Supervised Machine Learning

Reading

- Chapters 2.1-2.3
 - Hastie, Trevor, et al. The elements of statistical learning: data mining, inference, and prediction. Vol. 2. New York: springer, 2009.
 - Available online: https://hastie.su.domains/Papers/ESLII.pdf
- Chapters 2.1
 - James, Gareth, et al. An introduction to statistical learning.
 Vol. 112. New York: springer, 2013.
 - Available online: https://www.statlearning.com/

Supervised learning from a statistical point of view

Supervised Learning Overview

• For the first few weeks of the course, we will assume we are given a training set of N labeled examples: $(x_1, y_1), ..., (x_N, y_N)$

- The x_i variables are called features or inputs
 - Also called predictors or independent variables in the statistical community
- The y_i variables are called labels or outputs
 - Also called responses or dependent variables in the statistical community
 - Can be discrete (e.g., $y_i \in \{cat, dog\}$) as in classification
 - Typically represented numerically, e.g., cat = 0, dog = 1
 - Can be continuous (e.g., $y_i \in [0,1]$) as in regression

High-level Supervised Learning Task

- We are given a labeled training set $(x_1, y_1), ..., (x_N, y_N)$
 - E.g., the features x_i are images (i.e., pixel values) and the labels are digits (from 0 to 9)
- We are assuming that the features contain sufficient information in order to determine the label
 - One could also consider hidden state problems, where some features are not observed, but out of scope for this course
- A (naïve) goal is to find a function f such that $y_i = f(x_i)$
- What is an example function f that satisfies the above?
 - -A table of rules: if you see x_i , then output y_i
 - What is wrong with that function?

Statistical Learning

- We don't want to learn a function that only works on the training data
 - -The point of learning is to *generalize* to new data
- How do we think about new data?
- Each data point is a realization of random variables X, Y
- The variables X, Y follow an unknown distribution \mathcal{D}
 - -written $(X, Y) \sim \mathcal{D}$
- Each example (x_i, y_i) has some unknown probability under \mathcal{D}
- One way to approach the problem is to try to learn ${\mathcal D}$
 - —Then for new x_i , we output the y_i that has highest probability according to \mathcal{D}
 - However, learning an arbitrary unknown ${\mathcal D}$ is hard

Statistical Learning, cont'd

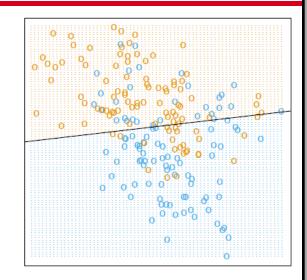
- The standard approach to the learning problem is to pick a class of functions $\mathcal F$ and choose the *best* $f\in\mathcal F$
 - We'll talk later about ways of defining *best*
- Examples of \mathcal{F} are all linear functions, all polynomials, all neural networks, etc.
- Ideally, ${\mathcal F}$ captures underlying relationship between ${\pmb X}$ and ${\pmb Y}$
 - If that is the case, the *best* f will satisfy Y = f(X)
 - for all (or most) pairs $(X, Y) \sim \mathcal{D}$
 - —The *best* f is selected based on the training set

Statistical Learning, cont'd

- This statistical formulation is nice, but why is it any good?
 - There is a lot of great theory for the above setting
- If data points are IID (independent and identically distributed)
 - For many classes \mathcal{F} , one can show that performing well on the training set implies good performance on unseen data
 - One can show that performing well on a test set (unseen during training) implies good performance on new data
- We will explore some of these notions in this class

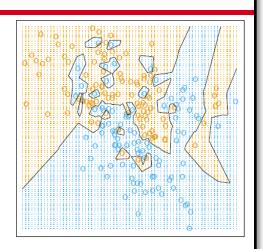
A Simple, but Effective, Classifier: Nearest Neighbor

- Suppose we are given a 2D training dataset with two classes as shown here
 - Orange and blue
 - (Ignore the line for now)
- If we see a new point, what's an easy way to choose a label?
 - Nearest neighbor!
 - Choose the label of the nearest point from the training set
- What is a problem with this strategy?

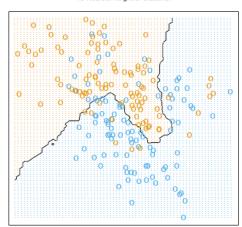


An Improved Version: K-Nearest Neighbors

- 1-Nearest neighbor similar to a table classifier
 - Decision space looks very fragmented
 - Overadapting to training data is overfitting
 - Unlikely to generalize well in some parts
- What's a simple way to improve?
- K-nearest neighbors!
- Instead of choosing the label of nearest neighbor, pick the majority label of the K-nearest neighbors
 - Much smoother decision boundary
 - More robust to noisy data







Formalizing KNN

- What is the "closest neighbor"?
 - Need a distance metric!
- Many distance metrics have been used in ML
- Most standard is Euclidean distance, or L_2 distance, or 2-norm
- The Euclidean distance between vectors $x, y \in \mathbb{R}^p$ is

$$||x - y||_2 = \sqrt{(x_1 - y_1)^2 + \dots + (x_p - y_p)^2}$$

= $\sqrt{(x - y)^T (x - y)}$

Formalizing KNN, cont'd

- Let $\mathcal{N}_K(x)$ be the indices of K closest neighbors for a point x
- Assume the labels are $y_i = 0$ for BLUE and $y_i = 1$ for ORANGE
- Then the following is the number of ORANGE votes by \mathcal{N}_K :

$$Or(\mathbf{x}) = \sum_{i \in \mathcal{N}_K(\mathbf{x})} y_i$$

• Finally, classify x as Orange if Or(x) > K/2:

$$f(x) = 1 \quad if \quad Or(x) > K/2$$

$$f(x) = 0 \quad if \quad Or(x) < K/2$$

Pick K to be odd to avoid ties