

Evaluation of Classification Models

Data Science Dojo

Agenda

- Evaluation of Classification Models
 - Confusion Matrix
 - Accuracy, Precision, Recall, F1 measure
- Building Robust Machine Learning Models
 - Bias/Variance Tradeoff
- Methods of Evaluation
 - Cross Validation
 - ROC Curve

The Limitations of Accuracy

- Consider a 2-class problem:
 - Number of Class 0 examples = 9990
 - Number of Class 1 examples = 10
- If the model predicts everything to be class 0, accuracy is $9990/10000 = 99.9\%$
 - Accuracy is misleading!

METRICS FOR EVALUATION

Confusion Matrix

| ACTUAL CLASS | PREDICTED CLASS | |
|-----------------|-----------------|----------|
| | | |
| | Class=Yes | Class=No |
| Class=Yes | a | b |
| | c | d |

a: TP (true positive)

b: FN (false negative)

c: FP (false positive)

d: TN (true negative)

Confusion Matrix

| | PREDICTED CLASS | | |
|--|-----------------|-----------|-----------|
| | | Class=Yes | Class=No |
| | ACTUAL CLASS | | |
| | Class=Yes | a (TP) | b (FN) |
| | Class=No | c (FP) | d (TN) |

$$\text{Accuracy} = \frac{TP + TN}{TP + TN + FP + FN} = \frac{a + d}{a + b + c + d}$$

Precision

$$p = \frac{TP}{TP + FP} = \frac{a}{a + c}$$

| ACTUAL CLASS | PREDICTED CLASS | |
|-----------------|-----------------|-----------|
| | | |
| | Class=Yes | Class=No |
| Class=Yes | a (TP) | b (FN) |
| | c (FP) | d (TN) |

Recall/Sensitivity

$$r = \frac{TP}{TP + FN} = \frac{a}{a + b}$$

| ACTUAL CLASS | PREDICTED CLASS | |
|-----------------|-----------------|-----------|
| | | |
| | Class=Yes | Class=No |
| Class=Yes | a (TP) | b (FN) |
| | c (FP) | d (TN) |

F1-Score

$$F1 = \frac{2rp}{r + p} = \frac{2a}{2a + b + c}$$

Harmonic mean of precision and recall

| | PREDICTED CLASS | |
|--------------|-----------------|-----------|
| | Class=Yes | Class=No |
| | Class=Yes | Class=No |
| ACTUAL CLASS | a (TP) | b (FN) |
| | c (FP) | d (TN) |

WILL MY MODEL BETRAY ME?

Is My Model Really Good?

- My model shows an accuracy of 90% in the **training environment**
- Would the model be 90% accurate in **production environment**?

Generalization

- A machine learning model should be able to handle any data set coming from the same distribution as the training set.
- **Generalization** refers to a models ability to handle any random variations of training data

Overfitting (Lack of generalization)

- The gravest and most common sin of machine learning
- Overfitting: learning so much from your data that you memorize it.
 - You do well on training data
 - But don't do well (or even fail miserably) on test data

Perils of Overfitting



Data Science Dojo

@DataScienceDojo

Perils of **#overfitting** @kaggle restaurant revenue prediction Pos 1 drops to 2041 in final ranking.



| | | |
|------|-------|--|
| 2041 | ↑7 | Cheng Jiang |
| 2042 | ↓2041 | BAYZ, M.D.  |
| 2043 | ↓81 | Alberto |

