

Statistical Learning Theory

Machine Learning Summer School, Kyoto, Japan

Alexander (Sasha) Rakhlin

University of Pennsylvania, The Wharton School
Penn Research in Machine Learning (PRiML)

August 27-28, 2012

References

Parts of these lectures are based on

- ▶ O. Bousquet, S. Boucheron, G. Lugosi:
“Introduction to Statistical Learning Theory”, 2004.
- ▶ MLSS notes by O. Bousquet
- ▶ S. Mendelson: “A Few Notes on Statistical Learning Theory”
- ▶ Lecture notes by S. Shalev-Shwartz
- ▶ Lecture notes (S. R. and K. Sridharan)
http://stat.wharton.upenn.edu/~rakhlin/courses/stat928/stat928_notes.pdf

Prerequisites: a basic familiarity with Probability is assumed.

Outline

Introduction

Statistical Learning Theory

The Setting of SLT

Consistency, No Free Lunch Theorems, Bias-Variance Tradeoff

Tools from Probability, Empirical Processes

From Finite to Infinite Classes

Uniform Convergence, Symmetrization, and Rademacher Complexity

Large Margin Theory for Classification

Properties of Rademacher Complexity

Covering Numbers and Scale-Sensitive Dimensions

Faster Rates

Model Selection

Sequential Prediction / Online Learning

Motivation

Supervised Learning

Online Convex and Linear Optimization

Online-to-Batch Conversion, SVM optimization

Example #1: Handwritten Digit Recognition

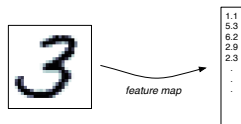
Imagine you are asked to write a computer program that recognizes postal codes on envelopes. You observe the huge amount of variation and ambiguity in the data:



One can try to hard-code all the possibilities, but likely to fail. It would be nice if a program looked at a large corpus of data and learned the distinctions!

Example #1: Handwritten Digit Recognition

Need to represent data in the computer. Pixel intensities is one possibility, but not necessarily the best one. Feature representation:



We also need to specify the “label” of this example: “3”. The *labeled example* is then

$$\left(\begin{array}{c} 1.1 \\ 5.3 \\ 6.2 \\ 2.9 \\ 2.3 \\ \vdots \\ \vdots \\ \vdots \end{array} , 3 \right)$$

After looking at many of these examples, we want the program to *predict* the label of the next hand-written digit.

Why Machine Learning?

- ▶ Impossible to hard-code all the knowledge into a computer program.
- ▶ The systems need to be adaptive to the changes in the environment.

Examples:

- ▶ Computer vision: face detection, face recognition
- ▶ Audio: voice recognition, parsing
- ▶ Text: document topics, translation
- ▶ Ad placement on web pages
- ▶ Movie recommendations
- ▶ Email spam detection

Machine Learning

(Human) learning is the process of acquiring knowledge or skill.

Quite vague. How can we build a mathematical theory for something so imprecise?

*Machine Learning is concerned with the design and analysis of algorithms that improve performance after observing **data**.*

That is, the acquired knowledge comes from data.

We need to make mathematically precise the following terms: *performance*, *improve*, *data*.

Learning from Examples

How is it possible to conclude something general from specific examples?

Learning is inherently an ill-posed problem, as there are many alternatives that could be consistent with the observed examples.

Learning can be seen as the process of *induction* (as opposed to *deduction*): “extrapolating” from examples.

Prior knowledge is how we make the problem well-posed.

Memorization is not learning, not induction. Our theory should make this apparent.

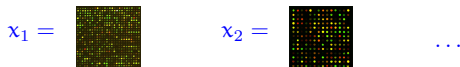
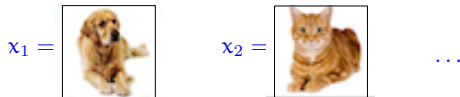
Very important to delineate assumptions. Then we will be able to prove mathematically that certain learning algorithms perform well.

Data

Space of inputs (or, predictors): \mathcal{X}

▷ e.g. $x \in \mathcal{X} \subset \{0, 1, \dots, 2^{16}\}^{64}$ is a string of pixel intensities in an 8×8 image.

▷ e.g. $x \in \mathcal{X} \subset \mathbb{R}^{33,000}$ is a set of gene expression levels.



$x_1 = \begin{bmatrix} 5 \\ 1 \\ 22 \\ \vdots \end{bmatrix}$ $x_2 = \begin{bmatrix} 1 \\ 0 \\ 17 \\ \vdots \end{bmatrix}$ # cigarettes/day
drinks/day
BMI

Data

Sometimes the space \mathcal{X} is uniquely defined for the problem. In other cases, such as in vision/text/audio applications, many possibilities exist, and a good feature representation is key to obtaining good performance.

This important part of machine learning applications will not be discussed in this lecture, and we will assume that \mathcal{X} has been chosen by the practitioner.

