

Classification Model Evaluation

Machine Learning



Hayley Bresina
One AI Client Enablement

Topics Covered

- Classification model evaluation overview
- An introduction to common evaluation metrics
 - Accuracy, Precision, Recall & F1 Score
 - Class Balance Chart
 - Confusion Matrix
 - ROC Curve & AUC
- Feature importance
- Interpretation + iteration

Learning Outcomes

You will:

- Grasp the importance of evaluating classifications to determine their accuracy & effectiveness in making predictions
- Understand how each evaluation metric provides unique insights; using them in conjunction is key
- See the importance of considering the context in which the model will be applied during evaluation
- Recognize that model building is an iterative process, involving refinement based on evaluation, domain knowledge, & stakeholder feedback



Classification Evaluation Overview



Overview

- Evaluating performance to determine if your model is producing accurate results
 - Involves comparing the model's predictions with actual outcomes
 - Allows you to compare different iterations and select the best one
 - Builds trust with stakeholders
- Critical to consider the broader context that the model will be applied
 - What are we using this model for?
 - Characteristics of the dataset



Classification Evaluation Metrics



Classification Evaluation Metrics

- Cross validation allows the accurate measurement of model performance
 - Before running the model, the dataset is split into 2 subsets: training & validation folds
 - Once the model is trained, performance is evaluated on the validation fold
- There are many evaluation tools - each offers unique insight into various aspects of the model
 - Best practice to use several multiple metrics together to see the broader picture of performance

Accuracy

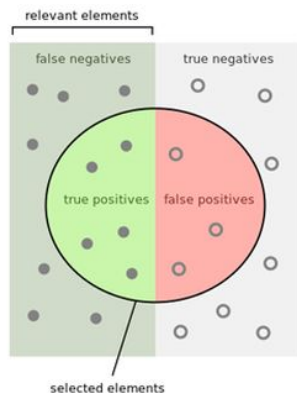
- **Accuracy** measures the proportion of correctly predicted instances out of the total number of predictions in the dataset
 - How often the model was right
 - **Misleading for imbalanced datasets**



$$\frac{TP + TN}{TP + TN + FP + FN}$$

Precision & Recall

- **Precision** quantifies the accuracy of positive predictions made by the model
 - The number of true positives divided by the total number of positive predictions (true & false positives)
- **Recall** measures how often the model correctly identifies true positives from all the actual positive samples in the dataset



How many selected items are relevant?

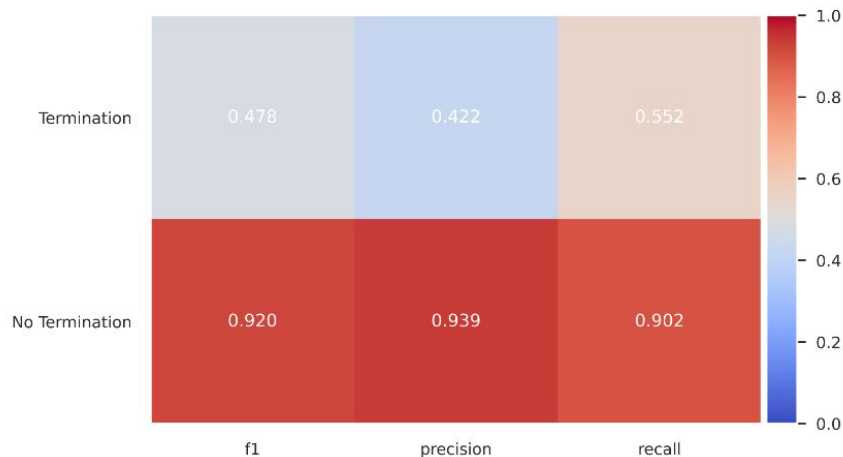
$$\text{Precision} = \frac{\text{true positives}}{\text{true positives} + \text{false positives}}$$

How many relevant items are selected?

$$\text{Recall} = \frac{\text{true positives}}{\text{true positives} + \text{false negatives}}$$

