

1 How To Use This Document

Highly regulated industries, such as banking and insurance, must comply with government regulations for model validation before a model can be put into production. This includes creating robust model development documentation. DataRobot automates the generation of model documentation, expediting the process required for regulatory compliance and following best practice for reducing model risk.

This document is split into two components: those sections that are automatically produced by DataRobot and those that require further input by the user. The sections in *blue italicized font* include specific instructions for the documenter and require additional user input of organization-specific information, such as business use cases, data sources, and implementation details. Once the sections are complete, remove the instructions. The remaining sections in non-blue italicized font are automatically populated by DataRobot and require no further input.

Copyright ©2023, DataRobot, Inc.

The logo for DataRobot, featuring the word "Data" in a bold, black, sans-serif font, followed by "Robot" in a bold, blue, sans-serif font.

Table of Contents

- 1 How To Use This Document
- 2 Model Performance Overview
- 3 Model Features Summary
 - 3.1 Model Features and Summary Statistics
 - 3.2 Data Quality Handling Report
- 4 Model Summary and Description
- 5 Feature Effects
- 6 Feature Impact Chart
- 7 Feature Impact Table
- 8 Validation Testing and Stability
 - 8.1 Cross Validation Scores
- 9 Model Results
- 10 Bias and Fairness

2 Model Performance Overview

As an additional layer of model validity, DataRobot not only evaluated the statistical metrics underlying the model, but also performed testing on in-sample records.

The performance metric used for this project was LogLoss. The model performance results are presented below for in-sample testing:

Scoring Type	Score (LogLoss)
cross_validation	0.3616*
holdout	0.0229*
validation	0.3633*

3 Model Features Summary

Below are two tables. The first contains a list of the final set of model feature variables, as well as summary statistics for the eXtreme Gradient Boosted Trees Classifier with Early Stopping (learning rate =0.02) model. The second table contains a detailed analysis of missing values.

The Model Features and Summary Statistics table provides a brief overview of the summary statistics of model features. This includes Feature Name, variable type (Var Type), number of unique values (Unique), Number of missing values (Missing), Mean, Standard Deviation (Std Dev), Median, Minimum Value (Min), Maximum Value (Max) and Assessment of target leakage risk (Target Leakage).

3.1 Model Features and Summary Statistics

Feature Name	Var Type	Unique	Missing	Mean	Std Dev	Median	Min	Max	Target Leakage
water_tech_clean	Categorical	19	6743	N/A	N/A	N/A	N/A	N/A	Low
water_source_clean	Categorical	10	4224	N/A	N/A	N/A	N/A	N/A	Low
management_clean	Categorical	9	576	N/A	N/A	N/A	N/A	N/A	Low
age_in_years	Numeric	4898	26	9.96	10.93	5.5	-10.0043	86.83	Low
distance_to_city	Numeric	19622	0	31890.56	32610.72	17324.42	82.604	140207.46	Low
local_population	Numeric	6529	2967	7313.56	13104.7	2514.0	0.0	78465.0	Low
population_1km	Numeric	11283	0	8227.35	14024.97	3059.017	-0.0	69122.89	Low
women_of_reproductive_age_ratio_1km	Numeric	2569	0	0.26	0.014	0.26	-0.0	0.28	Low
youth_ratio_1km	Numeric	2420	0	0.203	0.021	0.19	-0.013	0.23	Low

women_of_reproductive_age_ratio_10km	Numeric	5588	0	0.26	0.014	0.26	0.22	0.28	Low
youth_ratio_10km	Numeric	5477	0	0.203	0.0205	0.19	0.17	0.23	Low
children_under_five_ratio_10km	Numeric	5455	0	0.14	0.021	0.16	0.12	0.17	Low
men_ratio_100km	Numeric	656	0	0.49	0.0064	0.49	0.47	0.5	Low
women_of_reproductive_age_ratio_100km	Numeric	655	0	0.26	0.0077	0.26	0.23	0.27	Low
youth_ratio_100km	Numeric	655	0	0.2	0.012	0.19	0.17	0.22	Low
elderly_ratio_100km	Numeric	656	0	0.068	0.0042	0.07	0.0603	0.076	Low
children_under_five_ratio_100km	Numeric	656	0	0.15	0.013	0.16	0.12	0.17	Low
years_since_report	Numeric	4003	0	8.76	11.062	4.59	-10.0043	86.83	Low

The last column in this table is an assessment of target leakage risk. DataRobot automatically tests for target leakage on a per-feature basis during the Autopilot process. Target leakage, sometimes called data leakage, occurs when a model is trained using a dataset that includes information that would not be available at the time of prediction. This can produce overly optimistic model performance results during training, given a feature will near-completely describe the target (e.g., the number of late payments on a loan as a predictor for loan default at loan application date.)

