

## Module 25 – Precision and Recall



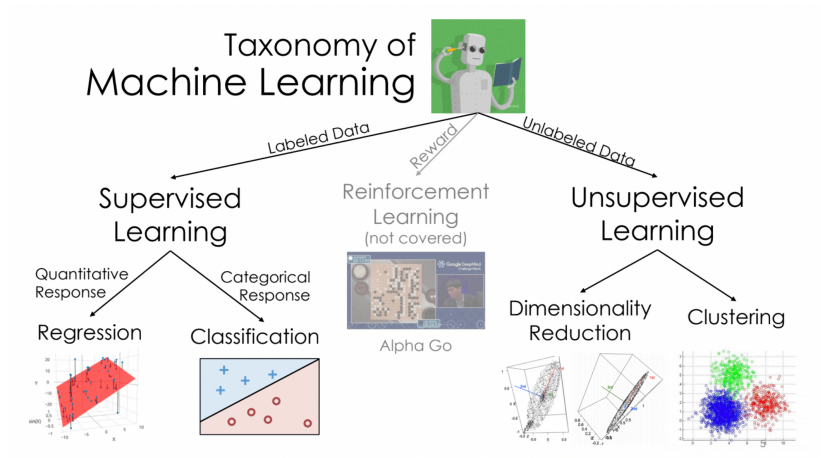
DSC 40A, Summer 2023

# Agenda

- ▶ Measuring quality of classification

## **Measuring quality of classification**

# Taxonomy of machine learning



# Classification problems

- ▶ In the classification problem, we make predictions based on data (called **training data**) for which we know the value of the **categorical** response variable.
- ▶ Example classification problems:
  - ▶ Deciding whether a patient has kidney disease.
  - ▶ Identifying handwritten digits.
  - ▶ Determining whether an avocado is ripe.
  - ▶ Predicting whether credit card activity is fraudulent.

# Assessing the quality of a classifier

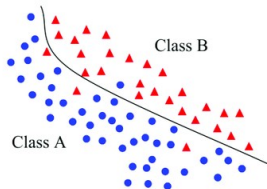
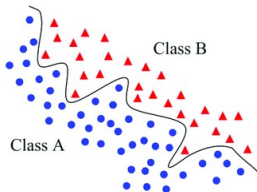
- ▶ Naive Bayes is one classification algorithm, or **classifier**, but there are many others.
- ▶ Is Naive Bayes any good? How do we measure how good of a job a classifier does?

## Discussion Question

Think back to regression (supervised learning with a quantitative response variable). How did we measure the quality of our predictions? Can we adopt a similar strategy?

## Unseen data

- ▶ A natural way to measure the quality of our classifications is to see how often we predict the right category.
- ▶ We want to make good predictions on **unseen data**. So we'll measure how often we classify examples correctly for a new set of **test data**.
- ▶ This provides us an indication about whether or not we are overfitting, and also helps us estimate general out of sample performance.



# Accuracy

- ▶ Classification **accuracy** is the proportion of examples in the test set that are correctly classified.
- ▶ Accuracy is measured on a 0 to 1 scale.



# Accuracy

- ▶ We can think of accuracy as an estimate for the probability of making a correct classification on an unseen example.

- ▶ Parameter:

$P(\text{successful classification})$

- ▶ Estimate:

$$\text{accuracy} = \frac{\text{\# correctly classified examples in test set}}{\text{size of test set}}$$

## Imbalanced classes

Alagille syndrome is a rare genetic condition that affects 1 in 40,000 people. We want to classify people as having this condition (**unhealthy**) or not having this condition (**healthy**).

### Discussion Question

Consider a classifier that classifies everyone as **healthy**.

1. What is the accuracy of this classifier?
2. What are the ethical repercussions of using this classifier?

## High accuracy is not enough

- ▶ We want to avoid overdiagnosis (telling someone they have the condition when they don't).
- ▶ We also want to avoid underdiagnosis (telling someone they're healthy when they're not).
- ▶ It's easy to avoid either one of these. It's hard to avoid both of these simultaneously, yet a good classifier should do exactly that.

## Different types of errors

	Actually <b>unhealthy</b>	Actually <b>healthy</b>
Classified as <b>unhealthy</b>		
Classified as <b>healthy</b>		

## Avoid overdiagnosis

	Actually <b>unhealthy</b>	Actually <b>healthy</b>
Classified as <b>unhealthy</b>	True positive	False positive
Classified as <b>healthy</b>	False negative	True negative

- ▶ How often does our prediction of the condition mean a person actually has the condition?
- ▶ Parameter:

$$P(\text{actually **unhealthy** | classified as **unhealthy**})$$

- ▶ Estimate:

$$\text{precision} = \frac{\# \text{ people in test set **correctly** classified as **unhealthy**}}{\# \text{ people in test set classified as **unhealthy**}}$$

## Avoid underdiagnosis

	Actually <b>unhealthy</b>	Actually <b>healthy</b>
Classified as <b>unhealthy</b>	True positive	False positive
Classified as <b>healthy</b>	False negative	True negative

- ▶ How often do we identify those that actually have the condition?
- ▶ Parameter:

$$P(\text{classified as } \mathbf{unhealthy} | \text{actually } \mathbf{unhealthy})$$

- ▶ Estimate:

$$\mathbf{recall} = \frac{\# \text{ people in test set } \mathbf{correctly} \text{ classified as } \mathbf{unhealthy}}{\# \mathbf{unhealthy} \text{ people in test set}}$$

# Precision vs. recall

	Actually <b>unhealthy</b>	Actually <b>healthy</b>
Classified as <b>unhealthy</b>	True positive	False positive
Classified as <b>healthy</b>	False negative	True negative

- ▶ Precision:

$$\begin{aligned}\text{precision} &= \frac{\text{\# people in test set **correctly** classified as **unhealthy**}}{\text{\# people in test set classified as **unhealthy**}} \\ &= \frac{\text{true positives}}{\text{true positives} + \text{false positives}}\end{aligned}$$

- ▶ Recall:

$$\begin{aligned}\text{recall} &= \frac{\text{\# people in test set **correctly** classified as **unhealthy**}}{\text{\# **unhealthy** people in test set}} \\ &= \frac{\text{true positives}}{\text{true positives} + \text{false negatives}}\end{aligned}$$

# Precision vs. recall

## Discussion Question

Consider a marketing model where you are optimizing who to show targeted ads to. In this case, our labels are “converted” vs “not converted.” In our training sample, we only have .0001 positive records (i.e. a severe class imbalance).

1. Suppose you select the top 500 ranked records as an arbitrary decision threshold. What are some reasonable values that you might expect for precision and recall?
2. If you were to increase the selection to the top 500,000, would you expect precision to go up or down? What about recall?



# Combining precision and recall

- ▶ We want high precision and high recall, but it's hard to have both.
- ▶ Let's combine them into a single measurement.
- ▶ Does the average of precision and recall work well?

$$\frac{P + R}{2}$$

- ▶ Compare:
  - ▶ Classifier A ( $P = 0, R = 1$ )
  - ▶ Classifier B ( $P = 0.5, R = 0.6$ )

# Combining precision and recall

- ▶ **Key insight:** Two moderate values are better than two extremes. Use the product, which shrinks when either term in the product is small.
- ▶ New way of combining precision and recall: **F-score**

$$\frac{2PR}{P + R}$$

- ▶ Compare:
  - ▶ Classifier A ( $P = 0, R = 1$ )
  - ▶ Classifier B ( $P = 0.5, R = 0.6$ )

## F-score

- ▶ The **F-score** combines the precision and recall of a classifier in a single measurement.

$$\frac{2PR}{P + R}$$

- ▶ Higher F-score  $\Rightarrow$  better classifier.

### Discussion Question

What would be the F-score of a “perfect classifier”?

## Summary

## Summary

- ▶ Accuracy is a simple way of measuring the quality of a classifier, but it can be misleading when classes are imbalanced.
- ▶ Precision and recall are two other ways of measuring the quality of a classifier, but they can be hard to achieve simultaneously.
- ▶ The F-score combines precision and recall into a single measurement that assesses the quality of a classifier on a 0 to 1 scale.

$$\frac{2PR}{P + R}$$