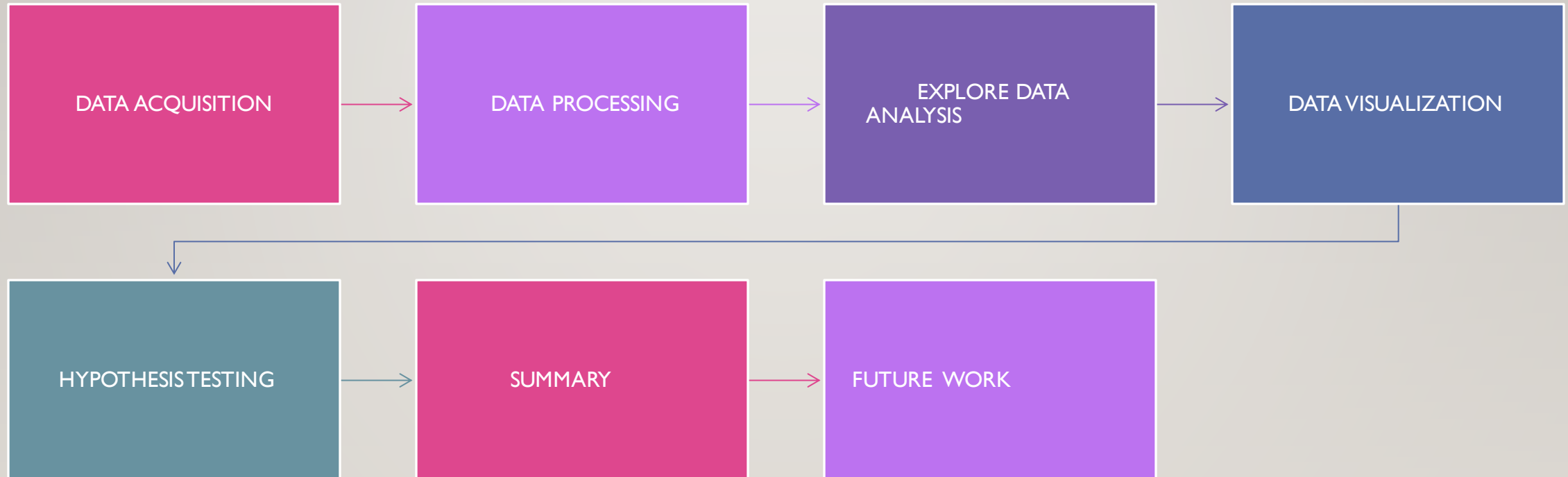
How far they live from their workplace.

Will the expense of transportation affect their absenteeism?

What are the reasons for absenteeism?