Data

Space of outputs (or, responses): \mathcal{Y}

- ▷ e.g. $y \in \mathcal{Y} = \{0, 1\}$ is a binary label (1 = “cat”)
- ▷ e.g. $y \in \mathcal{Y} = [0, 200]$ is life expectancy

A pair (x, y) is a *labeled* example.

▷ e.g. (x, y) is an example of an image with a label $y = 1$, which stands for the presence of a face in the image x

Dataset (or *training data*): examples $\{(x_1, y_1), \dots, (x_n, y_n)\}$

▷ e.g. a collection of images labeled according to the presence or absence of a face

The Multitude of Learning Frameworks

Presence/absence of labeled data:

- ▶ Supervised Learning: $\{(\mathbf{x}_1, \mathbf{y}_1), \dots, (\mathbf{x}_n, \mathbf{y}_n)\}$
- ▶ Unsupervised Learning: $\{\mathbf{x}_1, \dots, \mathbf{x}_n\}$
- ▶ Semi-supervised Learning: a mix of the above

This distinction is important, as labels are often difficult or expensive to obtain (e.g. can collect a large corpus of emails, but which ones are spam?)

Types of labels:

- ▶ Binary Classification / Pattern Recognition: $\mathcal{Y} = \{0, 1\}$
- ▶ Multiclass: $\mathcal{Y} = \{0, \dots, K\}$
- ▶ Regression: $\mathcal{Y} \subseteq \mathbb{R}$
- ▶ Structure prediction: \mathcal{Y} is a set of complex objects (graphs, translations)

The Multitude of Learning Frameworks

Problems also differ in the protocol for obtaining data:

- ▶ Passive
- ▶ Active

and in assumptions on data:

- ▶ Batch (typically i.i.d.)
- ▶ Online (i.i.d. or worst-case or some stochastic process)

Even more involved: Reinforcement Learning and other frameworks.

Why Theory?

“... theory is the first term in the Taylor series of practice”
– Thomas M. Cover, “1990 Shannon Lecture”

Theory and Practice should go hand-in-hand.

Boosting, Support Vector Machines – came from theoretical considerations.

Sometimes, theory is suggesting practical methods, sometimes practice comes ahead and theory tries to catch up and explain the performance.

This tutorial

First $2/3$ of the tutorial: we will study the problem of *supervised learning* (with a focus on binary classification) with an i.i.d. assumption on the data.

The last $1/3$ of the tutorial: we will turn to online learning without the i.i.d. assumption.

Outline

Introduction

Statistical Learning Theory

The Setting of SLT

Consistency, No Free Lunch Theorems, Bias-Variance Tradeoff

Tools from Probability, Empirical Processes

From Finite to Infinite Classes

Uniform Convergence, Symmetrization, and Rademacher Complexity

Large Margin Theory for Classification

Properties of Rademacher Complexity

Covering Numbers and Scale-Sensitive Dimensions

Faster Rates

Model Selection

Sequential Prediction / Online Learning

Motivation

Supervised Learning

Online Convex and Linear Optimization

Online-to-Batch Conversion, SVM optimization

Outline

Introduction

Statistical Learning Theory

The Setting of SLT

Consistency, No Free Lunch Theorems, Bias-Variance Tradeoff

Tools from Probability, Empirical Processes

From Finite to Infinite Classes

Uniform Convergence, Symmetrization, and Rademacher Complexity

Large Margin Theory for Classification

Properties of Rademacher Complexity

Covering Numbers and Scale-Sensitive Dimensions

Faster Rates

Model Selection

Sequential Prediction / Online Learning

Motivation

Supervised Learning

Online Convex and Linear Optimization

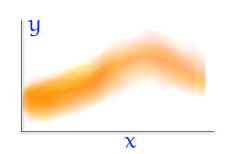
Online-to-Batch Conversion, SVM optimization

Statistical Learning Theory

The variable x is related to y , and we would like to learn this relationship from data.

The relationship is encapsulated by a distribution P on $\mathcal{X} \times \mathcal{Y}$.

Example: $x = [\text{weight}, \text{blood glucose}, \dots]$ and y is the risk of diabetes. We assume there is a relationship between x and y : it is less likely to see certain x co-occur with “low risk” and unlikely to see some other x co-occur with “high risk”. This relationship is encapsulated by $P(x, y)$.



This is an assumption about the *population* of all (x, y) . However, what we see is a *sample*.

Statistical Learning Theory

Data denoted by $\{(x_1, y_1), \dots, (x_n, y_n)\}$, where n is the sample size.

The distribution P is unknown to us (otherwise, there is no learning to be done).

The observed data are sampled independently from P (the *i.i.d. assumption*)

It is often helpful to write $P = P_x \times P_{y|x}$. The distribution P_x on the inputs is called the *marginal distribution*, while $P_{y|x}$ is the *conditional distribution*.

Statistical Learning Theory

Upon observing the training data $\{(x_1, y_1), \dots, (x_n, y_n)\}$, the learner is asked to summarize what she had learned about the relationship between x and y .

