

---

# AI-ADDIN *Plugin for Interpretable Machine Learning*

---

A PREPRINT

**Mengying Wang**  
Computer System Engineering  
Northeastern University  
NUID: 001357559  
wang.mengyin@husky.neu.edu

**Ruisi Gu**  
Information System  
Northeastern University  
NUID: 001816641  
gu.ru@husky.neu.edu

April 28, 2019

## ABSTRACT

Despite widespread of automates analytical model building, most of the methods are black box models, which may lead to the situation that people hardly trust a machine learning prediction therefore not be able to make any change due to the result. However, the aim of utilizing machine learning method to predict a situation is to make some improvements and get a better result.

Our goal is to build a prototype which suitable for different types of datasets, this prototype is able to interpret models and data in a more understandable way. We provide several models to fit the dataset with the help of H2O libraries, including generalized linear model, logistic regression, decision tree and gradient boosting model. After model training, we demonstrate the interpretability by three different plots, they are variable importance, partial dependence plot and individual conditional expectation (ICE). Finally, we compare our model by matrix AUC and select the best model. Based on the interpretable plots from the best model, we are able to interpret the whole dataset. With the completion of prototype based on Amazon Reviews, we test this framework on a different dataset, Yelp Reviews, in order to validate the feasibility of prototype[1].

**Keywords** Machine Learning Interpretability · PDP · ICE · H2O

## 1 Introduction

Machine Learning is rapidly growing area, new methods and research coming up every day. However, with machine learning continuously infiltrate into people's lives, we have to deal with the problem that how to let people trust a machine learning prediction. Another problem we are facing is that how people could figure out a suitable method if all they can get from a machine learning model is a simple prediction. For example, if you go to a doctor for obesity, it's hard to believe him if the doctor only gives you the conclusion that you will gain more weight. Beside the prove for this conclusion, you may also want to know the reason behind and also the method to prevent from gaining weight[2].

On the aspect of training model, we take the advantage of H2O library and run the leaderboard of models which can fit the dataset best. Beside some advanced models such as GBM and XGBoost, we also train some basic models like GLM and Logistic Regression. These models may seem easy to train and use, but very useful and powerful regards fitting various types of datasets. Separating data into training and testing set would allow us to compare the models by calculating AUC scores.

For model interpretability, we demonstrate our dataset in three different ways, which are variable importance, partial dependence plot and individual conditional expectation (ICE). Variable importance plot provides a list of the most significant variables in descending order by a mean decrease in Gini. The top variables contribute more to the model than the bottom ones and also have high predictive power in classifying default and non-default customers[3]. A partial dependence plot can show whether the relationship between the target and a feature is linear, monotonous or more complex[4]. Individual Conditional Expectation (ICE) plots display one line per instance that shows how the instance's prediction changes when a feature changes[4] [5].

For each model, we generate three plots that mentioned above. Based on two matrixes, we select the best model and draw our conclusion by analyzing the plots we obtain.

## 2 Methods

### 2.1 Data Processing

The tokenize function will split the reviews into words and remove any stop words, small words, or words with numbers in them. Then we group the vectors of similar words together in vector space by using Word2vec. Word2vec can make highly accurate guesses about a word's meaning based on past appearances. Combining aggregated word embeddings columns with the original dataset.

### 2.2 Training Models

In order to get a general idea of our dataset within limit parameters and minimize processing runtime, we can take advantage of Auto ML in H2O. With the help of Auto ML, we are able to find the “best” model from building large number of models without any prior knowledge. H2O's AutoML can be used for automating the machine learning workflow, which includes automatic training and tuning of many models within a user-specified time-limit[6].

**Generalized Linear Model** Since we already change our target into binomial class variable, we take family='binomial', model\_id='glm\_surrogate' as our first model. We use H2OGeneralizedLinearEstimator to initialize model and train the model based on our X (predictors), y (response) and training dataset.

**Logistic Regression** For logistic regression model, we initialize by using model\_id="glm\_logistic" in H2OGeneralizedLinearEstimator. And based on the same X, y and training dataset.

**GBM I** In Gradient Boosting Model I, we generate the model based on hyperparameter as follow: ntrees=1, sample\_rate=1, col\_sample\_rate=1, max\_depth=3. However, model performance of GBM I is poor. That's why we adjust our model hyperparameter and generate another GBM model, GBM II.

