# Predict bill_result
# Performance modeling of bill_result on Challenge_8_AK_AZ_WA_prepared

| Version | Author | Date |
|---------|--------|------|
| 1.0 | mcclenahan.kevin | 2023-04-10 19:31:05 |

# TABLE OF CONTENTS

# EXECUTIVE SUMMARY

A Binary classification Machine Learning model was built using Dataiku DSS Visual ML. Its goal is to predict **bill_result** given a total of 15 features. Using a dataset of 4738 rows, the process led to the selection of the Random Forest algorithm.

## Methodology

To ensure a good generalization capability for the ML model, a test strategy was set up. Data on which ML candidate models were not trained on was used for this purpose. The testing strategy was the following:

| Policy | Split the dataset |
|---|---|
| Use time ordering | No |
| Sampling method | First records |
| Record limit | 100000 |
| Split mode | Random |
| Use K-fold cross-testing | No |
| Train ratio | 0.8 |
| Random seed | 1337 |

See section 0 for detailed explanations about these options.

Before being tested, the ML candidate models had been tuned to find the best combination of hyperparameters according to the ROC AUC metric. This optimal hyperparameter search, based on assessing performance on a validation set, was done using the following methodology:

| Search strategy | |
|---|---|
| Strategy | Random search |
| Search parameters | |
| Random state (hyperparameter search) | 1337 |
| Max number of iterations | 24 |
| Max search time | 0 (no limit) |
| Parallelism | 4 |

Performance modeling of bill_result on Challenge_8_AK_AZ_WA_prepared

| Cross-validation | |
|---|---|
| Cross-validation strategy | K-fold |
| Number of folds | 5 |
| Random state (cross-validation split) | 1337 |
| Stratified | Yes |

See section **Error! Reference source not found.** for detailed explanations about these options.

## Results

The Random Forest algorithm was selected. The evaluation metric used to tune the hyperparameters was ROC AUC computed on the validation dataset. After the best hyperparameter combination was found, the same metric was also computed on the test dataset. The final value was 0.904.

## METHODOLOGY

This section deals with the methodological details:

- The *Problem Definition* consists of selecting the target (**bill_result**) and the type of problem (Binary classification).
- *Data Ingestion* analyzes each feature in order to maximize its prediction potential.
- *Model and Feature Tuning* describes the tested algorithms and the way to find the best hyperparameter set for each of them.
- The *Model Evaluation and Selection* strategy indicates how to compute the metrics that allow for comparison between the best-tuned algorithms so that the user can select the best algorithm according to one of the computed metrics (Here Random Forest).

## Problem Definition

A Binary classification Machine-Learning model was built using Dataiku DSS. Its goal is to predict **bill_result** given a total of 15 features.

The proportion of the target classes is shown on the graph below:

## Data Ingestion

During the data ingestion phase, the features are transformed into numerical features without missing values so as to be ingestible by the Machine Learning algorithm. The table below summarizes the processing applied to each of them.

| Feature Name | Status | Type | Processing |
|---|---|---|---|
| identifier | Input | Category | Dummy encoding |
| jurisdiction.classification | Rejected | Category | |
| session | Input | Category | Dummy encoding |
| classification | Input | Category | Dummy encoding |
| from_organization.classification | Input | Category | Dummy encoding |
| bill_result | Target | Numeric | |
| bill_subject_3 | Input | Category | Dummy encoding |
| bill_subject_4 | Input | Category | Dummy encoding |
| bill_subject_5 | Input | Category | Dummy encoding |
| state_party_affiliation | Input | Category | Dummy encoding |
| bill_subject_1 | Input | Category | Dummy encoding |
| from_organization.name | Input | Category | Dummy encoding |
| bill_subject_2 | Input | Category | Dummy encoding |
| jurisdiction.name | Input | Category | Dummy encoding |
| bill_party_affiliation | Input | Category | Dummy encoding |

Performance modeling of bill_result on Challenge_8_AK_AZ_WA_prepared

## Model and Feature Tuning

### Pre-processings

Once each feature has been processed, it is possible to combine them to generate new features:

- Pairwise linear feature generation: Disabled
- Pairwise polynomial feature generation (A*B) for all pairs of features: Disabled

### Tested Algorithms

A selection of algorithms (candidate models) was then trained on the Machine Learning dataset, with various combinations of hyperparameters. The section below details the tested algorithms and the space of hyperparameters for each of them. It begins with the selected algorithm and its hyperparameter selection and continues with the other tested algorithms.

*Selected Model*

The Random Forest algorithm has been finally selected.

> A Random Forest is made of many decision trees. Each tree in the forest predicts a record, and each tree "votes" for the final answer of the forest.
> The forest chooses the class having the most votes.
> A decision tree is a simple algorithm which builds a decision tree. Each node of the decision tree includes a condition on one of the input features.
> When "growing" (ie, training) the forest:
> - for each tree, a random sample of the training set is used;
> - for each decision point in the tree, a random subset of the input features is considered.
> Random Forests generally provide good results, at the expense of "explainability" of the model.

The settings for this algorithm are given below. For hyperparameters, the possible values or ranges are listed:

| Number of trees | Min: 80 |
| | Max: 200 |
| | Uniform distribution |
| **Feature sampling strategy** | Fixed proportion |
| **Proportion of features to sample** | Min: 0.1 |
| | Max: 0.7 |
| | Uniform distribution |
| **Maximum depth of tree** | Min: 6 |
| | Max: 20 |
| | Uniform distribution |
| **Minimum samples per leaf** | Min: 1 |
| | Max: 20 |
| | Uniform distribution |
| **Parallelism** | 4 |

*Alternative Models*

Other algorithms are also tested. They are listed below, along with their settings:

### Logistic Regression

Despite its name, Logistic Regression is a classification algorithm, using a linear model (i.e., it computes the target feature as a linear combination of input features).

Logistic Regression minimizes a specific cost function (called logit or sigmoid function), which makes it appropriate for classification.

A simple Logistic regression algorithm is prone to overfitting and sensitive to errors in the input dataset. To address these issues, it is possible to use a penalty (or regularization term) to the weights.

This implementation can use either 'L1' or 'L2' regularization terms.

| Regularization | Try with L1 regularization: Yes |
| | Try with L2 regularization: No |
| **C** | Min: 0.01 |
| | Max: 100 |
| | Log uniform distribution |

## LightGBM

LightGBM is a tree-based gradient boosting library designed to be distributed and efficient. This algorithm provides fast training speed, low memory usage, good accuracy and is capable of handling large scale data.

For more information on gradient tree boosting, see the "Gradient tree boosting" algorithm.

| Boosting type | Try Gradient Boosting Decision Tree: Yes<br>Try Gradient One-Side sampling: No |
|---|---|
| **Maximum number of trees** | Min: 50<br>Max: 200<br>Uniform distribution |
| **Maximum depth of trees** | -1 |
| **Number of leaves** | Min: 20<br>Max: 500<br>Uniform distribution |
| **Learning rate** | Min: 0.1<br>Max: 0.6<br>Uniform distribution |
| **Minimum split gain** | Min: 0<br>Max: 1<br>Uniform distribution |
| **Minimum child weight** | Min: 0.001<br>Max: 1<br>Uniform distribution |
| **Minimum leaf samples** | Min: 1<br>Max: 100<br>Uniform distribution |
| **Columns subsample ratio for trees** | Min: 0.5<br>Max: 1<br>Uniform distribution |
| **L1 regularization** | Min: 0<br>Max: 1<br>Uniform distribution |
| **L2 regularization** | Min: 0<br>Max: 1<br>Uniform distribution |
| **Use bagging** | Yes |
| **Subsample ratio** | 0.75 |
| **Subsample frequency** | 2 |
| **Early stopping** | Yes |

| | |
|---|---|
| **Early stopping rounds** | 4 |
| **Random state** | 1337 |
| **Parallelism** | 4 |

## Hyperparameter Search

The hyperparameter search is done for each algorithm separately. It consists of finding the combination of hyperparameters that results in the best-trained model according to the validation metric (ROC AUC) computed on the validation dataset.

