# LAB MID CAT

Name : Yashika kudesia    Registration No: 23MEC011

# DATA ROBOT AN AI TOOL FOR ML
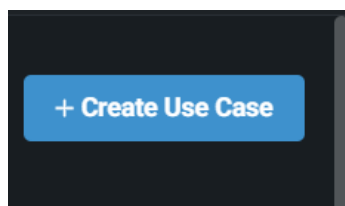
## Q:- WHAT IS DATA ROBOT ?

**Ans :-** **DataRobot is an AI platform used to build the predictive models enterprises depend on to accelerate their growth. The platform draws from a number of open-source machine learning R and Python-based libraries, including scikit-learn, H2O, TensorFlow, Vowpal Wabbit, Spark ML, and XGBoost. But with DataRobot's simple, drag-and-drop web-based interface, building and deploying sophisticated predictive models is a breeze...even for business analysts with little-to-no knowledge of machine learning or programming. By automating the selection of the ideal features, algorithms, and parameter values for building each model, the software supports best practices for new users. Meanwhile, the platform remains both flexible and extensible.**

## Q:- HOW IT WORKS ?

**Ans :-** **Here are the following steps that involve in**

**Step 1 :** **Get started with Creating the project**


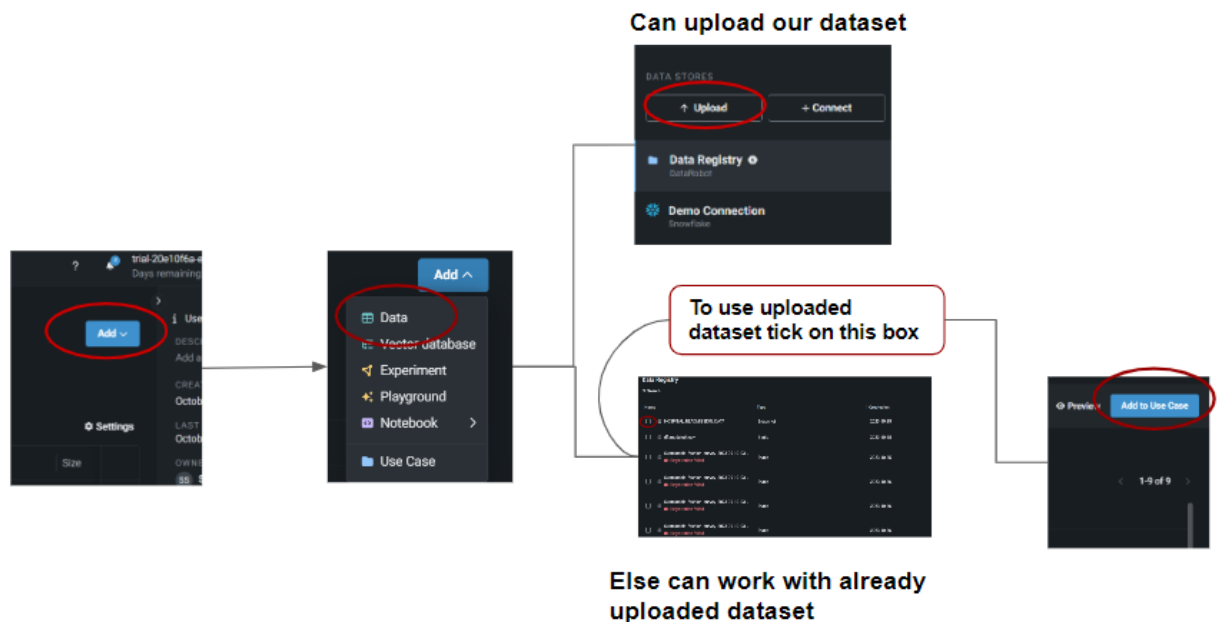