# DATA PIPELINE

---



# DATA OVERVIEW

Data was collected from 2015 July 06 to 2018 May 31

---

700  
observations

12 features

No missing  
values

1 object  
feature

10 integer  
features

1 float  
feature

# 12 FEATURES

ID

Reason for Absence

Date

Transportation Expense

Distance to Work

Age

Daily Workload Average

Body Mass Index

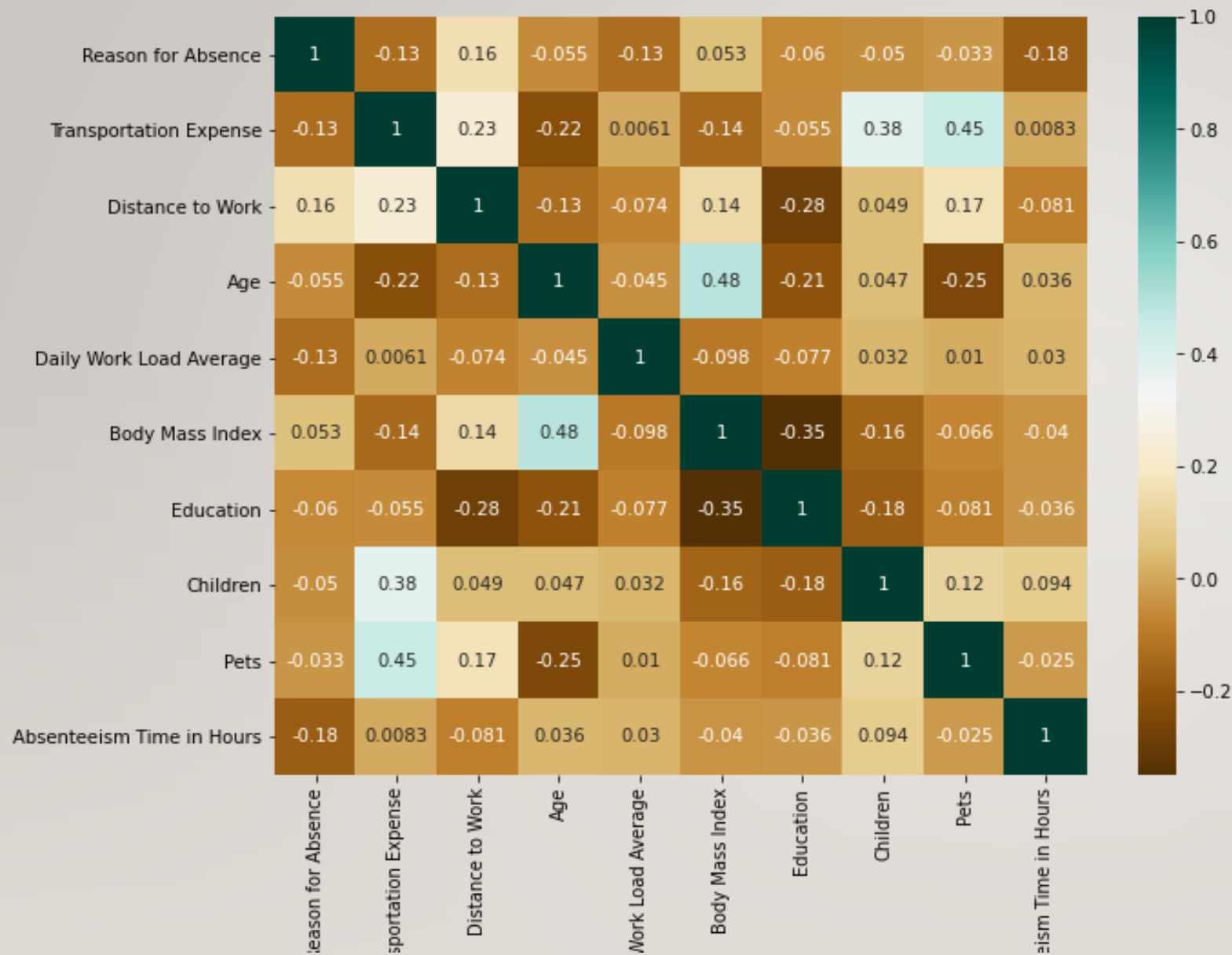
Children

Education

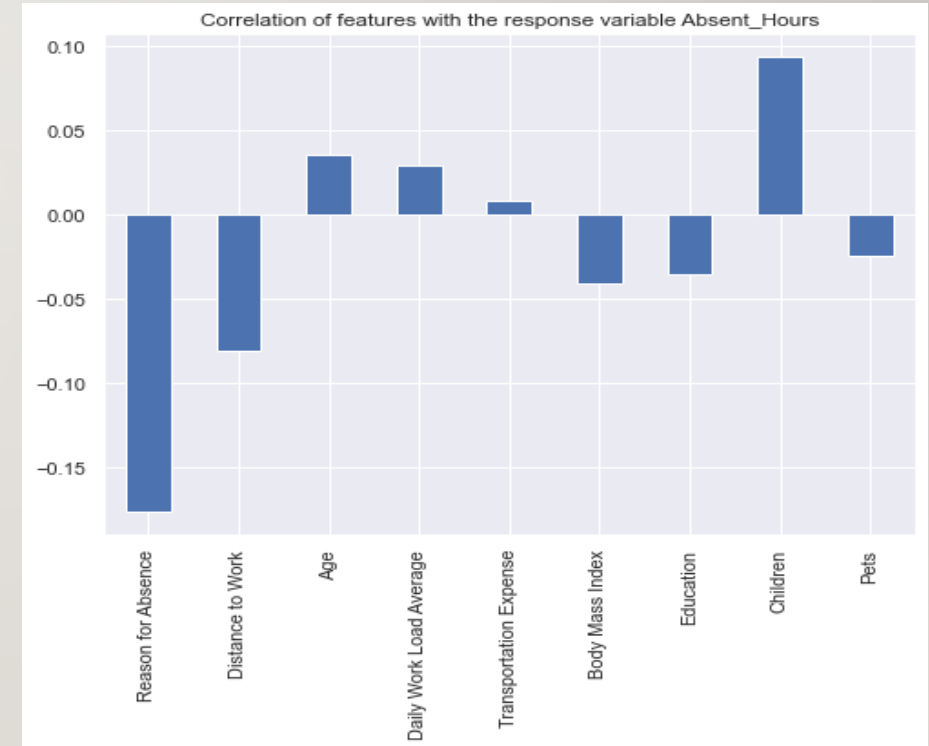
Pets

Absenteeism Time in Hours

Feature	Description
1	<b>ID:</b> Individual Identification
2	<b>Reason for Absence:</b> (Reasons 1 -21 are registered in the International Classification of Diseases (ICU), reasons 22-28 are not) 1: Certain infectious and parasitic diseases 2: Neoplasms 3: Disease of the blood and bleeding-forming organs and certain disorders involving in the immune mechanism 4: Endocrine, nutritional and metabolic diseases 5: Mental and behavioral disorders 6: Diseases of the nervous system 7: Diseases of the eye and adnexa 8: Diseases of the ear and mastoid process 9: Diseases of the circulatory system 10: Diseases of the respiratory system 11: Diseases of the digestive system 12: Diseases of the skin and subcutaneous tissue 13: Diseases of the musculoskeletal and connective tissue 14: Diseases of the genitourinary system 15: Pregnancy, childbirth and puerperium 16: Certain conditions originating in the perinatal period 17: Congenital malformations, deformations and chromosomal abnormalities 18: Symptoms, signs and abnormal clinical and laboratory findings, not elsewhere classified. 19: Injury, poisoning and certain other consequences of external causes 20: External causes of morbidity and mortality 21: Factors influencing health status and contact with health services 22: Patient follow-up 23: Medical consultation 24: Blood donation 25: Laboratory examination 26: Unjustified absence 27: Physiotherapy 28: Dental consultation
3	<b>Date:</b> date of absence
4	<b>Transportation Expense:</b> costs related to business travel such as fuel, parking and meals
5	<b>Distance to Work:</b> measured in kilometers
6	<b>Age:</b> years of age
7	<b>Daily Workload Average:</b> measured in minutes
8	<b>Body Mass Index</b>
9	<b>Education:</b> a categorical variable and presenting different levels of education
10	<b>Children:</b> number of children in the family
11	<b>Pets:</b> number of pets in the family
12	<b>Absenteeism Time in Hours:</b>



