

# Classification Metrics



**DATA SCIENCE BOOTCAMP**

# Classification Metrics

- How do I measure the performance of a (binary) classification algorithm?
- How do I compare different algorithms to see which is the best one?
- How do I know if my current performance is “good enough”?

# Types of Classification

Can consider 3 levels of responses to a binary classification problem:

- 1) **Hard classification** - For each test case, answer yes or no as to whether it is “in class” or not.
- 2) **Ranking classification** - Rank all the test cases from most likely to least likely to be “in class”.
- 3) **Probability estimation** - For each test case, estimate the probability that that the case is “in class”.

# Observations

- Can think of the three types of classification in order of increased sophistication.
- A probability estimation yields a ranking.
- A ranking yields a hard classification **for each possible threshold.**
- Each kind of classification has different metrics associated with its performance.

# Metrics for Classification Algorithms

- Hard Classification
  - Sensitivity (aka Recall)
  - Pos. Pred. Value (aka Precision)
  - Specificity
  - Accuracy
  - Others...
- Ranking Classification
  - AUC (Area Under ROC Curve, aka C-Statistic)
  - Examine Precision-Recall curve
- Probability Estimation
  - Log-likelihood of the test data (test-Deviance)
  - Brier Score
  - AIC, BIC, other penalized likelihood

# The Fish / Seaweed Analogy

A good net catches most of the fish, but allows the seaweed to pass through.

We can “cast a wider net” (lower the threshold) to catch more fish, but will probably catch more seaweed in the process.

TP = Fish in the net (good)

FP = Seaweed in the net (bad)

TN = Seaweed not in the net (good)

FN = Fish not in the net (bad)

# Metrics in Fish / Seaweed Framework

**Sensitivity (aka Recall / TPR):**

What % of the fish did I catch in my net?  
(Denominator: Fish)

**PPV/Precision:**

What % of my net is fish (vs. seaweed)?  
(Denominator: Net)

**Specificity:**

What % of the seaweed did I leave out of the net?  
(Denominator: Seaweed)

**NPV:**

What % of the non-net stuff is seaweed?  
(Denominator: non-Net)

**Accuracy:**

What % of stuff is in the right place?  
(Denominator: Everything)

# Balanced and Imbalanced Classes

- Some classification problems are **balanced**. (“in-class” and “out-of-class” are each about 50%)
- Many (most?) are **imbalanced**. Some extremely so.
  - Fraud detection (from .01% to .0001%)
  - Disease detection (<1%)
- In balanced problems, there is symmetry (and intuition) that falls apart on imbalanced problems.
- Suggestion: Keep an imbalanced problem in your head as your “canonical” example.
- Take-home: The more imbalanced the problem, the more careful one must be interpreting metrics.
- Main culprits: Accuracy, Specificity, False Positive Rate, and anything derived from them (AUC).



# A (Roughly) Balanced Problem

	Fish	Seaweed
Net	80	25
Non-Net	20	75

$$\text{Sens} = \text{TP} / (\text{TP} + \text{FN}) = 80 / 100 = .8$$

$$\text{Spec} = \text{TN} / (\text{TN} + \text{FP}) = 75 / 100 = .75$$

$$\text{Acc} = (\text{TP} + \text{TN}) / (\text{TP} + \text{TN} + \text{FP} + \text{FN}) = .775$$

$$\text{PPV} = \text{TP} / (\text{TP} + \text{FP}) = 80 / 105 = .762$$

$$\text{NPV} = \text{TN} / (\text{TN} + \text{FN}) = 75 / 95 = .789$$

Note: Seaweed = Fish = 100

$$\text{so: Acc} = (\text{Sens} + \text{Spec}) / 2$$

Net  $\sim$  Non-Net (105 vs 95)