**GBM II** In Gradient Boosting Model II, instead of setting GBM hyperparameters ourselves, we use the model\_id="gbm.hex" initializing by H2OGradientBoostingEstimator. Beside that, we set some stop rules as below, stopping\_metric = "AUC", stopping\_tolerance = 0.001, stopping\_rounds = 5, score\_tree\_interval = 10.

### 2.3 Model Interpretability

In this section, we selected three basic Model-Agnostic Methods - Variable Importance or standardized coefficient plot, Partial Dependence Plot (PDP) and Individual Conditional Expectation (ICE) to make interpretability analysis for our model.

**Variable Importance/ Standardized Coefficient** which feature affect the predict results mostly deeply?

This plot shows the top ten most important features in turn for the model. We called varimp\_plot() API of H2O to output a bar type plot for them, the degree of importance is the length of the bar, the first one with a longest bar is the most important feature in this model, it means that predict results is deeply influenced by this feature, and the changing of this factor would result in the change of response variables.

**Partial Dependence Plot (PDP)** How the feature affect the predict results?

The partial dependence plot (short PDP or PD plot) shows the marginal effect one or two features have on the predicted outcome of a machine learning model[7]. A mean status can be gotten in this plot, so it shows the general tendency between this feature and predict result.

**Individual Conditional Expectation (ICE)** How each instances affect the predict results?

The difference with the partial dependence plot (PDP) is that, the Individual Conditional Expectation (ICE) plots the tendency of every instance, so that we can get the difference between different instances and the scope of all instances. ecause H2o doesn't have the package to plot it, so we use the package PyCEbox for panda's dataframe, and convert it to suit for our h2o Frame. In order to apply ICE plot under H2O environment, we disassembled the package in pandas and convert into H2O version, then save it into a py file which is easy to refer. View the code in Appendix A.

### 3 Results

#### 3.1 Leaderboard

Figure 1 and Figure 2 are the output of H2O's AutoML, it can be used for automating the machine learning workflow, which includes automatic training and tuning of many models within a user-specified time-limit. With the help of Auto ML LeaderBoard, we are able to find the "best" model from building large number of models without any prior knowledge[6].

model_id	auc	logloss	mean_per_class_error	rmse	mse
XGBoost_1_AutoML_20190426_155539	0.896382	0.379217	0.204441	0.343743	0.118159

Figure 1: leaderboard of Prototype dataset

model_id	auc	logloss	mean_per_class_error	rmse	mse
XGBoost_1_AutoML_20190426_140523	0.846889	0.436274	0.275104	0.374472	0.140229
GLM_grid_1_AutoML_20190426_140523_model_1	0.84419	0.44428	0.271968	0.376596	0.141824
XGBoost_grid_1_AutoML_20190426_140523_model_8	0.843375	0.440441	0.260929	0.377569	0.142559
DeepLearning_grid_1_AutoML_20190426_140523_model_1	0.841085	0.466155	0.280023	0.382047	0.14596
XGBoost_2_AutoML_20190426_140523	0.840194	0.448026	0.32472	0.380326	0.144648
XGBoost_grid_1_AutoML_20190426_140523_model_10	0.839452	0.450697	0.317931	0.380994	0.145157
XGBoost_grid_1_AutoML_20190426_140523_model_3	0.83757	0.455414	0.298487	0.383773	0.147282
XGBoost_3_AutoML_20190426_140523	0.837124	0.45012	0.307592	0.381667	0.14567
GBM_5_AutoML_20190426_140523	0.836641	0.452982	0.327585	0.383299	0.146918
XGBoost_grid_1_AutoML_20190426_140523_model_2	0.836577	0.455647	0.303127	0.383852	0.147342

Figure 2: leaderboard of testing dataset

#### 3.2 Prototype

##### 3.2.1 Variable Importance/ Standardized Coefficient

Figure 3 has four sub figures for different models' variable importance or standardized coefficient plots in Amazon Review. For each plot, the first one with the longest bar is the most important feature in its model. So we can get that, for Generalized Linear Model, "HelpfulnessNumerator" is the most important, in the same way, C2 for Logistic Regression, Summary\_C91 for Gradient Boosting Model I and C45 for Gradient Boosting Model II.

