IBM Watson OpenScale

# German Credit Risk Model - Prod Evaluation Report

May 23, 2020

# Overview

Deployed model:

## German Credit Risk Model - Prod

Total red breaches
**2**

## Report Details

| | |
|---|---|
| Evaluated by: | IBM Workshop One (ibmworkshop1001@gmail.com) |
| Report generated by: | IBM Workshop One (ibmworkshop1001@gmail.com) |
| Report generated on: | May 23, 2020 17:52:40 UTC |

## Model Details

| | |
|---|---|
| Deployment ID: | e8d51043-effa-4300-acc0-16d4233a2f42 |
| Model name: | German Credit Risk Model - Prod |
| Model ID: | abf2a265-b88f-4e89-a0a0-2d94cfb0c4f5 |
| Data type: | Numeric/Categorical |
| Algorithm type: | Binary classification |
| Number of explanations: | 0 |

## Training data details

| | |
|---|---|
| Storage location: | Cloud Object Storage |
| Url: | https://s3-api.us-geo.objectstorage.softlayer.net |
| Resource instance id: | crn:v1:bluemix:public:cloud-object-storage:global:a/ 5488166adaf44c648817b27b5fcf636e: 733b6db2-9cdb-4c78-91e3-64818328ddb9:: |
| Filename: | german_credit_data_biased_training.csv |
| Bucket: | training-data-location |
| Label column: | Risk |
| Deployment prediction: | prediction |
| Training features: | Age, CheckingStatus, CreditHistory, CurrentResidenceDuration, Dependents, EmploymentDuration, ExistingCreditsCount, ExistingSavings, ForeignWorker, Housing, InstallmentPercent, InstallmentPlans, Job, LoanAmount, LoanDuration, LoanPurpose, OthersOnLoan, OwnsProperty, Sex, Telephone |

# Metrics

## Metric details

### Summary

| Deployed model | Model ID |
|---|---|
| German Credit Risk Model - Prod | abf2a265-b88f-4e89-a0a0-2d94cfb0c4f5 |

Metric

## Drift

### Summary

| | |
|---|---|
| Threshold violation: | N/A |
| Drop in accuracy: | 6% |
| Drop in data consistency: | 6% |
| Estimated accuracy: | 78% |
| Base accuracy: | 83% |
| Drift threshold: | 10% |
| Minimum sample size: | 100 |

Metric

## Fairness

Score
### 80%
**RED BREACH**

### Summary

| | |
|---|---|
| Threshold violation: | 17% |
| Fairness score: | 80% |
| Fairness threshold: | 98% |
| Favorable outcome: | No Risk |
| Unfavorable outcome: | Risk |
| Minimum sample size: | 100 |

### Sex

| | |
|---|---|
| Fairness score: | 80% |
| Fairness threshold: | 98% |
| Monitored group: | female |
| Reference group: | male |

### Age

| | |
|---|---|
| Fairness score: | 88% |
| Fairness threshold: | 98% |
| Monitored group: | 44-67 |

# Metrics

Reference group: 19-43

Metric

## Quality

Score
## 79%
**RED BREACH**

### Summary

| | |
|---|---|
| Quality score: | 0.79 |
| Quality threshold: | 0.8 |
| Threshold violation: | 0.01 |
| Minimum sample size: | 100 |

### Statistics

| | |
|---|---|
| True positive rate (TPR): | 0.62 |
| Area under ROC: | 0.79 |
| Precision: | 0.84 |
| F1-Measure: | 0.71 |
| Accuracy: | 0.85 |
| Logarithmic loss: | 0.37 |
| False positive rate (FPR): | 0.05 |
| Area under PR: | 0.73 |
| Recall: | 0.62 |

## Test summary

**Tests passed**
**1**
**Tests failed**
**2**

Number of evaluated records
0

# Appendix

| Quality Measures | Area under ROC |
| --- | --- |
| | Area under PR |
| | Accuracy |
| | True positive rate (TPR) |
| | False positive rate (FPR) |
| | Recall |
| | Precision |
| | F1-measure |
| | Logarithmic loss |
| Fairness measures | Fairness |
| Drift measures | Drop in accuracy |
| | Drop in data consistency |
| | Estimated accuracy |
| | Base accuracy |
| Performance measures | Throughput |

# Appendix

Quality measures

## Area under ROC

The Area under ROC is plotted parametrically as the True positive rate versus the False positive rate with respect to a threshold T.

## Area under PR

Area under Precision Recall gives the total for both Precision + Recall. Precision (P) is defined as the number of true positives (Tp) over the number of true positives plus the number of false positives (Fp)

Formula

$$\text{Precision} = \frac{\text{number of true positives}}{(\text{number of true positives} + \text{number of false positives})}$$

Recall (R) is defined as the number of true positives (Tp) over the number of true positives plus the number of false negatives (Fn).

$$\text{Recall} = \frac{\text{number of true positives}}{(\text{number of true positives} + \text{number of false negatives})}$$

# Appendix

Quality measures

## Accuracy

Base accuracy is calculated from the training data. It is the percentage of predictions that the model got correct when tested against the training data.

