

Bigdata Analytics

Final Report

INE5015 – 22057

빅데이터 애널리틱스 8조

| 최준희 | 이강산 | 장혜연 | 황태영 |

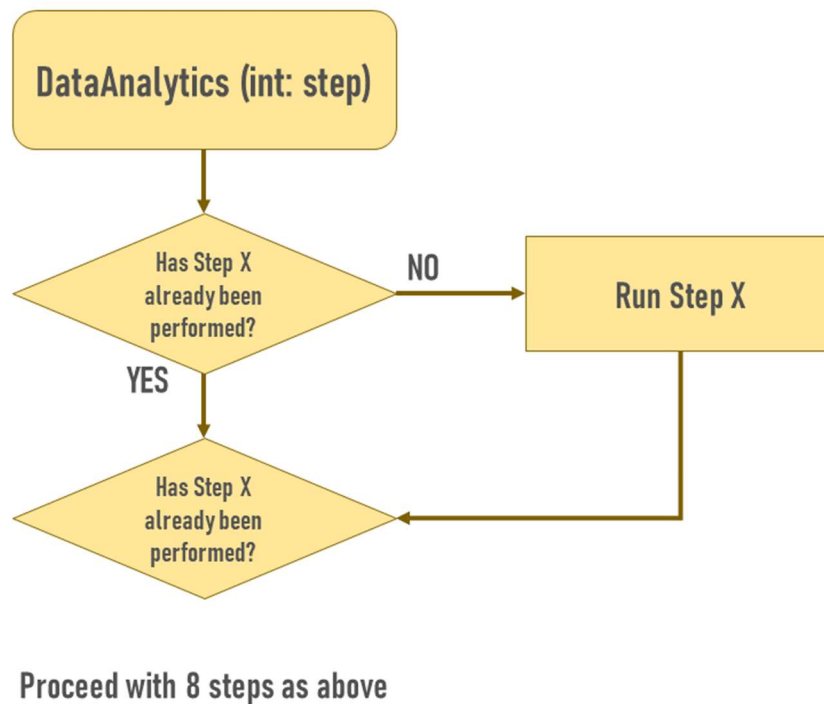
List of Contents

- Logic of Data Preprocessing
- Each step-by-step process
 - Step 0 ~ Step 2 : raw file refining
 - Data Cleaning Overview (3 Steps)
 - ◆ Step 3 : correlation check and correction
 - ◆ Step 4 : Missing Value Imputation
 - ◆ Step 5 : Outlier corrections
 - Step 6 : Data Scaling
- Feature Selection
- Data over/down Sampling overview
- Performance
 - The performance of two samplings
 - Prediction accuracy according to each algorithm
 - ◆ Logistic Regression
 - ◆ Decision Tree
 - ◆ Random Forest
 - ◆ Boosting

GITHUB : https://github.com/ChoiJunhee/INE5015_bigdata_analytics

Logic of our Data Processing

전반적인 Preprocessing 과정은 기능, 목적에 따라 Step으로 구별했고, 매 Step마다 시각적, 통계적 데이터를 확인하면서 진행하였다. 아래는 Step에 대한 구성 Logic이다.



DataAnalytics 함수에서 필요한 함수를 호출하여 사용하는 구조이며, 각 단계마다 csv 파일을 저장하도록 하여 변화를 직관적으로 확인하기 용이하도록 디자인하였다.