Feature	Importance
C2	1.0
C4	0.7
C9	0.68
C12	0.65
C10	0.6
C13	0.58
C1	0.58
C8	0.55
C3	0.55
C5	0.5

Variable	Importance
summary_C10	1.0
summary_C18	0.3
Q1	0.25
summary_C12	0.2
summary_C16	0.15
summary_C15	0.1
ProductID	0.0
UserID	0.0
HighPulseDisseminator	0.0
HighPulseDisseminator	0.0

Variable	Importance (approx.)
C05	0.98
C03	0.88
summary_C09	0.55
summary_C05	0.42
C09	0.38
summary_C40	0.35
summary_C13	0.32
summary_C45	0.32
C15	0.32
C07	0.30

Figure 3: Variable Importance/ Standardized Coefficient Plots for Amazon Reviews

### 3.2.2 Partial Dependence Plot

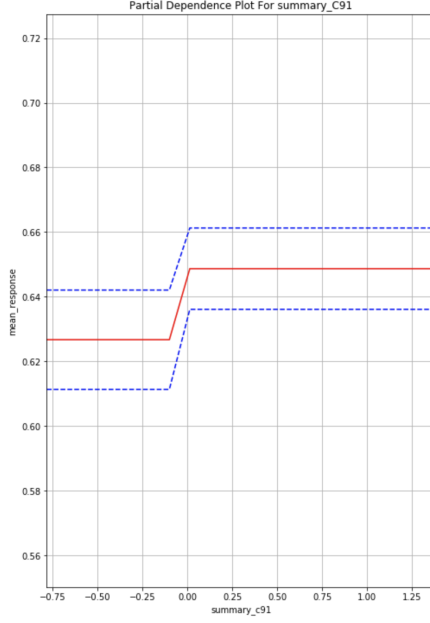
Partial Dependence Plot For HelpfulnessNumerator

The plot shows the mean response (y-axis, ranging from 0.3 to 1.1) as a function of the HelpfulnessNumerator (x-axis, ranging from 0 to 800). The solid red line represents the estimated mean response, which starts at approximately 0.6 for a HelpfulnessNumerator of 0 and increases sharply, reaching a plateau of 1.0 around a HelpfulnessNumerator of 100. The dashed blue lines represent the confidence interval, which is widest at low values of HelpfulnessNumerator and narrows as the value increases, converging to the mean response of 1.0 for values above 500.

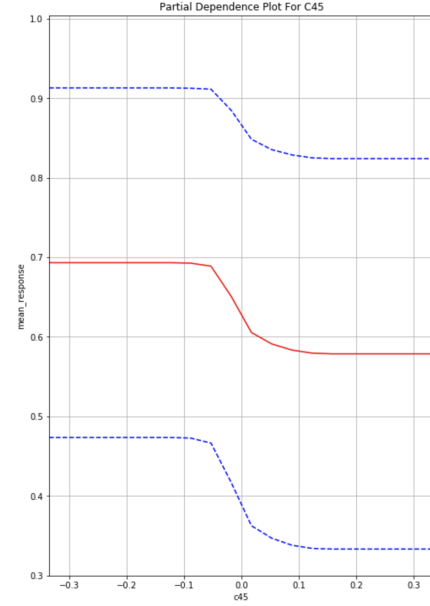
A partial dependence plot for variable  $c_2$ . The x-axis is labeled  $c_2$  and ranges from -0.4 to 0.4. The y-axis is labeled 'mean response' and ranges from 0.2 to 1.0. Three lines are plotted: a solid red line (top), a dashed blue line (middle), and a dotted blue line (bottom). All three lines show a negative linear relationship between  $c_2$  and the mean response.

$c_2$	Red Line (Solid)	Blue Line (Dashed)	Blue Line (Dotted)
-0.4	0.83	0.68	1.00
-0.3	0.77	0.58	0.97
-0.2	0.71	0.48	0.94
-0.1	0.65	0.38	0.91
0.0	0.59	0.28	0.88
0.1	0.53	0.18	0.85
0.2	0.47	0.08	0.82
0.3	0.41	0.02	0.79
0.4	0.35	0.00	0.76

4



(c) Gradient Boosting I

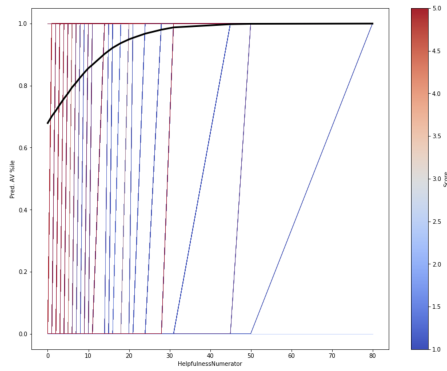


(d) Gradient Boosting II

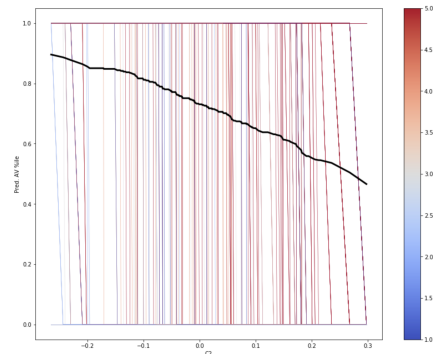
Figure 4: Partial Dependence Plots (PDP) for Amazon Reviews