## True positive rate (TPR)

The True positive rate is calculated by the following formula:

Formula

$$TPR = \frac{\text{number of true positives}}{(\text{number of true positives} + \text{number of false negatives})}$$

## False positive rate (FPR)

The false positive rate is calculated as the total number of false positives divided by the number of false positives and the number of true negatives.

$$FPR = \frac{\text{number of false positives}}{(\text{number of false positives} + \text{number of true negatives})}$$

# Appendix

Quality measures

## Recall

Recall (R) is defined as the number of true positives (Tp) over the number of true positives plus the number of false negatives (Fn).

Formula

$$\text{Recall} = \frac{\text{number of true positives}}{(\text{number of true positives} + \text{number of false negatives})}$$

## Precision

Precision (P) is defined as the number of true positives (Tp) over the number of true positives plus the number of false positives (Fp).

Formula

$$\text{Precision} = \frac{\text{number of true positives}}{(\text{number of true positives} + \text{number of false positives})}$$

# Appendix

Quality measures

## F1-Measure

The F1-Measure is the weighted harmonic average, or mean, of precision and recall.

Formula

$$F1 = 2 * \frac{(\text{precision} * \text{recall})}{(\text{precision} + \text{recall})}$$

## Logarithmic loss

For a binary model, Logarithmic loss is calculated by using the following formula:

Formula

$$-(y \log (p) + (1-y) \log (1-p))$$

where p = true label and y = predicted probability

For a multi-class model, Logarithmic loss is calculated by using the following formula:

$$-\sum_{c=1}^{M} Y_{o,c} \log(P_{o,c})$$

where M > 2, p = true label, and y = predicted probability

# Appendix

Fairness measures

## Fairness

The fairness metric used in Watson OpenScale is disparate impact, which is a measure of how the rate at which an unprivileged group receives a certain outcome or result compares with the rate at which a privileged group receives that same outcome or result.

Formula

$$\text{Disparate impact} = \frac{(num\_positives(privileged=False)/num\_instance(privileged=False)}{(num\_positives(privileged=True)/num\_instance(privileged=True)}$$

# Appendix

Drift measures

## Drop in accuracy

Watson OpenScale analyzes each transaction to estimate if the model prediction is accurate. If the model prediction is inaccurate, the transaction is marked as drifted. The Estimated accuracy is then calculated as the fraction of non-drifted transactions to the total number of transactions analyzed. The Base accuracy is the accuracy of the model on the test data. Watson OpenScale calculates the extent of the drift in accuracy as the difference between Base accuracy and Estimated accuracy. Further, Watson OpenScale analyzes all the drifted transactions; and then, groups transactions based on the similarity of each feature's contribution to the drift in accuracy. In each cluster, Watson OpenScale also estimates the important features that played a major role in the drift in accuracy and classifies their feature impact as large, some, and small.

## Drop in data consistency

Watson OpenScale analyzes each transaction for data inconsistency, by comparing the transaction content with the training data patterns. If a transaction violates one or more of the training data patterns, the transaction is marked as drifted. Watson OpenScale then estimates the magnitude of data inconsistency as the fraction of drifted transactions to the total number of transactions analyzed. Further, Watson OpenScale analyzes all the drifted transactions; and then, groups transactions that violate similar training data patterns into different clusters. In each cluster, Watson OpenScale also estimates the important features that played a major role in the data inconsistency and classifies their feature impact as large, some, and small.

# Appendix

Drift measures

## Estimated accuracy

Estimated accuracy is the accuracy score at runtime estimated by Watson OpenScale. As part of drift monitor configuration, Watson OpenScale trains a drift detection model that identifies when the original model is likely to provide an incorrect response to a transaction. As the original model receives a new transaction, the transaction is evaluated by the drift model. If the drift model believes that the model likely provided an incorrect response, the transaction is identified as a drifted transaction. The Estimated accuracy is then calculated as the fraction of non-drifted transactions to the total number of transactions analyzed.

Formula

$$\text{Estimated Accuracy} = \frac{\text{Number of non-drifted transactions*}}{\text{Total number of transactions}}$$

*determined by the Watson OpenScale drift model

## Base Accuracy

This is calculated from the training data. It is the percentage of predictions that the model got correct when tested against the training data.

# Appendix

Performance measures

## Throughput

Throughput measures the average scoring requests per minute.

Formula

$$\frac{\text{Number of transactions received in 1 hour}}{60 \text{ minutes}}$$