DataRobot tests for target leakage risk using Alternating Conditional Expectation (ACE) to measure the association between each feature and the target; the ACE score is normalized using the project optimization metric so that its value is in the range [0,1]. If above a certain threshold (see below), DataRobot will create a new feature list with those features flagged and possibly removed, and the user is notified by a banner in the user interface during modeling. Notably, because the definition of target leakage is directly tied with prediction time and not strength of association between a feature and the target, it's possible for DataRobot to not identify all sources of target leakage. Therefore, to reduce the risk for potential target leakage in the feature list, it's important to apply subject matter expertise.

The thresholds for target leakage risk are based on a normalized ACE score:

- High risk: > 0.975, flagged and removed
- Moderate risk: > 0.85, flagged but not removed
- Low risk: < 0.85, no action

The following table provides a summary of missing values. It includes the name of the feature, its type, a summary of the missing value count (both number of rows and as a percentage), and information on the type of imputation applied to the feature.

3.2 Data Quality Handling Report

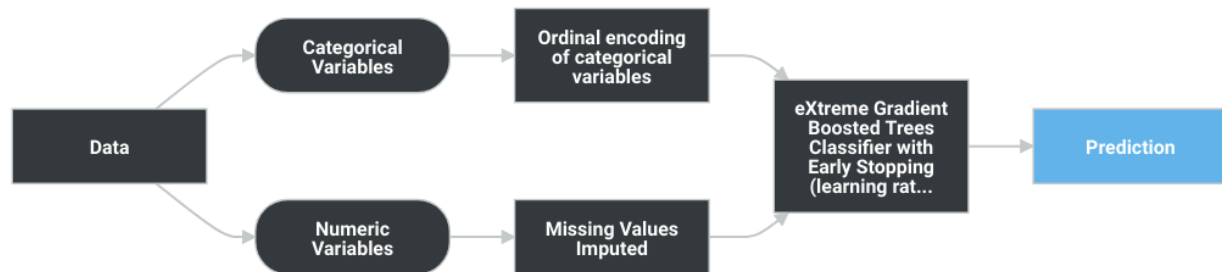
Feature Name	Var Type	Missing Count	Missing Percentage	Imputation Name	Imputation Description
water_tech_clean	Categorical	8611	26	Ordinal encoding of categorical variables	Imputed value: -2

water_source_clean	Categorical	7123	22	Ordinal encoding of categorical variables	Imputed value: -2
local_population	Numeric	3650	11	Missing Values Imputed	Imputed value: -9999
management_clean	Categorical	579	2	Ordinal encoding of categorical variables	Imputed value: -2
age_in_years	Numeric	43	0	Missing Values Imputed	Imputed value: -9999
distance_to_city	Numeric	0	0	Missing Values Imputed	Imputed value: 11852.693
population_1km	Numeric	0	0	Missing Values Imputed	Imputed value: 3905.7463
women_of_reproductive_age_ratio_1km	Numeric	0	0	Missing Values Imputed	Imputed value: 0.2628
youth_ratio_1km	Numeric	0	0	Missing Values Imputed	Imputed value: 0.1928
women_of_reproductive_age_ratio_10km	Numeric	0	0	Missing Values Imputed	Imputed value: 0.2624
youth_ratio_10km	Numeric	0	0	Missing Values Imputed	Imputed value: 0.1927
children_under_five_ratio_10km	Numeric	0	0	Missing Values Imputed	Imputed value: 0.1562
men_ratio_100km	Numeric	0	0	Missing Values Imputed	Imputed value: 0.4854
women_of_reproductive_age_ratio_100km	Numeric	0	0	Missing Values Imputed	Imputed value: 0.2634
youth_ratio_100km	Numeric	0	0	Missing Values Imputed	Imputed value: 0.208
elderly_ratio_100km	Numeric	0	0	Missing Values Imputed	Imputed value: 0.0718
children_under_five_ratio_100km	Numeric	0	0	Missing Values Imputed	Imputed value: 0.1372
years_since_report	Numeric	0	0	Missing Values Imputed	Imputed value: 2.7571