The learner's summary takes the form of a function $\hat{f}_n : \mathcal{X} \mapsto \mathcal{Y}$. The *hat* indicates that this function depends on the training data.

Learning algorithm: a mapping $\{(x_1, y_1), \dots, (x_n, y_n)\} \mapsto \hat{f}_n$.

The quality of the learned relationship is given by comparing the response $\hat{f}_n(x)$ to y for a pair (x, y) independently drawn from the same distribution \mathbb{P} :

$$\mathbb{E}_{(x,y)} \ell(\hat{f}_n(x), y)$$

where $\ell : \mathcal{Y} \times \mathcal{Y} \mapsto \mathbb{R}$ is a *loss function*. This is our **measure of performance**.

Loss Functions

- ▶ Indicator loss (classification): $\ell(\mathbf{y}, \mathbf{y}') = \mathbf{I}_{\{\mathbf{y} \neq \mathbf{y}'\}}$
- ▶ Square loss: $\ell(\mathbf{y}, \mathbf{y}') = (\mathbf{y} - \mathbf{y}')^2$
- ▶ Absolute loss: $\ell(\mathbf{y}, \mathbf{y}') = |\mathbf{y} - \mathbf{y}'|$

Examples

Probably the simplest learning algorithm that you are probably familiar with is *linear least squares*:

Given $(\mathbf{x}_1, \mathbf{y}_1), \dots, (\mathbf{x}_n, \mathbf{y}_n)$, let

$$\hat{\beta} = \arg \min_{\beta \in \mathbb{R}^d} \frac{1}{n} \sum_{i=1}^n (y_i - \langle \beta, \mathbf{x}_i \rangle)^2$$

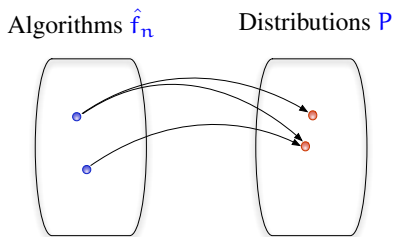
and define

$$\hat{f}_n(\mathbf{x}) = \langle \hat{\beta}, \mathbf{x} \rangle$$

Another basic method is *regularized least squares*:

$$\hat{\beta} = \arg \min_{\beta \in \mathbb{R}^d} \frac{1}{n} \sum_{i=1}^n (y_i - \langle \beta, \mathbf{x}_i \rangle)^2 + \lambda \|\beta\|^2$$

Methods vs Problems



Expected Loss and Empirical Loss

The *expected loss* of any function $f: \mathcal{X} \mapsto \mathcal{Y}$ is

$$\mathbf{L}(f) = \mathbb{E}\ell(f(x), y)$$

Since \mathbf{P} is unknown, we cannot calculate $\mathbf{L}(f)$.

However, we can calculate the *empirical loss* of $f: \mathcal{X} \mapsto \mathcal{Y}$

$$\hat{\mathbf{L}}(f) = \frac{1}{n} \sum_{i=1}^n \ell(f(x_i), y_i)$$

... again, what is random here?

Since data $(x_1, y_1), \dots, (x_n, y_n)$ are a random i.i.d. draw from \mathcal{P} ,

- ▶ $\hat{\mathbf{L}}(f)$ is a random quantity
- ▶ \hat{f}_n is a random quantity (a random function, output of our learning procedure after seeing data)
- ▶ hence, $\mathbf{L}(\hat{f}_n)$ is also a random quantity
- ▶ for a given $f: \mathcal{X} \rightarrow \mathcal{Y}$, the quantity $\mathbf{L}(f)$ is *not* random!

It is important that these are understood before we proceed further.

The Gold Standard

Within the framework we set up, the smallest expected loss is achieved by the *Bayes optimal* function

$$f^* = \arg \min_f \mathbf{L}(f)$$

where the minimization is over all (measurable) prediction rules $f: \mathcal{X} \mapsto \mathcal{Y}$.

The value of the lowest expected loss is called the *Bayes error*:

$$\mathbf{L}(f^*) = \inf_f \mathbf{L}(f)$$

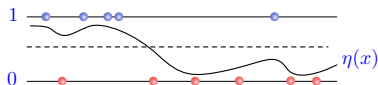
Of course, we cannot calculate any of these quantities since \mathbf{P} is unknown.

Bayes Optimal Function

Bayes optimal function f^* takes on the following forms in these two particular cases:

- ▶ Binary classification ($\mathcal{Y} = \{0, 1\}$) with the indicator loss:

$$f^*(x) = \mathbf{I}_{\{\eta(x) \geq 1/2\}}, \quad \text{where} \quad \eta(x) = \mathbb{E}[Y|X = x]$$

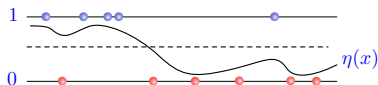


Bayes Optimal Function

Bayes optimal function f^* takes on the following forms in these two particular cases:

- ▶ Binary classification ($\mathcal{Y} = \{0, 1\}$) with the indicator loss:

$$f^*(x) = \mathbf{I}_{\{\eta(x) \geq 1/2\}}, \quad \text{where} \quad \eta(x) = \mathbb{E}[Y|X = x]$$



- ▶ Regression ($\mathcal{Y} = \mathbb{R}$) with squared loss:

$$f^*(x) = \eta(x), \quad \text{where} \quad \eta(x) = \mathbb{E}[Y|X = x]$$

The big question: is there a way to construct a learning algorithm with a guarantee that

$$\mathbf{L}(\hat{f}_n) - \mathbf{L}(f^*)$$

is small for large enough sample size n ?