Train/Test partition is not enough

Labelled Data

Training Data

**Blind Holdout
Data**

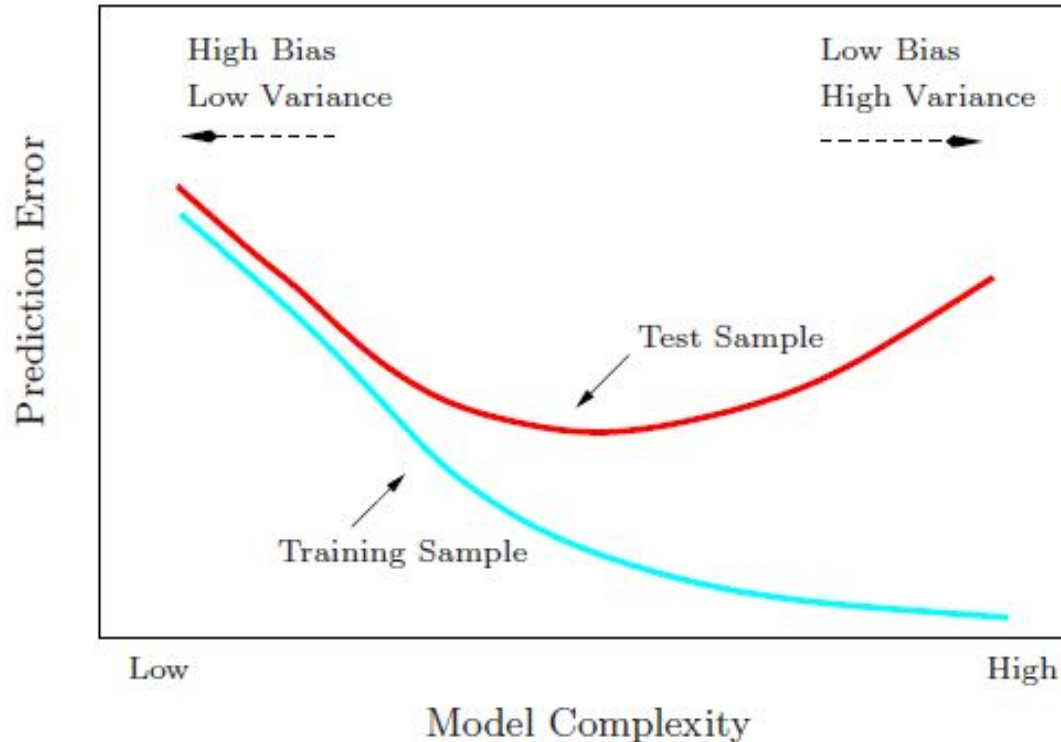
70%

30%

Blind Holdout Dataset

- The person building the model has no access to the blind holdout data set
 - Why do we need to lock it away?
- Even in presence of a 70/30 split, you may end up with a model that is not **generalized**

Bias/Variance Tradeoff

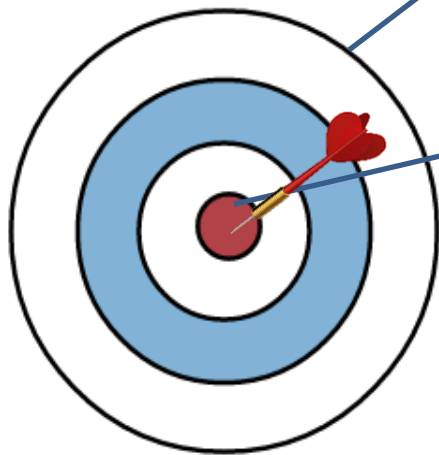


You can beat your data to confession



The generation of random numbers is too important to be left to chance.

