Submission Date: 15 Oct 2023 11.59 PM

Weightage: 30%

INDIAN INDUSTRIES ARE GIVING DIRECT EMPLOYMENT TO MILLIONS OF WORKERS. ACQUIRING A NEW TALENT IS ALWAYS A KEY CONCERN OF THE ORGANIZATIONS, ESPECIALLY WHEN THEY DON'T HAVE ADEQUATE EXPERIENCE IN IDENTIFYING THE RIGHT TALENT AND APPOINTING HIM/HER TO THE APPROPRIATE ROLE. IN MANY CASES, IT IS DIFFICULT TO FIND AN EXACT MATCH FOR THE JOB SPECIFIED. IF AN OFFER IS DENIED, THEN THE HUMAN RESOURCE (HR) DEPARTMENT HAS TO REPEAT THE ENTIRE RECRUITMENT PROCESS RESULTING IN ADDITIONAL EFFORT.

THAT'S WHY MANY OF THE ORGANIZATIONS DEPEND ON THE THIRD PARTY SPECIALIST SERVICE PROVIDERS WHO HAVE WIDE EXPERIENCE IN HIRING PROCESSES ACROSS DIFFERENT INDUSTRIES. "HRMAGIC" IS ONE OF THE KNOWN NAMES IN THIS SPACE HAVING A HUGE CLIENT BASE SPREAD ACROSS DIFFERENT SECTORS LIKE IT, MANUFACTURING, FINANCE ETC. OVER THE YEARS THEY HAVE REALIZED THE PAIN ORGANIZATIONS FACE WITH RESPECT TO THE HIRING PROCESS. AT SAME TIME THEY HAVE COLLECTED A HUGE DATASET OF THE RECRUITMENT DRIVES THEY HAVE CONDUCTED FOR THE VARIOUS ORGANIZATIONS. WHILE WORKING ON DIFFERENT RECRUITMENT DRIVES, THEY HAVE OBSERVED THAT 30% OF THE CANDIDATES WHO ACCEPT THE JOB OFFER, DO NOT JOIN THE COMPANY. THIS LEADS TO HUGE LOSS OF REVENUE AND TIME AS THE COMPANIES NEED TO INITIATE THE RECRUITMENT PROCESS AGAIN TO FILL IN THE WORKFORCE DEMAND. WHILE LOOKING FOR THE SOLUTIONS FOR THIS PROBLEM, THE TOP MANAGEMENT HAS COME ACROSS THE USE OF PREDICTIVE ANALYTICS TECHNIQUES APPLIED BY FOREIGN TALENT MANAGEMENT FIRMS. THEY FOUND THOSE TECHNIQUES QUITE INTERESTING AND HAVE ASKED THE TEAM TO EXPLORE MORE ON THE SAME.

"HRMAGIC " TEAM HAS DECIDED TO PROCEED WITH A PROTOTYPE USING WHICH THEY WANT TO FIND OUT IF A MODEL CAN BE BUILT TO PREDICT THE LIKELIHOOD OF A CANDIDATE JOINING THE COMPANY. IF THE LIKELIHOOD IS HIGH, THEN ONLY THE COMPANY WILL GO AHEAD AND OFFER THE JOBS TO THE CANDIDATES. THEY HAVE ALREADY EXTRACTED THE PORTION OF THE DATASET FROM THEIR HUGE PILE OF DATA WHICH CONTAINS SEVERAL ATTRIBUTES ABOUT CANDIDATES ALONG WITH A COLUMN THAT INDICATES IF THE CANDIDATE FINALLY JOINED THE COMPANY OR NOT. THE DETAILED DESCRIPTION OF THE OTHER ATTRIBUTES IS PROVIDED ALONG WITH THE DATASET FILE.

FOR THAT PURPOSE THEY HAVE REACHED OUT TO SEVERAL GROUPS WHO CAN HELP THE "HRMAGIC" TEAM TO BUILD THE REQUIRED MODEL. YOURS IS ONE OF THEM. AFTER CAREFUL EVALUATION OF THE MODELS RECEIVED FROM SEVERAL GROUPS, THE MANAGEMENT WILL TAKE A CALL ON WHETHER TO GO AHEAD AND INVEST MORE INTO THIS PREDICTIVE MODELING STRATEGY OR NOT.

