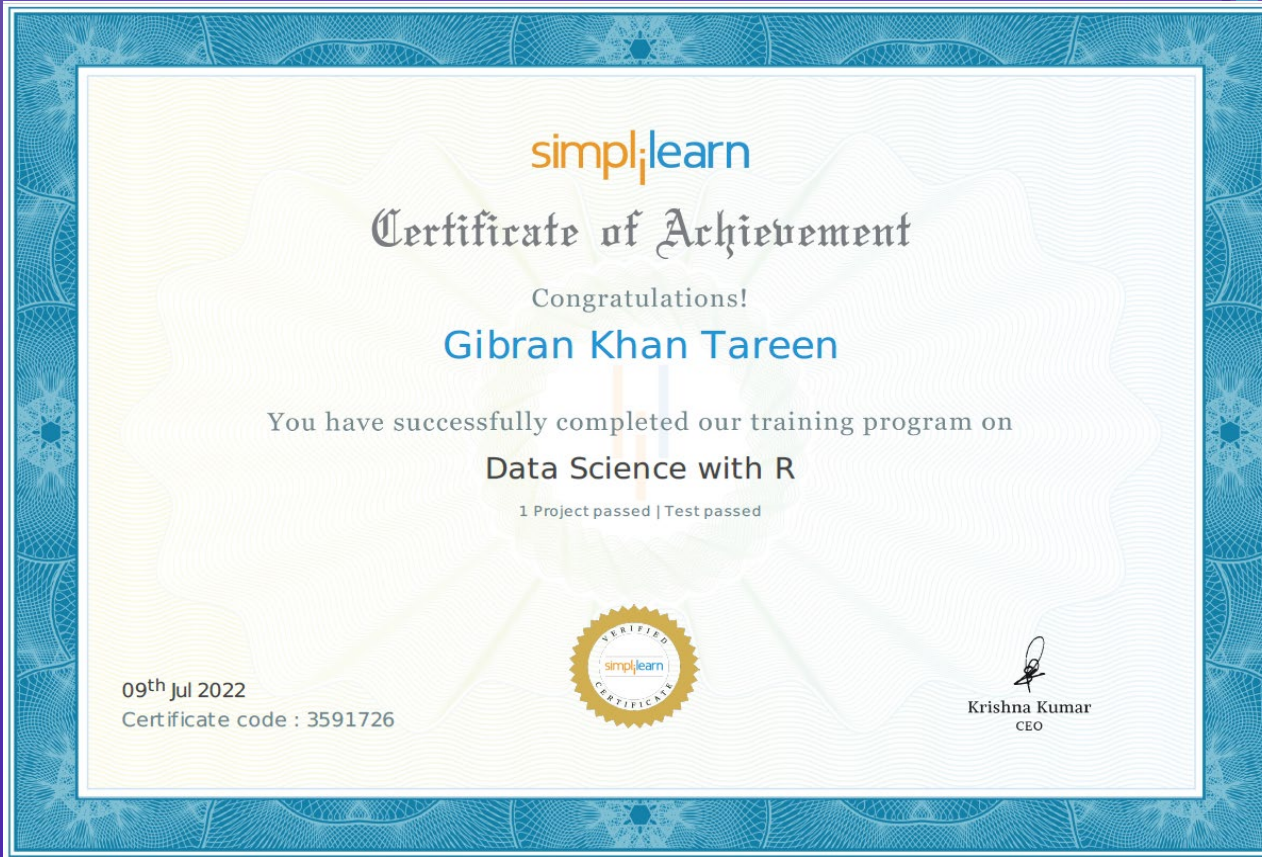# Data Science With R

## Summer Training
### ETP VIVA

Made by Gibran Khan Tareen
Registration Number: 12100173

# Certificate from The Institution

# Table of contents

**01** A Brief Introduction
About the main objectives of our project and language used

**02** Important Concepts
Discuss the important concepts used in the project

**03** Few Glimpses
Of important outcomes which we found during the Project

**04** Presenting Project Report
Finally I will show my project report, source code and execution

"Torture the data, and it will **confess** to anything."