4 Model Summary and Description

The particular model referenced in this document: eXtreme Gradient Boosted Trees Classifier with Early Stopping (learning rate =0.02). This model was developed in a project created with v7686f5bec5c9e5ac of DataRobot. This model is denoted within DataRobot by the Project ID: 63ff9a625bbdbb837a8d124d and the Model ID: 640966ced23c55a6d6fdfa3d. The project was created on 2023-03-01 18:33:06.

The model development workflow process (i.e., the model blueprint) is detailed in the figure below.



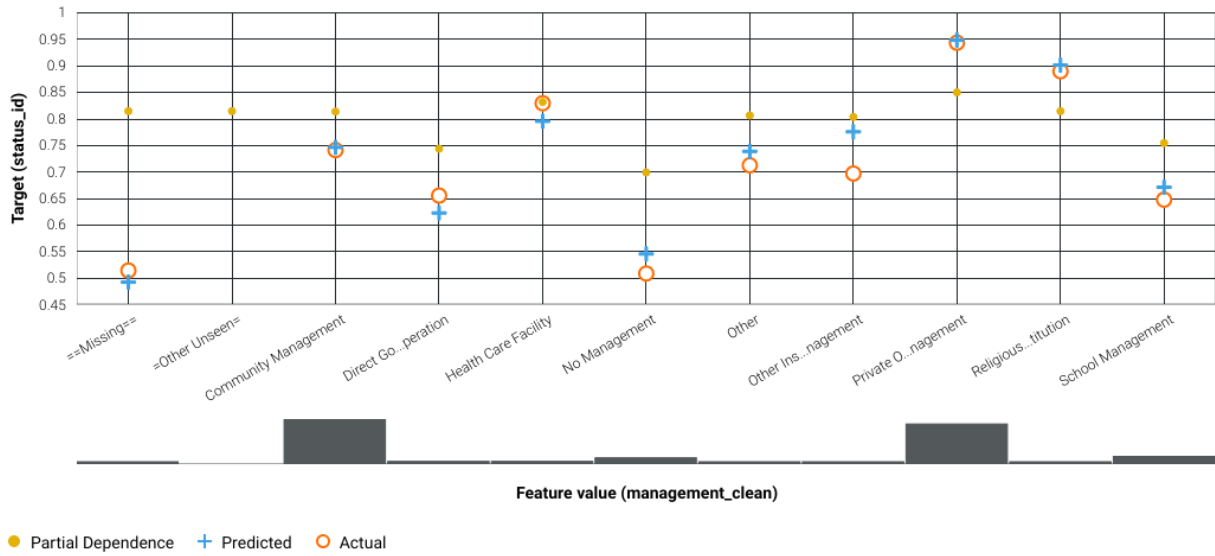
A Blueprint represents the high-level end-to-end procedure for fitting the model, including any preprocessing steps, algorithms, and post-processing. It illustrates the many steps involved in transforming input predictors and targets into a model. Each element (or, “node”) in a blueprint can represent multiple steps.

The following elements connect to create the blueprint:

