

# Low Level Design (LLD)

## Insurance Premium Prediction

Revision Number – 1.0

Last Date of Revision: 16/05/2024

## Document Version Control

Date	Version	Description	Author
10/05/2024	1.0	Introduction Architecture	Sakshi Choube
12/05/2024	1.1	Data Preprocessing	Sakshi Choube
16/05/2024	1.2	Deployment	Sakshi Choube

# Contents

## Introduction

1. Why this Low-Level Design Document?	5
Architecture	5
2. Architecture Description	5
2.1 Data Gathering	5
2.2 Data Description	5
2.3 Data Ingestion	8
2.4 Data Transformation	8
2.5 Modelling	8
2.6 Batch Prediction	8
2.7 Training & Prediction Pipeline	8
2.8 UI Integration	9
2.9 Data from User	9
2.10 Data Validation	9
2.11 Rendering Result	9
3. Deployment	9
3.1 Unit Test	10

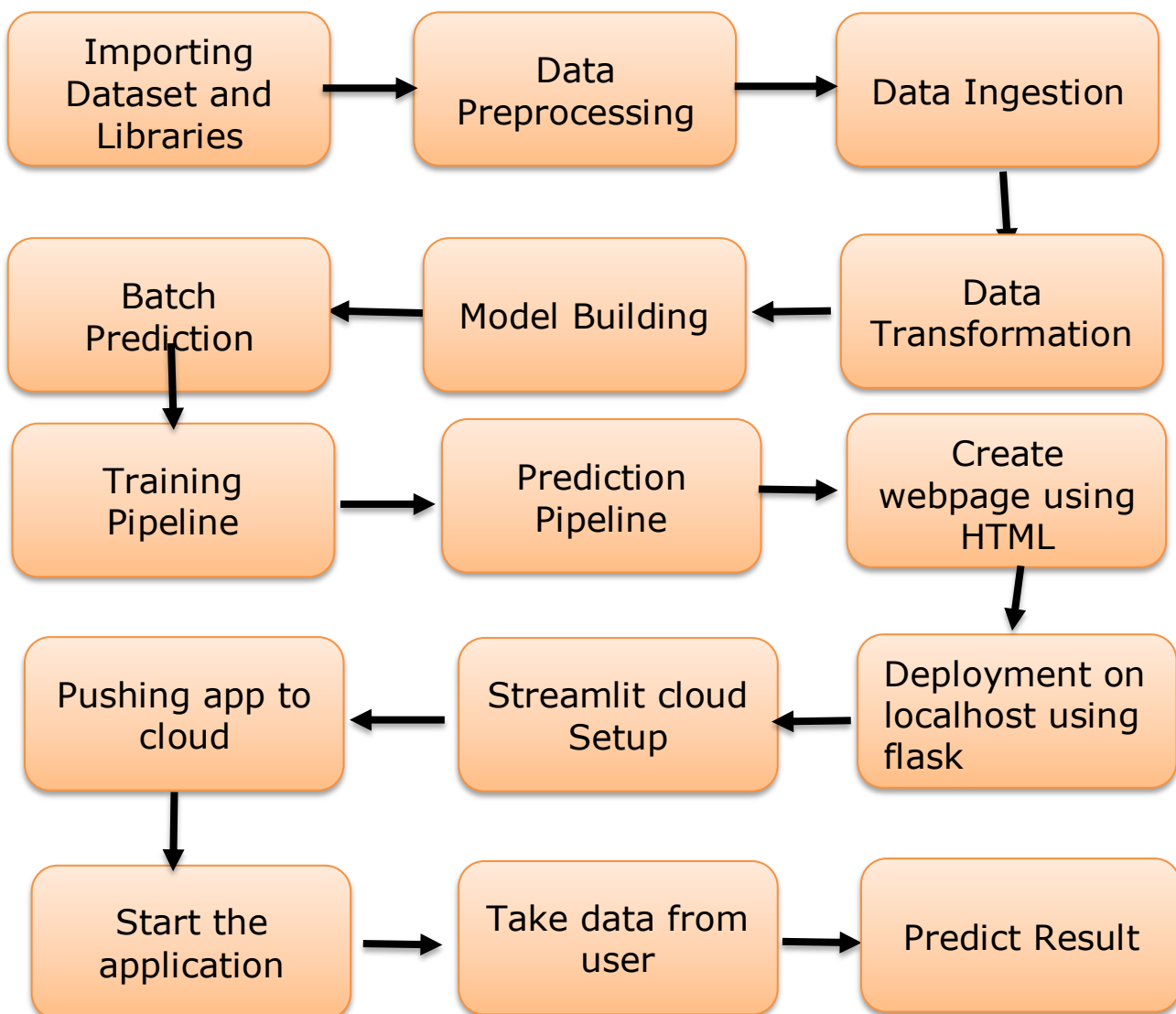
## Introduction

### 1. Why this Low-Level Design Document?

The main purpose of this LLD documentation is to feature the required details of the project and supply the outline of the machine learning model and also the written code. This additionally provides the careful description on however the complete project has been designed end-to-end.

T

## Architecture



## 2. Architecture Description

### 2.1. Data Gathering

The data for project is available on Kaggle dataset. This dataset serves as the primary source for our project's data analysis and premium prediction tasks.