Outline

Introduction

Statistical Learning Theory

The Setting of SLT

Consistency, No Free Lunch Theorems, Bias-Variance Tradeoff

Tools from Probability, Empirical Processes

From Finite to Infinite Classes

Uniform Convergence, Symmetrization, and Rademacher Complexity

Large Margin Theory for Classification

Properties of Rademacher Complexity

Covering Numbers and Scale-Sensitive Dimensions

Faster Rates

Model Selection

Sequential Prediction / Online Learning

Motivation

Supervised Learning

Online Convex and Linear Optimization

Online-to-Batch Conversion, SVM optimization

Consistency

An algorithm that ensures

$$\lim_{n \rightarrow \infty} \mathbf{L}(\hat{f}_n) = \mathbf{L}(f^*) \quad \text{almost surely}$$

is called *consistent*. Consistency ensures that our algorithm is approaching the best possible prediction performance as the sample size increases.

The good news: consistency is possible to achieve.

- ▶ easy if \mathcal{X} is a finite or countable set
- ▶ not too hard if \mathcal{X} is infinite, and the underlying relationship between \mathbf{x} and \mathbf{y} is “continuous”

The bad news...

In general, we cannot prove anything “interesting” about $\mathbf{L}(\hat{f}_n) - \mathbf{L}(f^*)$, unless we make further assumptions (incorporate *prior knowledge*).

What do we mean by “nothing interesting”? This is the subject of the so-called “No Free Lunch” Theorems. Unless we posit further assumptions,

The bad news...

In general, we cannot prove anything “interesting” about $\mathbf{L}(\hat{f}_n) - \mathbf{L}(f^*)$, unless we make further assumptions (incorporate *prior knowledge*).

What do we mean by “nothing interesting”? This is the subject of the so-called “No Free Lunch” Theorems. Unless we posit further assumptions,

- ▶ For any algorithm \hat{f}_n , any n and any $\epsilon > 0$, there exists a distribution \mathbf{P} such that $\mathbf{L}(f^*) = 0$ and

$$\mathbb{E}\mathbf{L}(\hat{f}_n) \geq \frac{1}{2} - \epsilon$$

The bad news...

In general, we cannot prove anything “interesting” about $\mathbf{L}(\hat{f}_n) - \mathbf{L}(f^*)$, unless we make further assumptions (incorporate *prior knowledge*).

What do we mean by “nothing interesting”? This is the subject of the so-called “No Free Lunch” Theorems. Unless we posit further assumptions,

- ▶ For any algorithm \hat{f}_n , any n and any $\epsilon > 0$, there exists a distribution \mathbf{P} such that $\mathbf{L}(f^*) = 0$ and

$$\mathbb{E}\mathbf{L}(\hat{f}_n) \geq \frac{1}{2} - \epsilon$$

- ▶ For any algorithm \hat{f}_n , and any sequence α_n that converges to 0, there exists a probability distribution \mathbf{P} such that $\mathbf{L}(f^*) = 0$ and for all n

$$\mathbb{E}\mathbf{L}(\hat{f}_n) \geq \alpha_n$$

Reference: (Devroye, Györfi, Lugosi: *A Probabilistic Theory of Pattern Recognition*),
(Bousquet, Boucheron, Lugosi, 2004)

is this really “bad news”?

Not really. We always have some domain knowledge.

Two ways of incorporating prior knowledge:

- ▶ Direct way: assume that the distribution \mathbf{P} is not arbitrary (also known as a modeling approach, generative approach, statistical modeling)
- ▶ Indirect way: redefine the goal to perform as well as a reference set \mathcal{F} of predictors:

$$\mathbf{L}(\hat{f}_n) - \inf_{f \in \mathcal{F}} \mathbf{L}(f)$$

This is known as a discriminative approach. \mathcal{F} encapsulates our *inductive bias*.

Pros/Cons of the two approaches

Pros of the **discriminative** approach: we never assume that P takes some particular form, but we rather put our prior knowledge into “what are the types of predictors that will do well”. Cons: cannot really interpret \hat{f}_n .

Pros of the **generative** approach: can estimate the model / parameters of the distribution (*inference*). Cons: it is not clear what the analysis says if the assumption is actually violated.

Both approaches have their advantages. A machine learning researcher or practitioner should ideally know both and should understand their strengths and weaknesses.

In this tutorial we only focus on the discriminative approach.