AS A PART OF THIS PROJECT, YOU WILL BE MAKING USE OF DATASETS PROVIDED AND HELP BUILDING THE MODEL. THE MAIN OBJECTIVE OF THIS PROJECT IS TO GIVE YOU REAL LIFE EXPERIENCE WHILE DOING DATA ACQUISITION, DATA INTEGRATION, DATA CLEANING AND DATA TRANSFORMATION BEFORE ATTEMPTING ANY OF THE ANALYTICAL ACTIVITY ON THE DATA.

THE VARIOUS TASKS THAT YOU WILL BE DOING AS A PART OF THIS EXERCISE WILL BE AS FOLLOWS:

THE FIRST STEP IN BUILDING A MODEL IS TO COLLECT OR EXTRACT DATA ON THE DEPENDENT VARIABLE AND INDEPENDENT VARIABLES FROM DIFFERENT DATASETS PROVIDED. DATA COLLECTION / EXTRACTION IS A TIME CONSUMING AND EXPENSIVE PROCESS.

- AS A TEAM YOU HAVE TO MERGE THE CANDIDATE DATA PROVIDED BASED ON THE COMMON FACTORS.

- SIMULTANEOUSLY YOU NEED TO SEARCH WHAT ARE THE ATTRIBUTES THOSE ARE TAKEN INTO CONSIDERATION WHILE CARRYING OUT RECRUITMENT ACTIVITY.

- THE RESEARCH WILL PROVIDE YOU MORE INSIGHT ON DETERMINING THE USEFUL ATTRIBUTES PRESENT IN THE DATASET OR GIVE POINTERS TO THE VARIABLES THAT NEED TO BE DERIVED. MAKE A LIST OF SUCH POTENTIAL ATTRIBUTES.

- THE OUTCOME OF THIS STEP WILL BE A MERGED DATA SET CONTAINING THE 200 CANDIDATES DATA.

- DOCUMENT ALL YOUR EFFORTS APPROPRIATELY IN THE JUPYTER NOTEBOOKS WITH DESCRIPTION AND CODE.

- THE WEIGHTAGE FOR THIS TASK WILL BE 5 MARKS

## (B) PRE-PROCESS DATA

BEFORE THE MODEL IS BUILT, IT IS ESSENTIAL TO ENSURE THE QUALITY OF THE DATA FOR ISSUES SUCH AS RELIABILITY, COMPLETENESS, USEFULNESS, ACCURACY, MISSING DATA AND OUTLIERS.

DATA IMPUTATION TECHNIQUES MAY BE USED TO DEAL WITH MISSING DATA. USE OF DESCRIPTIVE STAT AND VISUALIZATION MAY BE USED TO IDENTIFY THE EXISTENCE OF OUTLIERS AND VARIABILITY IN THE DATASET. MANY NEW VARIABLES CAN BE DERIVED AND ALSO USED IN MODEL BUILDING. CATEGORICAL DATA HAS TO BE PRE-PROCESSED USING DUMMY VARIABLES, BEFORE IT IS USED TO MODEL BUILDING.

- YOU HAVE TO CLEANSE YOUR DATASETS TO REMOVE ALL SUCH DAUNTING ISSUES.

- NARRATE ALL THE ISSUES WHICH YOU ENCOUNTER DURING THIS EXERCISE CLEARLY WITH APPROPRIATE EXPLANATION AND CODE.

- THE WEIGHTAGE FOR THIS TASK WILL BE 7 MARKS.

## (C) PERFORM DESCRIPTIVE ANALYTICS ON DATA

IT IS ALWAYS GOOD TO PERFORM DESCRIPTIVE ANALYTICS BEFORE MOVING TO BUILDING A PREDICTIVE ANALYTICS MODEL. IT WILL HELP TO UNDERSTAND THE VARIABILITY IN THE MODEL. IT ALSO HELPS IN IDENTIFYING THE RELATIONSHIP BETWEEN THE VARIABLES PRESENT IN THE DATASET.

- USE THE EXPLORATORY DATA ANALYSIS TECHNIQUE ON THE DATASET IN ORDER TO FIND OUT THE INTERESTING INSIGHTS THAT ARE HIDDEN WITHIN THE DATA CAPTURED.

- DESCRIBE ALL EDA STEPS THOSE ARE DONE WITH THE OBSERVATIONS OBTAINED OUT OF IT WITH THE HELP OF PYTHON CODE IN JUPYTER NOTEBOOK.

- THE WEIGHTAGE FOR THIS TASK WILL BE 7 MARKS.