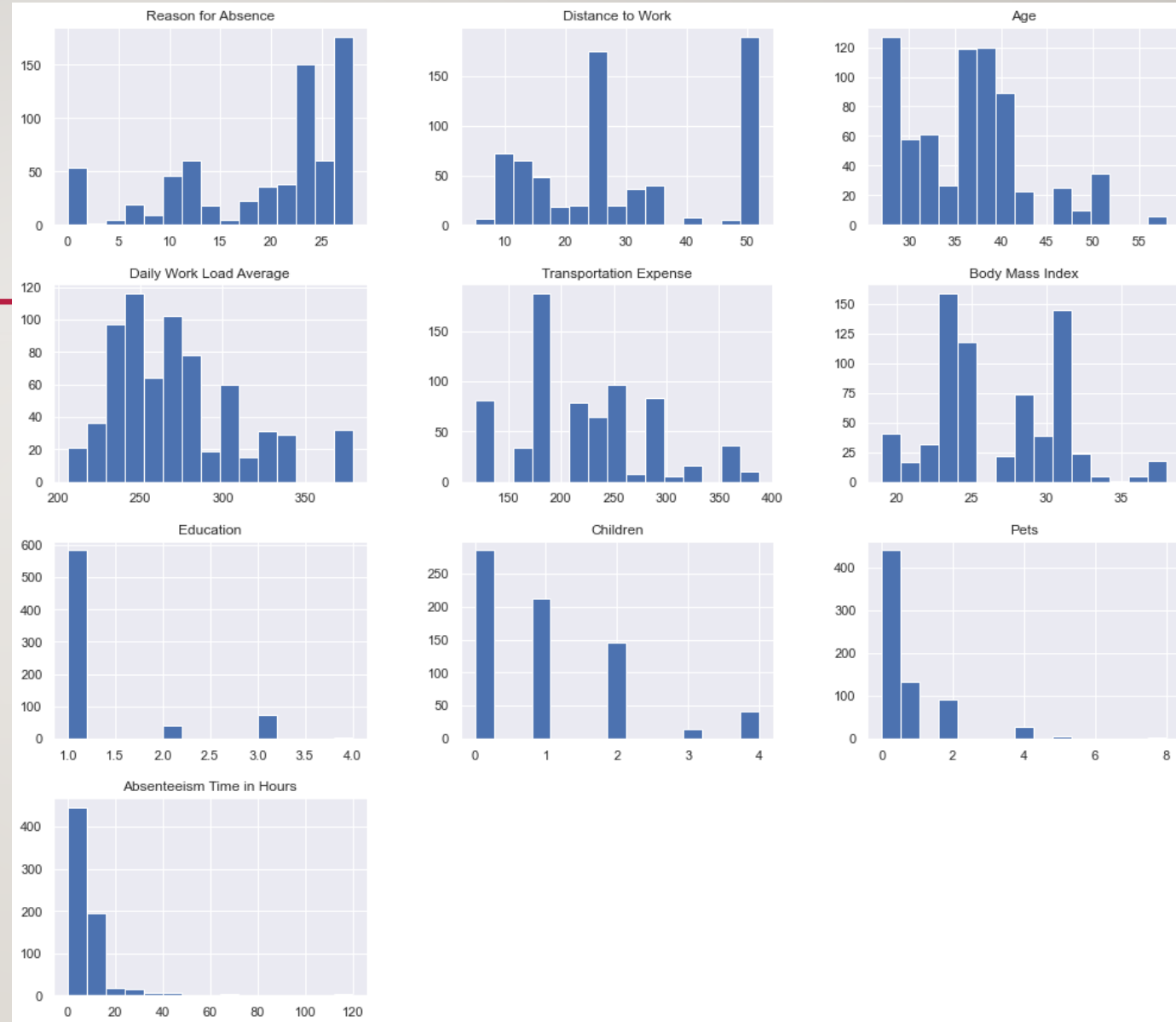
# DATA CORRELATION

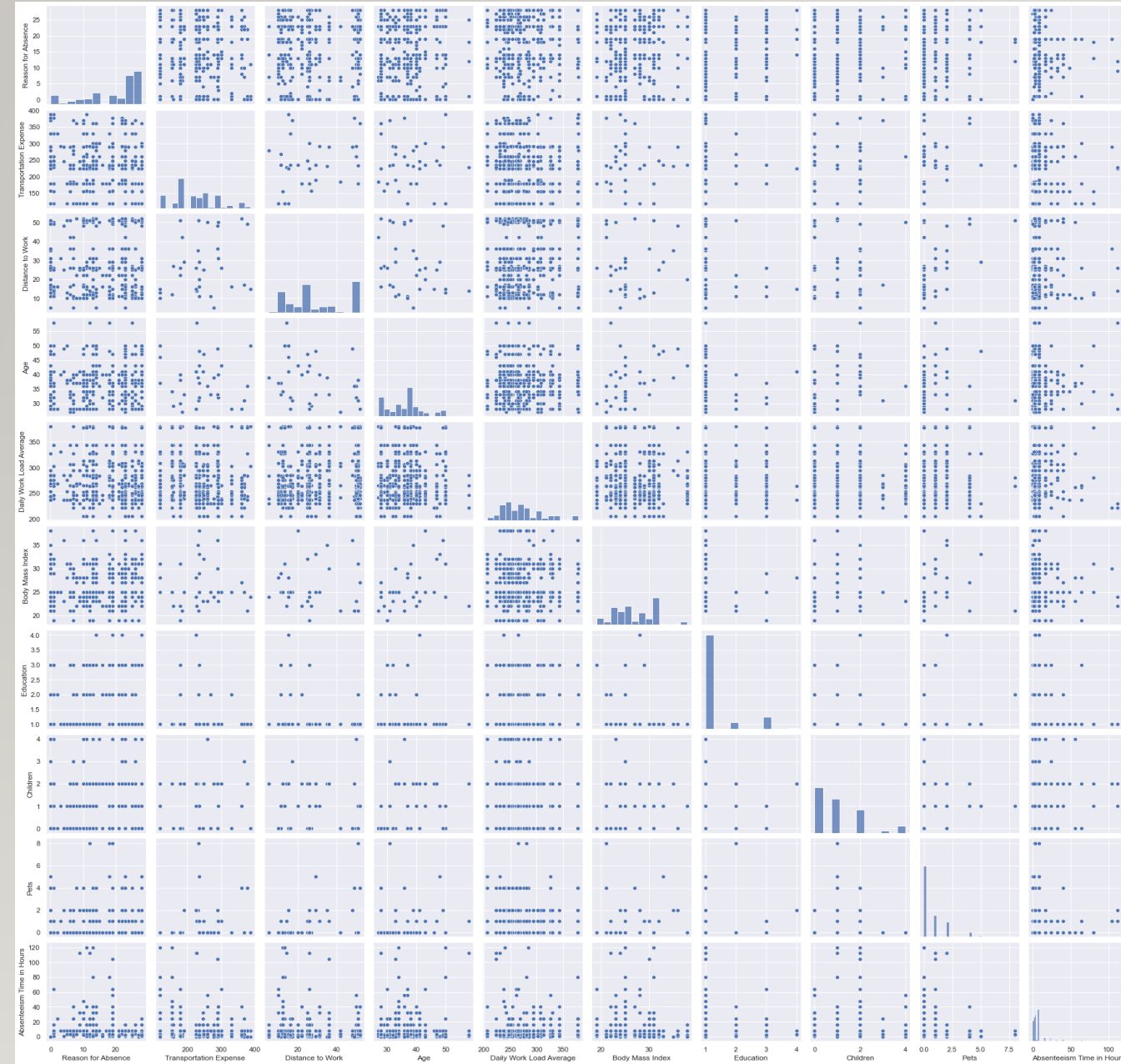




# DISTRIBUTION OF NUMERIC FEATURES

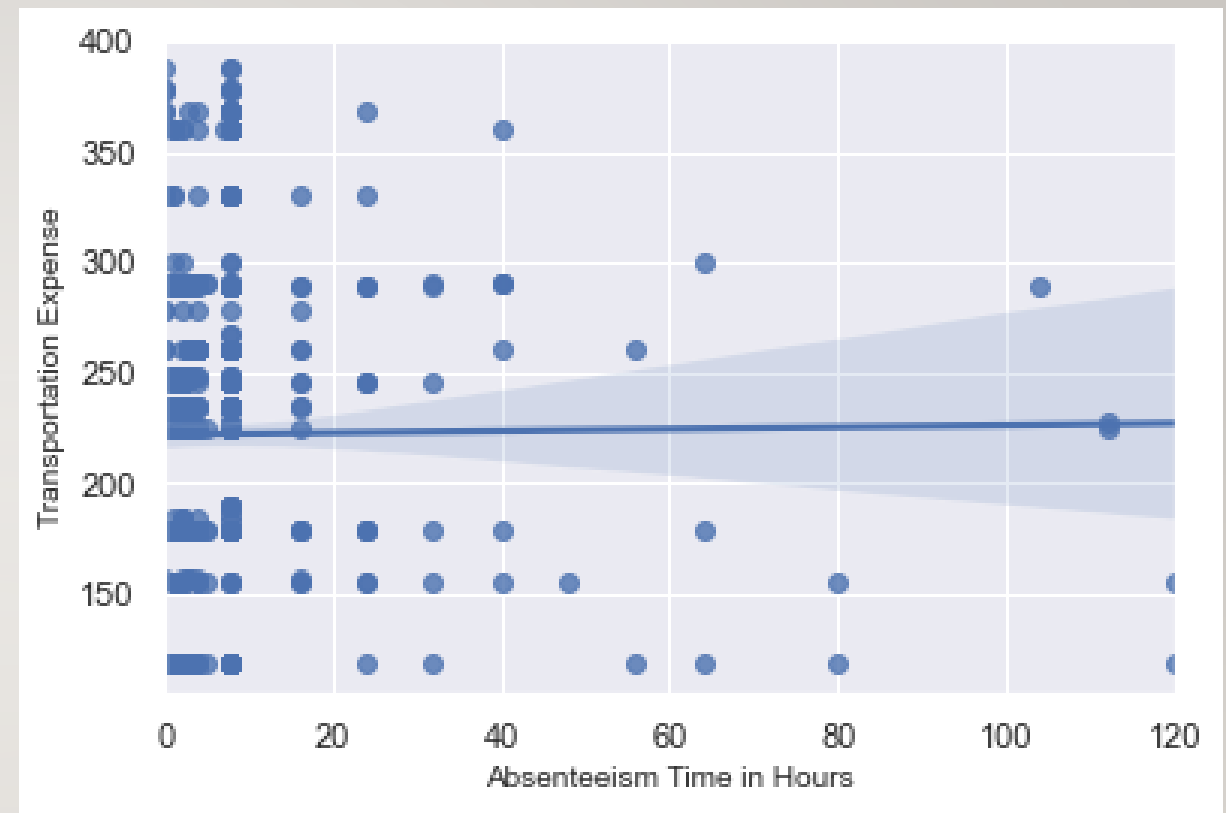
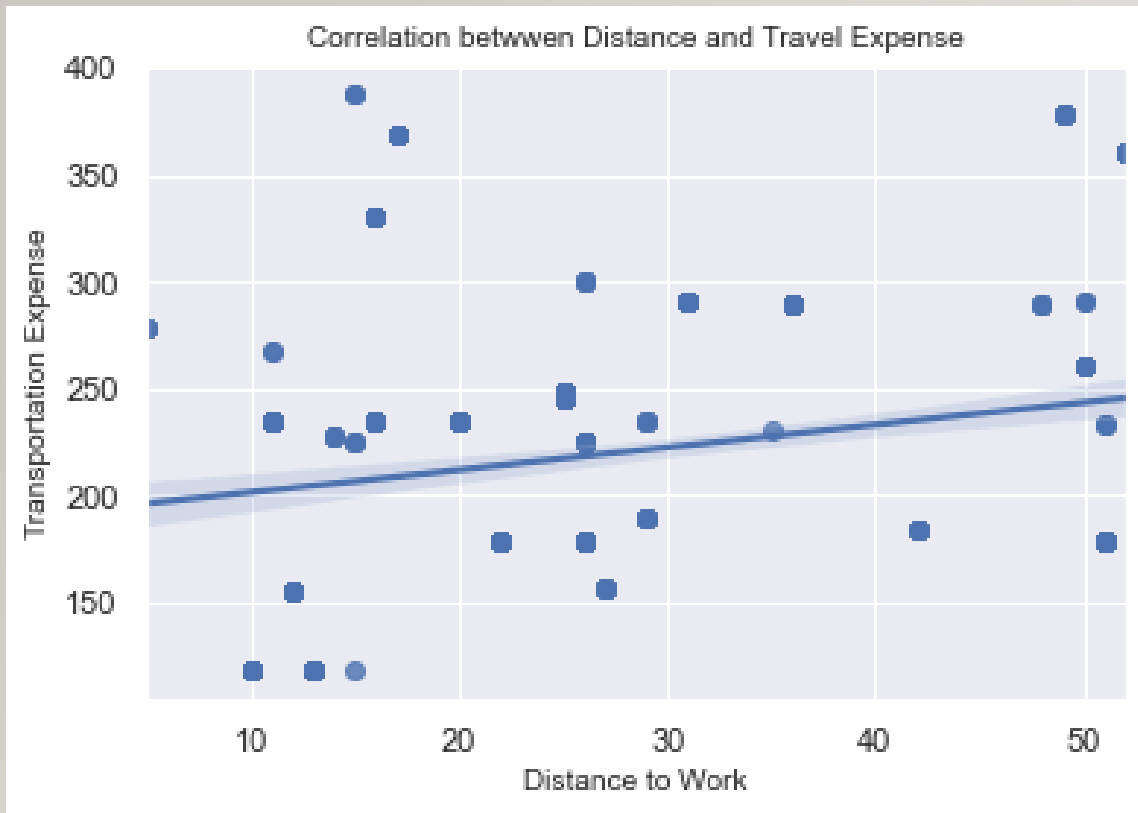
- Click to add text