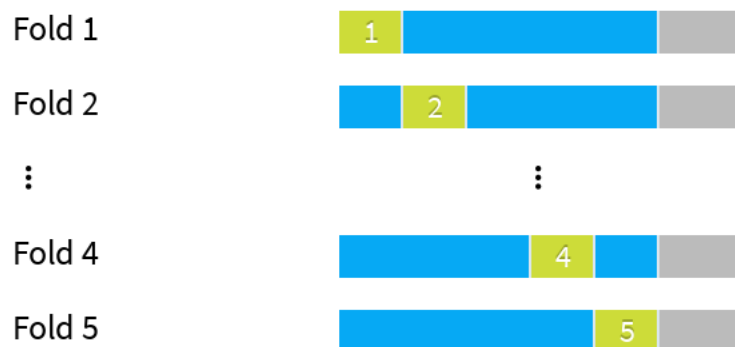The actual search settings for all the tested algorithms, including the selected one, are the following:

| | |
|---|---|
| **Search strategy** | |
| **Strategy** | Random search |
| **Search parameters** | |
| **Random state (hyperparameter search)** | 1337 |
| **Max number of iterations** | 24 |
| **Max search time** | 0 (no limit) |
| **Parallelism** | 4 |
| **Cross-validation** | |
| **Cross-validation strategy** | K-fold |
| **Number of folds** | 5 |
| **Random state (cross-validation split)** | 1337 |
| **Stratified** | Yes |

> **Legend**
> - *Randomize grid search:* If true, the grid was shuffled before the search.
> - *Max number of iterations:* This parameter sets the number of points of the grid that have been evaluated.
> - *Max search time:* Maximum search time in minutes.
> - *Parallelism:* -1 for automatic. It sets the number of hyperparameter searches that are performed simultaneously.
> - *Stratified:* If true, the same target distribution is kept in all the splits.

Illustration:

Performance modeling of bill_result on Challenge_8_AK_AZ_WA_prepared

## 5-fold cross-validation

Fold 1

Fold 2

⋮                ⋮

Fold 4

Fold 5

The metrics used to rank hyperparameter points are computed by cross-validation.

In **K-fold cross-validation** the dataset is partitioned into k equally sized subsets. Then, k-1 subsets are used as 🔵 folded train sets while the remaining subset is retained to validate the model.

This process is then repeated k times, once for each fold defined by the subset used as 🟡 validation set.

*Note:* A grey area appears on the graphic to illustrate the data that is used for the test dataset.

Weighting Strategy
Each row weight is defined as inversely proportional to the cardinality of its target class.

## Evaluation and Selection
The last part of the methodology consists of comparing the performance of each algorithm trained using the best hyperparameter combination. The policy can consist in either:

- Splitting the dataset by setting apart a test dataset, also called the hold-out dataset, for this performance evaluation. The train ratio indicates the amount of the dataset used in training, the remaining being used for evaluation.
- Performing a K-fold evaluation. It allows a more precise performance evaluation, at the expense of increased computation time.

This is indicated by the policy and the split mode in the table below.

When the original dataset is very big, the required computational resources may be too large compared to the expected benefit of training algorithms on it. As a result, the training, validation, and testing may be performed on a subset of the dataset. The sampling method given in the table below defines how it is built.

| Policy | Split the dataset |
|---|---|
| Use time ordering | No |
| Sampling method | First records |
| Record limit | 100000 |
| Split mode | Random |
| Use K-fold cross-testing | No |
| Train ratio | 0.8 |
| Random seed | 1337 |

Illustration:

## First 100000 records & 0.8 train ratio

Sampling

Evaluation

The metrics used to rank models obtained by different algorithms are computed on the ● test set. The final model is trained on the ● train set.

> **Legend**
> - Policy:
>   - *Split the dataset:* Split a subset of the dataset.
>   - *Explicit extracts from the dataset:* Use two extracts from the dataset, one for the train set, one for the test set.
>   - *Explicit extracts from two datasets:* Use two extracts from two different datasets, one for the train set, one for the test set.
>   - *Split another dataset:* Split a subset of another dataset, compatible with the dataset.

- o *Explicit extracts from another dataset:* Use two extracts from another dataset, one for the train set, one for the test set.
- *Sampling method:* A subset may have been extracted in order to limit the computational resources required by the evaluation and selection process. The *Record limit* gives its size.
  - o *No sampling (whole data)*: the complete dataset has been kept.
  - o *First records*: The first N rows of the dataset have been kept (or all the dataset if it has fewer rows. The current dataset has 4738 rows). It may result in a very biased view of the dataset.
  - o *Random (approx. ratio)*: Randomly selects approximately X% of the rows.
  - o *Random (approx. nb. records)*: Randomly selects approximately N rows.
  - o *Column values subset (approx. nb. records)*: Randomly selects a subset of values and chooses all rows with these values, in order to obtain approximately N rows. This is useful for selecting a subset of customers, for example.
  - o *Class rebalance (approx. nb. records)*: Randomly selects approximately N rows, trying to rebalance equally all modalities of a column. It does not oversample, only undersamples (so some rare modalities may remain under-represented). Rebalancing is not exact.
  - o *Class rebalance (approx. ratio)*: Randomly selects approximately X% of the rows, trying to rebalance equally all modalities of a column. It does not oversample, only undersamples (so some rare modalities may remain under-represented). Rebalancing is not exact.
- Partitions:
  - o *All partitions:* Use all partitions of the dataset.
  - o *Select partitions:* Use an explicitly selected list of partitions.
  - o *Latest partition:* Use the "latest" partition currently available in the dataset. "Latest" is only defined for single-dimension time-based partitioning.
- *Time variable:* By enabling time-based ordering, DSS checks that the train and the test sets are consistent with the time variable. Moreover, DSS guarantees that:
  - o The train set is sorted according to the selected variable.
  - o The hyperparameter search is done with training sets and validation sets consistent with the ordering induced by the time variable.
- *Split mode:* If "*K-fold cross-test*" is selected, it gives error margins on metrics, but strongly increases training time.
- *Train ratio:* Proportion of the sample that goes to the train set. The rest goes to the test set.
- *Number of folds:* Number of folds K to divide the dataset into.
- *Random seed:* Using a fixed random seed allows for reproducible results.

