

Data Science Pipeline Application – Mobot

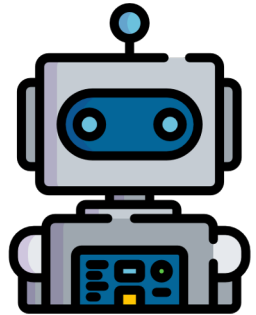
Team Members :
Jane, Andy, Henry

Agenda

- *What is our idea ?*
- *Application Structure*
- *Data Flow*

What is our idea ?

*We want to build a **Data Science Pipeline Application** that could be beneficial to the future students who are interested in developing models for data analysis project.*



Mobot

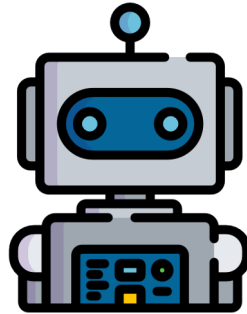
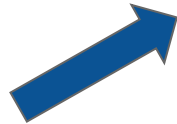
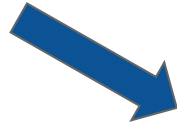
The Structure of the Application



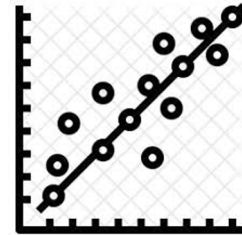
**Raw
Data**



**Executi
on Plan**



Mobot



**The
Best
Model**

The Structure of the Application

Mission	Main Objective	Action	Detail
Data Collection	Gather all the data needed.	Put data to correct folder	The user has to collect the datasets and put it in the same folder.
Data Preprocess	Convert multiple original datasets into a single ready-to-run dataset.	Missing Data Imputation	The user has to decide a way to fill NA. (Mean, Median, Zero and KDE)
		Data Transformation	The user has to decide a way to transorm data. (Root Square, Log)
		Create a Flat Table	The user has to pick a column name as key to merge all the original datasets.
Training / Testing Data Split	Split processed dataset into training/testing datasets.	Train / Test Split	The user has to decide a way to train/test data split. (Percentage, Column Value)
Model Training	Run the processed training dataset with various kinds of model.	Run data with models	The user has to decide types of model to run. (AIC, Simple Linear Regression, Stepwise)
Model Evaluation	Evaulate the performance of different model running the preprocessed data set.	Estimate performance of models run.	The user has to decide a way to estimate the performance of model. (RMSE...etc)

Data Flow



Source



Preprocessed



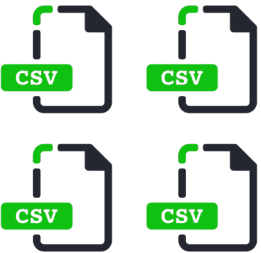
Split



Model



Estimate



Original
Data Sets



Preprocessed
Original
Dataset



Training



Testing



Models
Summary
Data

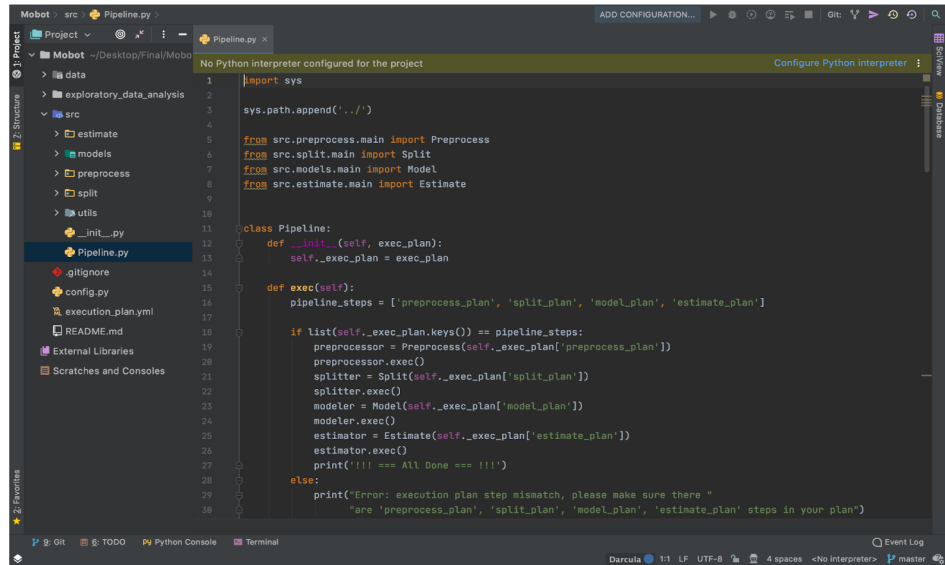


Models
Summary
Data with
Evaluation

Final Output Data Set Example

	model_name	response	predictors	criteria	rmse
0	AIC	['Recovery Rate']	Physicians.density..physicians.1.000.population./Life expectancy at birth (years)/GDP - per capita (PPP) (US\$)		0.259886543
1	SimpleLm	['Recovery Rate']	Health.expenditures....of.GDP./Literacy..../Physicians.density..physicians.1.000.population./Obesity - adult prevalence rate (%)/Life expectancy at birth (years)/H_bed_density/Imigrate_Rate/Pop_Density/GDP - per capita (PPP) (US\$)/Unemployment rate (%)		0.26027653
2	StepWise_0.1	['Recovery Rate']	Physicians.density..physicians.1.000.population./Life expectancy at birth (years)/GDP - per capita (PPP) (US\$)	p_value 0.1	0.259886543
3	StepWise_0.0 1	['Recovery Rate']	Life expectancy at birth (years)	p_value 0.01	0.259581175
4	StepWise_0.0 2	['Recovery Rate']	Life expectancy at birth (years)	p_value 0.02	0.259581175

Code Demo



```
1 import sys
2
3 sys.path.append('../')
4
5 from src.preprocess.main import Preprocess
6 from src.split.main import Split
7 from src.models.main import Model
8 from src.estimate.main import Estimate
9
10
11 class Pipeline:
12     def __init__(self, exec_plan):
13         self._exec_plan = exec_plan
14
15     def exec(self):
16         pipeline_steps = ['preprocess_plan', 'split_plan', 'model_plan', 'estimate_plan']
17
18         if list(self._exec_plan.keys()) == pipeline_steps:
19             preprocessor = Preprocess(self._exec_plan['preprocess_plan'])
20             preprocessor.exec()
21             splitter = Split(self._exec_plan['split_plan'])
22             splitter.exec()
23             modeler = Model(self._exec_plan['model_plan'])
24             modeler.exec()
25             estimator = Estimate(self._exec_plan['estimate_plan'])
26             estimator.exec()
27             print('!!! == All Done == !!!')
28         else:
29             print("Error: execution plan step mismatch, please make sure there "
30                   "are 'preprocess_plan', 'split_plan', 'model_plan', 'estimate_plan' steps in your plan")
```