**Step 2 :** **Prepare the Dataset**

**To prep the data using DataRobot we can import a local dataset or can connect to an external data source.**
**To complete the quickstart, you first log in to DataRobot Data Prep.**
**Once you log in, complete these steps:**
- ➢ **Click on ADD**
- ➢ **Then Data**
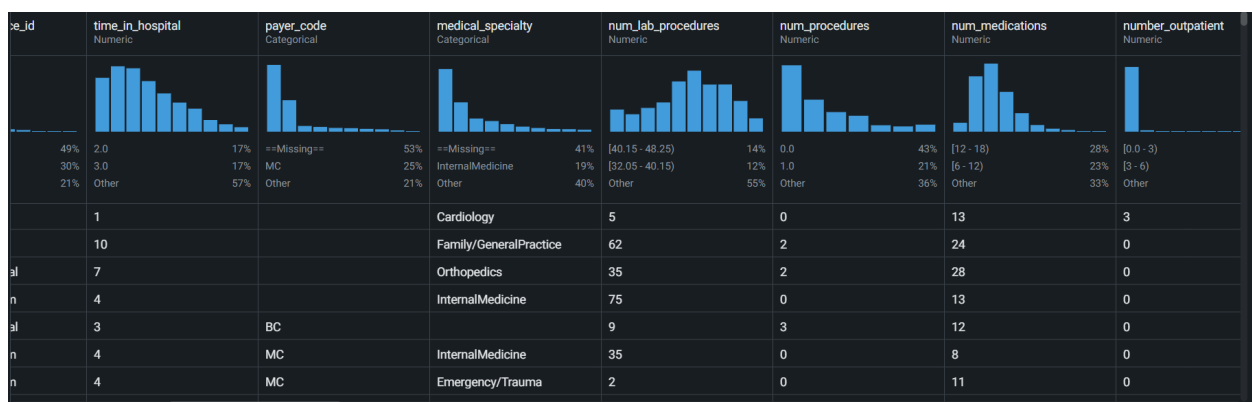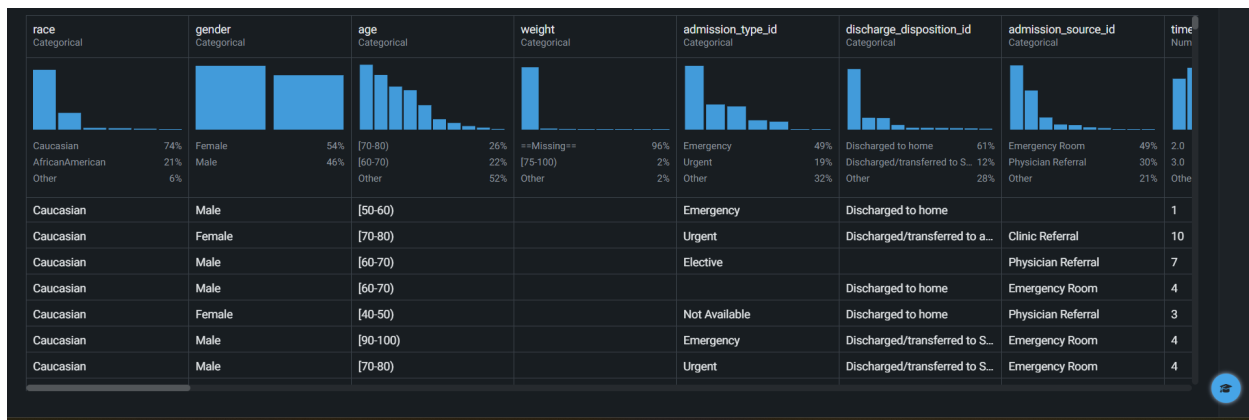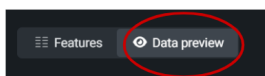- ➢ **Upload(can publish your dataset)  or Use uploaded dataset.**

**The hospital readmission dataset for analysis**

Can upload our dataset

To use uploaded dataset tick on this box

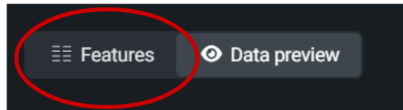Else can work with already uploaded dataset

## Step 3 : Data visualization

**To analyze the dataset, can directly click the dataset.**
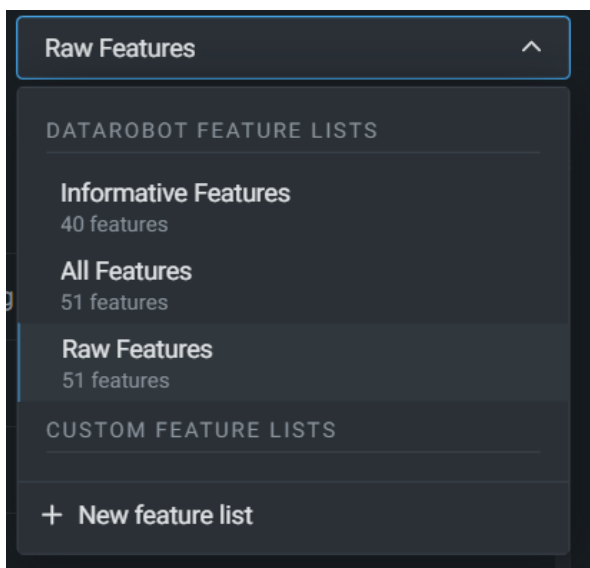**There can check the histogram and detailed values of the data**