## EXPERIMENT RESULTS

The methodology detailed in the previous section has been run. The obtained results are presented in this section.
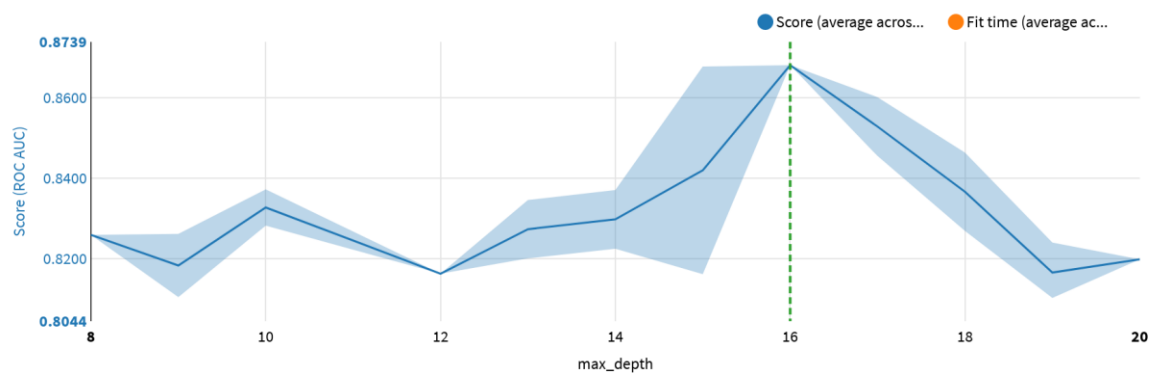
## Selected Model

Random Forest was finally selected by the user with the optimal set of hyperparameters given in the table below:
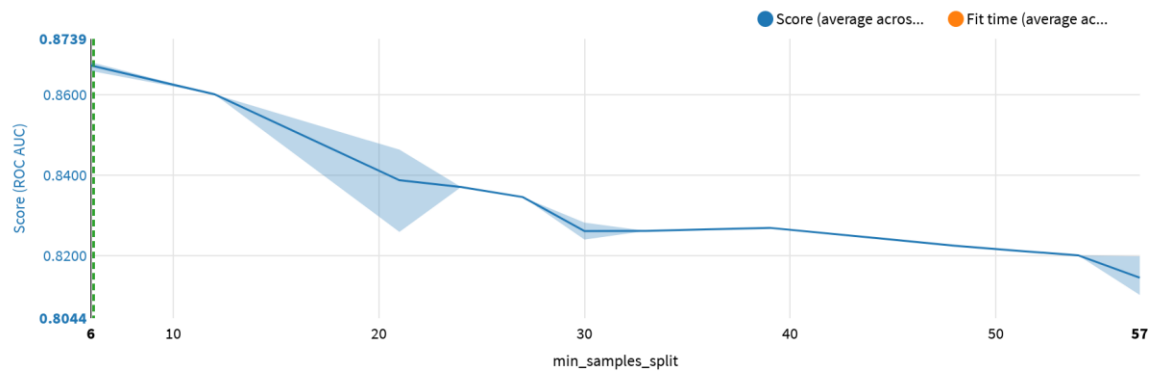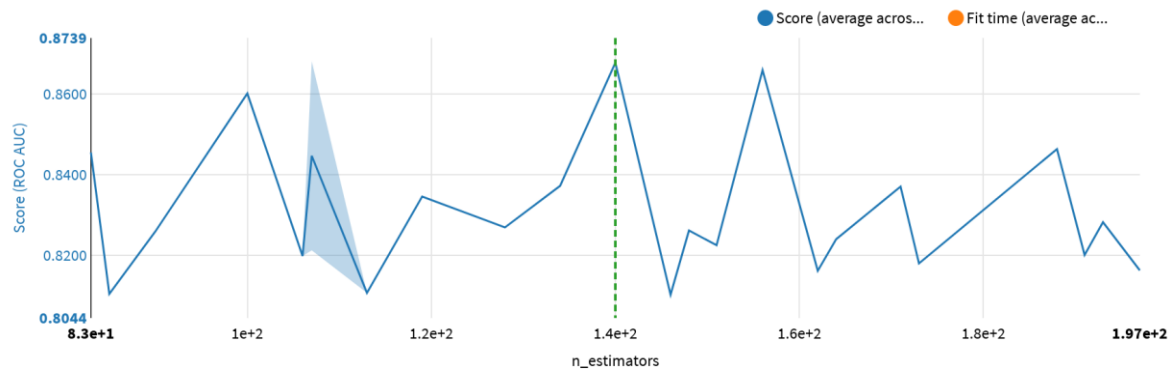
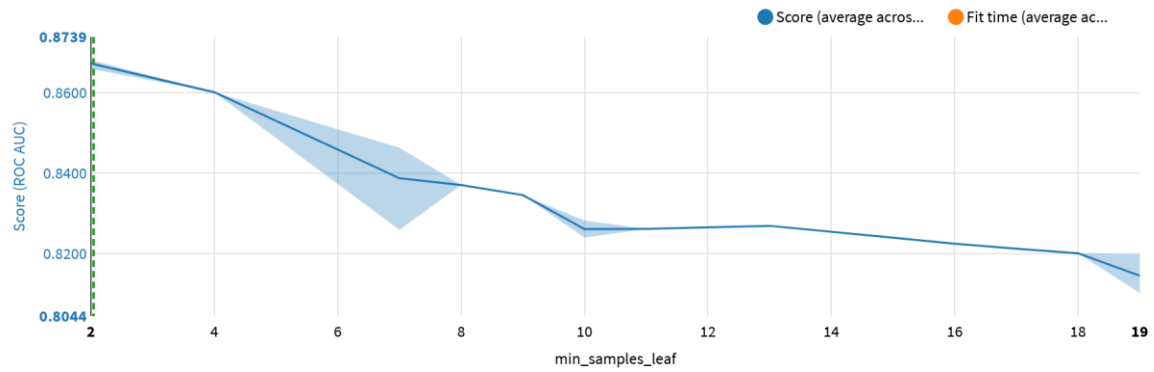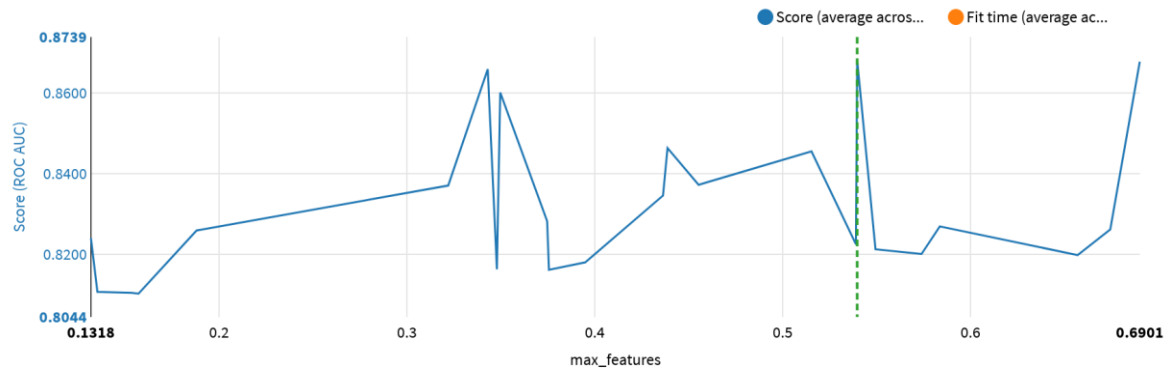| Algorithm | Random forest classification | Split quality criterion | Gini |
|---|---|---|---|
| Number of trees | 107 | Use bootstrap | Yes |
| Max trees depth | 16 | Feature sampling strategy | prop |
| Min samples per leaf | 2 | Used features | 54% |
| Min samples to split | 6 | | |

See section 0 for detailed explanations on the algorithm and its hyperparameters.