$$\text{so: Acc} \sim (\text{NPV} + \text{PPV}) / 2$$

# Imbalanced Problem (100x seaweed)

	Fish	Seaweed
Net	80	2500
Non-Net	20	7500

$$\text{Sens} = \text{TP} / (\text{TP} + \text{FN}) = 80 / 100 = .8$$

$$\text{Spec} = \text{TN} / (\text{TN} + \text{FP}) = 7500 / 10000 = .75$$

$$\text{Acc} = (\text{TP} + \text{TN}) / (\text{TP} + \text{TN} + \text{FP} + \text{FN}) = 7580 / 10100 = .7504$$

$$\text{PPV} = \text{TP} / (\text{TP} + \text{FP}) = 80 / 2580 = .031$$

$$\text{NPV} = \text{TN} / (\text{TN} + \text{FN}) = 7500 / 7520 = .997$$

Note: Seaweed = 100 x Fish

so:  $\text{Acc} \approx \text{Spec}$

and:  $\text{Net} \approx 1/3 * \text{Non-Net}$

so:  $\text{Acc} \approx (3 * \text{NPV} + \text{PPV}) / 4$

and NPV is inflated!

# Imbalanced Problem (Even more seaweed!)

	Fish	Seaweed
Net	80	2500
Non-Net	20	17500

Add “easily filterable” seaweed to our ocean.

$$\text{Sens} = \text{TP} / (\text{TP} + \text{FN}) = 80 / 100 = .8$$

$$\text{Spec} = \text{TN} / (\text{TN} + \text{FP}) = 17500 / 20000 = .875$$

$$\text{Acc} = (\text{TP} + \text{TN}) / (\text{TP} + \text{TN} + \text{FP} + \text{FN}) = 17580 / 20100 = .8746$$

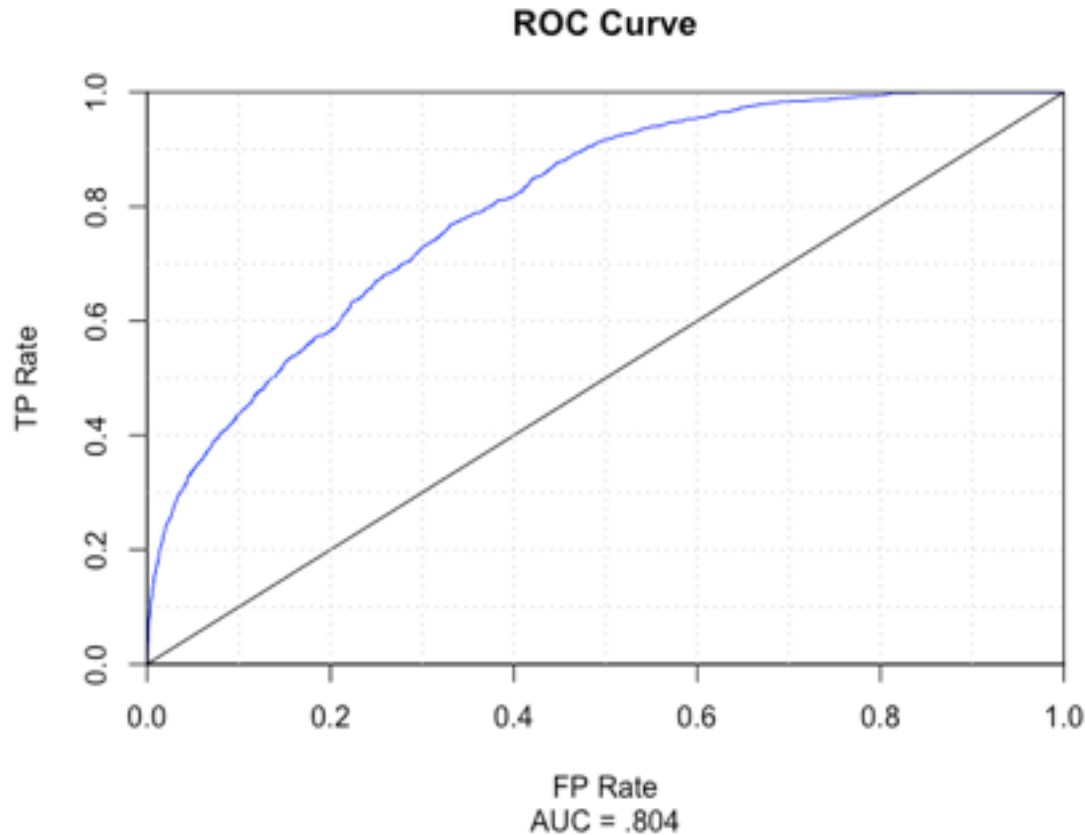
$$\text{PPV} = \text{TP} / (\text{TP} + \text{FP}) = 80 / 2580 = .031$$