### 3.2.3 Individual Conditional Expectation

Figure 5 has four sub figures for different models' Individual Conditional Expectation (ICE) in Amazon Review. Based on all these four plots, the instances in Figure 5(a), Figure 5(b) and Figure 5(c), especially the Logistic Regression. On the contrary, the instances in Figure 5(c) are focused in a small range. We also can get that the general tendency of ICE change accords with PDP in Figure 4 above.



(a) Generalized Linear



(b) Logistic Regression

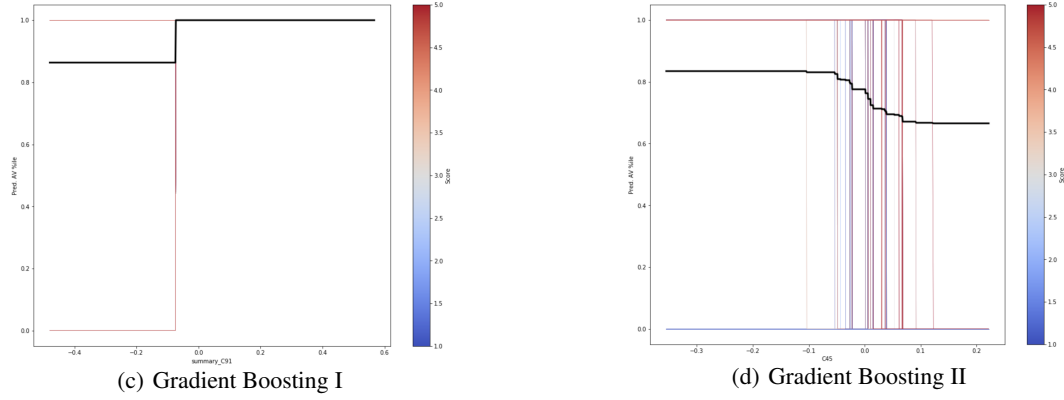


Figure 5: Individual Conditional Expectation (ICE) for Yelp Reviews

### 3.3 Test on new dataset

#### 3.3.1 Variable Importance/ Standardized Coefficient

Figure 6 has four sub figures for different models' variable importance or standardized coefficient plots in Yelp Review. We can get that, for Generalized Linear Model, "user\_id" is the most important, in the same way, C40 for Logistic Regression, C2 for Gradient Boosting Model I and Gradient Boosting Model II.



Figure 6: Variable Importance/ Standardized Coefficient Plots for Yelp Reviews

#### 3.3.2 Partial Dependence Plot

Figure 7 has four sub figures for different models' Partial Dependence Plots (PDP) in Yelp Review. We use the most important feature which we got in Figure 6 as x axis, and the mean of response variables is the y axis. In Generalized Linear Model, We can get that the positive reviews always have a "cool" larger than 20, and in Logistic Regression, the

larger the C40 is, the larger probability to get a positive review. In the last two Gradient Boosting Models, the mean of response variables change smoothly with the change of important feature C2.