| race Categorical | gender Categorical | age Categorical | weight Categorical | admission_type_id Categorical | discharge_disposition_id Categorical | admission_source_id Categorical | time Num |
|---|---|---|---|---|---|---|---|
| Caucasian 74% | Female 54% | [70-80] 26% | ==Missing== 96% | Emergency 49% | Discharged to home 61% | Emergency Room 49% | 2.0 |
| AfricanAmerican 21% | Male 46% | [60-70] 22% | [75-100] 2% | Urgent 19% | Discharged/transferred to S... 12% | Physician Referral 30% | 3.0 |
| Other 6% | | Other 52% | Other 2% | Other 32% | Other 28% | Other 21% | Othe |
| Caucasian | Male | [50-60] | | Emergency | Discharged to home | | 1 |
| Caucasian | Female | [70-80] | | Urgent | Discharged/transferred to a... | Clinic Referral | 10 |
| Caucasian | Male | [60-70] | | Elective | | Physician Referral | 7 |
| Caucasian | Male | [60-70] | | | Discharged to home | Emergency Room | 4 |
| Caucasian | Female | [40-50] | | Not Available | Discharged to home | Physician Referral | 3 |
| Caucasian | Male | [90-100] | | Emergency | Discharged/transferred to S... | Emergency Room | 4 |
| Caucasian | Male | [70-80] | | Urgent | Discharged/transferred to S... | Emergency Room | 4 |



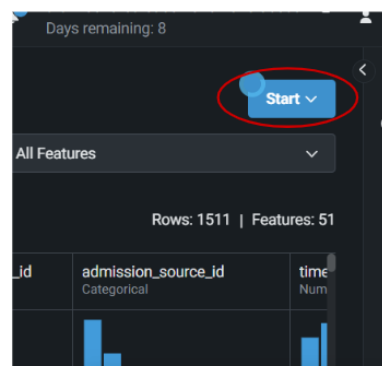| ce_id | time_in_hospital Numeric | payer_code Categorical | medical_specialty Categorical | num_lab_procedures Numeric | num_procedures Numeric | num_medications Numeric | number_outpatient Numeric |
|---|---|---|---|---|---|---|---|
| 49% | 2.0 17% | ==Missing== 53% | ==Missing== 41% | [40.15 - 48.25] 14% | 0.0 43% | [12 - 18] 28% | [0.0 - 3) |
| 30% | 3.0 17% | MC 25% | InternalMedicine 19% | [32.05 - 40.15) 12% | 1.0 21% | [6 - 12) 23% | [3 - 6) |
| 21% | Other 57% | Other 21% | Other 40% | Other 55% | Other 36% | Other 33% | Other |
| | 1 | | Cardiology | 5 | 0 | 13 | 3 |
| | 10 | | Family/GeneralPractice | 62 | 2 | 24 | 0 |
| al | 7 | | Orthopedics | 35 | 2 | 28 | 0 |
| n | 4 | | InternalMedicine | 75 | 0 | 13 | 0 |
| al | 3 | BC | | 9 | 3 | 12 | 0 |
| n | 4 | MC | InternalMedicine | 35 | 0 | 8 | 0 |
| n | 4 | MC | Emergency/Trauma | 2 | 0 | 11 | 0 |