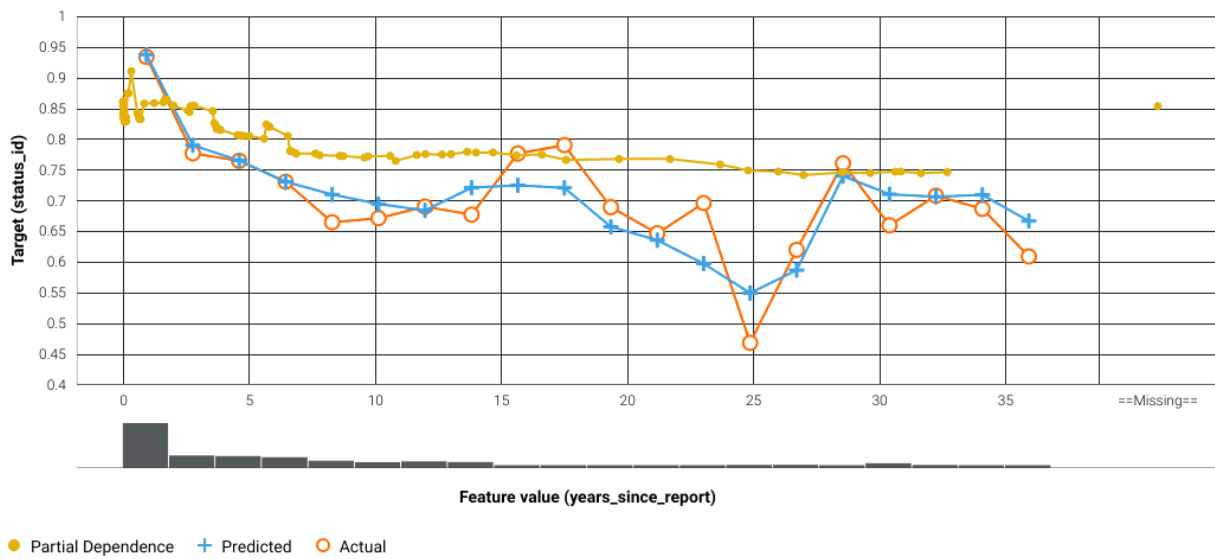
- Ordinal encoding of categorical variables
- Missing Values Imputed
- eXtreme Gradient Boosted Trees Classifier with Early Stopping (learning rate =0.02)