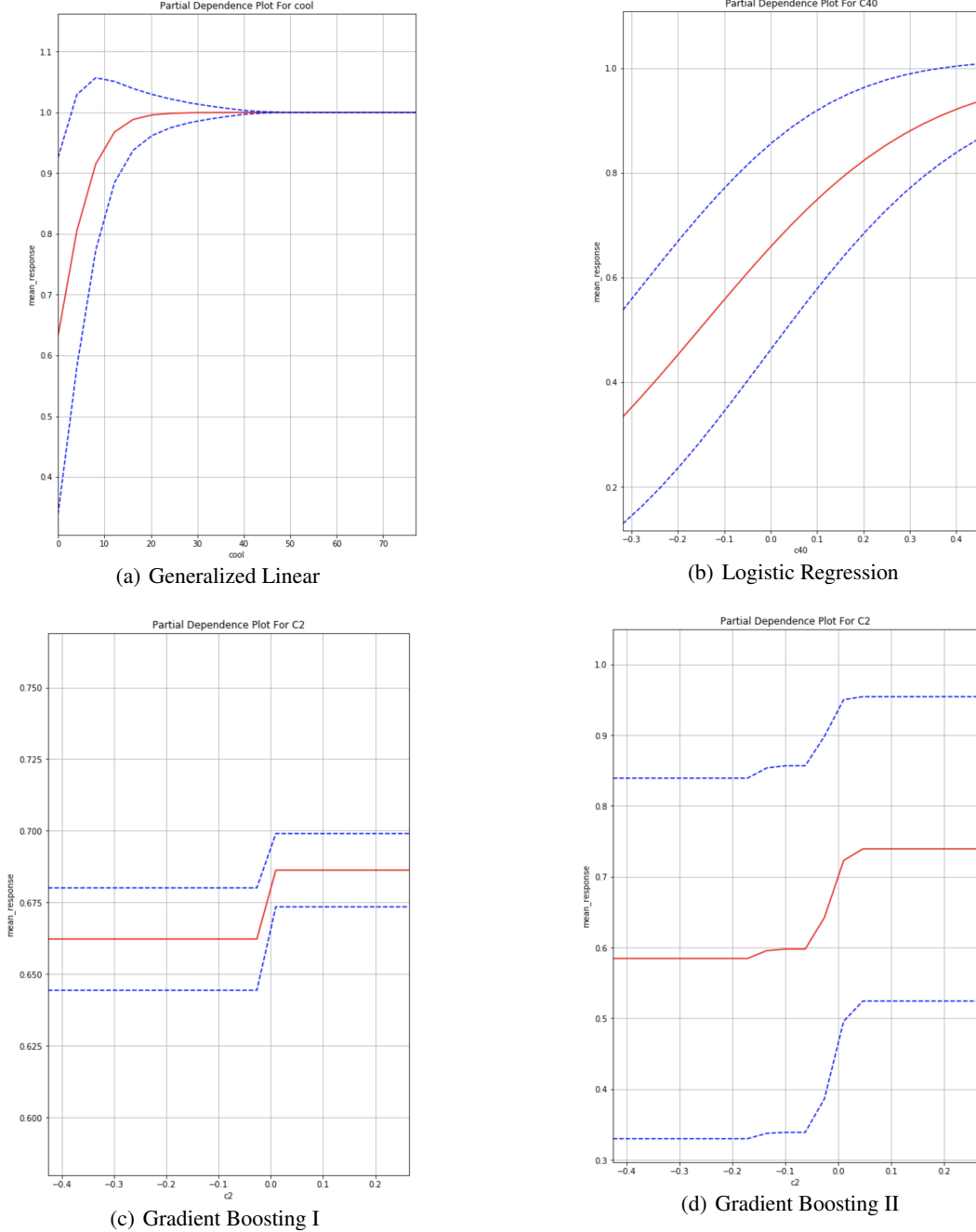


Figure 7: Partial Dependence Plots (PDP) for Yelp Reviews

### 3.3.3 Individual Conditional Expectation

Figure 8 has four sub figures for different models' Individual Conditional Expectation (ICE) in Yelp Review. We can get that the instances' distribution in each plot are similar to the Figure 5 which are ICE for Amazon Review. On the other hand, we can also get the information between different instances in each model, for example the Figure 8(b), c40 is uniformly distributed, so the the mean of the probability rises everywhere uniformly, but C2 in Figure 8(d) it is opposite.