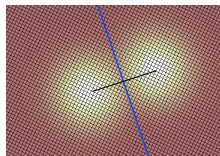
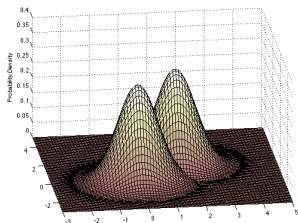
Example: Linear Discriminant Analysis

Consider the classification problem with $\mathcal{Y} = \{0, 1\}$. Suppose the *class-conditional densities* are multivariate Gaussian with the same covariance $\Sigma = \mathbf{I}$:

$$p(\mathbf{x}|\mathbf{y} = 0) = (2\pi)^{-k/2} \exp\left\{-\frac{1}{2}\|\mathbf{x} - \mu_0\|^2\right\}$$

and

$$p(\mathbf{x}|\mathbf{y} = 1) = (2\pi)^{-k/2} \exp\left\{-\frac{1}{2}\|\mathbf{x} - \mu_1\|^2\right\}$$



The “best” (Bayes) classifier is $f^* = \mathbf{I}_{\{p(\mathbf{y}=1|\mathbf{x}) \geq 1/2\}}$ which corresponds to the half-space defined by the decision boundary $p(\mathbf{x}|\mathbf{y} = 1) \geq p(\mathbf{x}|\mathbf{y} = 0)$. This boundary is *linear*.

Example: Linear Discriminant Analysis

The (linear) optimal decision boundary comes from our generative assumption on the form of the underlying distribution.

Alternatively, we could have indirectly postulated that we will be looking for a linear discriminant between the two classes, without making distributional assumptions. Such linear discriminant (classification) functions are

$$\mathbf{I}_{\{(w, x) \geq b\}}$$

for a unit-norm w and some bias $b \in \mathbb{R}$.

Quadratic Discriminant Analysis: If unequal correlation matrices Σ_1 and Σ_2 are assumed, the resulting boundary is quadratic. We can then define classification function by

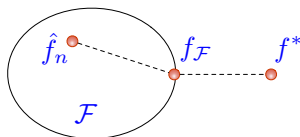
$$\mathbf{I}_{\{q(x) \geq 0\}}$$

where $q(x)$ is a quadratic function.

Bias-Variance Tradeoff

How do we choose the inductive bias \mathcal{F} ?

$$\mathbf{L}(\hat{f}_n) - \mathbf{L}(f^*) = \underbrace{\mathbf{L}(\hat{f}_n) - \inf_{f \in \mathcal{F}} \mathbf{L}(f)}_{\text{Estimation Error}} + \underbrace{\inf_{f \in \mathcal{F}} \mathbf{L}(f) - \mathbf{L}(f^*)}_{\text{Approximation Error}}$$



Clearly, the two terms are at odds with each other:

- ▶ Making \mathcal{F} larger means smaller approximation error but (as we will see) larger estimation error
- ▶ Taking a larger sample n means smaller estimation error and has no effect on the approximation error.
- ▶ Thus, it makes sense to trade off size of \mathcal{F} and n . This is called *Structural Risk Minimization*, or *Method of Sieves*, or *Model Selection*.

Bias-Variance Tradeoff

We will only focus on the estimation error, yet the ideas we develop will make it possible to read about model selection on your own.

Note: if we guessed correctly and $f^* \in \mathcal{F}$, then

$$\mathbf{L}(\hat{f}_n) - \mathbf{L}(f^*) = \mathbf{L}(\hat{f}_n) - \inf_{f \in \mathcal{F}} \mathbf{L}(f)$$

For a particular problem, one hopes that prior knowledge about the problem can ensure that the approximation error $\inf_{f \in \mathcal{F}} \mathbf{L}(f) - \mathbf{L}(f^*)$ is small.

Occam's Razor

Occam's Razor is often quoted as a principle for choosing the simplest theory or explanation out of the possible ones.

However, this is a rather philosophical argument since simplicity is not uniquely defined. We will discuss this issue later.

What we will do is to try to understand “complexity” when it comes to behavior of certain stochastic processes. Such a question will be well-defined mathematically.

Looking Ahead

So far: represented prior knowledge by means of the class \mathcal{F} .

Looking forward, we can find an algorithm that, after looking at a dataset of size n , produces \hat{f}_n such that

$$L(\hat{f}_n) - \inf_{f \in \mathcal{F}} L(f)$$

decreases (in a certain sense which we will make precise) at a non-trivial rate which depends on “richness” of \mathcal{F} .

This will give a *sample complexity* guarantee: how many samples are needed to make the error smaller than a desired accuracy.