$$\text{NPV} = \text{TN} / (\text{TN} + \text{FN}) = 17500 / 17520 = .9989$$

Specificity, Accuracy and NPV all improved! Why?  
Sensitivity and PPV did not.

# Metrics for Ranking Classification

ROC Curve - Plot Sensitivity vs FPR (FPR = 1-Spec)  
(ROC = Receiver Operating Characteristic)



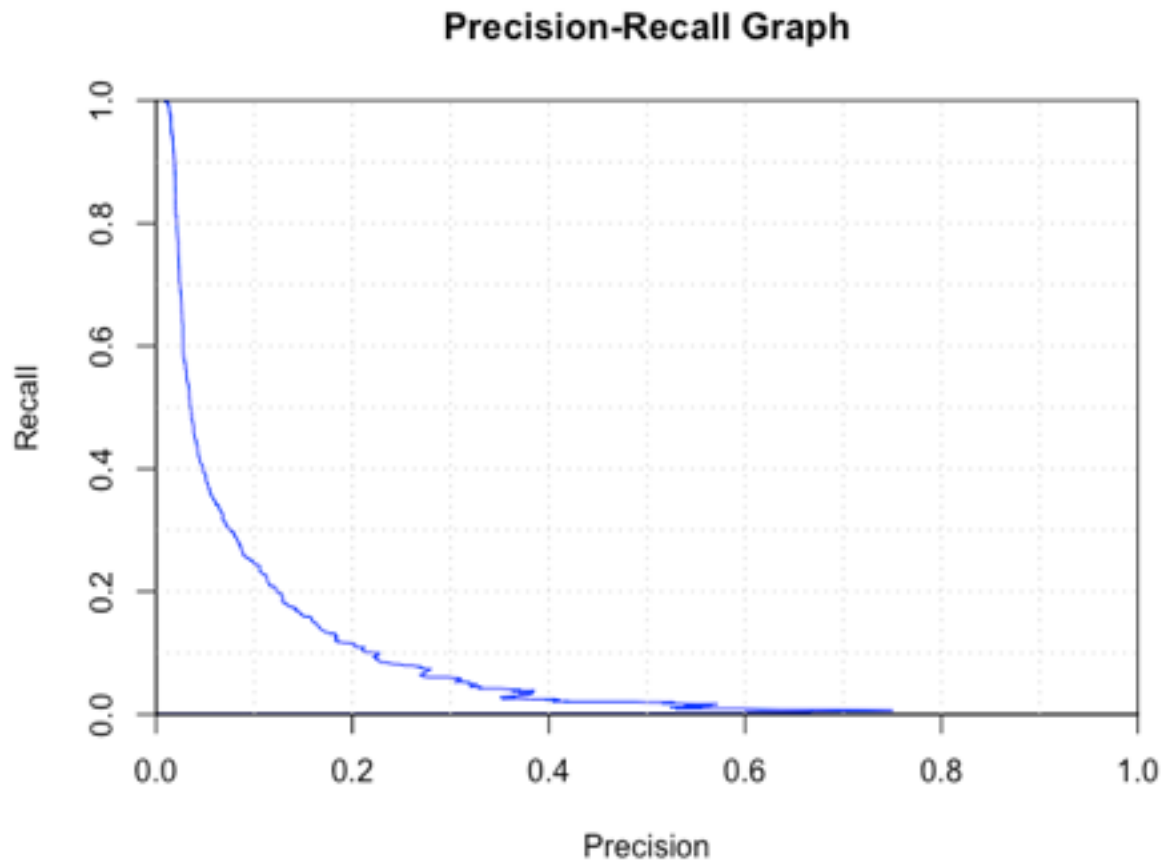
# Metrics for Ranking Classification

## AUC - Area Under the ROC Curve

- Average specificity across all sensitivity values
- AUC is the probability that a random “in-class” case is ranked higher than a random “out-of-class” case
- Can game AUC (and Spec., Acc.) by adding obvious negatives to population.
- Not a useful measure to compare models unless they are on the exact same test set.
- Evaluates entire range of thresholds, even though you may only operate in a narrow range