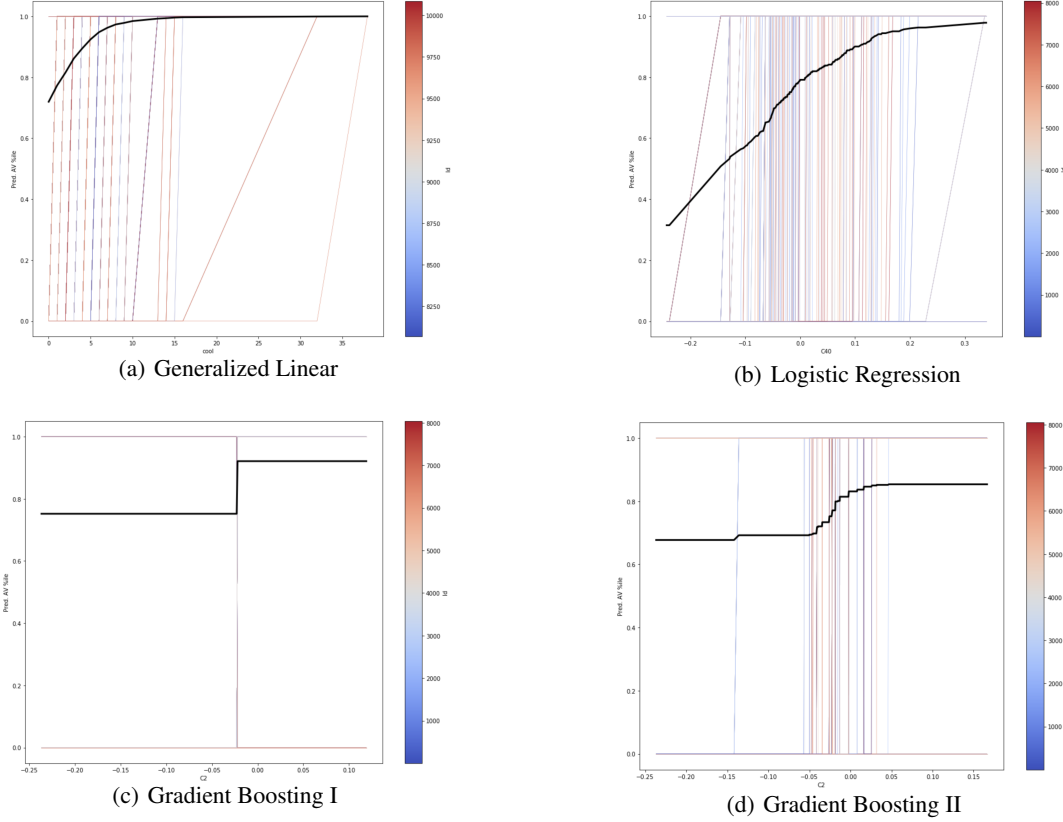


Figure 8: Individual Conditional Expectation (ICE) for Yelp Reviews

## 4 Discussion

As we can see the plots of prototype and new dataset in Result part, we can draw the conclusion that our interpretable prototype does well at least in review type of dataset, with the prove that we only slightly adjust the code to fit the new dataset and it still runs well.

### 4.1 Amazon (prototype)

#### 4.1.1 Compare and Select the Best Model

We can easily select the model which has the best performance based on matrix AUC score. In this case, Generalized Linear Model has the highest AUC score, which is 0.878. Thus, all interpretable explanation would be based on model Generalized Linear Model.

#### 4.1.2 Interpretability

As for the GLM model, the importance of predictors can be detected by standard coefficients value. Thus, “Helpfulness-Numerator” has the highest absolute value of standard coefficients which means it is the most important predictor in this case. It makes sense since in the real world, people intend to trust the comments when others willing to stand for. And from the partial dependence plot, we can draw the conclusion that people more likely to give a positive review when the number of Helpfulness is greater than 50. In ICE plots, it’s obvious that the lower the review score, the higher the number of “HelpfulnessNumerator” needs in order to get a good review.



## 4.2 Yelp (testing)

### 4.2.1 Compare and Select the Best Model

With the help of tables of AUC score, we are able to select the best model for the Yelp Reviews dataset, that is Generalized Linear Model. The AUC of GBM is 0.876, and the three interpretable plots for this model will be the one we choose to further explain.

### 4.2.2 Interpretability

As we can see from the standardized coefficient plot, predictor “user\_id” take the lead. We can understand this circumstance as some people willing to give average high score while others has a strict standards for the items they grade. From the partial dependence plot of predictor “cool”, we can conclude that the number of “cool” is making huge different when range 0 to 20. However, once the number of “cool” excess 20, it won’t be that important for the score of reviews. With the help of ICE plots, we find out that most of the “cool” number is range from 0 to 15 and evenly distributed.

## 4.3 Future

- For the prototype, we will enrich our model interpretability. In order to make our interpretability more easily and clearly, we plan to use pandas library to generate LIME, ALE and other interpretable plots.
- For testing data, since all the data we interpret so far is Reviews data. In the future, we will try some different types of dataset and try to improve our prototype also fit these types of dataset.
- In addition, we also need to create a user interface for customers. From this interface, customers only need to provide limit predictors like original dataset and would be able to obtain the interpretable plots and explanations as well.
- We used text formatting datasets in final project, in the future, we will add image samples to make our model suit for more situations.