Outline

Introduction

Statistical Learning Theory

The Setting of SLT

Consistency, No Free Lunch Theorems, Bias-Variance Tradeoff

Tools from Probability, Empirical Processes

From Finite to Infinite Classes

Uniform Convergence, Symmetrization, and Rademacher Complexity

Large Margin Theory for Classification

Properties of Rademacher Complexity

Covering Numbers and Scale-Sensitive Dimensions

Faster Rates

Model Selection

Sequential Prediction / Online Learning

Motivation

Supervised Learning

Online Convex and Linear Optimization

Online-to-Batch Conversion, SVM optimization

Types of Bounds

In expectation **vs** in probability (control the mean vs control the tails):

$$\mathbb{E} \left\{ \mathbf{L}(\hat{f}_n) - \inf_{f \in \mathcal{F}} \mathbf{L}(f) \right\} < \psi(n) \quad \text{vs} \quad \mathbb{P} \left(\mathbf{L}(\hat{f}_n) - \inf_{f \in \mathcal{F}} \mathbf{L}(f) \geq \epsilon \right) < \psi(n, \epsilon)$$

Types of Bounds

In expectation **vs** in probability (control the mean **vs** control the tails):

$$\mathbb{E} \left\{ \mathbf{L}(\hat{f}_n) - \inf_{f \in \mathcal{F}} \mathbf{L}(f) \right\} < \psi(n) \quad \text{vs} \quad \mathbb{P} \left(\mathbf{L}(\hat{f}_n) - \inf_{f \in \mathcal{F}} \mathbf{L}(f) \geq \epsilon \right) < \psi(n, \epsilon)$$

The in-probability bound can be inverted as

$$\mathbb{P} \left(\mathbf{L}(\hat{f}_n) - \inf_{f \in \mathcal{F}} \mathbf{L}(f) \geq \phi(\delta, n) \right) < \delta$$

by setting $\delta := \psi(\epsilon, n)$ and solving for ϵ .

In this lecture, we are after the function $\phi(\delta, n)$. We will call it “the rate”.

“With high probability” typically means logarithmic dependence of $\phi(\delta, n)$ on $1/\delta$. Very desirable: the bound grows only modestly even for high confidence bounds.

Sample Complexity

Sample complexity is the sample size required by the algorithm \hat{f}_n to guarantee $\mathbf{L}(\hat{f}_n) - \inf_{f \in \mathcal{F}} \mathbf{L}(f) \leq \epsilon$ with probability at least $1 - \delta$. Of course, we just need to invert a bound

$$\mathbb{P} \left(\mathbf{L}(\hat{f}_n) - \inf_{f \in \mathcal{F}} \mathbf{L}(f) \geq \phi(\delta, n) \right) < \delta$$

by setting $\epsilon := \phi(\delta, n)$ and solving for n . In other words, $n(\epsilon, \delta)$ is sample complexity of the algorithm \hat{f}_n if

$$\mathbb{P} \left(\mathbf{L}(\hat{f}_n) - \inf_{f \in \mathcal{F}} \mathbf{L}(f) \geq \epsilon \right) \leq \delta$$

as soon as $n \geq n(\epsilon, \delta)$.

Hence, “rate” can be translated into “sample complexity” and vice versa.

Easy to remember: rate $O(1/\sqrt{n})$ means $O(1/\epsilon^2)$ sample complexity, whereas rate $O(1/n)$ is a smaller $O(1/\epsilon)$ sample complexity.

Types of Bounds

Other distinctions to keep in mind: We can ask for bounds (either in expectation or in probability) on the following random variables:

$$\mathbf{L}(\hat{f}_n) - \mathbf{L}(f^*) \quad (\text{A})$$

$$\mathbf{L}(\hat{f}_n) - \inf_{f \in \mathcal{F}} \mathbf{L}(f) \quad (\text{B})$$

$$\mathbf{L}(\hat{f}_n) - \hat{\mathbf{L}}(\hat{f}_n) \quad (\text{C})$$

$$\sup_{f \in \mathcal{F}} \{\mathbf{L}(f) - \hat{\mathbf{L}}(f)\} \quad (\text{D})$$

$$\sup_{f \in \mathcal{F}} \{\mathbf{L}(f) - \hat{\mathbf{L}}(f) - \text{pen}_n(f)\} \quad (\text{E})$$

Let's make sure we understand the differences between these random quantities!

Types of Bounds

Upper bounds on (D) and (E) are used as *tools* for achieving the other bounds. Let's see why.

Obviously, for any algorithm that outputs $\hat{f}_n \in \mathcal{F}$,

$$\mathbf{L}(\hat{f}_n) - \hat{\mathbf{L}}(\hat{f}_n) \leq \sup_{f \in \mathcal{F}} \{\mathbf{L}(f) - \hat{\mathbf{L}}(f)\}$$

and so a bound on (D) implies a bound on (C).

How about a bound on (B)? Is it implied by (C) or (D)? It depends on what the algorithm does!

Denote $f_{\mathcal{F}} = \arg \min_{f \in \mathcal{F}} \mathbf{L}(f)$. Suppose (D) is small. It then makes sense to ask the learning algorithm to minimize or (approximately minimize) the empirical error (why?)

