

DSC 40A

Theoretical Foundations of Data Science I

In This Video

We've looked at mean error and mean squared error. How do both of these ways of measuring the quality of a prediction fit into a general framework?

Recommended Reading

Course Notes: Chapter 1, Section 2

A General Framework

- ▶ We started with the **mean error**:

$$R(h) = \frac{1}{n} \sum_{i=1} |y_i - h|$$

- ▶ Then we introduced the **mean squared error**:

$$R_{\text{sq}}(h) = \frac{1}{n} \sum_{i=1} (y_i - h)^2$$

- ▶ They have the same form: both are averages of some measurement that represents how different h is from the data.

A General Framework

- ▶ Definition: A **loss function** $L(h, y)$ takes in a prediction h and a right answer, y , and outputs a number measuring how far h is from y (bigger = further).
- ▶ The **absolute loss**:

$$L_{\text{abs}}(h, y) = |y - h|$$

- ▶ The **square loss**:

$$L_{\text{sq}}(h, y) = (y - h)^2$$

A General Framework

- ▶ Suppose that y_1, \dots, y_n are some data points, h is a prediction, and L is a loss function. The **empirical risk** is the average loss on the data set:

$$R_L(h) = \frac{1}{n} \sum_{i=1}^n L(h, y_i)$$

- ▶ The goal of learning: find h that minimizes R_L . This is called **empirical risk minimization (ERM)**.

Designing a learning algorithm using ERM

1. Pick a loss function.
 2. Pick a way to minimize the average loss on the data (empirical risk).
- ▶ **Key Idea:** The choice of loss function determines the properties of the result and the difficulty of computing it.

Example: 0-1 Loss

1. Pick as our loss function the **0-1 loss**:

$$L_{0,1}(h, y) = \begin{cases} 0, & \text{if } h = y \\ 1, & \text{if } h \neq y \end{cases}$$

2. Minimize empirical risk:

$$R_{0,1}(h) = \frac{1}{n} \sum_{i=1}^n L_{0,1}(h, y_i)$$

Example: 0-1 Loss

1. Pick as our loss function the **0-1 loss**:

$$L_{0,1}(h, y) = \begin{cases} 0, & \text{if } h = y \\ 1, & \text{if } h \neq y \end{cases}$$

2. Minimize empirical risk:

$$R_{0,1}(h) = \frac{1}{n} \sum_{i=1}^n L_{0,1}(h, y_i)$$

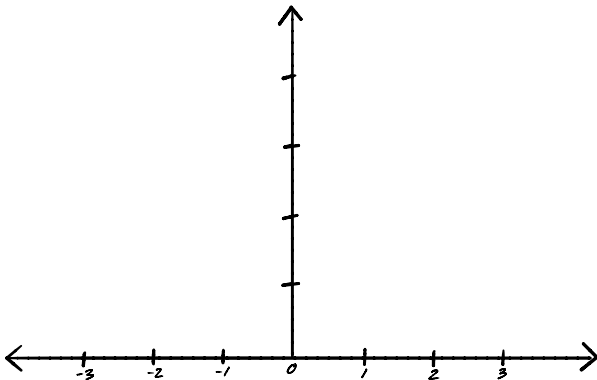
Question

Suppose y_1, \dots, y_n are all distinct. What is the value of $R_{0,1}(y_1)$?

- a) 0 b) $\frac{1}{n}$ c) $\frac{n-1}{n}$ d) 1

Minimizing Empirical Risk

$$R_{0,1}(h) = \frac{1}{n} \sum_{i=1}^n \begin{cases} 0, & \text{if } h = y_i \\ 1, & \text{if } h \neq y_i \end{cases}$$



Different Loss Functions Lead to Different Predictions

Loss	Minimizer	Outliers	Differentiable	Algorithm
L_{abs}	median	insensitive	no	not simple
L_{sq}	mean	sensitive	yes	simple, fast
$L_{0,1}$	mode	insensitive	no	simple, fast

- ▶ The optimal predictions are all **summary statistics** that measure the **center** of the data set in different ways.

Summary

- ▶ The mean error and the mean squared error fit into a general framework of **empirical risk minimization**.
- ▶ By changing the loss function, we change which prediction is considered the best.
- ▶ The optimal predictions each measure the **center** of the data set.
- ▶ **Next Time:** We'll design a more complicated loss function.