# Metrics for Ranking Classification

The Precision-Recall graph gives a different perspective



# Probability Estimation vs. Ranking

- For some applications, it may be important to estimate an accurate probability of being in-class.
- May be desirable to use a metric that rewards accuracy and calibration as well as stratification.
- Simple idea: evaluate the likelihood of the test data  
 $p_1 = .6, p_2 = .8, p_3 = .9$  and  $y_1 = 1, y_2 = 0, y_3 = 1$   
Likelihood =  $.6 * (1-.8) * .9$
- Test-Deviance (DIC) =  $-2 * (\log\text{-likelihood of test data})$  (smaller deviance is better)
- Brier Score: Essentially the Mean Squared Error

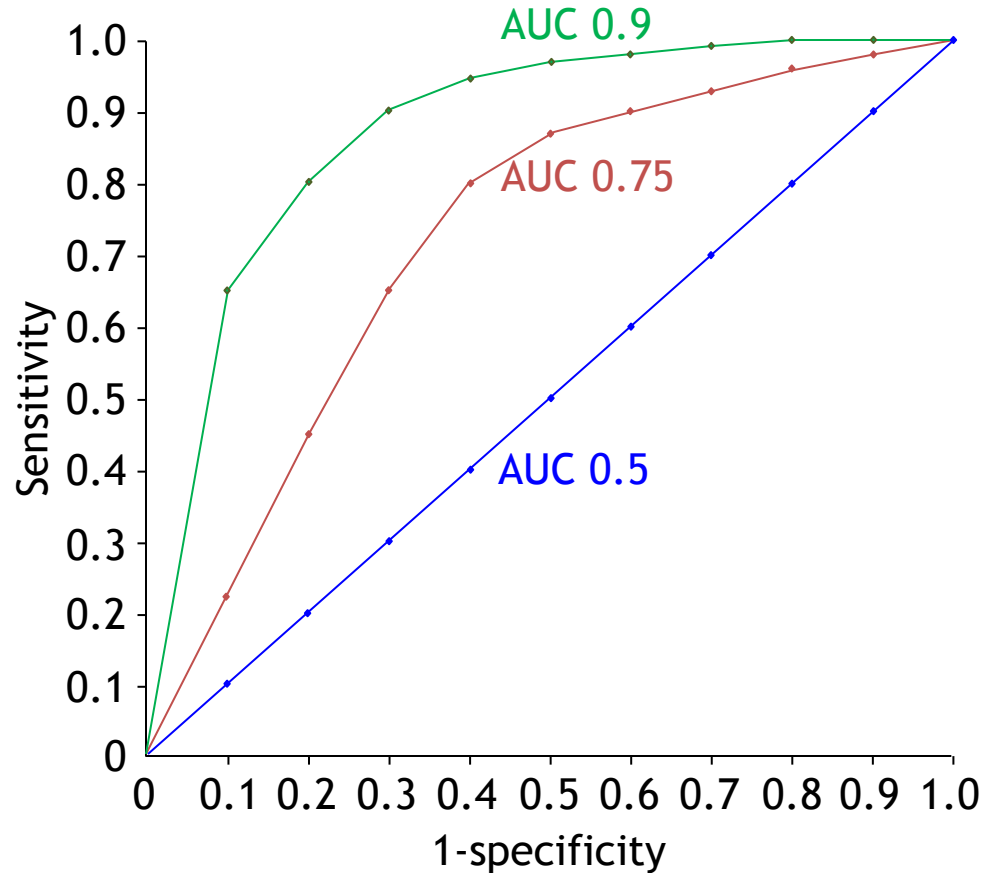
$$BS = \frac{1}{n} \sum_{i=1}^n (\hat{p}_i - y_i)^2$$