**From here we can check different type of feature in our dataset**
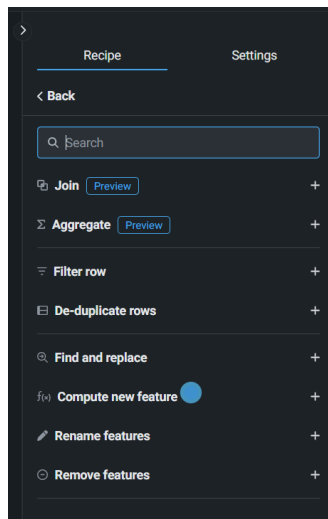


## Step 4 : Modeling

**When we click on start we will get the two option**
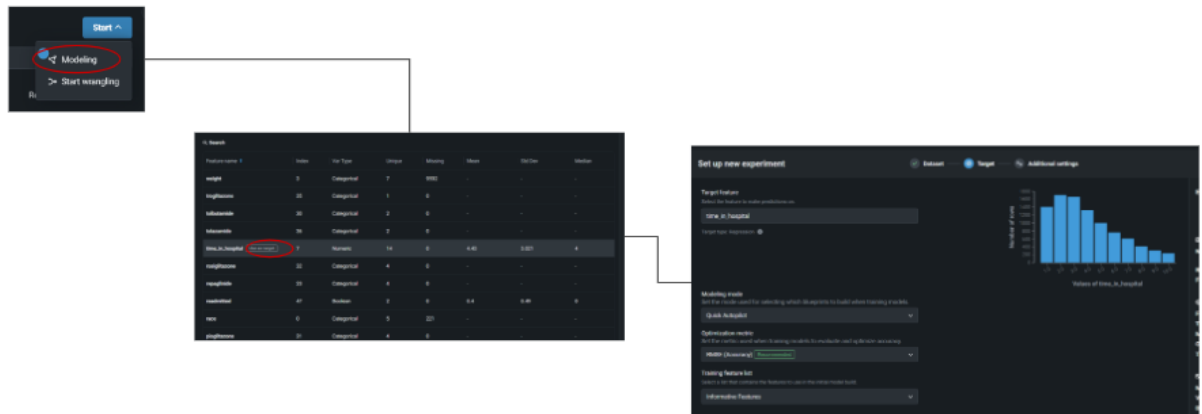
➢ **Modeling**

➢ **Wrangling**

**On clicking wrangling :** Start wrangling to pull a random sample of data from the data source and begin transformation operations.

**Click Add operation :** to build a wrangling "recipe." Each new operation updates the live sample to reflect the transformation. Note that if you wrangle your training dataset, you will want to apply the same operations to your scoring dataset to ensure you have the same columns.



**On clicking Modeling : we will get the following function**
  ➢ **Click on modeling**
  ➢ **Set the target**
  ➢ **Change the mode of model(if needed)**



  ➢ **Click on additional settings**
  ➢ **In data partitioning**

- Select the partitioning method
- Select the validation type according to your need
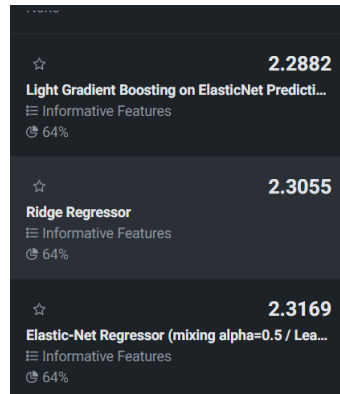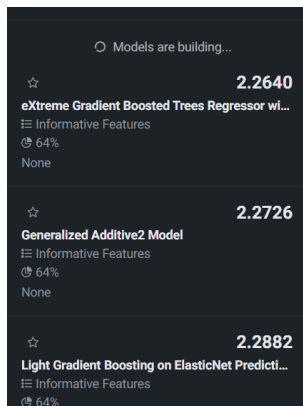


➤ **In additional settings**



➤ **Click start modeling**
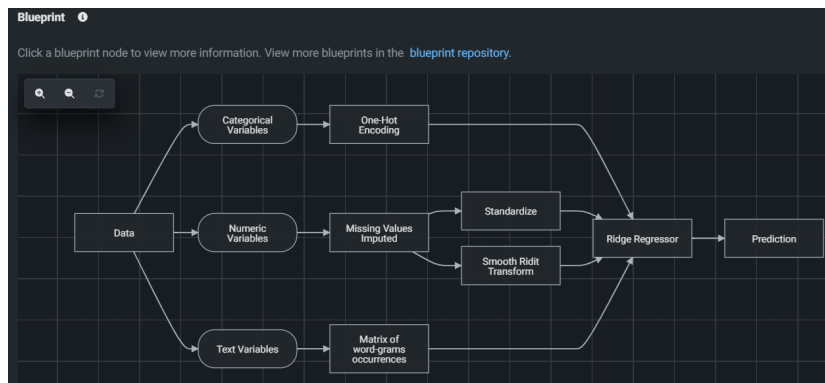
## Step 5 : Model and Outputs