# PAIRPLOT SHOWING CORRELATION ALL NUMERIC FEATURES

---

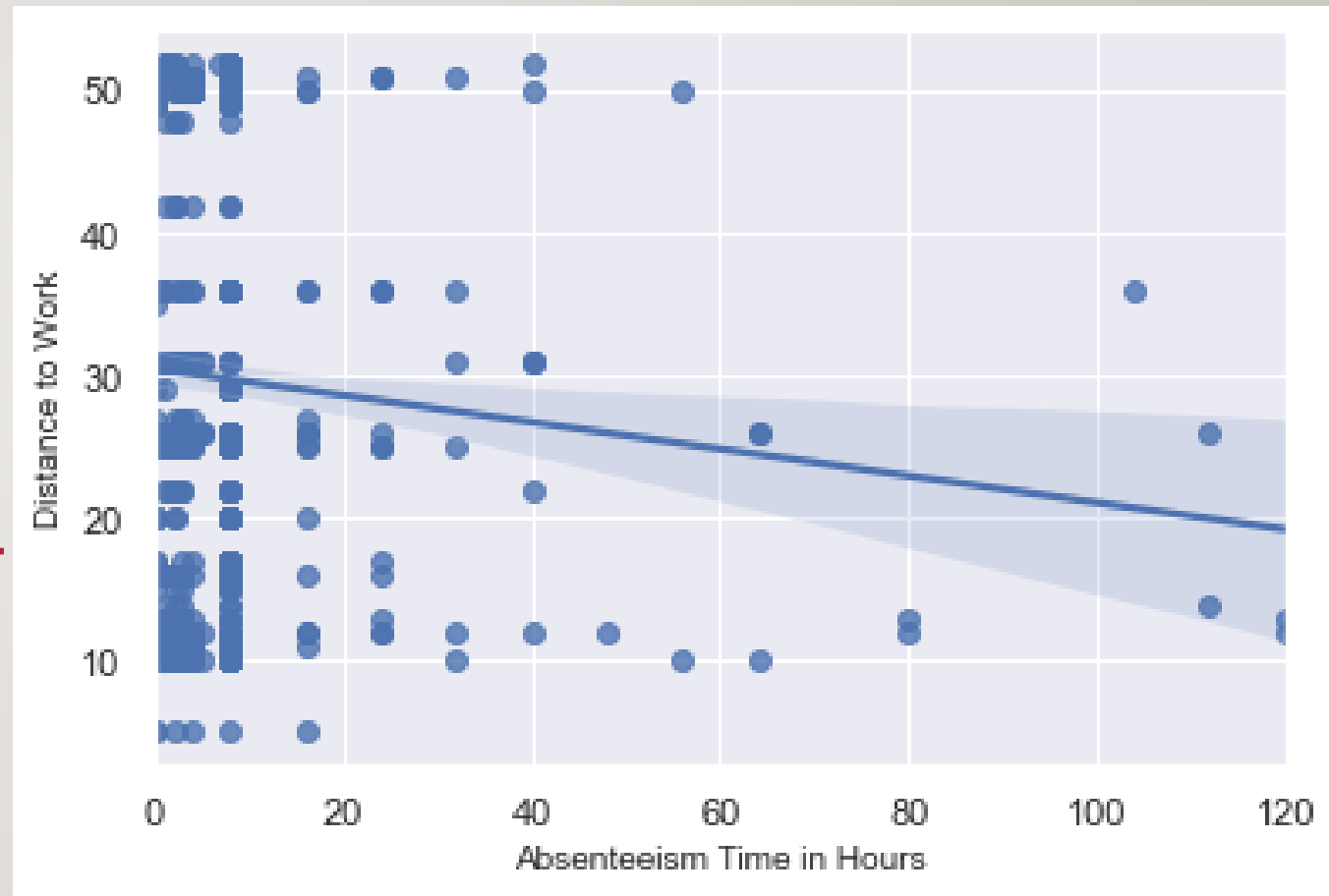


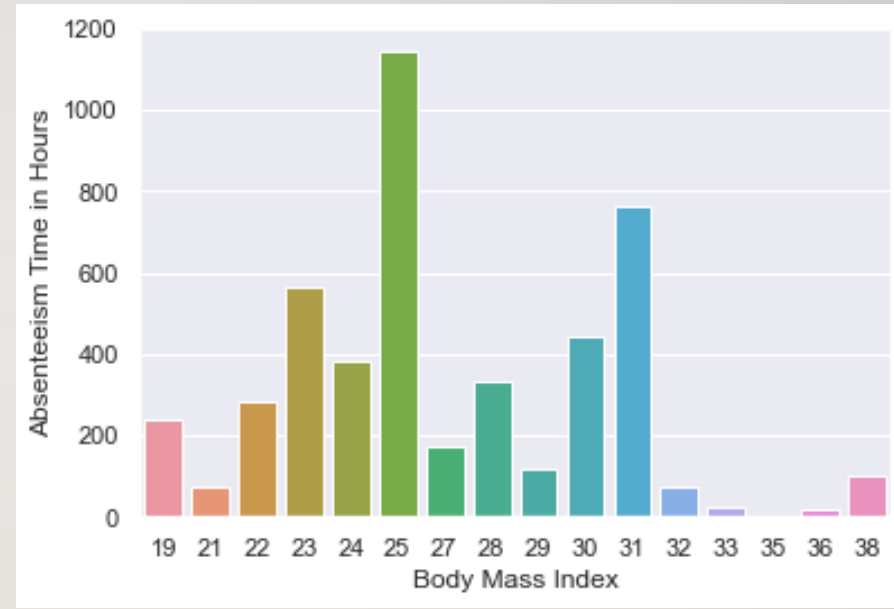
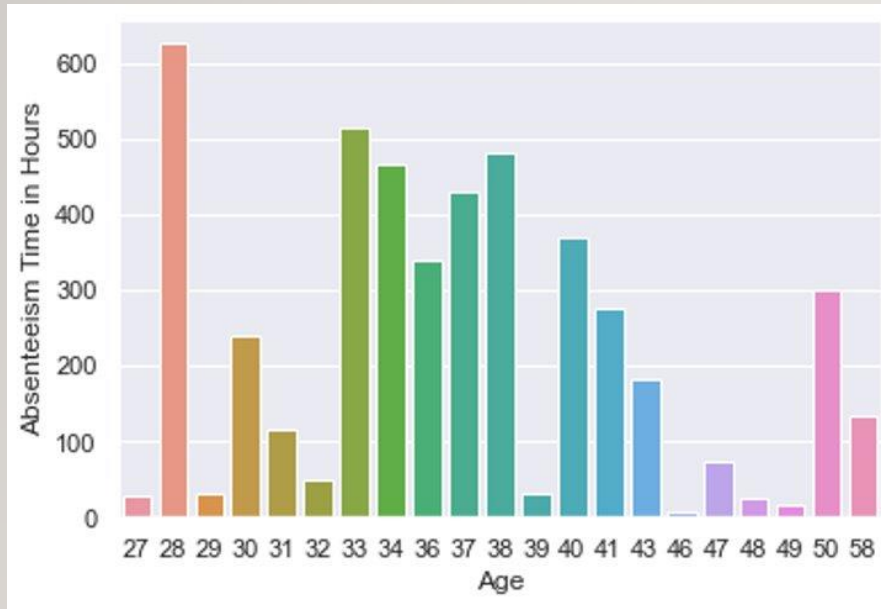
# EDA