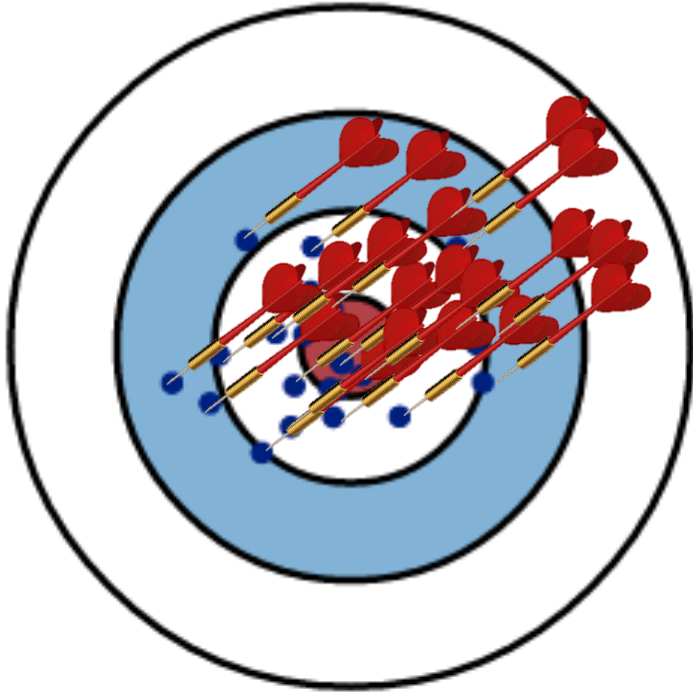
Bias/Variance Trade-off



Each dartboard represents a model

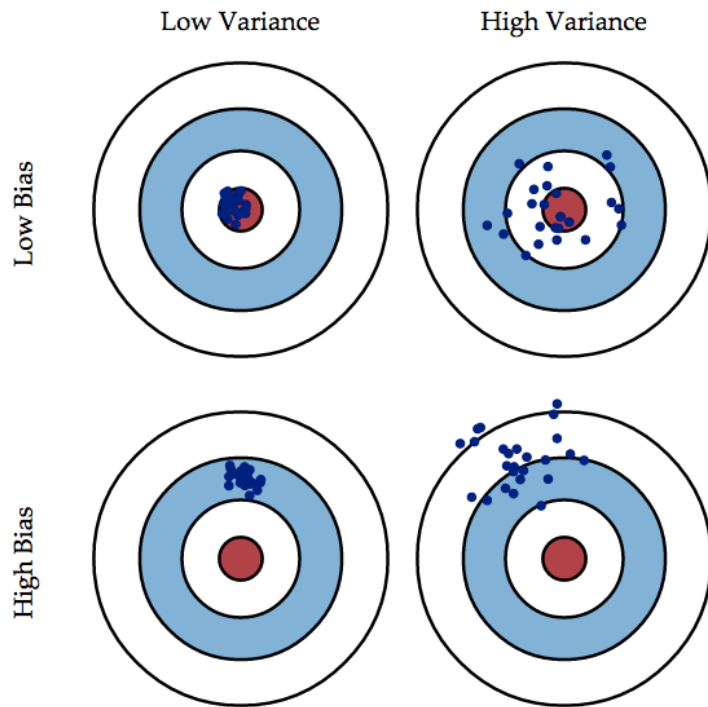
Bullseye is the theoretical best performance (accuracy, precision, recall or something else)

Bias/Variance Trade-off



- Test your model on several variations of the dataset
- Each dot represents a random variation of the test data set

Bias/Variance Trade-off



METHODS OF EVALUATION

Cross validation

- Split data into k disjoint partitions
- Train on $k-1$ partitions and test on 1
- Repeat k times

Cross validation (k=10)



Holdout Set

- 70% for training, 30% for testing
- 60/40 or 50/50 also possible
- Repeated holdout: Apply 70/30, 60/40 or 50/50 many times.

Stratified Sampling

- Use when class distribution is skewed
- Ensures that all partitions have fixed ratio of classes
 - Same ratio as training set
 - If training set is 5% class 1, 95% class 2, so is each partition

ROC CURVE

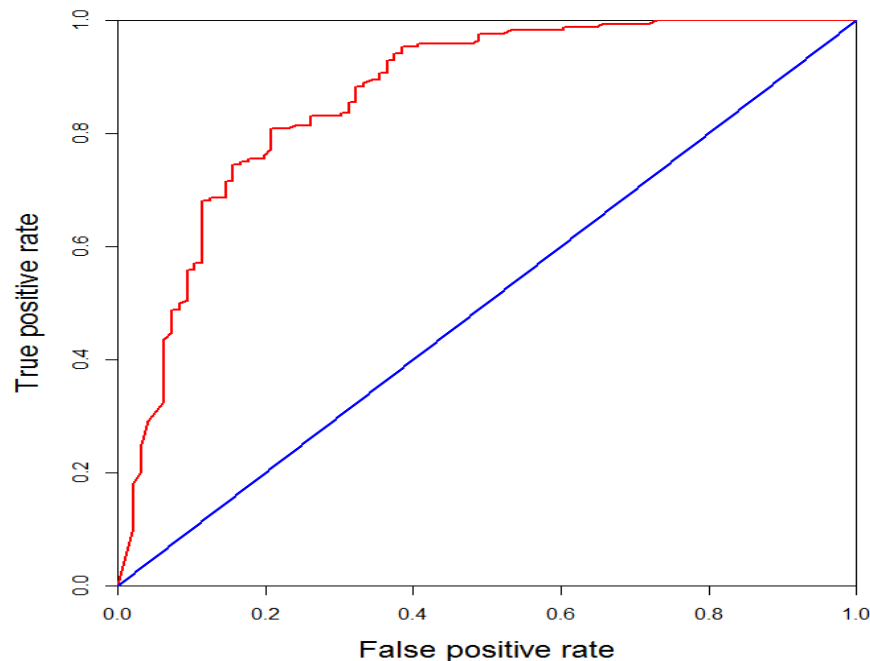
Controlling Precision and Recall

- What if probabilities are reported?
- Threshold
 - The probability value which separates positive predictions from negative predictions
 - Adjusts class label metrics

| Pid | Prediction | T=0.5 | T=0.25 | T=0.75 |
|-----|------------|----------|----------|----------|
| 2 | .95 | Survived | Survived | Survived |
| 3 | .86 | Survived | Survived | Survived |
| 5 | .02 | Dead | Dead | Dead |
| 7 | .15 | Dead | Dead | Dead |
| 13 | .48 | Dead | Survived | Dead |
| 14 | .35 | Dead | Survived | Dead |
| 21 | .12 | Dead | Dead | Dead |
| 24 | .01 | Dead | Dead | Dead |
| 34 | .74 | Survived | Survived | Dead |
| 54 | .63 | Survived | Survived | Dead |