Once modeling starts, Workbench begins to construct a model Leaderboard. Ultimately, DataRobot will select and retrain the most accurate model and mark it as prepared for deployment. Since the process takes some time, click on any completed model and familiarize yourself with the insights available for model evaluation.
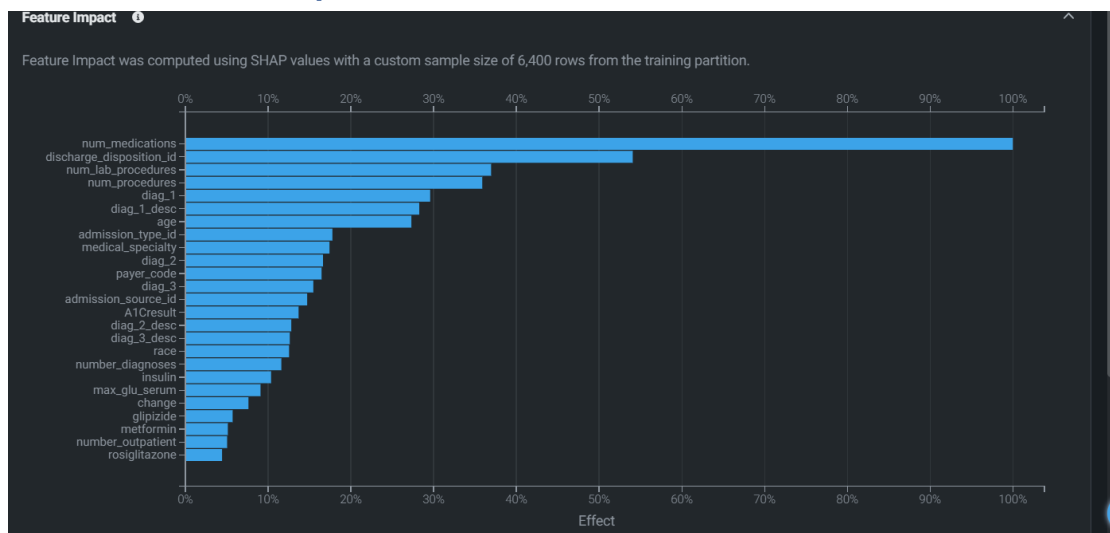
## ➢ Select the model
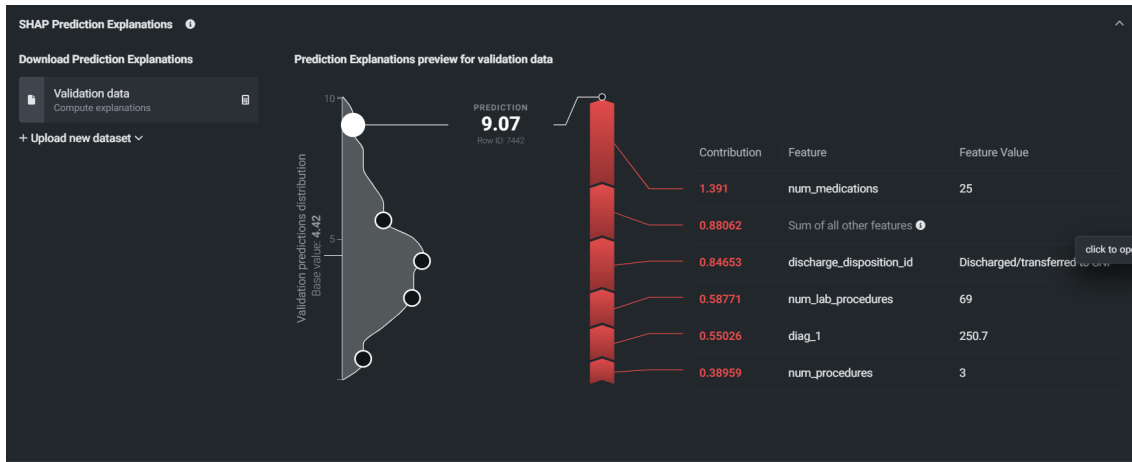