Canonical Algorithms

Empirical Risk Minimization (ERM) algorithm:

$$\hat{f}_n = \arg \min_{f \in \mathcal{F}} \hat{\mathbf{L}}(f)$$

Regularized Empirical Risk Minimization algorithm:

$$\hat{f}_n = \arg \min_{f \in \mathcal{F}} \hat{\mathbf{L}}(f) + \text{pen}_n(f)$$

We will deal with the regularized ERM a bit later. For now, let's focus on ERM.

Remark: to actually *compute* $f \in \mathcal{F}$ minimizing the above objectives, one needs to employ some optimization methods. In practice, the objective might be optimized only approximately.

Performance of ERM

If \hat{f}_n is an ERM,

$$\begin{aligned} \mathbf{L}(\hat{f}_n) - \mathbf{L}(f_{\mathcal{F}}) &\leq \underbrace{\{\mathbf{L}(\hat{f}_n) - \hat{\mathbf{L}}(\hat{f}_n)\}}_{(C)} + \{\hat{\mathbf{L}}(\hat{f}_n) - \hat{\mathbf{L}}(f_{\mathcal{F}})\} + \{\hat{\mathbf{L}}(f_{\mathcal{F}}) - \mathbf{L}(f_{\mathcal{F}})\} \\ &\leq \underbrace{\{\mathbf{L}(\hat{f}_n) - \hat{\mathbf{L}}(\hat{f}_n)\}}_{(C)} + \{\hat{\mathbf{L}}(f_{\mathcal{F}}) - \mathbf{L}(f_{\mathcal{F}})\} \\ &\leq \underbrace{\sup_{f \in \mathcal{F}} \{\mathbf{L}(f) - \hat{\mathbf{L}}(f)\}}_{(D)} + \{\hat{\mathbf{L}}(f_{\mathcal{F}}) - \mathbf{L}(f_{\mathcal{F}})\} \end{aligned}$$

because the second term is negative. So, (C) also implies a bound on (B) when \hat{f}_n is ERM (or “close” to ERM). Also, (D) also implies a bound on (B).

What about this extra term $\hat{\mathbf{L}}(f_{\mathcal{F}}) - \mathbf{L}(f_{\mathcal{F}})$? Central Limit Theorem says that for i.i.d. random variables with bounded second moment, the average converges to the expectation. Let’s quantify this.

Hoeffding Inequality

Let W, W_1, \dots, W_n be i.i.d. such that $\mathbb{P}(a \leq W \leq b) = 1$. Then

$$\mathbb{P}\left(\mathbb{E}W - \frac{1}{n} \sum_{i=1}^n W_i > \epsilon\right) \leq \exp\left(-\frac{2n\epsilon^2}{(b-a)^2}\right)$$

and

$$\mathbb{P}\left(\frac{1}{n} \sum_{i=1}^n W_i - \mathbb{E}W > \epsilon\right) \leq \exp\left(-\frac{2n\epsilon^2}{(b-a)^2}\right)$$

Let $W_i = \ell(f_{\mathcal{F}}(x_i), y_i)$. Clearly, W_1, \dots, W_n are i.i.d. Then,

$$\mathbb{P}\left(|\mathbf{L}(f_{\mathcal{F}}) - \hat{\mathbf{L}}(f_{\mathcal{F}})| > \epsilon\right) \leq 2 \exp\left(-\frac{2n\epsilon^2}{(b-a)^2}\right)$$

assuming $a \leq \ell(f_{\mathcal{F}}(x), y) \leq b$ for all $x \in \mathcal{X}, y \in \mathcal{Y}$.

Wait, Are We Done?

Can't we conclude directly that (C) is small? That is,

$$\mathbb{P}\left(\mathbb{E}\ell(\hat{f}_n(x), \mathbf{y}) - \frac{1}{n} \sum_{i=1}^n \ell(\hat{f}_n(x_i), y_i) > \epsilon\right) \leq 2 \exp\left(-\frac{2n\epsilon^2}{(b-a)^2}\right) \quad ?$$

Wait, Are We Done?

Can't we conclude directly that (C) is small? That is,

$$\mathbb{P}\left(\mathbb{E}\ell(\hat{f}_n(x), y) - \frac{1}{n} \sum_{i=1}^n \ell(\hat{f}_n(x_i), y_i) > \epsilon\right) \leq 2 \exp\left(-\frac{2n\epsilon^2}{(b-a)^2}\right) \quad ?$$

No! The random variables $\ell(\hat{f}_n(x_i), y_i)$ are **not** necessarily independent and it is possible that

$$\mathbb{E}\ell(\hat{f}_n(x), y) = \mathbb{E}W \neq \mathbb{E}\ell(\hat{f}_n(x_i), y_i) = \mathbb{E}W_i$$

The expected loss is “out of sample performance” while the second term is “in sample”.

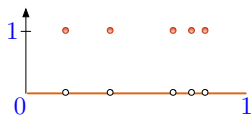
We say that $\ell(\hat{f}_n(x_i), y_i)$ is a *biased estimate* of $\mathbb{E}\ell(\hat{f}_n(x), y)$.

How bad can this bias be?