## Alternative Models

For the selected algorithm, the following other hyperparameter combinations were tried and led to lower performance. As an example, the plot below shows the evolution of the performance for each hyperparameter:

Performance modeling of bill_result on Challenge_8_AK_AZ_WA_prepared

The table below lists all the performed trainings:

| max_-dept h | max_-features | min_-samples_lea f | n_estimator s | min_-samples_spli t | Score | Score StdDev | Fit Time | Fit Time StdDev | Score Time | Score Time StdDev |
|---|---|---|---|---|---|---|---|---|---|---|
| 8 | 0.1880151778946672 | 7 | 90 | 21 | 0.82589 2 | 0.01828 7 | 4.489200 | 1.54666 6 | 1.23560 0 | 0.94948 4 |
| 9 | 0.6744155342027564 | 11 | 148 | 33 | 0.82613 7 | 0.01841 7 | 21.99840 0 | 1.51526 5 | 3.60980 0 | 0.91009 1 |
| 9 | 0.1534270403214963 | 19 | 85 | 57 | 0.81040 4 | 0.01496 6 | 4.757600 | 0.91546 7 | 1.33320 0 | 0.42374 2 |
| 10 | 0.4552678209332060 6 | 7 | 134 | 21 | 0.83722 8 | 0.01889 2 | 16.25140 0 | 0.88082 9 | 2.88180 0 | 0.51327 6 |
| 10 | 0.3746275505124097 | 10 | 193 | 30 | 0.82820 4 | 0.01657 4 | 18.62120 0 | 0.95392 5 | 4.96000 0 | 0.32337 5 |
| 12 | 0.3478847808547953 5 | 19 | 197 | 57 | 0.81624 5 | 0.01533 6 | 15.86060 0 | 0.34897 8 | 5.12440 0 | 0.42173 6 |
| 13 | 0.4363178483784972 | 9 | 119 | 27 | 0.83455 2 | 0.01813 0 | 16.24000 0 | 0.66154 9 | 1.36100 0 | 0.71453 8 |
| 13 | 0.5739371275638874 | 18 | 191 | 54 | 0.82006 2 | 0.01651 8 | 23.19900 0 | 1.05950 0 | 4.24220 0 | 0.99069 7 |
| 14 | 0.3219712593606489 6 | 8 | 171 | 24 | 0.83704 1 | 0.01681 3 | 15.95920 0 | 0.93997 6 | 5.36640 0 | 0.73384 8 |
| 14 | 0.5390140412005745 | 16 | 151 | 48 | 0.82247 3 | 0.01628 1 | 18.59800 0 | 0.88736 2 | 3.73020 0 | 0.48499 7 |
| 15 | 0.3755901323287399 5 | 19 | 162 | 57 | 0.81612 3 | 0.01441 8 | 16.24140 0 | 0.74457 1 | 2.46720 0 | 0.57353 9 |
| 15 | 0.6901291630861847 | 2 | 140 | 6 | 0.86778 2 | 0.01438 6 | 27.80100 0 | 1.16501 1 | 4.92080 0 | 0.61788 1 |
| 15 | 0.3430125288649570 3 | 2 | 156 | 6 | 0.86593 0 | 0.01570 1 | 16.74380 0 | 0.40010 8 | 5.17820 0 | 0.51540 6 |
| 15 | 0.3948889178988518 3 | 19 | 173 | 57 | 0.81798 2 | 0.01515 5 | 15.31380 0 | 1.71663 6 | 3.26500 0 | 0.90656 3 |
| 16 | 0.5396887316142894 | 2 | 107 | 6 | 0.86811 6 | 0.01432 5 | 18.59460 0 | 0.90444 0 | 2.88140 0 | 0.33476 4 |
| 17 | 0.3496623637425988 | 4 | 100 | 12 | 0.86014 6 | 0.01683 0 | 11.15480 0 | 1.29519 9 | 3.36440 0 | 0.60693 1 |
| 17 | 0.5153530858925176 | 7 | 83 | 21 | 0.84551 6 | 0.01758 5 | 12.13720 0 | 0.87221 4 | 2.36200 0 | 0.50918 2 |
| 18 | 0.583659121407126 | 13 | 128 | 39 | 0.82690 0 | 0.01625 4 | 15.23720 0 | 0.92266 5 | 3.11100 0 | 0.49443 1 |
| 18 | 0.4387113876847230 6 | 7 | 188 | 21 | 0.84634 0 | 0.01809 0 | 22.94280 0 | 2.71066 9 | 5.18080 0 | 1.18581 7 |
| 19 | 0.1570886759243958 2 | 19 | 146 | 57 | 0.81021 2 | 0.01515 5 | 6.896800 | 0.69243 0 | 3.40300 0 | 0.63702 2 |
| 19 | 0.1352075953480675 3 | 19 | 113 | 57 | 0.81066 7 | 0.01399 5 | 5.441600 | 1.70290 8 | 2.20800 0 | 0.82933 2 |

Performance modeling of bill_result on Challenge_8_AK_AZ_WA_prepared

| 19 | 0.1317641885691451 | 10 | 164 | 30 | 0.823986 | 0.015232 | 8.059200 | 0.908996 | 3.957200 | 1.111253 |
| 19 | 0.5494346253691668 | 17 | 107 | 51 | 0.821226 | 0.015818 | 11.418800 | 1.162745 | 2.725200 | 0.524342 |
| 20 | 0.6570442788259865 | 19 | 106 | 57 | 0.819792 | 0.015824 | 12.935800 | 1.099623 | 2.902600 | 0.334616 |

The selected algorithm was compared with other algorithms. The table below gives the performance obtained with the combination of hyperparameters that optimizes the ROC AUC metric:

SESSION 1

Random forest (s1) 🏆 0.904 ☆

Logistic Regression (… 0.874 ☆

LightGBM (s1) 0.890 ☆

Complete performance results obtained with the other evaluated metrics are given below:

| Accuracy | SESSION 1 |
| --- | --- |
| | Random forest (s1) 🏆 0.884 ☆ |
| | Logistic Regression (… 0.857 ☆ |
| | LightGBM (s1) 0.860 ☆ |

Performance modeling of bill_result on Challenge_8_AK_AZ_WA_prepared

| | |
|---|---|
| Precision | **SESSION 1** <br><br> 🔴 Random forest (s1)  🏆 0.903  ☆ <br><br> 🟢 Logistic Regression (...  0.868  ☆ <br><br> 🔵 LightGBM (s1)  0.892  ☆ |
| Recall | **SESSION 1** <br><br> 🔴 Random forest (s1)  0.963  ☆ <br><br> 🟢 Logistic Regressi...  🏆 0.974  ☆ <br><br> 🔵 LightGBM (s1)  0.945  ☆ |
| F1 Score | **SESSION 1** <br><br> 🔴 Random forest (s1)  🏆 0.932  ☆ <br><br> 🟢 Logistic Regression (...  0.918  ☆ <br><br> 🔵 LightGBM (s1)  0.918  ☆ |

Performance modeling of bill_result on Challenge_8_AK_AZ_WA_prepared

| | |
|---|---|
| Cost Matrix Gain | **SESSION 1**<br><br>● Random forest (s1) 🏆 0.768 ☆<br><br>● Logistic Regression (… 0.766 ☆<br><br>● LightGBM (s1) 0.750 ☆ |
| Log Loss | **SESSION 1**<br><br>● Random forest (s1) 🏆 0.374 ☆<br><br>● Logistic Regression (… 0.430 ☆<br><br>● LightGBM (s1) 0.425 ☆ |
| ROC AUC | **SESSION 1**<br><br>● Random forest (s1) 🏆 0.904 ☆<br><br>● Logistic Regression (… 0.874 ☆<br><br>● LightGBM (s1) 0.890 ☆ |