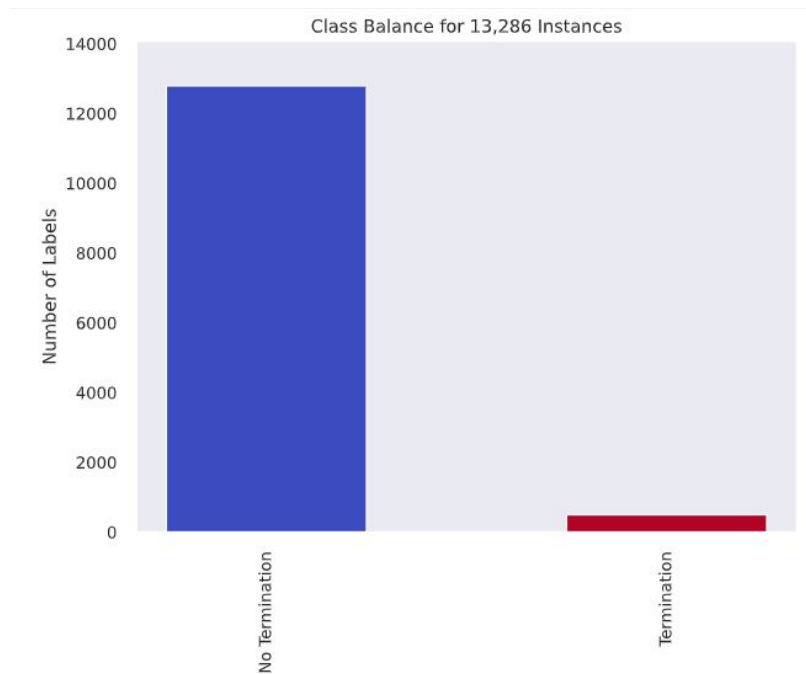
F1 Score

- The **F1 Score** is harmonic mean of precision & recall
 - Considers both false positives & false negatives
 - Helps quantify the value of trade off between precision & recall



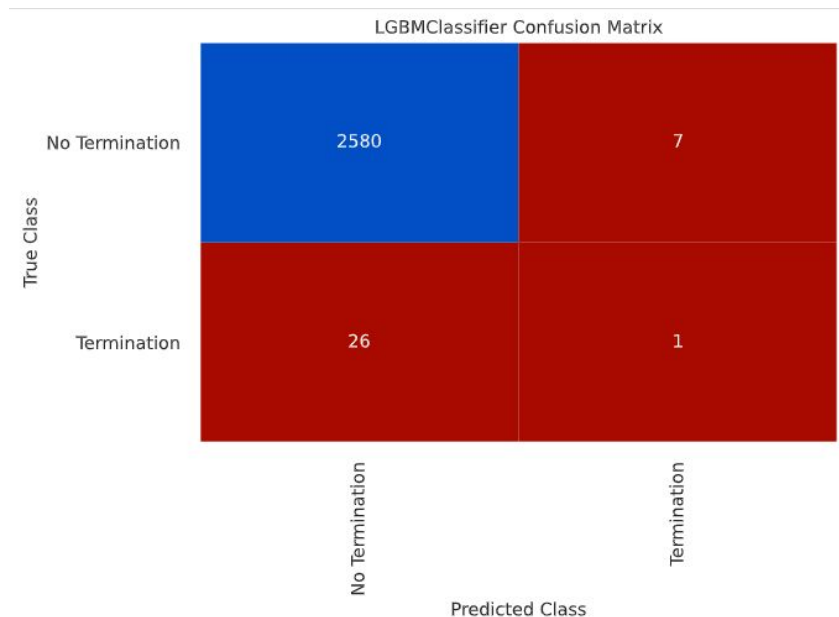
Class Balance

- **Class Balance:** chart offering a visual representation of the distribution of predicted labels
 - Helps identify potential dataset **imbalances**



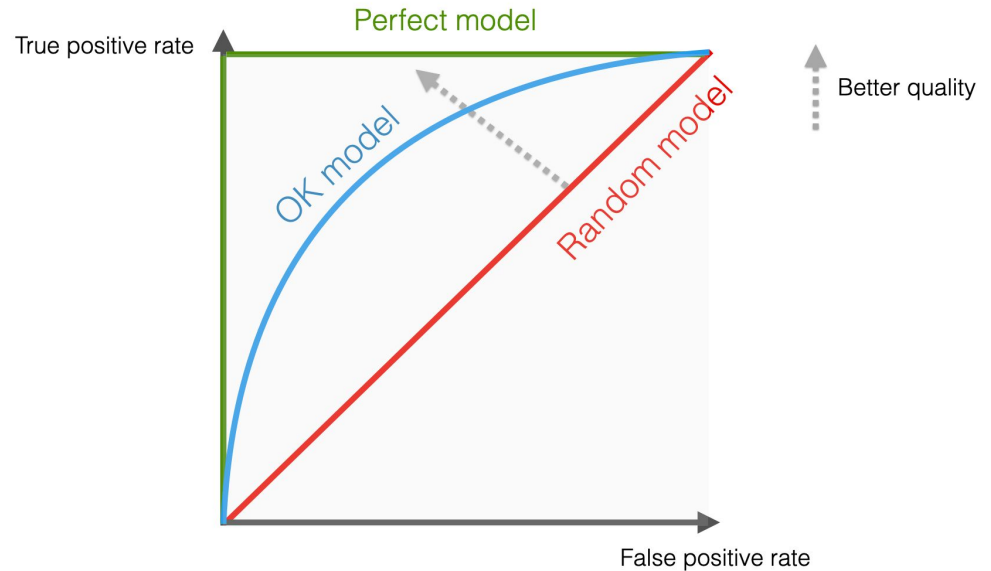
Confusion Matrix

- **Confusion Matrix:** table summarizing the model's predictions compared to the actual labels
 - Shows counts of true positives, true negatives, false positives, & false negatives