## References

- [1] Generalized Linear Models and Mixed-Effects in Agriculture. <https://www.r-bloggers.com/generalized-linear-models-and-mixed-effects-in-agriculture/>
- [2] Machine Learning for Particle Data When You are Not a Physicist. <https://towardsdatascience.com/machine-learning-for-particle-data-when-you-are-not-a-physicist-dad77beb90e0>
- [3] Pratap Dangeti . In *Statistics for Machine Learning*, July 2017.
- [4] Christoph Molnar. Fast classification of handwritten on-line arabic characters. In *Interpretable Machine Learning*, April 2019.
- [5] Interpretable Machine Learning with Python. [https://github.com/jphall663/interpretable\\_machine\\_learning\\_with\\_python](https://github.com/jphall663/interpretable_machine_learning_with_python)
- [6] H2O AutoML. Automatic Machine Learning. <http://docs.h2o.ai/h2o/latest-stable/h2o-docs/automl.html#automl-automatic-machine-learning>
- [7] Friedman, Jerome H Greedy function approximation: A gradient boosting machine. In *Annals of statistics*, 2001: 1189-1232.

A

```

1  from __future__ import division
2
3  import six
4
5  from matplotlib import colors, cm
6  from matplotlib import pyplot as plt
7  import numpy as np
8  import pandas as pd
9  import h2o
10
11 def _get_grid_points(x, num_grid_points):
12     if num_grid_points is None:
13         return x.unique()
14     else:
15         # unique is necessary, because if num_grid_points is too much
16         # larger than x.shape[0], there will be duplicate quantiles (even
17         # with interpolation)
18         return x.quantile(np.linspace(0, 1, num_grid_points)).unique()
19
20
21 def _get_point_x_iloc(grid_index, data_index):
22     data_level = 'data_{}'.format(grid_index.name)
23
24     return (np.abs(np.subtract
25                   .outer(grid_index,
26                         data_index.get_level_values(data_level)))
27           .argmin(axis=0))
28
29
30 def _get_quantiles(x):
31     return np.greater.outer(x, x).sum(axis=1) / x.size
32
33
34 def ice(data, column, model, num_grid_points=None):
35     """
36     Generate individual conditional expectation (ICE) curves for a model.
37
38     :param data: the sample data from which to generate ICE curves
39     :type data: pandas DataFrame
40
41     :param column: the name of the column in data that will be varied to
42                   generate ICE curves
43     :type column: str
44
45     :param predict: the function that generates predictions from the model.
46                     Must accept a DataFrame with the same columns as data.
47     :type predict: callable
48
49     :param num_grid_points: the number of grid points to use for the
50                             independent variable of the ICE curves. The independent variable
51                             values for the curves will be quantiles of the data.
52
53     If None, the values of the independent variable will be the unique
54     values of data[column].
55     :type num_grid_points: None or int
56

```

```

57 :return: A DataFrame whose columns are ICE curves. The row index is the
58         independent variable, and the column index is the original data point
59         corresponding to that ICE curve.
60 :rtype: pandas DataFrame
61 """
62 data = data.as_data_frame()
63 x_s = _get_grid_points(data[column], num_grid_points)
64 ice_data, orig_column = _to_ice_data(data, column, x_s)
65 hf = h2o.H2OFrame(ice_data)
66 hfd = model.predict(hf)
67 ice_data['ice_y'] = hfd.as_data_frame()['predict'].as_matrix()
68 ice_data['data_{}'.format(column)] = orig_column
69
70 other_columns = ['data_{}'.format(column)] +
71                 [col for col in data.columns if col != column]
72 ice_data = ice_data.pivot_table(values='ice_y',
73                                 index=other_columns, columns=column).T
74
75 return ice_data
76
77
78 def ice_plot(ice_data, frac_to_plot=1.,
79             plot_points=False, point_kwargs=None,
80             x_quantile=False, plot_pdp=False,
81             centered=False, centered_quantile=0.,
82             color_by=None, cmap=None, figsize=(14,11),
83             ax=None, pdp_kwargs=None, **kwargs):
84     """
85     Plot the ICE curves
86
87     :param ice_data: the ICE data generated by :func:pycebox.ice.ice
88     :type ice_data: pandas DataFrame
89
90     :param frac_to_plot: the fraction of ICE curves to plot. If less than
91                         one, randomly samples columns of ice_data to plot.
92     :type frac_to_plot: float
93
94     :param plot_points: whether or not to plot the original data points on
95                         the ICE curves. In this case, point_kwargs is passed as keyword
96                         arguments to plot.
97     :type plot_points: bool
98
99     :param x_quantile: if True, the plotted x-coordinates are the quantiles
100                       of ice_data.index
101     :type x_quantile: bool
102
103     :param plot_pdp: if True, plot the partial dependence plot. In this
104                     case, pdp_kwargs is passed as keyword arguments to plot.
105
106     :param centered: if True, each ICE curve is centered to zero at the
107                     percentile closest to centered_quantile.
108     :type centered: bool
109
110     :param color_by: If a string, color the ICE curve by that level of the
111                     column index.
112
113                     If callable, color the ICE curve by its return value when applied to
114                     a DataFrame of the column index of ice_data
115     :type color_by: None, str, or callable