5 Feature Effects

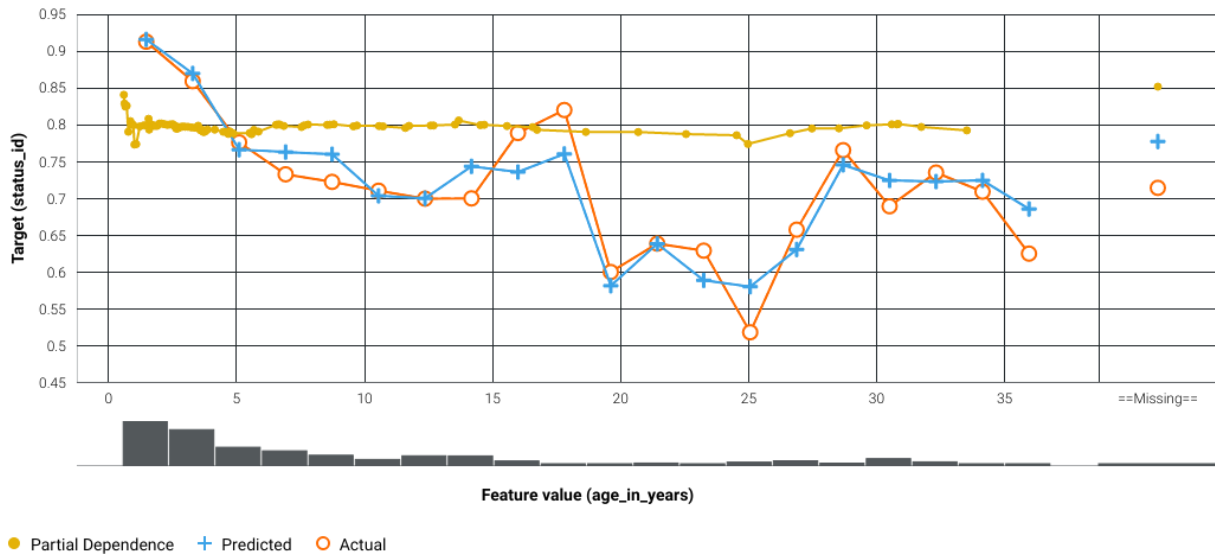
management_clean



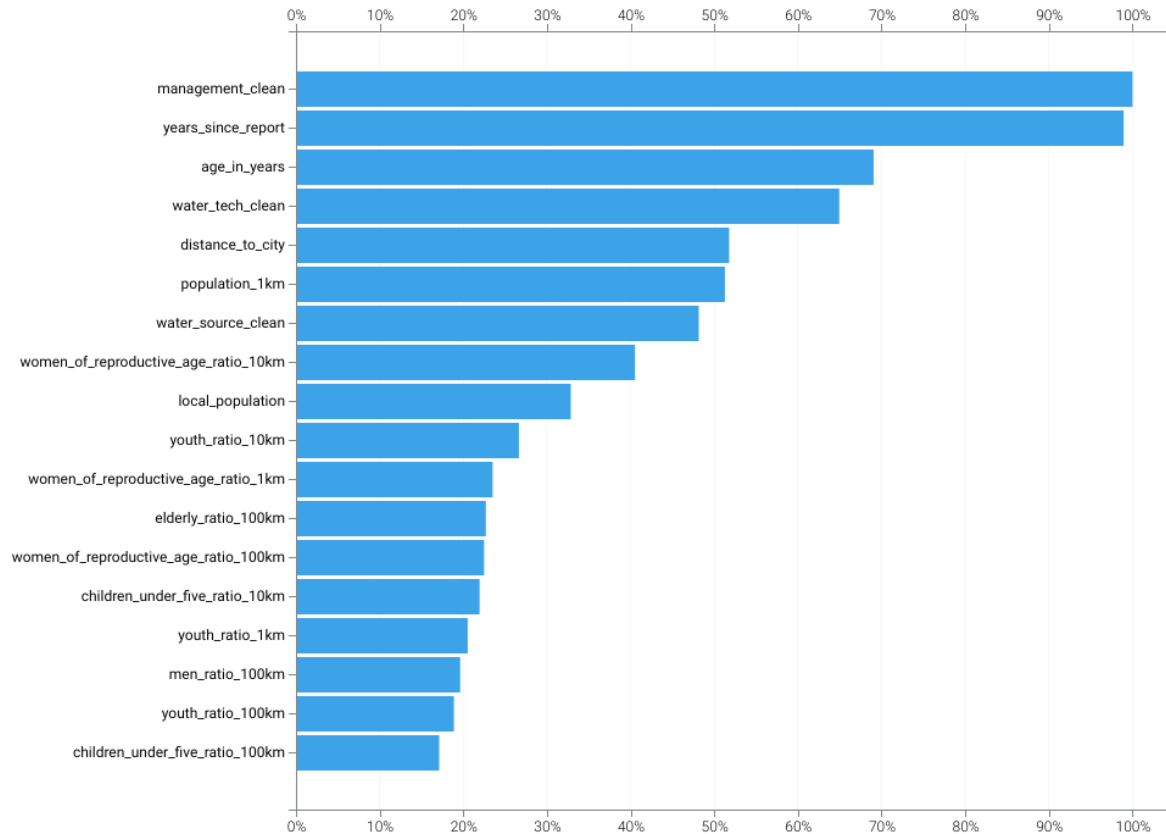
years_since_report



age_in_years



6 Feature Impact Chart



7 Feature Impact Table

Feature Name	Impact Normalized	Impact Unnormalized
management_clean	1.0	0.0702
years_since_report	0.9894	0.0695
age_in_years	0.6905	0.0485
water_tech_clean	0.6494	0.0456
distance_to_city	0.5178	0.0364
population_1km	0.5127	0.036
water_source_clean	0.4814	0.0338
women_of_reproductive_age_ratio_10km	0.4053	0.0285

local_population	0.3285	0.0231
youth_ratio_10km	0.2664	0.0187
women_of_reproductive_age_ratio_1km	0.235	0.0165
elderly_ratio_100km	0.2269	0.0159
women_of_reproductive_age_ratio_100km	0.2249	0.0158
children_under_five_ratio_10km	0.2194	0.0154
youth_ratio_1km	0.2052	0.0144
men_ratio_100km	0.1962	0.0138
youth_ratio_100km	0.1887	0.0133
children_under_five_ratio_100km	0.1711	0.012

8 Validation Testing and Stability

To find patterns in a dataset from which it can make predictions, an algorithm must first learn from a historical example – typically from a historical dataset that contains the output variable you want to predict. However, if a model is trained too closely on its training data then it may be overfit. Overfitting is a modeling error that occurs when a model is too closely fit to training data and therefore performs poorly on out-of-sample data (data that was not used to train the model). Overfitting generally results in an overly complex model that explains idiosyncrasies and random noise in the training data, rather than the underlying trends that the model was intended to capture. To avoid overfitting, the best practice is to evaluate model performance on out-of-sample data. If the model performs very well on in-sample data, (the training data) but poorly on out-of-sample data, that may be an indication that the model is overfit.