Performance modeling of bill_result on Challenge_8_AK_AZ_WA_prepared

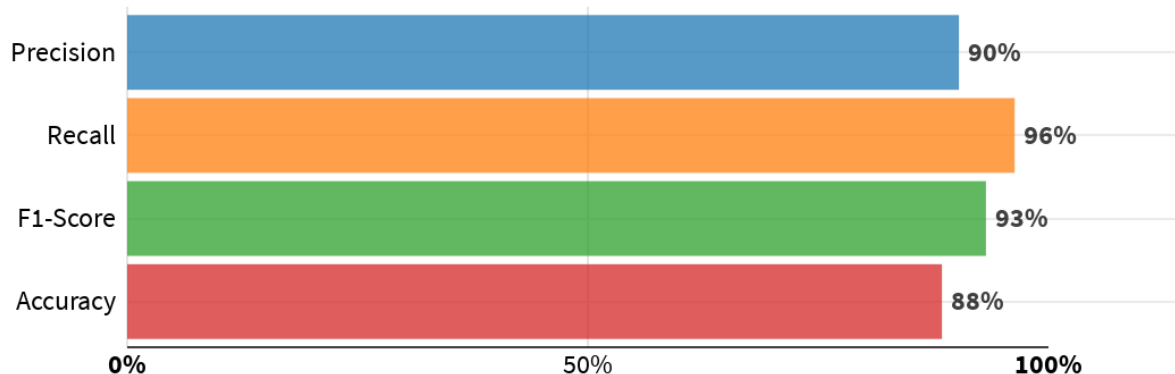| | |
|---|---|
| Calibration Loss | **SESSION 1**<br><br>⬤ Random forest (s1)  🏆 0.143  ☆<br><br>⬤ Logistic Regression (…  0.166  ☆<br><br>⬤ LightGBM (s1)  0.180  ☆ |
| Lift | **SESSION 1**<br><br>⬤ Random forest (s1)  🏆 1.206  ☆<br><br>⬤ Logistic Regression (…  1.198  ☆<br><br>⬤ LightGBM (s1)  🏆 1.206  ☆ |

## SELECTED MODEL RESULTS

## Selected Model Metrics

One way to assess the classification model performance is to use the "confusion matrix", which compares actual values (from the test dataset) to predicted values:

| | Predicted 1 | Predicted 0 | Total |
|---|---|---|---|
| Actually 1 | 920 | 35 | 955 |
| Actually 0 | 99 | 105 | 204 |
| Total | 1019 | 140 | 1159 |

A classifier produces a probability that a given object belongs to the "positive" class (**1**). The threshold (or "cut-off") is the number beyond which the prediction is considered "positive". If set too low, it may predict "negative" too often, if set too high, too rarely. The confusion matrix was obtained with a threshold set at 0.275. The optimal value according to the F1 is 0.275.

From this confusion matrix, several statistical metrics can be computed:



| | |
|---|---|
| Precision | 90% |
| Recall | 96% |
| F1-Score | 93% |
| Accuracy | 88% |

**Legend**
- *Precision*: Proportion of correct predictions among "positive" (**1**) predictions.
- *Recall*: Proportion of actually "positive" (**1**) records correctly predicted as "positive".
- *F1-score*: Harmonic mean of precision and recall.

Performance modeling of bill_result on Challenge_8_AK_AZ_WA_prepared

The confusion matrix also allows to evaluate the average gain of using the classifier thanks to the provided costs of good and bad classifications:

| If model predicts 1 | and value is 1 | the gain is | 1 | × | 920 | = | 920.00 |
| | but value is 0 | the gain is | -0.3 | × | 99 | = | -29.70 |
| | and value is 0 | the gain is | 0 | × | 105 | = | 0.00 |
| | but value is 1 | the gain is | 0 | × | 35 | = | 0.00 |
| **Average gain per record** | | | **0.77** | × | 1159 | = | 890.30 |

The detailed metrics obtained on the test dataset are given below.

Threshold independent

| Log loss | Error metric that takes into account the predicted probabilities (the lower the better) | 0.3739 |
|---|---|---|
| ROC - AUC Score | Area under the ROC; from 0.5 (random model) to 1 (perfect model) | 0.9039 |
| Calibration loss | Average distance between calibration curve and diagonal. From 0 (perfectly calibrated) up to 0.5. | 0.1431 |

Threshold dependent

| Accuracy | Proportion of correct predictions (positive and negative) in the test set | 0.8844 |
|---|---|---|
| Precision | Proportion of positive predictions that were indeed positive (in the test set) | 0.9028 |
| Recall | Proportion of actual positive values found by the classifier | 0.9634 |
| F1 Score | Harmonic mean between Precision and Recall | 0.9321 |
| Hamming loss | Fraction of labels that are incorrectly predicted (the lower the better) | 0.1156 |
| Cost matrix gain | Average gain per record that the test set (1159 rows) would yield given the specified gain for each outcome.Specified gains: TP = 1, TN = 0, FP = -0.3, FN = 0. | 0.7682 |
| Matthews Correlation Coefficient | Correlation coefficient between actual and predicted values. +1 = perfect, 0 = no correlation, -1 = perfect anti-correlation | 0.5587 |