## ➢ After selection the model can see the following feature
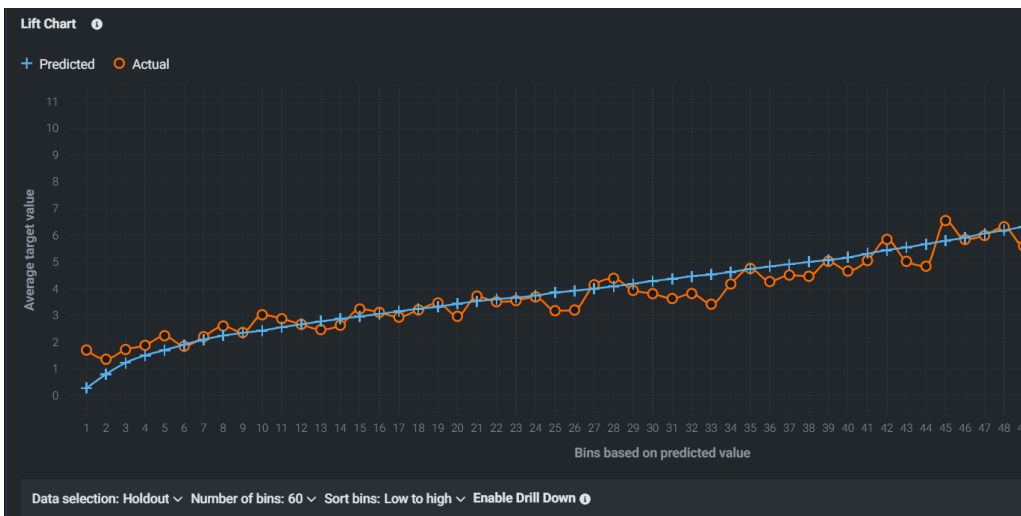
### ○ Blueprint of the model



### ○ Feature Impact

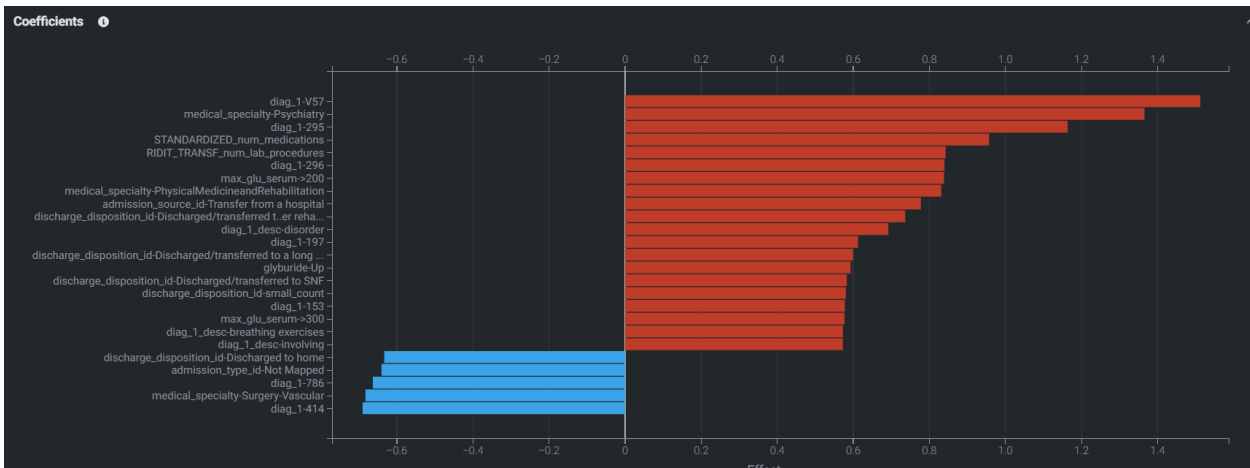- ○ **Prediction Explanation**



- ○ **Lift chart**



**The value of the number of bins, data selection and sort bin can be change**
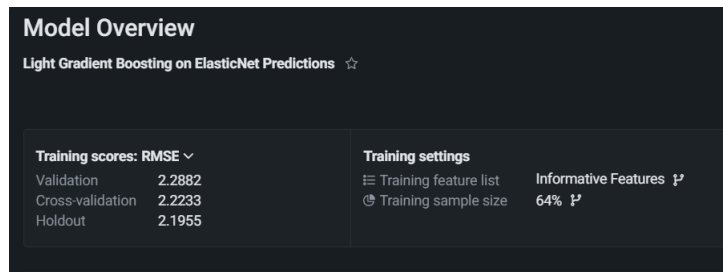
- ○ **Residuals**