DataRobot uses standard modeling techniques to validate model performance and ensure that overfitting does not occur. DataRobot used a robust model k-fold cross-validation framework to test the out-of-sample stability of a model's performance. In addition to cross-validation partitioning, DataRobot uses a holdout sample to further test out-of-sample model performance and ensure the model is not overfit.

The following procedure was used during development to insure that overfitting did not occur:

- All values of the feature "fold" represent separate partitions used for cross-validation

DataRobot calculates the Cross Validation scores for each of the training data partitions or folds. The project metric used to calculate the score is LogLoss.

8.1 Cross Validation Scores

Fold	Cross Validation Score (LogLoss)
------	----------------------------------

Fold 1	0.36334
Fold 2	0.36268
Fold 3	0.36475
Fold 4	0.35946
Fold 5	0.35799

9 Model Results

DataRobot runs performance testing during the model development process to evaluate model results and reliability. The validation, cross-validation, and holdout (if applicable) out-of-sample performance scores are presented below, as well as the number of observations for each partition. The performance metric used for this project was LogLoss and the project included a total of 32,866 observations. An asterisk (*) next to a score, whether validation or holdout, indicates that DataRobot used in-sample predictions to derive the score. (In-samples predictions are those that include data from the validation or holdout partitions due to sample size used to build the model.)

Scoring Type	Score (LogLoss)
cross_validation	0.3616*
holdout	0.0229*
validation	0.3633*

10 Bias and Fairness

The Bias & Fairness feature is not a legal compliance tool. Please carefully review the product Documentation to ensure that you fully understand the functionality of the feature before using it. DataRobot recommends that you take local legal advice where you wish to rely on the results of this feature for complying with any legal obligations. Product Documentation can be found here: <https://app.datarobot.com/docs/modeling/special-workflows/b-and-f/index.html>

DataRobot's Bias and Fairness testing identifies whether the model exhibits biased behavior towards any classes in the dataset's protected features, based on the selected definition of fairness. Protected features and the fairness metric are chosen before Autopilot is started. DataRobot also provides a workflow that guides you towards an appropriate definition of fairness for the specific use case.

DataRobot's Bias and Fairness feature includes two model-level insights:

- Per-Class Bias, which shows whether the model is treating certain protected groups differently as measured by the selected fairness metric. This identifies if there is biased behavior, and if so, how that bias manifests, but not why.
- Cross-Class Data Disparity, which shows how different protected classes differ in their data distribution. This offers deeper insight into why the model is treating groups differently.

Together, these insights can help identify potential mitigation strategies for bias in the dataset and model, such as improving data collection or data sampling for specific groups.

Bias and Fairness testing was used in this project. The selected protected features were `clean_adm1`. The favorable target outcome was No. The selected fairness metric was `FavorableAndUnfavorablePredictiveValueParity`. The fairness threshold was set at 0.8.

Favorable Predictive Value & Unfavorable Predictive Value Parity (also known as Precision and NPV Parity) measure fairness by Equal Error, and it is best suited for cases when the target is severely unbalanced. These metrics measure whether your model disproportionately discovers favorable cases or omits unfavorable cases correctly for any particular group. It is calculated relative to the model's predictions--the proportion of predicted favorable or predicted unfavorable cases that are classified correctly. Unlike True Favorable and Unfavorable Rate Parity, Favorable Predictive Value & Unfavorable Predictive Value Parity are calculated relative to the predicted class, rather than the ground truth, which allows it to capture more meaningful information when there are very few predictions for a certain class, such as in imbalanced datasets. Higher values for Favorable Predictive Value and Unfavorable Predictive Value mean that the model is more accurate for that class relative to that metric. It's important to measure both Favorable Predictive Value & Unfavorable Predictive Value Parity together, because poor performance in either one can lead to biased outcomes, often for different stakeholders. These metrics will help you investigate how your model balances both of these types of accuracy across your protected classes.