```

```

116
117 :param cmap:
118 :type cmap: matplotlib Colormap
119
120 :param figsize: size of the figure
121 :type figsize: tuple (width, height)
122
123 :param ax: the Axes on which to plot the ICE curves
124 :type ax: None or matplotlib Axes
125
126 Other keyword arguments are passed to plot
127 """
128 if not ice_data.index.is_monotonic_increasing:
129     ice_data = ice_data.sort_index()
130
131 if centered:
132     quantiles = _get_quantiles(ice_data.index)
133     centered_quantile_iloc = np.abs(quantiles - centered_quantile)
134     .argmin()
135     ice_data = ice_data - ice_data.iloc[centered_quantile_iloc]
136
137 if frac_to_plot < 1.:
138     n_cols = ice_data.shape[1]
139     icols = np.random.choice(n_cols, size=frac_to_plot * n_cols,
140                             replace=False)
141     plot_ice_data = ice_data.iloc[:, icols]
142 else:
143     plot_ice_data = ice_data
144
145
146 if x_quantile:
147     x = _get_quantiles(ice_data.index)
148 else:
149     x = ice_data.index
150
151 if plot_points:
152     point_x_iloc = _get_point_x_iloc(plot_ice_data.index,
153                                     plot_ice_data.columns)
154     point_x = x[point_x_iloc]
155     point_y = plot_ice_data.values[point_x_iloc,
156                                   np.arange(point_x_iloc.size)]
157
158
159 if ax is None:
160     _, ax = plt.subplots(figsize=figsize)
161
162 if color_by is not None:
163     if isinstance(color_by, six.string_types):
164         colors_raw = plot_ice_data.columns.get_level_values(color_by)
165         .values
166     elif hasattr(color_by, '__call__'):
167         col_df = pd.DataFrame(list(plot_ice_data.columns.values),
168                               columns=plot_ice_data.columns.names)
169         colors_raw = color_by(col_df)
170     else:
171         raise ValueError('color_by must be a string or function')
172
173     norm = colors.Normalize(colors_raw.min(), colors_raw.max())
174     m = cm.ScalarMappable(norm=norm, cmap=cmap)

```

```

175         for color_raw, (_, ice_curve) in zip(colors_raw,
176                                             plot_ice_data.iteritems()):
177             c = m.to_rgba(color_raw)
178             ax.plot(x, ice_curve, c=c, zorder=0, **kwargs)
179     else:
180         ax.plot(x, plot_ice_data, zorder=0, **kwargs)
181
182     if plot_points:
183         ax.scatter(point_x, point_y, zorder=10, **(point_kwargs or {}))
184
185     if plot_pdp:
186         pdp_kwargs = pdp_kwargs or {}
187         pdp_data = pdp(ice_data)
188         ax.plot(x, pdp_data, **pdp_kwargs)
189
190     return ax
191
192
193 def pdp(ice_data):
194     """
195     Calculate a partial dependence plot from ICE data
196
197     :param ice_data: the ICE data generated by :func:pycebox.ice.ice
198     :type ice_data: pandas DataFrame
199
200     :return: the partial dependence plot curve
201     :rtype: pandas Series
202     """
203     return ice_data.mean(axis=1)
204
205
206 def _to_ice_data(data, column, x_s):
207     """
208     Create the DataFrame necessary for ICE calculations
209
210     ice_data = pd.DataFrame(np.repeat(data.values, x_s.size, axis=0),
211                             columns=data.columns)
212     data_column = ice_data[column].copy()
213     ice_data[column] = np.tile(x_s, data.shape[0])
214
215     return ice_data, data_column

```