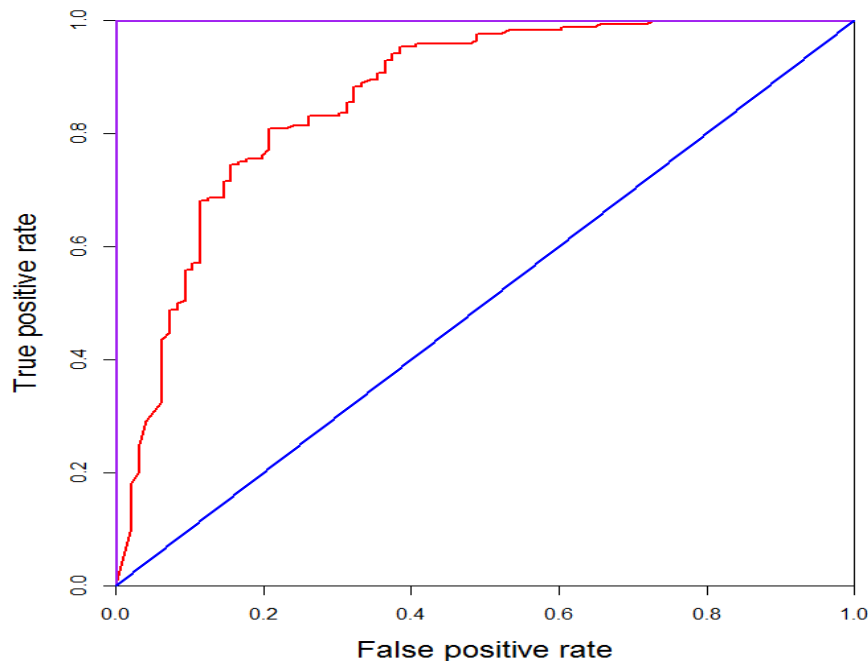
ROC(Receiver Operating Characteristic)

- Developed to analyze noisy signals
- TP on the y-axis vs FP on the x-axis
- Plot points for different threshold values
- Curve represents quality of model *independent* of threshold

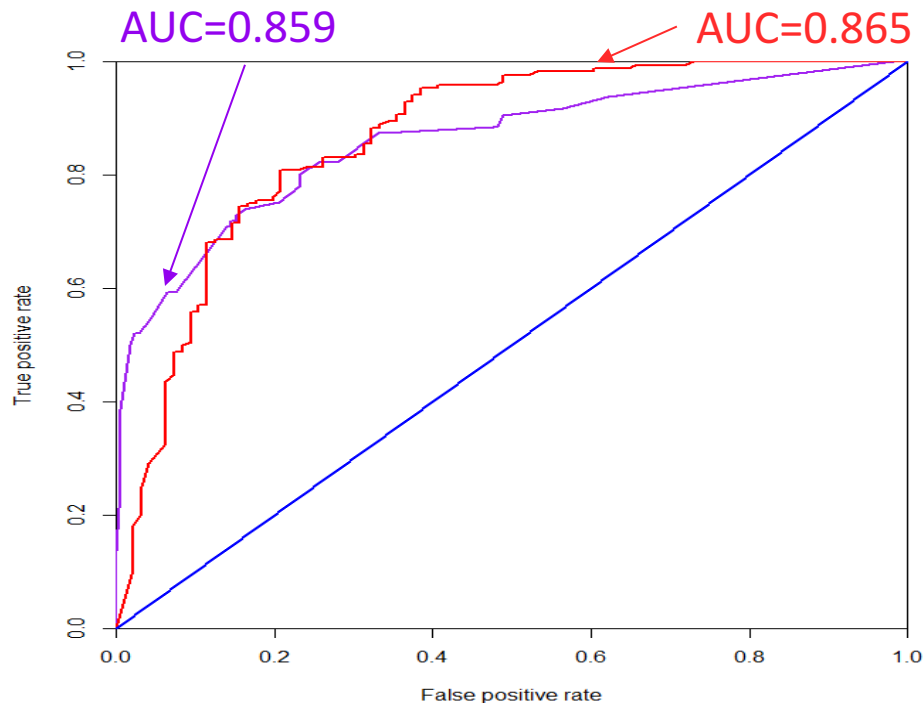


ROC Curve

- Ideal curve (purple)
 - 100% True Positives
 - 0% False Positives
- Random chance (blue)
 - Worst case
- Below diagonal line?
 - Prediction is opposite of the true class



Using ROC for Model Comparison



- No model consistently outperforms the other
 - Purple is better at low thresholds
 - Red is better at high thresholds
- Area Under ROC Curve (AUC)
 - Calculate the area under the curves
 - Compare models directly

QUESTIONS