- ○ **Coefficient**
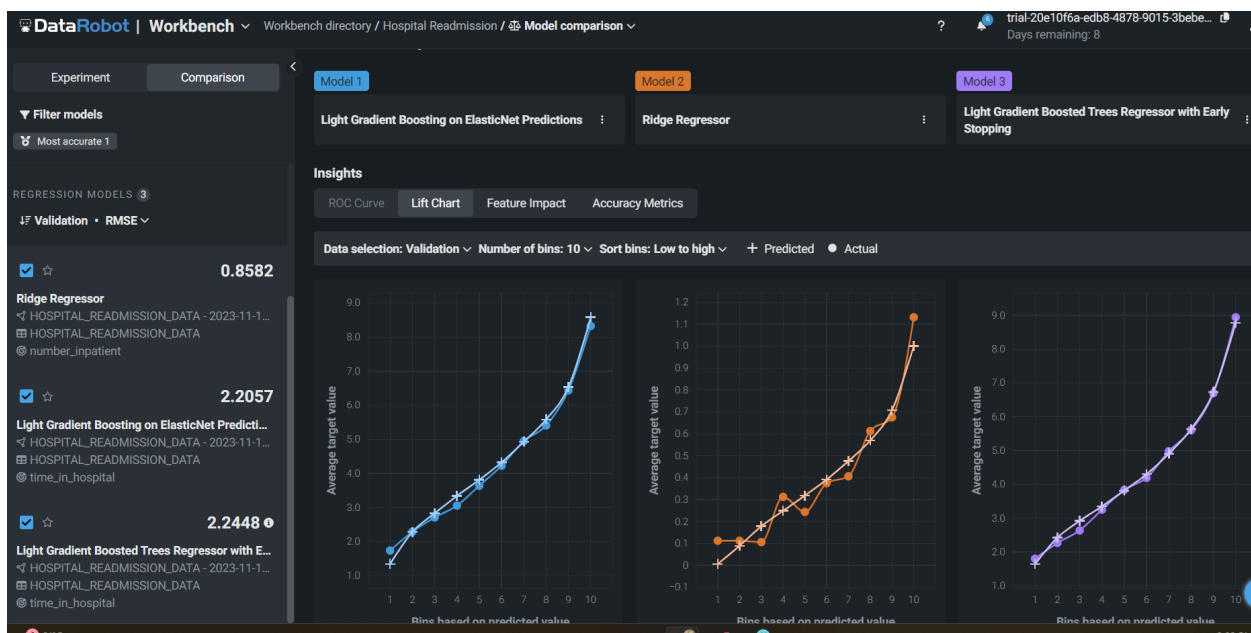


- ➢ **On the partition (validation, cross validation & holdout)**
  - ○ **Can select the matrix according to need**

○ **Model overview**

**Model Overview**

**Light Gradient Boosting on ElasticNet Predictions** ☆

| Training scores: RMSE ∨ | | Training settings | |
|---|---|---|---|
| Validation | 2.2882 | Training feature list | Informative Features |
| Cross-validation | 2.2233 | Training sample size | 64% |
| Holdout | 2.1955 | | |

# In this different type of model and their prediction can also be compared



**Lineage**

| Datasets | Experiments | Model blueprints | **Model info** |

Hide lineage values that are the same for selected models.

| | Model 1 | Model 2 | Model 3 |
|---|---|---|---|
| **Blueprint description** | One-Hot Encoding | Missing Values Imputed | Standardize | Matrix of word-grams occurrences | Ordinal encoding of categorical variables | Ridge Regressor | Light Gradient Boosting on ElasticNet Predictions **View blueprint** | One-Hot Encoding | Missing Values Imputed | Standardize | Smooth Ridit Transform | Matrix of word-grams occurrences | Ridge Regressor **View blueprint** | Ordinal encoding of categorical variables | Missing Values Imputed | Converter for Tex Mining | Auto-Tuned Word N-Gram Text Modeler using token occurrences | Light Gradient Boosted Trees Regressor with Earl Stopping **View blueprint** |
| **Blueprint family** | Light Gradient Boosting on ElasticNet Predictions | Ridge Regressor | Light Gradient Boosted Trees Regressor wit Early Stopping |
| **Model size** | 3.29 MB | 905 KB | 4.53 MB |
| **Sample size** | 64% 6400 of 10000 rows | 64% 6400 of 10000 rows | 100% 10000 of 10000 rows |
| **Time to predict 1,000 rows** ⓘ | 0.5016 ms | 0.5790 ms | 0.3502 ms |
| **Properties** | | | |