The threshold dependent metrics have been computed thanks to the confusion matrix while the others are based on predicted probabilities.
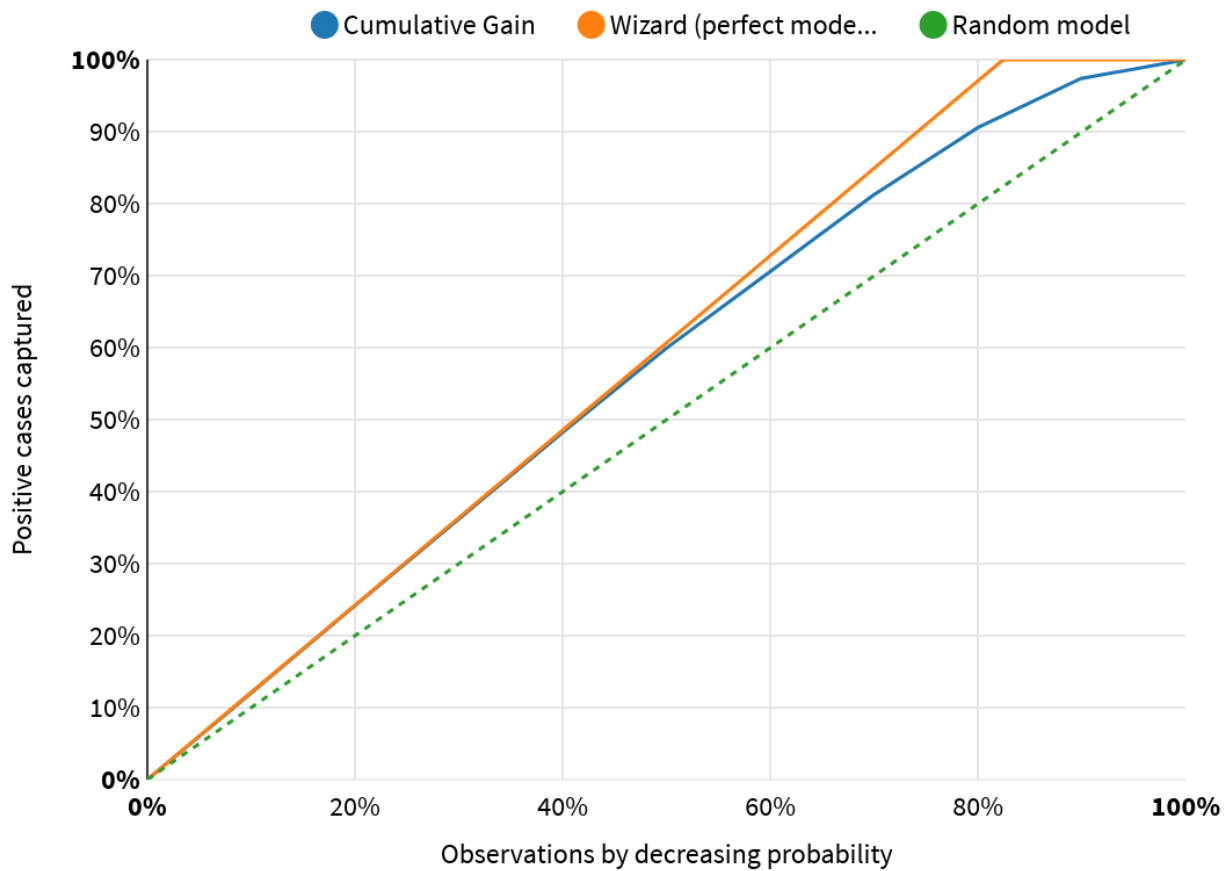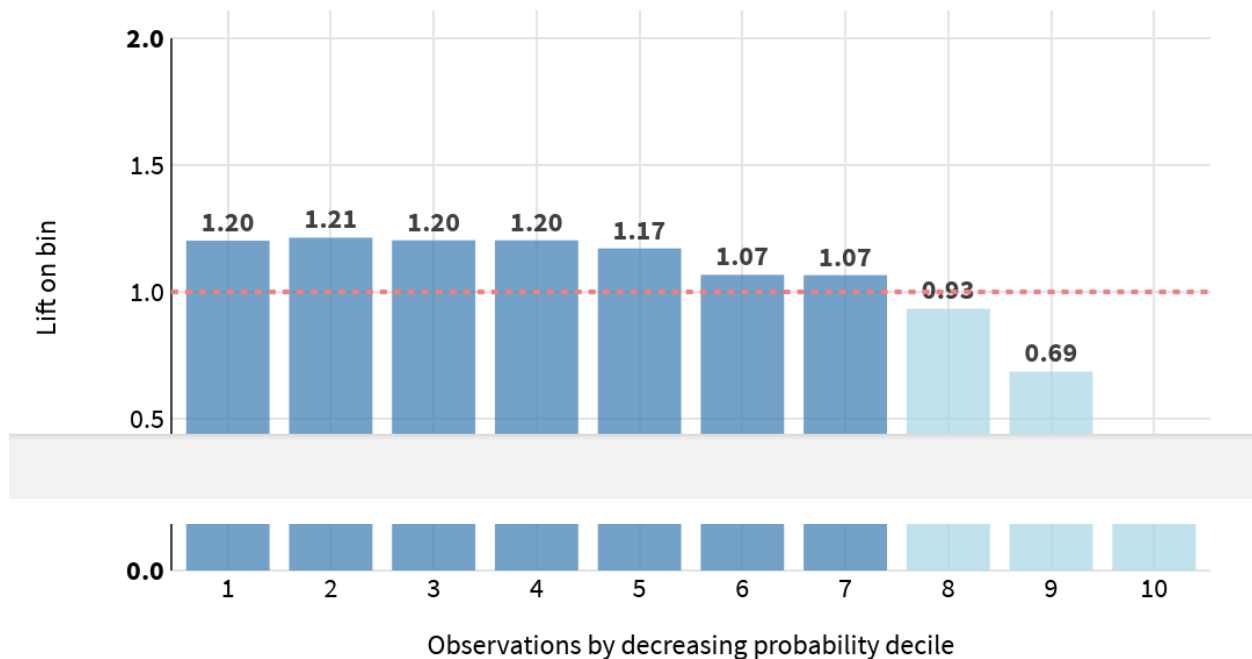
The ml assertions metrics are given below.

| Name | Criteria | Expected class | Expected valid ratio | Rows matching criteria | Rows dropped by the model | Valid ratio | Result |
|---|---|---|---|---|---|---|---|
| | No assertions defined in the settings | | | | | | |

## Selected Model Performance Charts

### Lift Charts

A binary classifier produces a probability that a given record is "positive" (Here **1**). The lift is the ratio between the results of this model and the results obtained with a random model. Lift curves are particularly useful for "targeting" kinds of problems (churn prevention, marketing campaign targeting...)

**Lift on bin** (y-axis)

Bar values by decile:
1: 1.20, 2: 1.21, 3: 1.20, 4: 1.20, 5: 1.17, 6: 1.07, 7: 1.07, 8: 0.93, 9: 0.69

Observations by decreasing probability decile

**Cumulative Lift Curve**

The curve displays the benefits of targeting a population subset with a model. On the horizontal axis, the percentage of the population which is targeted is shown. On the vertical axis, it is the percentage of found positive records (Here **1**).

- The dotted diagonal illustrates a *random model* (i.e., targeting 40% of the population will find 40% of the positive records).
- The *wizard* curve above shows a perfect model, i.e., a model that selects first all actually positive records.
- The cumulative gain curve   shows the actual percentage of actually positives found by this model. The steeper the curve, the better.

**Per-bin lift chart**

This chart sorts the observations by deciles of decreasing predicted probability. It shows the lift in each of the bins.
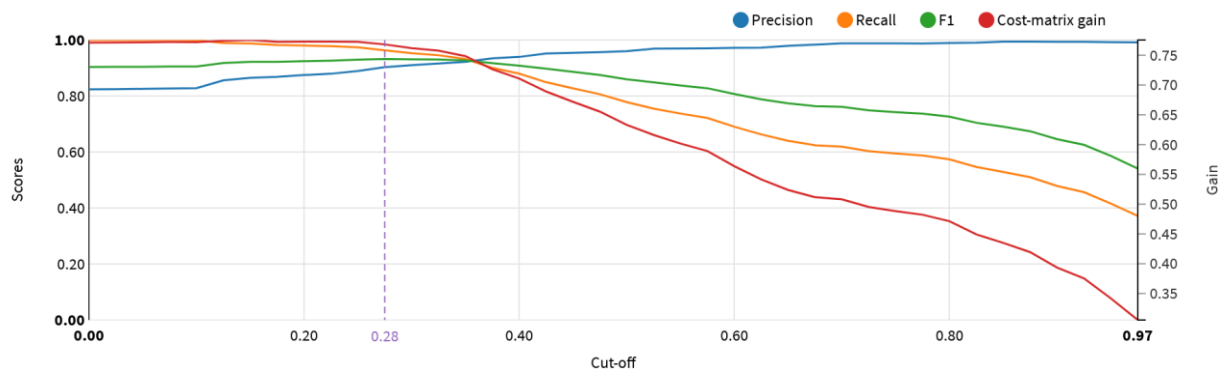
If there is 20% of positives (here **1**) in your test set, but 60% in the first bin of probability, then the lift of this first bin is 3. This means that targeting only the observations in this bin would yield 3

times as many positive results as a random sampling (equally sized bars at the level of the dotted line).