---

### 2.2. Data Description

We have train (1070) and test (268) data set, train data set has both input and output  
Columns are:

- Age: Age of the insured individuals, a key factor in premium calculation due to its influence on health risk.
- Sex: Gender of the insured individuals, impacting premiums based on gender-specific health risks.
- BMI (Body Mass Index): Measure of body weight relative to height, influencing premiums based on health implications.
- Children: Number of dependents covered, affecting policy costs due to family size.
- Smoker: Smoking status (yes/no), a significant factor in premium pricing due to health risks associated with smoking.
- Region: Geographic location of the insured, which can affect healthcare costs and insurance pricing.
- Expenses: Actual medical expenses incurred, providing insight into healthcare costs and utilization for premium calculation.

variable to be predicted.

age	sex	bmi	children	smoker	region	expenses
19	female	27.9	0	yes	southwest	16884.92
18	male	33.8	1	no	southeast	1725.55
28	male	33	3	no	southeast	4449.46
33	male	22.7	0	no	northwest	21984.47
32	male	28.9	0	no	northwest	3866.86
31	female	25.7	0	no	southeast	3756.62
46	female	33.4	1	no	southeast	8240.59
37	female	27.7	3	no	northwest	7281.51
37	male	29.8	2	no	northeast	6406.41
60	female	25.8	0	no	northwest	28923.14
25	male	26.2	0	no	northeast	2721.32
62	female	26.3	0	yes	southeast	27808.73
23	male	34.4	0	no	southwest	1826.84
56	female	39.8	0	no	southeast	11090.72
27	male	42.1	0	yes	southeast	39611.76
19	male	24.6	1	no	southwest	1837.24
52	female	30.8	1	no	northeast	10797.34
23	male	23.8	0	no	northeast	2395.17
56	male	40.3	0	no	southwest	10602.39
30	male	35.3	0	yes	southwest	36837.47
60	female	36	0	no	northeast	13228.85
30	female	32.4	1	no	southwest	4149.74

## 2.3. Data Ingestion

The cornerstone of project was established through a systematic process of data acquisition and ingestion. Utilizing Kaggle, a reputable platform renowned for its high-quality datasets, we identified and acquired the crucial data required for our Insurance price prediction project. This dataset, integral to our goal of accurate price forecasting, was meticulously downloaded and securely stored within our local system infrastructure. Subsequently, we initiated the data ingestion phase, where the dataset seamlessly integrated into our project's data pipeline. This meticulous approach ensures that our project is built upon a solid foundation, setting the stage for robust and precise price prediction models and analysis.

## 2.4. Data Transformation

- Steps performed in pre-processing are:
  - First read data from Artifact folder
  - Checking unnecessary columns
  - One column has product id which is unique for every product so I deleted that column.
  - Checked for null values
  - there are too many null values are present in two columns that's why I deleted them
  - Performed one-hot encoder on categorical columns.
  - Scaling is performed for needed information.
  - And, the info is prepared for passing to the machine learning formula

## 2.5. Modelling

The pre-processed information is then envisioned and every one the specified insights are being drawn. though from the drawn insights, the info is at random unfold however still modelling is performed with completely different machine learning algorithms to form positive we tend to cowl all the chances. and eventually, Gradient Boosting performed well .

## 2.6. Batch Prediction

In the pursuit of creating a comprehensive and efficient system, we have successfully executed batch prediction as a pivotal component of our project. Leveraging a meticulously designed data transformation pipeline, we have harnessed the power of our predictive model to generate accurate and timely batch predictions. This milestone signifies the culmination of our efforts in seamlessly processing and analyzing data, resulting in actionable insights that drive informed decision-making. As we prepare our Low-Level Design Document, this achievement underscores the significance of our data transformation pipeline and predictive model, which will be elaborately detailed to ensure clarity and scalability in our system architecture.

## 2.7. Training And Prediction Pipeline

The training pipeline serves as the backbone for developing our predictive models Meanwhile, the prediction pipeline enables us to seamlessly apply these trained models to new data, ensuring that our insights and forecasts remain consistently accurate and adaptable to real world scenarios. This dual pipeline approach embodies our commitment to providing a comprehensive, data-driven solution that empowers decision- makers with the most reliable and up-to-date information. As we delve into the creation of our Low-Level Design Document, we will intricately detail these pipelines, showcasing their sophistication and efficiency in our system architecture.



## 2.8. UI Integration

Both CSS and HTML files are being created and are being integrated with the created machine learning model. All the required files are then integrated to the app.py file and tested locally

## 2.9 Data from User

The data from the user is retrieved from the created HTML web page and Streamlit application.

## 2.10 Data Validation

The data provided by the user is then being processed by app.py and application.py (streamlit application) file and validated. The validated data is then sent for the prediction.

## 2.11 Rendering Result

The data sent for the prediction is then rendered to the web page.

## Deployment

The tested model is then deployed on local machine and Streamlit Cloud. So, users can access the project from any internet devices.

### 3.1 Unit Test

Test Case	Description Pre-Requisite	Expected Result
Verify whether the Application URL is accessible to the user	1. Application URL should be defined	Application URL should be accessible to the user
Verify whether the Application Loads completely for the user when the URL is accessed.	1. Application URL accessible 2. Application deployed	The Application should load completely for the user when the URL is accessed
Verify whether user is able to edit all input fields	1. Application is accessible 2. User is signed up to the application 3. User is logged in to the application	User should be able to edit all input fields
Verify whether user gets Submit button to submit the inputs	1. Application is accessible 2. User is signed up to the application 3. User is logged in to the application	User should get submit button to submit the inputs
Verify whether user is presented with Predicted results on clicking submit	1. Application is accessible 2. User is signed up to the application 3. User is logged in to the application	User should be present with Predicted results on clicking submit