# Penalized Likelihoods

- Typically used to compare models when an independent test set is **not** available
- Idea - adding more parameters (e.g. in logistic regression) will always reduce likelihood of the **training** data.
- Therefore, to compare models we consider the likelihood plus a penalty term for model complexity.
- Attempts to measure whether the gain in likelihood “justifies” the additional parameters.
- $AIC = -2(LL) + 2k$  (where  $k$  is number of parameters)
- $BIC = -2(LL) + k*(\ln n)$  (more conservative:  $\ln n > 2$ )

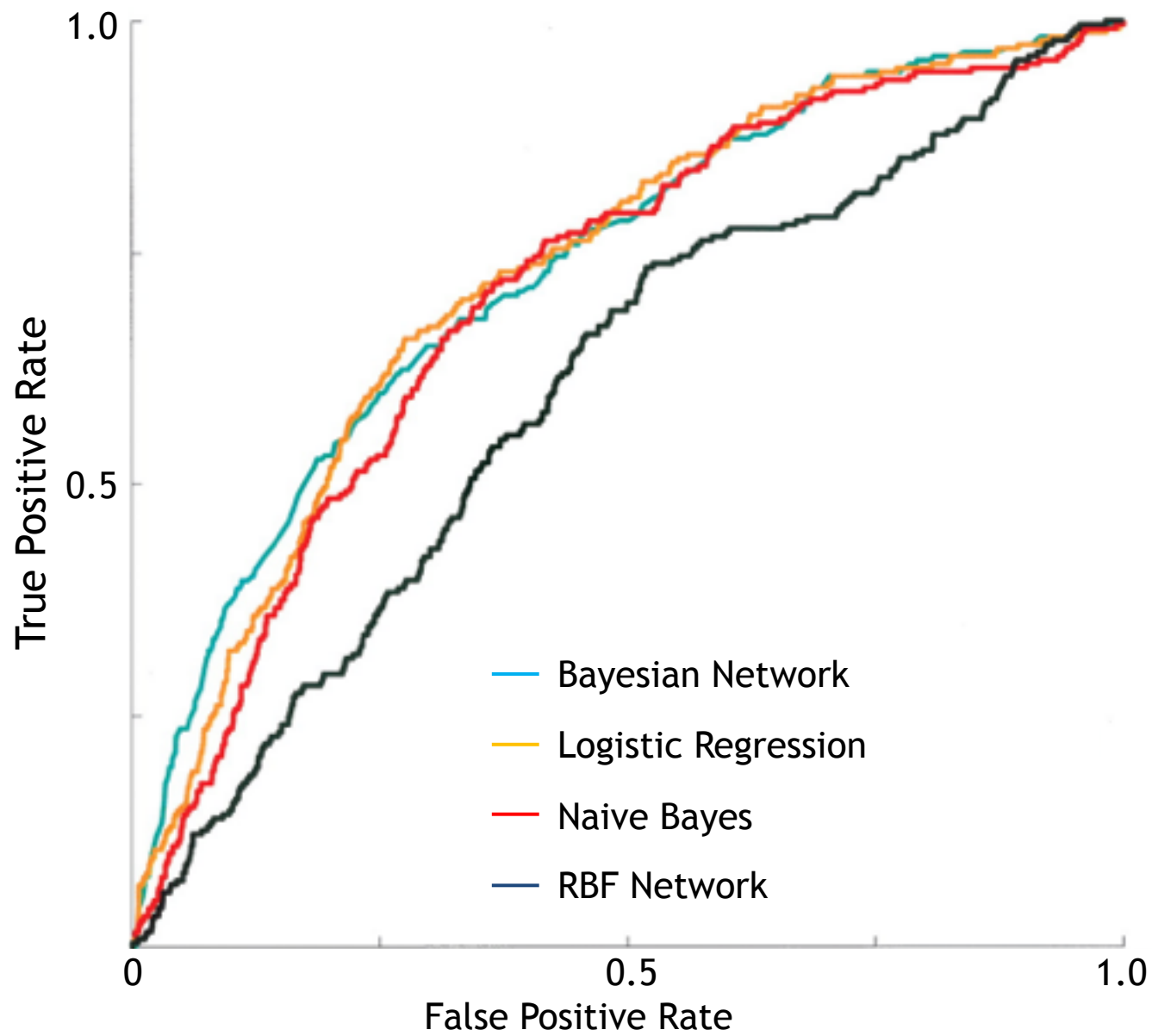


Some Extra Slides Follow



## Area under curve (AUC)

An evaluation of a classification algorithm (including all possible thresholds)



# from sklearn.metrics import .....

## Classification metrics

See the [Classification metrics](#) section of the user guide for further details.

<code>metrics.accuracy_score(y_true, y_pred[, ...])</code>	Accuracy classification score.
<code>metrics.auc(x, y[, reorder])</code>	Compute Area Under the Curve (AUC) using the trapezoidal rule
<code>metrics.average_precision_score(y_true, y_score)</code>	Compute average precision (AP) from prediction scores
<code>metrics.classification_report(y_true, y_pred)</code>	Build a text report showing the main classification metrics
<code>metrics.confusion_matrix(y_true, y_pred[, ...])</code>	Compute confusion matrix to evaluate the accuracy of a classification