## (D) FEATURE ENGINEERING

E ENGINEERING IS AN IMPORTANT STEP TO DEVELOP AND IMPROVE PERFORMANCE OF MACHINE LEARNING MODELS. THESE TECHNIQUES CAN ALSO BE USED TO IDENTIFY THE VARIABLES THAT IMPACTS THE OUTCOME OF THE MODEL. ITS BASICALLY PROCESS OF IDENTIFYING AND EXTRACTING THE USEFUL FEATURES FROM THE AVAILABLE DATA. THE PRIMARY GOAL IS TO DERIVE A SET OF FEATURES THAT BEST REPRESENT THE INSIGHTS HIDDEN IN THE DATA, WITH A SIMPLER MODEL THAT GENERALIZES WELL TO FUTURE (UNKNOWN) OBSERVATIONS.

- IN THIS STEP YOU HAVE TO USE KNOWLEDGE OF FEATURE SELECTION METHODS TO IDENTIFY THE VARIABLES THAT HAVE GREATER IMPACT ON THE OUTCOME.
- AN ELABORATE DESCRIPTION OF THE IMPACT OBSERVED IS EXPECTED AS AN OUTCOME OF THIS STEP.
- THE WEIGHTAGE FOR THIS TASK WILL BE 6 MARKS.

**(E) MODEL BUILDING AND DIAGNOSTICS**

IN THIS STAGE FIRST DATA IS DIVIDED INTO TRAIN AND TEST DATA. THE SUBSET CAN ALSO BE CREATED USING RANDOM / STRATIFIED SAMPLING PROCEDURE. THIS IS AN IMPORTANT STEP TO MEASURE THE PERFORMANCE OF A MODEL USING A DATASET NOT USED IN MODEL BUILDING. IT IS ALSO ESSENTIAL TO CHECK FOR ANY OVERFITTING OF THE MODEL. THE MODEL IS BUILT USING A TRAINING DATASET TO ESTIMATE THE MODEL PARAMETERS. THE METHOD OF CLASSIFICATIONS CAN BE UTILIZED FOR THE SAME.

- YOU HAVE TO PREPARE A LOGISTIC REGRESSION MODEL TO PREDICT THE PROBABILITY OF A CANDIDATE JOINING THE COMPANY. ASSUME "OFFERREJECTED" AS POSITIVE CASES AND "JOINED" AS NEGATIVE CASES.
- FIND THE SIGNIFICANT FEATURES FROM THE ABOVE MODEL AND BUILD ANOTHER LOGISTIC REGRESSION MODEL WITH ONLY THE SIGNIFICANT VARIABLES.
- COMPARE THE PERFORMANCE OF BOTH MODELS USING VARIOUS MODEL ATTRIBUTES AND RECOMMEND A MODEL THAT CAN BE USED BY "HRMAGIC"
- THE WEIGHTAGE FOR THIS TASK WILL BE 5 MARKS

**NOTES:**

- THIS IS A TAKE-HOME PROJECT TO BE CARRIED OUT BY A GROUP OF LEARNERS.
- AS PER THE NEED, THE DEMOS / VIVAS CAN BE ARRANGED FURTHER ON.
- WHEREVER REQUIRED YOU CAN MAKE APPROPRIATE ASSUMPTIONS BUT MAKE SURE THAT YOU HAVE SPELT THEM APPROPRIATELY IN THE SUBMITTED DOCUMENTS.
- THIS IS A PROGRAMMING EXERCISE - REQUIRING THE APPROACH OF APPROPRIATE MODEL BUILDING.
- YOU MAY CONSULT / DISCUSS WITH OTHER LEARNERS PERIPHERAL ASPECTS SUCH AS THE ENVIRONMENT BUT NOT ON SOLVING THE SPECIFIC PROBLEMS IN TERMS OF DESIGN OR IMPLEMENTATION.
- YOU HAVE TO WRITE THE APPROPRIATE PYTHON CODE IN JUPYTER NOTEBOOK TO SUPPORT YOU ANSWERS AND SUBMIT WITH FOLLOWING NOMENCLATURE
    - FE_PROJECT1_<GROUP_ID>.IPYNB
- IN CASE OF ANY FURTHER QUERIES, IF THOSE ARE GENERIC ONES, LEARNERS ARE ENCOURAGED TO USE DISCUSSION FORUMS, OTHERWISE THEY CAN REACH OUT TO ME AT PPAWAR@WILP.BITS-PILANI.AC.IN.
- MANAGE YOUR EFFORTS PROPERLY AS THERE IS NO SCOPE TO SHIFT THE DEADLINES ANNOUNCED ABOVE.

1) Logistics Regression
2) Logistics Regression Python tutorial
3) Logistics Regression explained