The bars should decrease progressively from left to right, and the higher the bars on the left, the better.
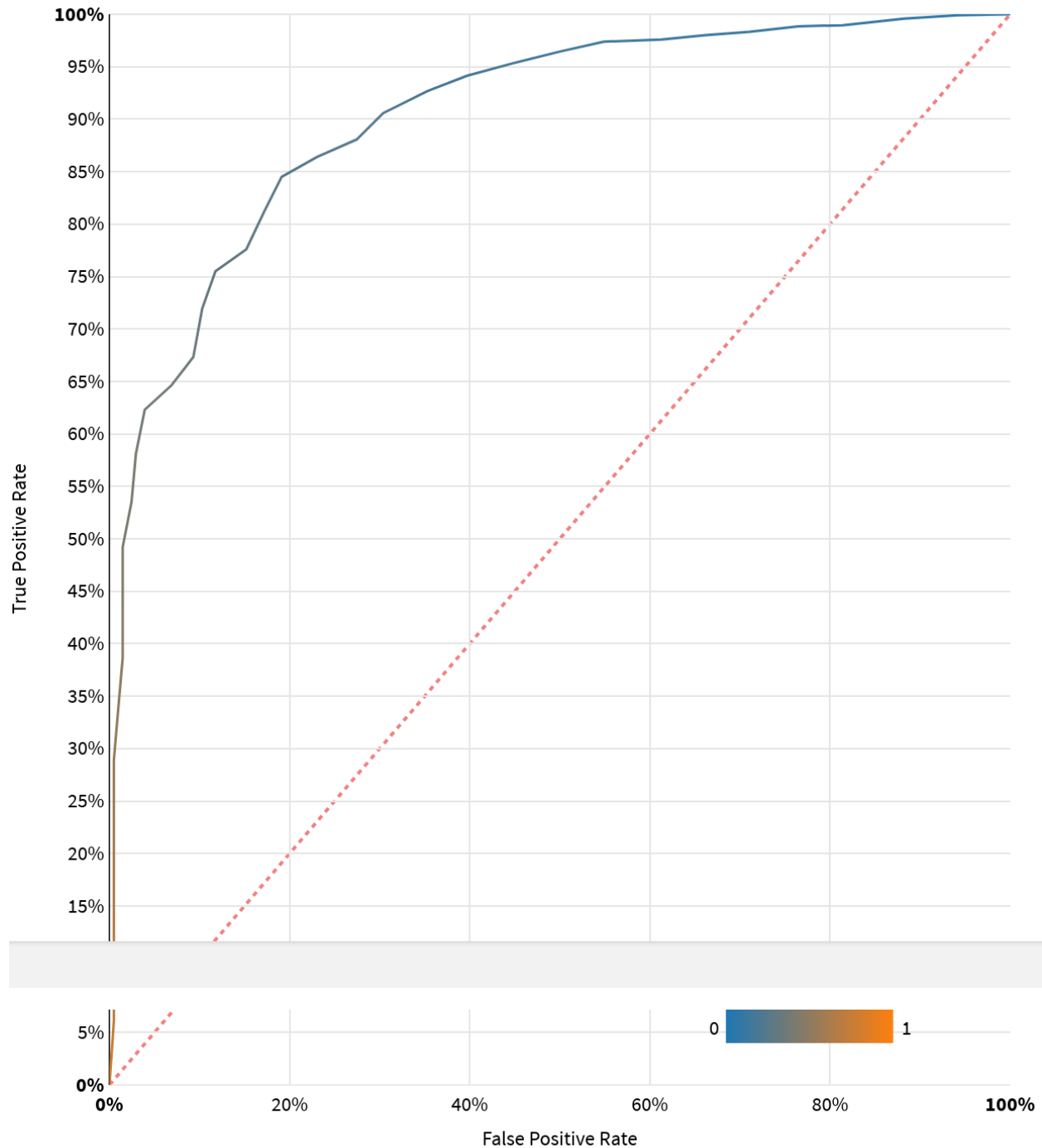
## Decision Chart
The chart below shows how the threshold-based performance metrics of the classifier vary depending on the threshold.



## ROC Curve
The Receiver Operating Characteristic (or ROC) curve shows the true positive rate versus the false positive resulting from different cutoffs in the predictive model. The "faster" the curve climbs, the better it is. On the contrary, a curve close to the diagonal line corresponds to a model with bad predictive power.
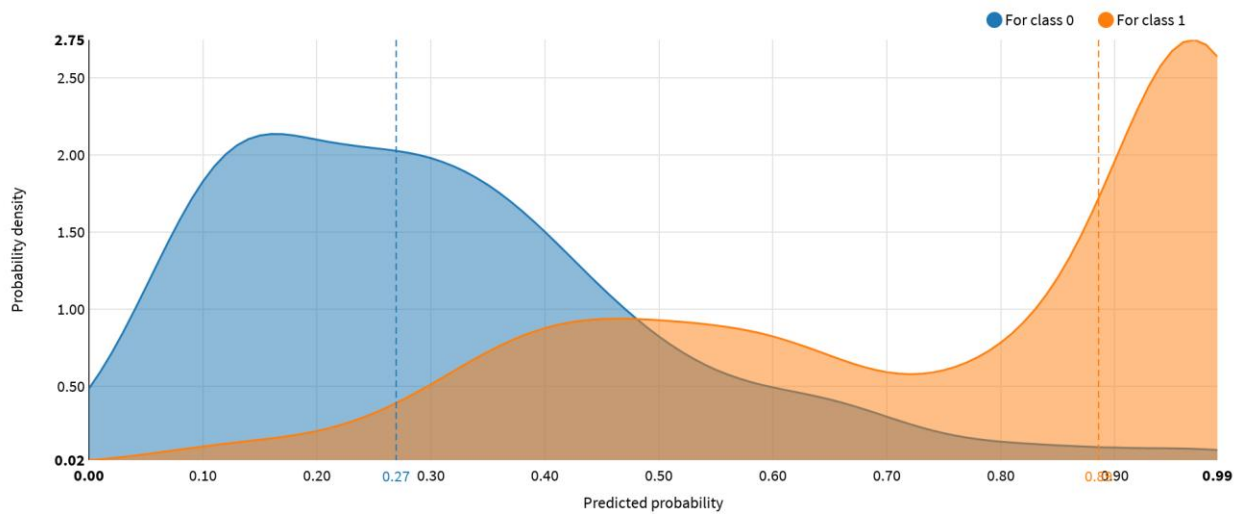
## Density Chart

The density chart illustrates how the model succeeds in recognizing (and separating) the classes (e.g., 1 and 0 for binary classification). It shows the probability distribution of the actual classes in the test set given the predicted probability of being of the "positive" class (Here **1**). The two density functions show the probability density of rows in the test set that actually belongs to the "positive" class vs. rows that do not.