—Ronald Coase, *Economics, Nobel prize Laureate*

**01**

# A Brief
# Introduction

What are the Project Objectives
and Language We Used In It

# Basic Structure of    Our Dataset

| Attributes | Description |
| --- | --- |
| AGE | Age of the patient discharged |
| FEMALE | A binary variable that indicates if the patient is female |
| LOS | Length of stay in days |
| RACE | Race of the patient (specified numerically) |
| TOTCH | Hospital discharge costs |
| APRDRG | All Patient Refined Diagnosis Related Groups |

# Main Objectives of  Our Project

### Record Patient Stats

Find which age category most frequently visit the hospital and has the maximum expenditure

### Check For Malpractice

Find if there is any malpractice going on by analyzing the race of the patient is related to the hospitalization costs.

### Allocation of resources

To properly utilize the costs, the agency has to analyze the severity of the hospital costs by age and gender.

### Find Main Cost Factors

To perform a complete analysis, the agency wants to find the variable(s) that mainly affect hospital costs.

# Language Used
## for our Project: R

R is the most popular language in the world of Data Science. It is heavily used in analyzing data that is both structured and unstructured. This has made R, the standard language for performing statistical operations. R allows various features that set it apart from other Data Science languages.

# 02

## Important Concepts

Now we will Discuss the important concepts used in the project

# Important   Concepts Used

## Hypothesis Test

Its is a method of statistical inference used to decide whether the data at hand sufficiently support a particular hypothesis.

## Anova Testing

Analysis of variance is used to investigate relations between categorical variables and continuous variable

## Linear Regression

It is a statistical approach for modeling the relationship between a dependent variable and a given set of independent variables.

# Where did I Implement These Concepts?

| | Description | Implemented In |
|---|---|---|
| Hypothesis Test | An act in statistics whereby an ana tests an assumption regarding a population parameter. It provides evidence concerning the plausibili the hypothesis, given the data. | We used concept of Hypothesis testing using ANNOVA Testing our Goal 3 |
| Annova Testing | It is a type of hypothesis testing fo population variance used to find th relations between categorical variables and continuous variable Programming | We implemented ANNOVA Tes in our Goal 3 |
| Linear Regression | It is a commonly used type of predictive analysis. It is a statistica method that allows us to summari and study relationships between t continuous (quantitative) variables | We implemented Linear Regression by using our own Model in Goal 4, 5 and 6 |

# 03

# Few Glimpses of Outcomes

We will now see some important outcomes which we found during the Project

# Screenshot: Patient Statistics

# Patient Statistics
# Maximum Patients
# Top 3 Age Groups

Age GROUP 0

HasMaximumPatients

**61%**

Age GROUP 1  **08%**

Age GROUP 2  **06%**

GROUP 0

# Screenshot: Patient Statistics

# ₹678,118

AGE group 0 has the maximum Hospital expenditure

# Screenshot of Check for Malpractice

# Check for Malpractice ANNOVA TESTING

Total No. of Races: **6**

Dependent Var.: **TOTCH**

## Change impact

| Low | Medium | High |
|-----|--------|------|

| DF | F | MEAN | SUM SQ |
|----|---|------|--------|

the "F value" is quite low (**0.244**). This clearly indicates that the variation between hospitalization costs among different races is very small as compared to the variation of hospitalization costs within each race.

We observed that we have more data for **RACE 1** (484 out of 500 patients) in comparison to all other races. **This makes the observations biased** We conclude by saying "There is Insufficient data to verify if a patient's race affects his expenditure."

The p-value (labeled Pr >F) is greater than Significance value ie. Alpha (0.05) and the "Residual values" (deviation of the observed values) was quite high, so both of these Observations indicate that there is **no relationship** between race and hospital costs, thereby accepting the Null hypothesis.

# Screenshot of Finding Main Cost Factor(s)

# Finding Main Variables Which Are Affecting    Cost Factors

| | Variable Factor | Impact | | | Level of Impact |
|---|---|---|---|---|---|
| 1 | AGE | ✔ | | | Very High |
| 2 | FEMALE | | | ✔ | Very Low |
| 3 | LOS (Length of Stay) | ✔ | | | Very High |
| 4 | RACE | | | ✔ | Very Low |
| 5 | APRDRG (All Patient Refined Diagnosis Related Groups) | | ✔ | | Medium |

# Thank You!

This Shall be it for the Presentation Part. I will now show the Project Report and Source Code

# 04

# Presenting
# Project Report

Project report, Source code
and execution