---

# EDA

---





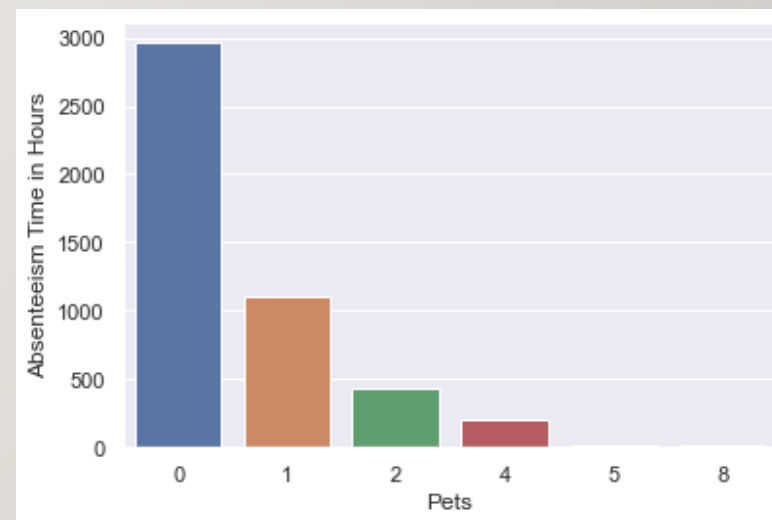
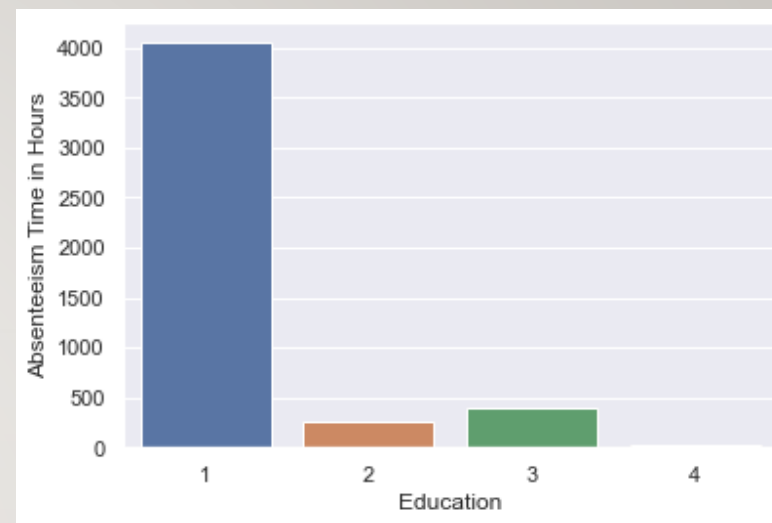
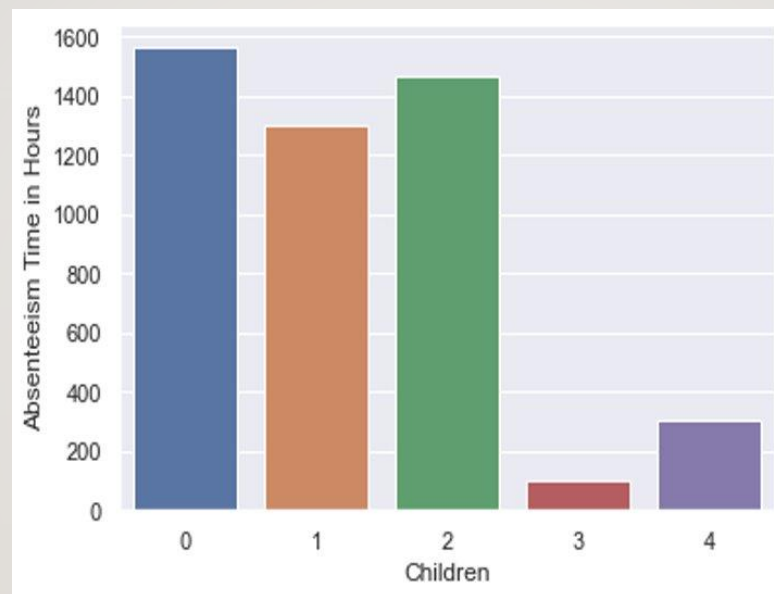
# EDA

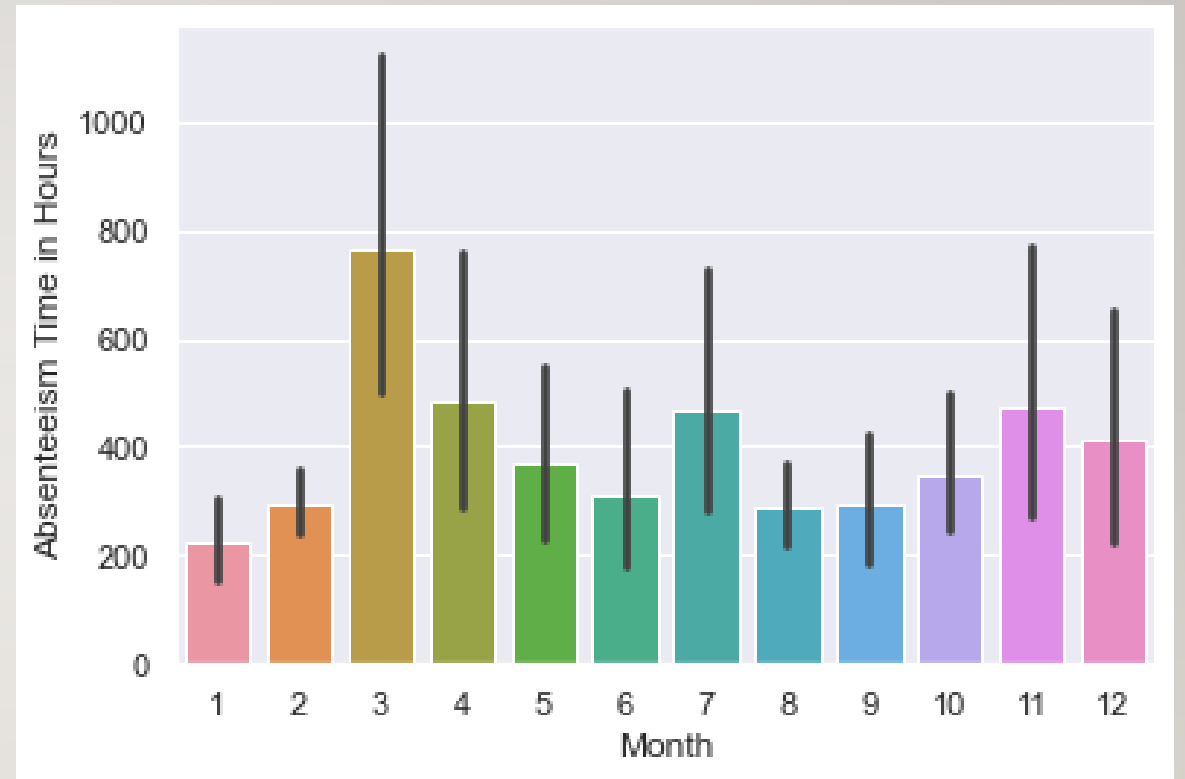
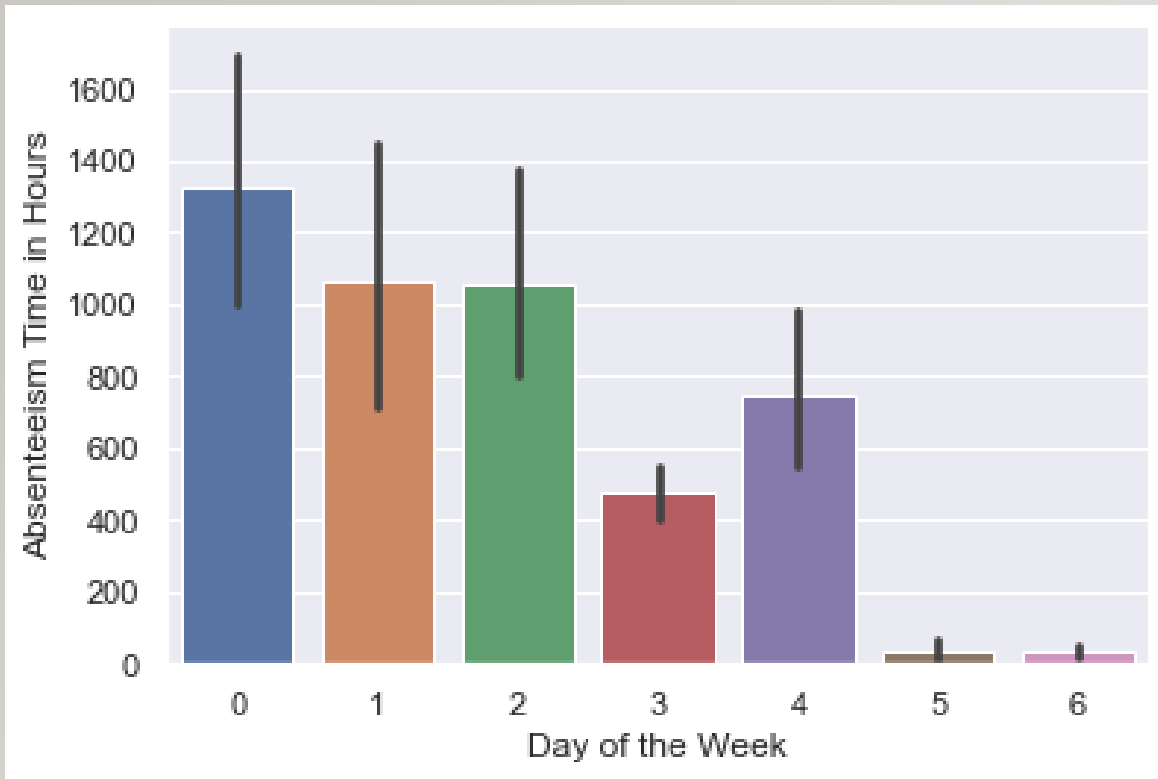
---



# EDA

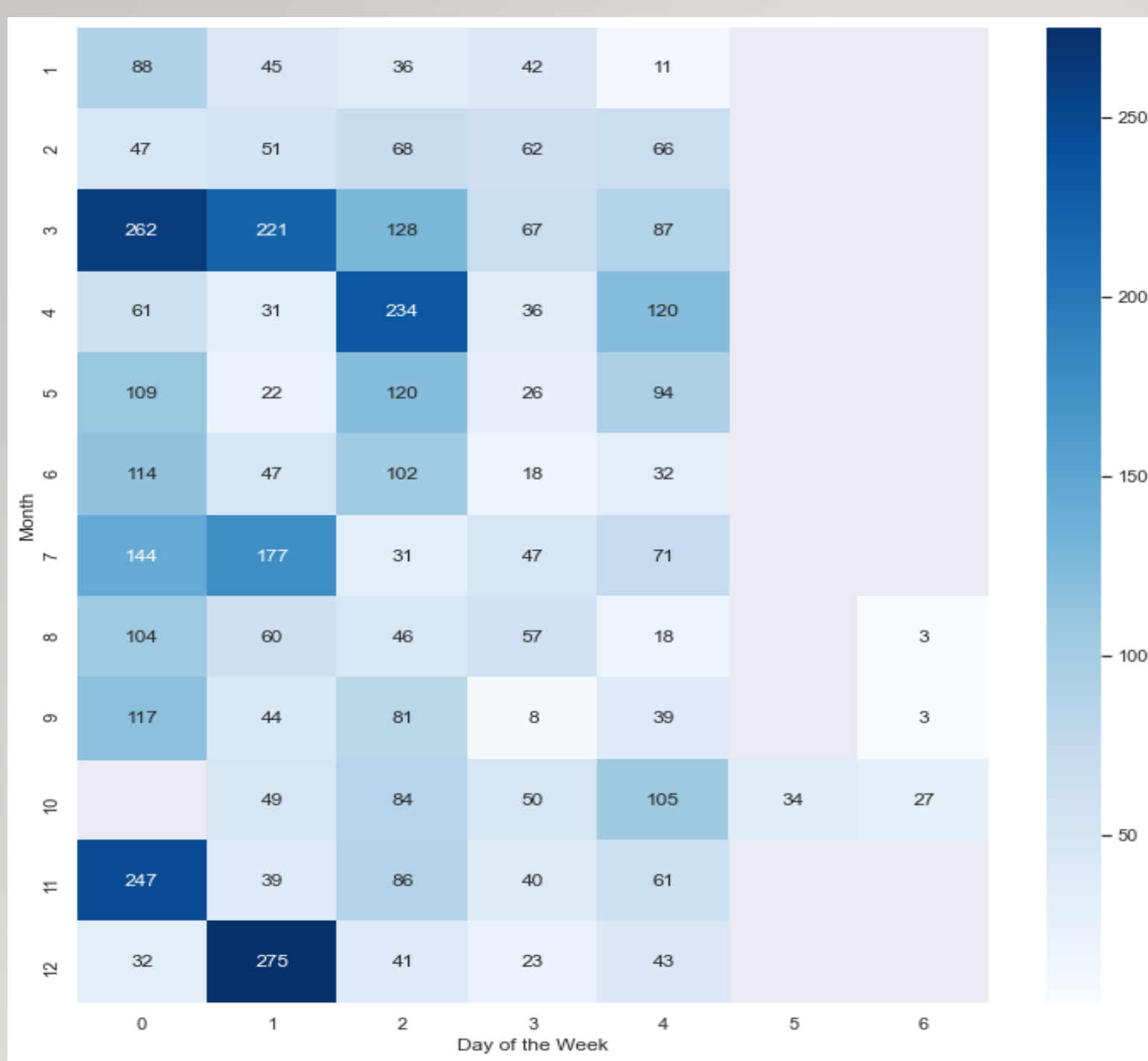
---





# EDA

---

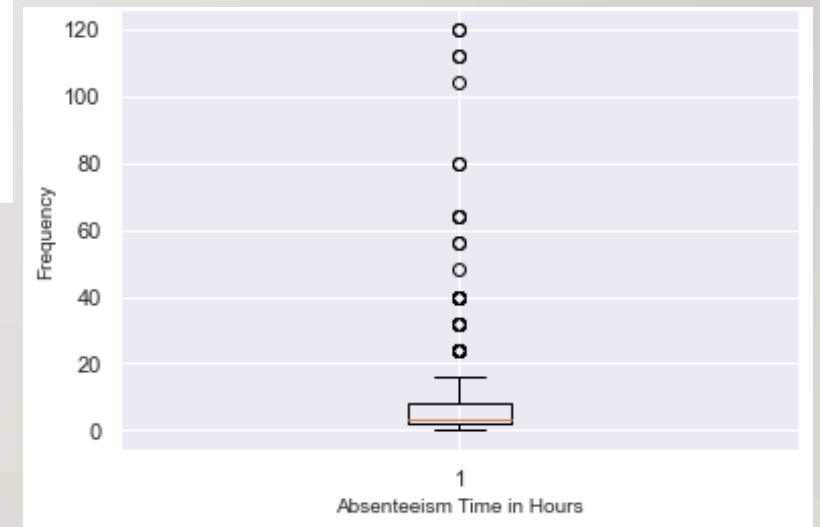
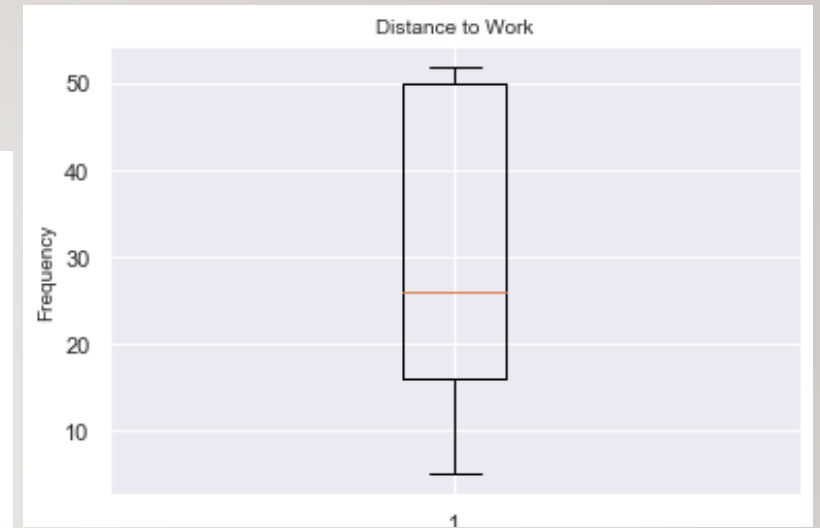
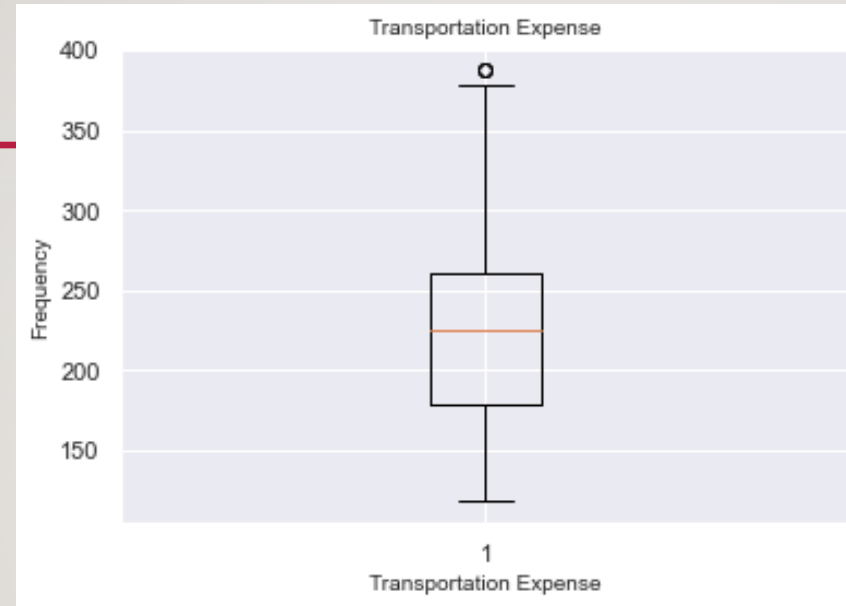