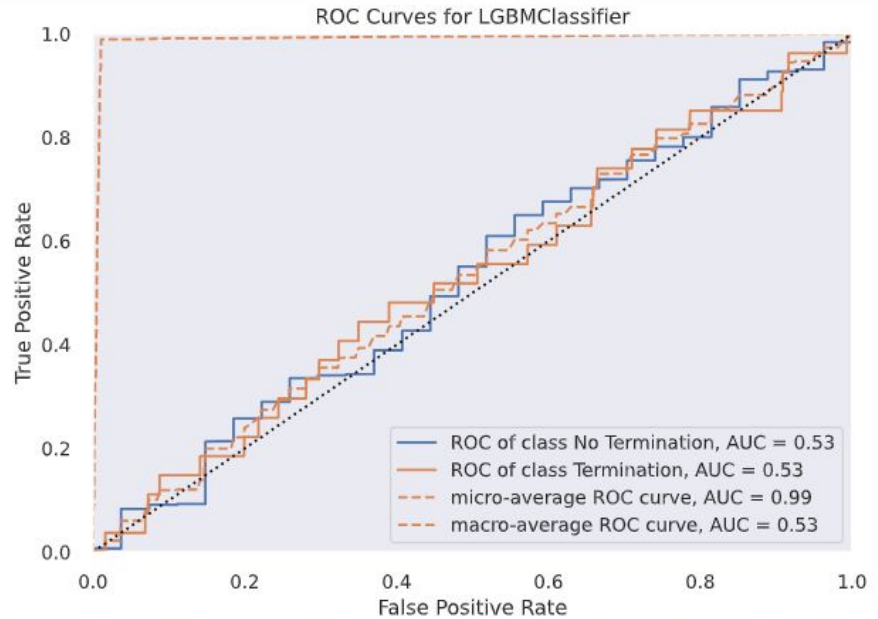
Receiver Operating Characteristics (ROC) Curve

- **ROC Curve:** graphical representation that helps convey how well a model can distinguish between positive & negative instances
 - Plots the true positive rate against the false positive rate across different threshold settings



Area Under the Curve (AUC)

- **AUC:** the area under the ROC curve; provides a value for the overall performance of the model.
 - Values closer to 1 indicate better model performance and better discrimination ability





Feature Importance



Feature Importance

- **Feature importance:** assessment of the contribution of individual features towards making accurate predictions
 - Ranks & scores the selected features to convey how important including each feature in the model is to making accurate predictions
 - Complements other performance metrics
 - Allows users to remove “bad” features

Feature Importances

Feature Name	Values
Average Tenure (EOP)_scaled	991
Department Headcount Null Filled_scaled	711
one.employee.department_scaled	628
one.employee.home_zip_scaled	279
one.employee.eeojob_scaled	168
one.employee.safety_incidents_scaled	121
one.employee.fte_scaled	102

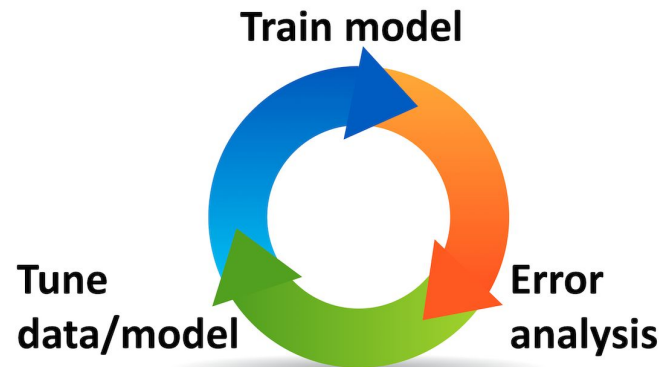


Interpretation + Iteration



Interpretation + Iteration

- Model building is an iterative process
 - Involves refining the model based on insights gained from evaluation results, domain knowledge, & feedback from stakeholders
 - Model should evolve to meet changing needs





Thanks for watching!