<code>metrics.f1_score(y_true, y_pred[, labels, ...])</code>	Compute the F1 score, also known as balanced F-score or F-measure
<code>metrics.fbeta_score(y_true, y_pred, beta[, ...])</code>	Compute the F-beta score
<code>metrics.hamming_loss(y_true, y_pred[, classes])</code>	Compute the average Hamming loss.
<code>metrics.hinge_loss(y_true, pred_decision[, ...])</code>	Average hinge loss (non-regularized)
<code>metrics.jaccard_similarity_score(y_true, y_pred)</code>	Jaccard similarity coefficient score
<code>metrics.log_loss(y_true, y_pred[, eps, ...])</code>	Log loss, aka logistic loss or cross-entropy loss.
<code>metrics.matthews_corrcoef(y_true, y_pred)</code>	Compute the Matthews correlation coefficient (MCC) for binary classes
<code>metrics.precision_recall_curve(y_true, ...)</code>	Compute precision-recall pairs for different probability thresholds
<code>metrics.precision_score(y_true, y_pred[, ...])</code>	Compute precision, recall, F-measure and support for each class
<code>metrics.recall_score(y_true, y_pred[, ...])</code>	Compute the precision
<code>metrics.roc_auc_score(y_true, y_score)</code>	Compute the recall
<code>metrics.roc_curve(y_true, y_score)</code>	Compute Area Under the Curve (AUC) from prediction scores
<code>metrics.scorer(y_true, y_pred)</code>	Compute Decision Evaluation Metric (DEM)

Always remember,

Fit the model to a **training set**,

Calculate performance

(**accuracy**, **precision**, **recall**, **f1**, **AUC**,  
etc.) on a **test set**

or (better) on a k-fold **cross validation**  
scheme