---

Heatmap showing absent  
hours throughout  
weekdays and months

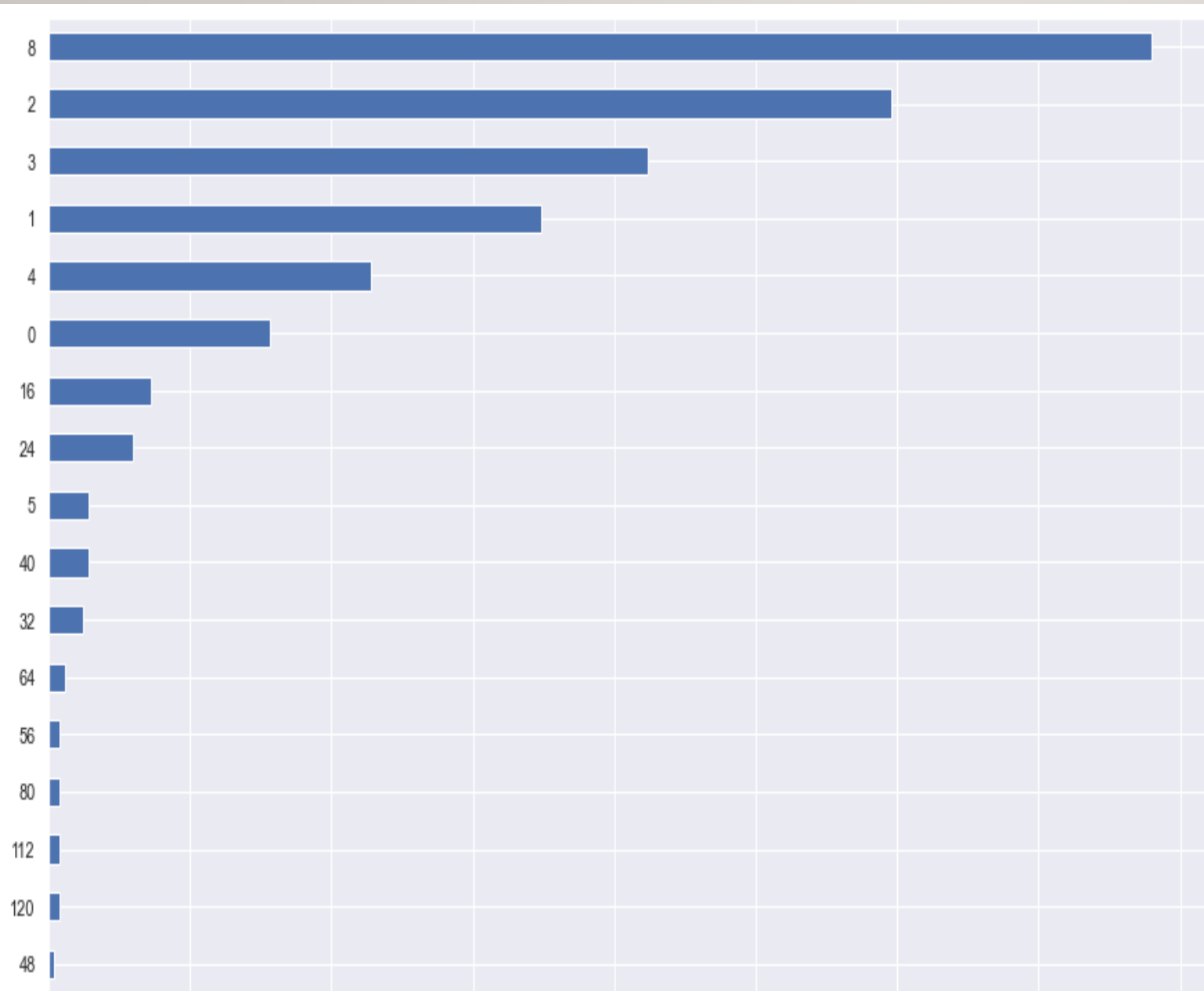


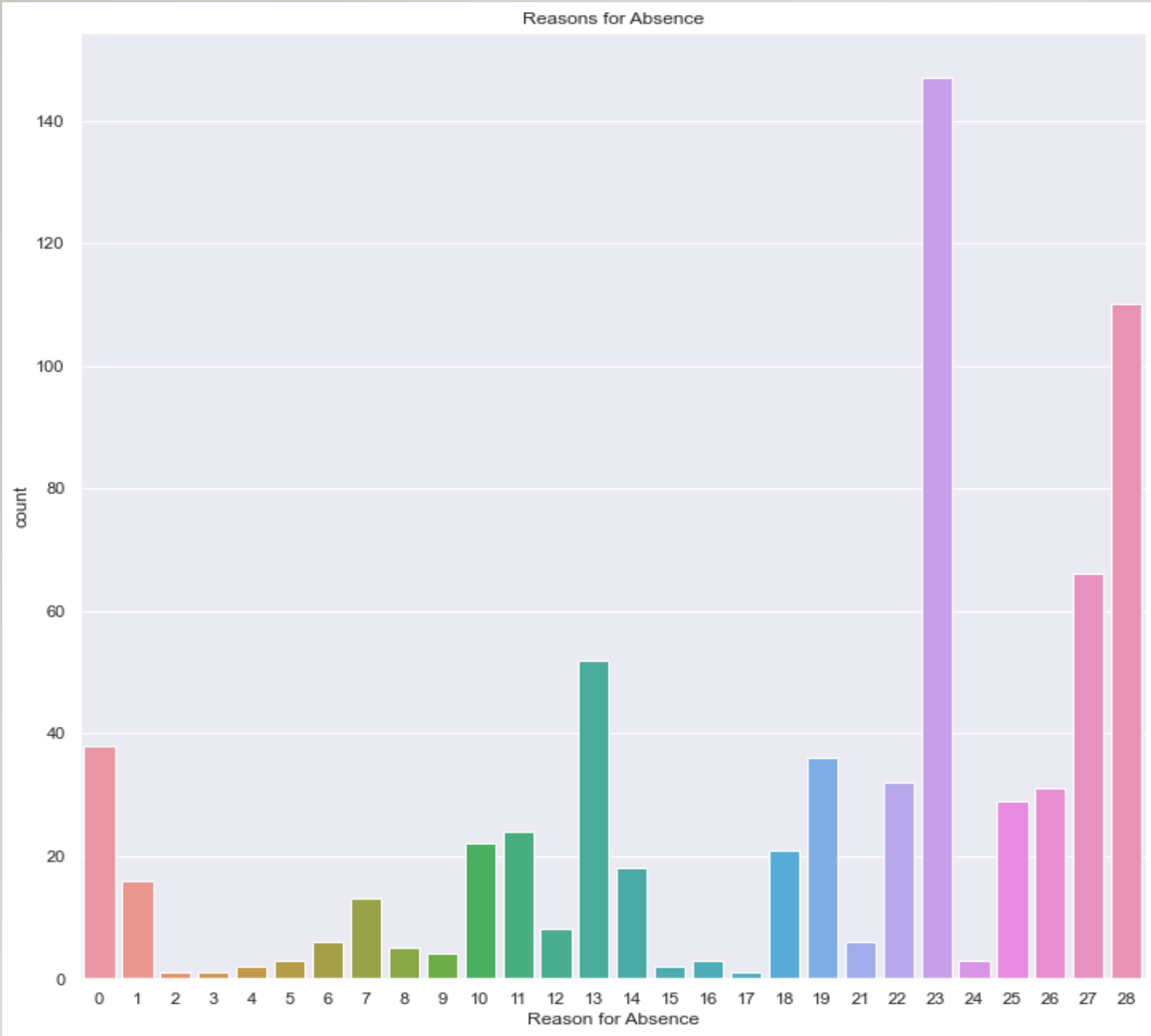
# EDA



# Absenteeism Time in Hours

---





# Reasons for Absent

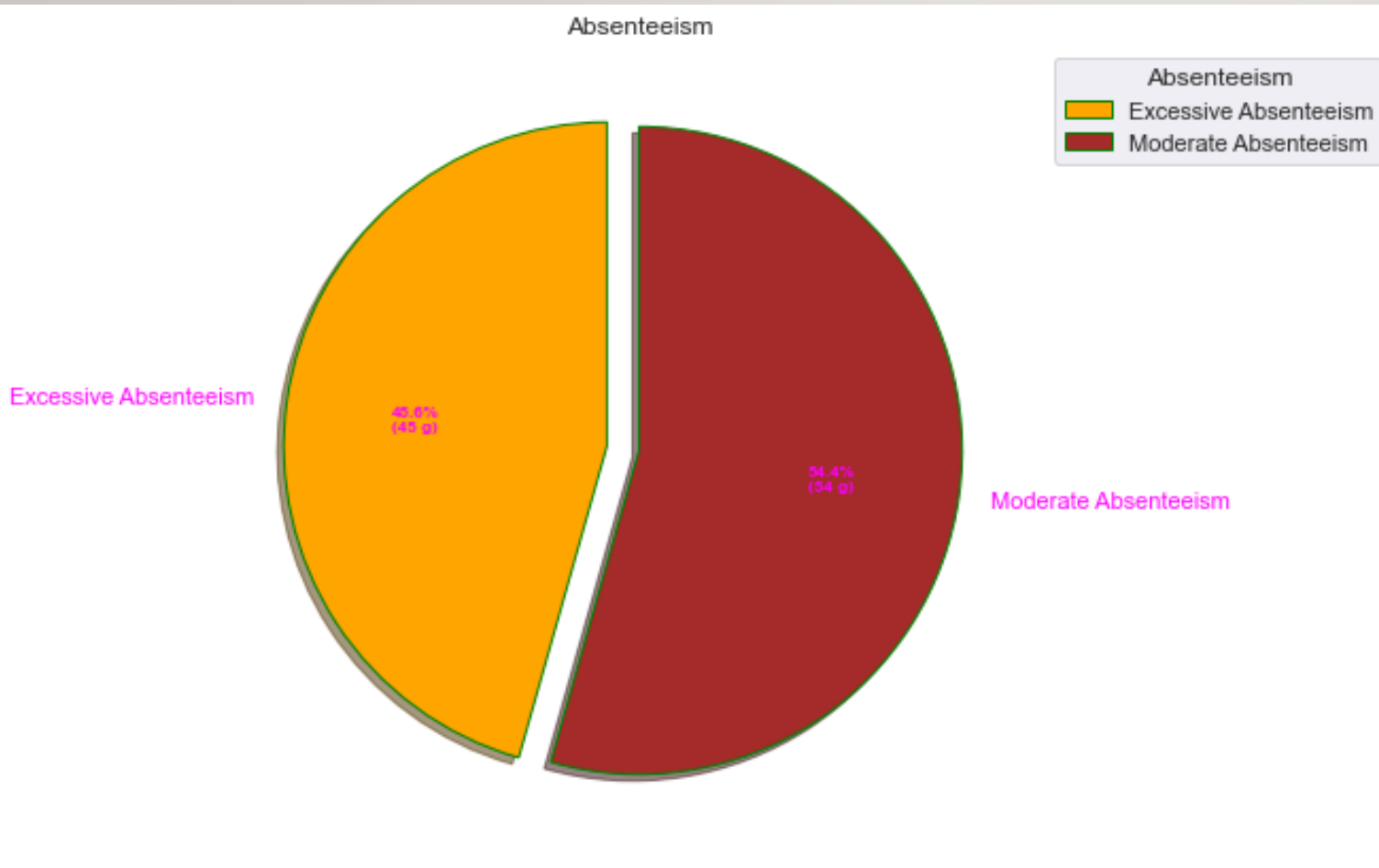
*23: Medical consultation*

*27: Physiotherapy*

*28: Dental consultation*

# PIE CHART SHOWING PERCENTAGE OF ABSENTEEISM HOUR

---



# HYPOTHESIS TESTING

ALPHA = 0.05

T = -1.971299347442617

P = 0.04961460049096922

WE REJECT OUR NULL HYPOTHESIS

Null Hypothesis (H<sub>0</sub>): There is no difference in the Absenteeism time between the near and far distance to workplace.

Alternative Hypothesis (H<sub>A</sub>): There is a difference in the Absenteeism time between the near and far distance to workplace.



# HYPOTHESIS TESTING

Null Hypothesis ( $H_0$ ): There is no difference in the Absenteeism time between the less and more expense to travel to workplace.

Alternative Hypothesis ( $H_A$ ): There is a difference in the Absenteeism time between the less and more expense to travel to workplace.

**alpha = 0.05**

$t = -6.2864120267018455$

$p = 8.547685104806746e-10$

We reject our null hypothesis.



# HYPOTHESIS TESTING

---

ALPHA = 0.05

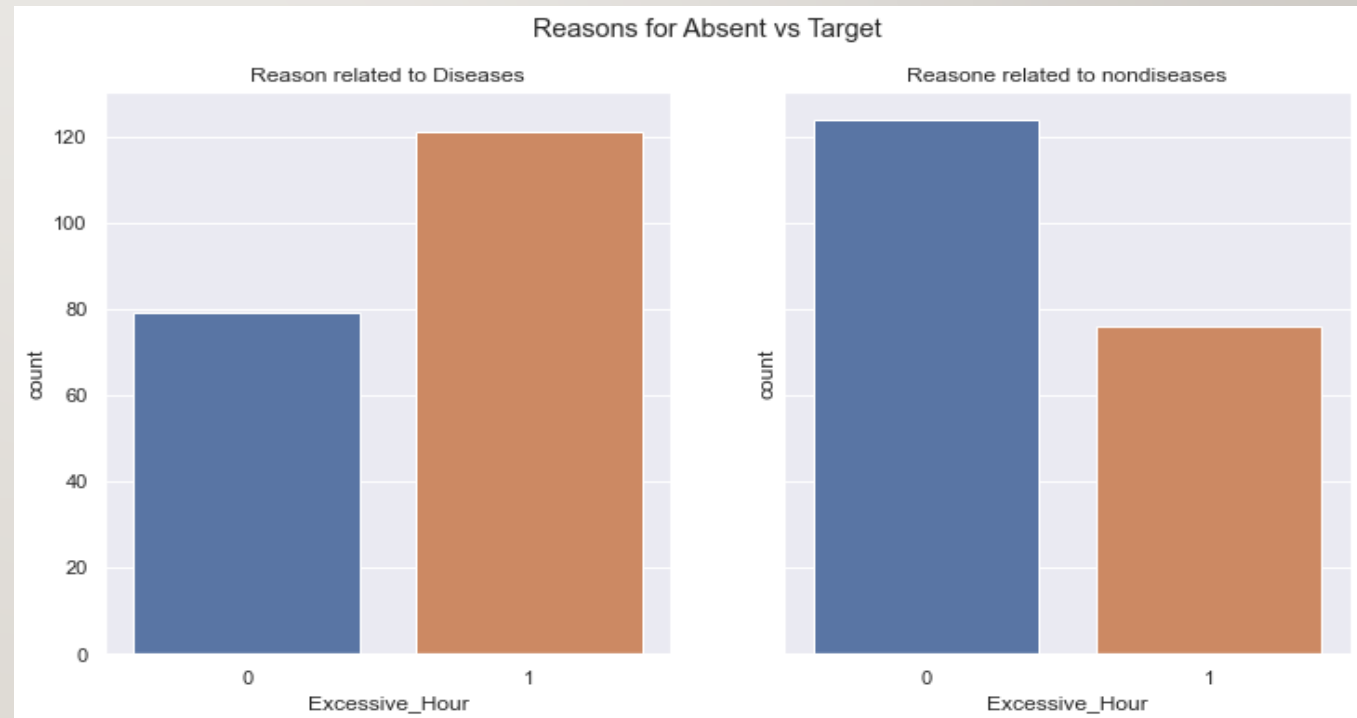
T = 4.6074075338187495

P = 5.497664640181139E-06

WE REJECT OUR NULL  
HYPOTHESIS.

Null Hypothesis (H<sub>0</sub>): There is no difference in the Absenteeism time between the disease reasons and non-disease reasons.

Alternative Hypothesis (H<sub>A</sub>): There is a difference in the Absenteeism time between the disease reasons and non-disease reasons.



# SUMMARY

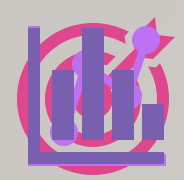
---

Employee is living far from workplace is more likely to absent from work

Employee who spend more expense for travel to work is also likely to absent from work

The reasons for absence related to diseases would spend more absent hour than reasons related to non-diseases.





# FUTURE WORK

---

Analysis other features to be more accuracy for predict absenteeism hours purposes.

Use logistic regression model to classify the target.

Collect more data to be able to get better analysis