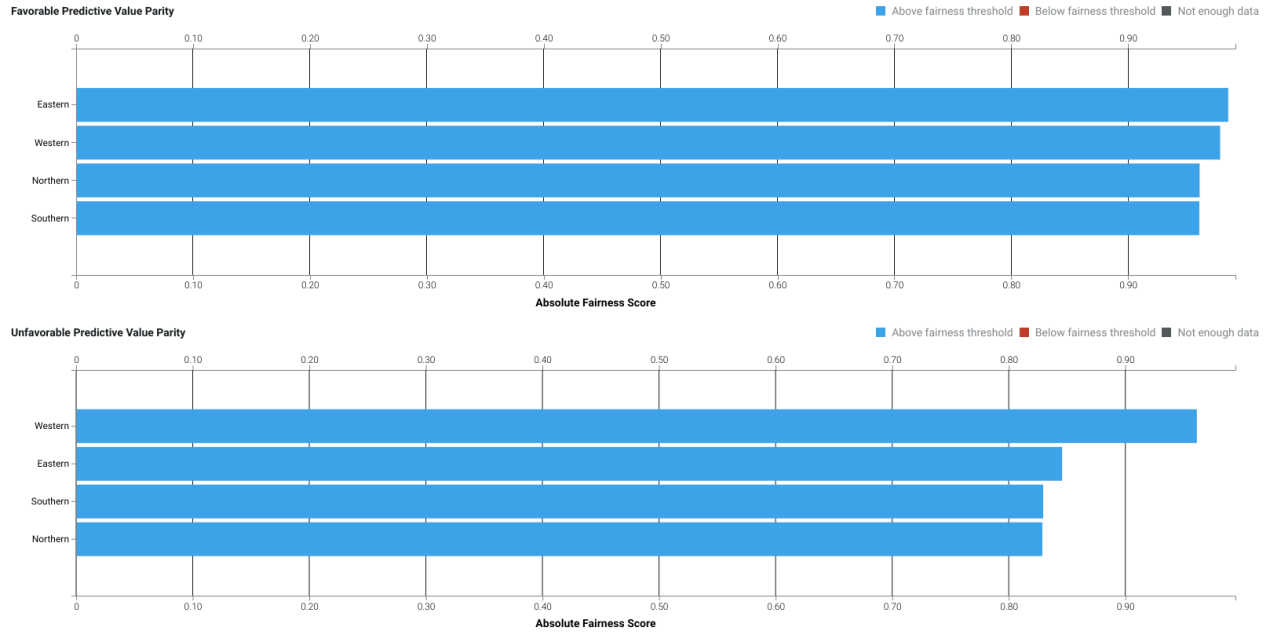
The Per-Class Bias graph shows whether the model exhibits biased behavior across protected features. The top fairness score across each protected class is scaled to 1.0, and the fairness scores for every other class are scaled relative to that value. If the fairness score for a class crosses the selected fairness threshold, the bar for that class is shown in red. If DataRobot used in-sample predictions to derive the model's performance scores (see **Overview of Model Results**), the fairness scores were calculated using in-sample validation data.

If there is not enough data for a class, its score is still calculated, but the bar for that class is shown in gray. The heuristic for whether a class does not have enough data is the following:

- If the class has <100 rows in the validation data, then it does not have enough data.
- If the class has between 100 and 1,000 rows in the validation data, but has fewer than <10% of the rows of the majority class, then it does not have enough data.
- If the class has >1,000 rows, then it has enough data.

The following figure is the Per-Class Bias graph for each protected feature in this project:

clean_adm1



The Cross-Class Data Disparity graph shows how different protected classes differ in their data distribution, in order to understand why the model treats each class differently based on the dataset. The X-axis depicts the feature importance of each feature in the dataset, while the Y-axis shows the Population Stability Index (PSI) for that feature compared across the two selected classes of the protected feature. The higher the feature importance, the more important that feature is to the model. The higher the PSI, the more differences there are for that feature across each of the two classes.