Performance modeling of bill_result on Challenge_8_AK_AZ_WA_prepared
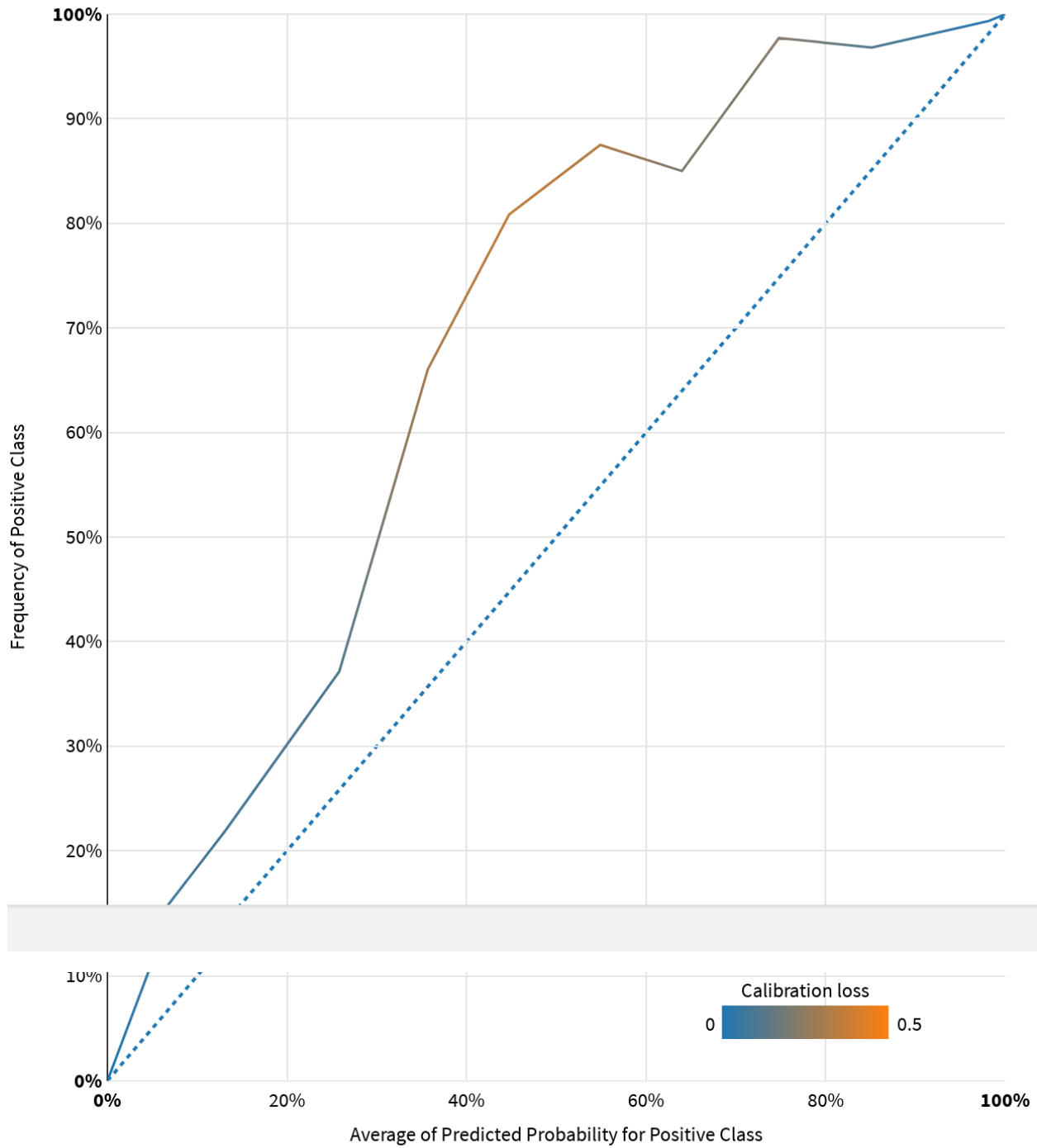
A perfect model entirely separates the density functions:

- The colored areas should not overlap.
- The density function of the "positive" class (**1**) should be entirely on the right.
- The density function of the "negative" class (**0**) should be entirely on the left.

The dotted vertical lines mark the medians.



### Calibration

Calibration denotes the consistency between predicted probabilities and their actual frequencies observed on a test dataset.
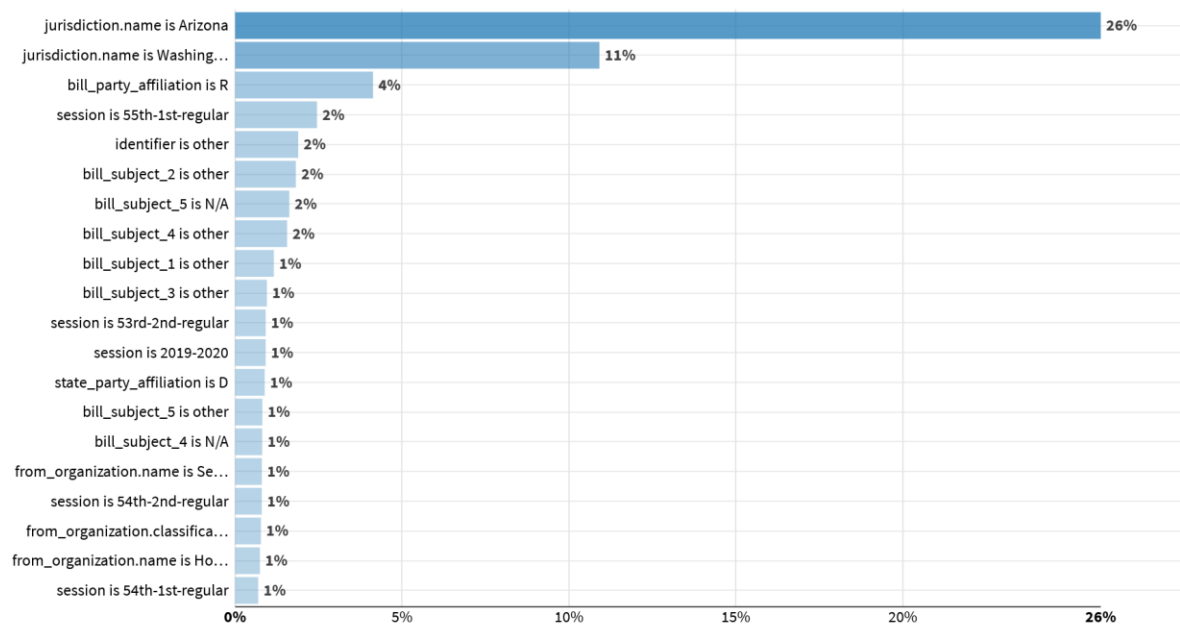
A perfectly calibrated model should have a calibration curve that is exactly on the diagonal line.

In reality, the calibration curve is often quite distinct from the diagonal line, and the average distance between the two measures the quality of the calibration.

Performance modeling of bill_result on Challenge_8_AK_AZ_WA_prepared

## Sensitivity Testing and Analysis

The selected algorithm has provided feature importance values that assess which features have a significant impact on its performance.



## Diagnostics

ML Diagnostics are designed to identify and help troubleshoot potential problems and suggest possible improvements at different stages of training and building machine learning models.

| Dataset sanity checks | Nothing to report |
|---|---|
| Modeling parameters | Nothing to report |

Performance modeling of bill_result on Challenge_8_AK_AZ_WA_prepared

| Training speed | Nothing to report |
|---|---|
| Overfit detection | Nothing to report |
| Leakage detection | Nothing to report |
| Model check | Nothing to report |
| ML assertions | Nothing to report |
| Abnormal predictions detection | Nothing to report |

## DEPLOYMENT AND MONITORING

## Implementation Details

- The backend used by the model is: Python (in memory)
- The model can be found here: https://dss-efcd7e91-005a5115-dku.us-east-1.app.dataiku.io/projects/GEORGIA_TECH_PRACTICUM_2023_SPRING/analysis/GsSL7IFi/ml/p/P72qF6zQ/A-GEORGIA_TECH_PRACTICUM_2023_SPRING-GsSL7IFi-P72qF6zQ-s1-pp1-m1/report
- The name of the generated file is: Dataiku Model Documentation - Predict bill_result on Performance modeling of bill_result on Challenge_8_AK_AZ_WA_prepared - Random Forest.docx
- The timing of the training was the following:

| Preprocessed in | 0.2s |
|---|---|
| Trained in | 557.5s |
| Loading train set | 0.0s |
| Loading test set | 0.0s |
| Collecting statistics | 0.1s |
| Preprocessing train set | 0.1s |
| Preprocessing test set | 0.0s |
| Hyperparameter searching | 551.2s |
| Fitting model | 5.2s |
| Saving model | 0.1s |
| Scoring model | 0.9s |
| Preparing explainability | 0.2s |

## Version Control

- The model was trained at 2023-04-09 07:52:13 (In the DSS server time zone).
- The model was trained with the following version of DSS: 11.3.2
- With the following code environment: DSS builtin environment

# ANNEXES

The first 3 levels of the decision tree are represented below: