

## Lecture 25 – Precision and Recall



DSC 40A, Spring 2023

# Announcements

- ▶ Homework 7 was due last night, but can still be turned in if you have a slip day remaining. This is the **last homework!**
- ▶ Midterm 2 is Monday during lecture.
- ▶ Final Exam is **Saturday, June 10**. Given in two separate parts, both of which are optional.
  - ▶ Part 1 is 9-9:50am. Can replace Midterm 1.
  - ▶ Part 2 is 10-10:50am. Can replace Midterm 2.
- ▶ Next week is review only. No new content.

## Midterm 2 is Monday during lecture

- ▶ You may use an unlimited number of handwritten note sheets for Midterm 2 (and Final Part 2). Start working on this now as you study!
- ▶ No calculators.
- ▶ Leave all answers **unsimplified** in terms of permutations, combinations, factorials, exponents, etc.
- ▶ Assigned seats will be posted on Campuswire.
- ▶ We will not answer questions during the exam. State your assumptions if anything is unclear.

## Midterm 2 is Monday during lecture

- ▶ The exam will definitely include short-answer questions such as multiple choice or filling in the numerical answer to a probability or combinatorics question. Short-answer questions will be graded on correctness only, so you don't need to show your work or provide explanation for these questions.
- ▶ The exam may also include long-answer homework-style questions, which would require explanation and be graded with partial credit.
- ▶ Midterm 2 covers all material that was not covered on Midterm 1. Clustering is in scope, but the vast majority will be probability and combinatorics. This week's lectures are also in scope.

# Agenda

- ▶ Recap: Text classification with Naive Bayes
- ▶ Measuring quality of classification

## Text classification

## Recap: Naive Bayes for spam classification

- ▶ To classify an email, we'll use Bayes' theorem to calculate the probability of it belonging to each class:

$$P(\text{spam} \mid \text{features}) = \frac{P(\text{spam}) \cdot P(\text{features} \mid \text{spam})}{P(\text{features})}$$

$$P(\text{ham} \mid \text{features}) = \frac{P(\text{ham}) \cdot P(\text{features} \mid \text{ham})}{P(\text{features})}$$

- ▶ We'll find the larger probability by comparing numerators, and predict that class.
- ▶ To compute the numerator, we make the naive assumption that the features are conditionally independent given the class.

## Concrete example

- ▶ Dictionary: “prince”, “money”, “free”, and “just”.
- ▶ Dataset of 5 emails (red are spam, green are ham):
  - ▶ **“I am the prince of UCSD and I demand money.”**
  - ▶ **“Tapioca Express: redeem your free Thai Iced Tea!”**
  - ▶ **“DSC 10: free points if you fill out CAPEs!”**
  - ▶ **“Click here to make a tax-free donation to the IRS.”**
  - ▶ **“Free career night at Prince Street Community Center.”**



## Concrete example

- ▶ What happens if we try to classify the email “just what’s your price, prince”?

# Smoothing

- ▶ **Without** smoothing:

$$P(x^{(i)} = 1 \mid \text{spam}) \approx \frac{\# \text{ spam containing word } i}{\# \text{ spam containing word } i + \# \text{ spam not containing word } i}$$

- ▶ **With** smoothing:

$$P(x^{(i)} = 1 \mid \text{spam}) \approx \frac{(\# \text{ spam containing word } i) + 1}{(\# \text{ spam containing word } i) + 1 + (\# \text{ spam not containing word } i) + 1}$$

- ▶ When smoothing, we add 1 to the count of every group whenever we're estimating a conditional probability.

## Concrete example with smoothing

- ▶ What happens if we try to classify the email “just what’s your price, prince”?



## Modifications and extensions

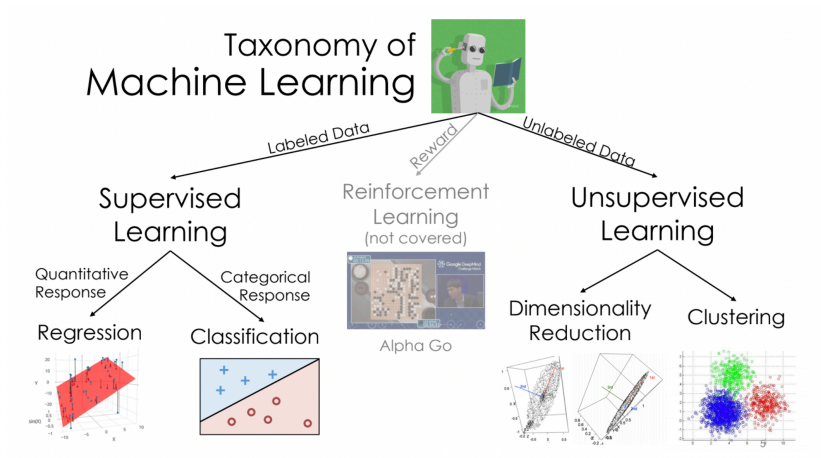
- ▶ **Idea:** Use pairs (or longer sequences) of words rather than individual words as features.
  - ▶ This better captures the dependencies between words.
  - ▶ It also leads to a much larger space of features, increasing the complexity of the algorithm.

## Modifications and extensions

- ▶ **Idea:** Use pairs (or longer sequences) of words rather than individual words as features.
  - ▶ This better captures the dependencies between words.
  - ▶ It also leads to a much larger space of features, increasing the complexity of the algorithm.
- ▶ **Idea:** Instead of recording whether each word appears, record how many times each word appears.
  - ▶ This better captures the importance of repeated words.

## **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

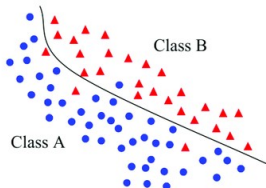
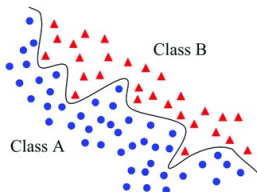
- ▶ Naive Bayes is one classification algorithm, or **classifier**, but there are many others.
- ▶ Is Naive Bayes any good? How do we measure how good of a job a classifier does?

## 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 avoids **overfitting**.



# 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

|                                | 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           |

### Discussion Question

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

1. What is the precision of this classifier?
2. What is the recall of this classifier?

## 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           |

### Discussion Question

Now consider a classifier that classifies everyone as **un-healthy**.

1. What is the precision of this classifier?
2. What is the recall of this classifier?

# 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.

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