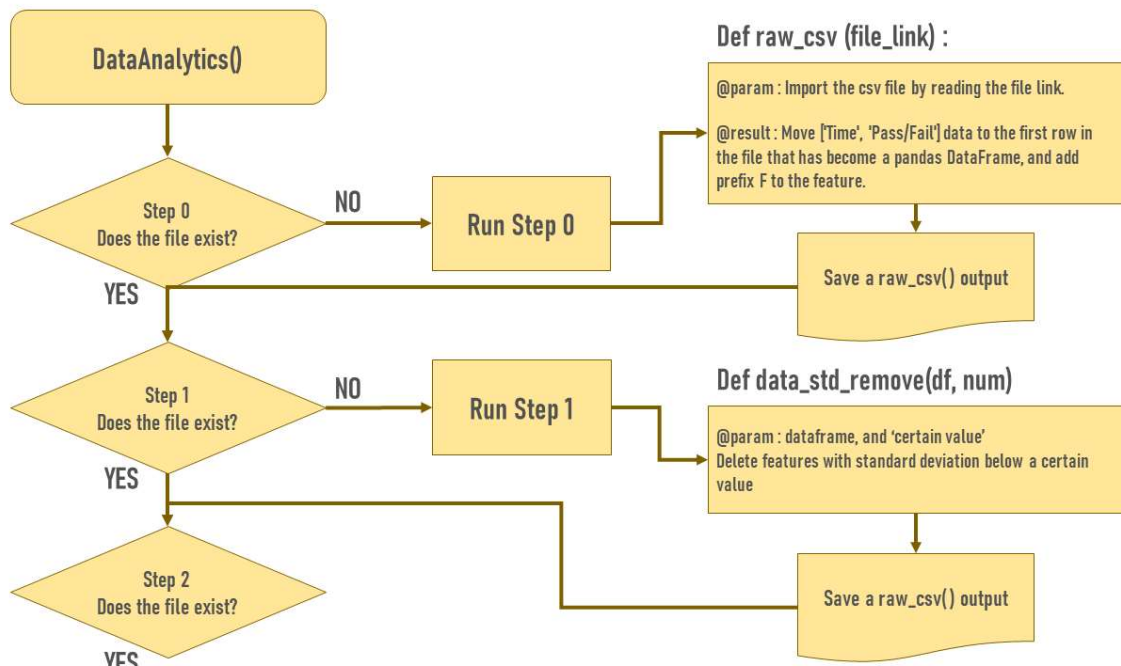
각 Step을 크게 나누어 보면 Raw data를 가공에 용이하도록 조정하는 부분, Pass/Fail 데이터를 분리해 총 3가지 데이터셋으로 나누어 진행할 수 있도록 하는 부분이 있다.

이후 각 Feature (독립변수)간 종속성을 확인하고, 제거하기 위해 진행하는 Correlation 과, 결측치 처리/보정 과정, 이상치 처리/보정 과정. 즉, 데이터 클리닝 과정이 있다.

마지막으로 Feature Selection과 샘플링을 통해 최종 데이터 셋 후보를 선정하고, 퍼포먼스를 확인하는 과정으로 마무리한다.

Step 0 ~ Step 2 : Ready to Preprocess

이 단계는 데이터 파일을 불러오고, 이후 보정에 있어 편의성을 높이기 위해 보정한다.



Explanation of Steps

STEP	STEP DESCRIPTION
STEP 0	raw_csv 함수를 통해, Pandas Read_CSV를 실행하여 데이터 셋을 받아온다. 이후 작업의 편의를 위해 행 위치 변경, Prefix 추가 등의 과정을 거친 Dataframe 파일을 받고, 저장한다.
STEP 1	일정 계수 미만의 표준편차를 가진, 특징 선택의 필요도가 낮은 Feature들을 제거하고, Dataframe 파일을 저장한다. 그리고 데이터 셋을 3개로 나누는 과정을 거친다.
STEP 2	데이터 셋은 Pass Data, Fail Data, Both Data로 나뉘며, 이렇게 데이터 셋을 나누는 것은 프로젝트 초기의 아이디어 중 유용한, 유의미한 결과를 가져오는 계기가 된다.

Step 0 부터 Step 2까지의 3 과정은 데이터 전 처리를 위한 준비 단계에 불과하다.

> 위 과정에서 1566개 (1463+104)의 테스트 케이스, 247개의 Feature Data가 남았다.

Step 3 ~ Step 5 : Data Cleaning

데이터 클리닝 개요

Step 3, 4, 5에 대한 개념 설명

구현 및 도입하게 된 계기 혹은 근거

2~3P 예상됨