Example

- ▶ $\mathcal{X} = [0, 1]$, $\mathcal{Y} = \{0, 1\}$
- ▶ $\ell(f(X_i), Y_i) = \mathbf{I}_{\{f(X_i) \neq Y_i\}}$
- ▶ distribution $\mathbf{P} = \mathbf{P}_x \times \mathbf{P}_{y|x}$ with $\mathbf{P}_x = \text{Unif}[0, 1]$ and $\mathbf{P}_{y|x} = \delta_{y=1}$
- ▶ function class

$$\mathcal{F} = \cup_{n \in \mathbb{N}} \{f = f_S : S \subset \mathcal{X}, |S| = n, f_S(x) = \mathbf{I}_{\{x \in S\}}\}$$



ERM \hat{f}_n **memorizes** (perfectly fits) the data, but has no ability to generalize. Observe that

$$0 = \mathbb{E}\ell(\hat{f}_n(x_i), y_i) \neq \mathbb{E}\ell(\hat{f}_n(x), y) = 1$$

This phenomenon is called *overfitting*.

Example

Not only is (C) large in this example. Also, uniform deviations (D) do not converge to zero.

For any $n \in \mathbb{N}$ and any $(x_1, y_1), \dots, (x_n, y_n) \sim P$

$$\sup_{f \in \mathcal{F}} \left\{ \mathbb{E}_{x, y} \ell(f(x), y) - \frac{1}{n} \sum_{i=1}^n \ell(f(x_i), y_i) \right\} = 1$$

Where do we go from here? Two approaches:

1. understand how to upper bound uniform deviations (D)
2. find properties of algorithms that limit in some way the bias of $\ell(\hat{f}_n(x_i), y_i)$. *Stability* and *compression* are two such approaches.

Uniform Deviations

We first focus on understanding

$$\sup_{f \in \mathcal{F}} \left\{ \mathbb{E}_{\mathbf{x}, \mathbf{y}} \ell(f(\mathbf{x}), \mathbf{y}) - \frac{1}{n} \sum_{i=1}^n \ell(f(\mathbf{x}_i), \mathbf{y}_i) \right\}$$

If $\mathcal{F} = \{f_0\}$ consists of a single function, then clearly

$$\sup_{f \in \mathcal{F}} \left\{ \mathbb{E} \ell(f(\mathbf{x}), \mathbf{y}) - \frac{1}{n} \sum_{i=1}^n \ell(f(\mathbf{x}_i), \mathbf{y}_i) \right\} = \left\{ \mathbb{E} \ell(f_0(\mathbf{x}), \mathbf{y}) - \frac{1}{n} \sum_{i=1}^n \ell(f_0(\mathbf{x}_i), \mathbf{y}_i) \right\}$$

This quantity is $O_P(1/\sqrt{n})$ by Hoeffding's inequality, assuming $a \leq \ell(f_0(\mathbf{x}), \mathbf{y}) \leq b$.

Moral: for “simple” classes \mathcal{F} the uniform deviations (D) can be bounded while for “rich” classes not. We will see how far we can push the size of \mathcal{F} .

A bit of notation to simplify things...

To ease the notation,

- ▶ Let $\mathbf{z}_i = (x_i, y_i)$ so that the training data is $\{\mathbf{z}_1, \dots, \mathbf{z}_n\}$
- ▶ $g(\mathbf{z}) = \ell(f(\mathbf{x}), \mathbf{y})$ for $\mathbf{z} = (x, y)$
- ▶ Loss class $\mathcal{G} = \{g : g(\mathbf{z}) = \ell(f(\mathbf{x}), \mathbf{y})\} = \ell \circ \mathcal{F}$
- ▶ $\hat{g}_n = \ell(\hat{f}_n(\cdot), \cdot)$, $g_{\mathcal{G}} = \ell(f_{\mathcal{F}}(\cdot), \cdot)$
- ▶ $g^* = \arg \min_g \mathbb{E}g(\mathbf{z}) = \ell(f^*(\cdot), \cdot)$ is Bayes optimal (loss) function

We can now work with the set \mathcal{G} , but keep in mind that each $g \in \mathcal{G}$ corresponds to an $f \in \mathcal{F}$:

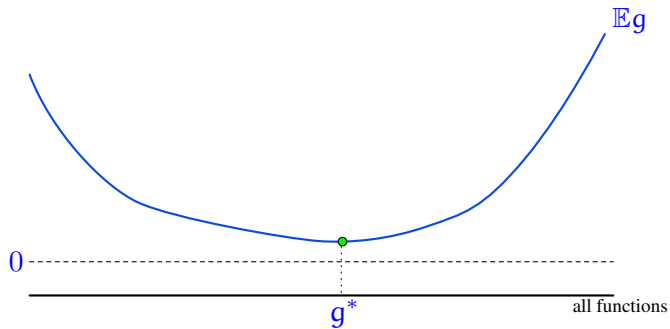
$$g \in \mathcal{G} \iff f \in \mathcal{F}$$

Once again, the quantity of interest is

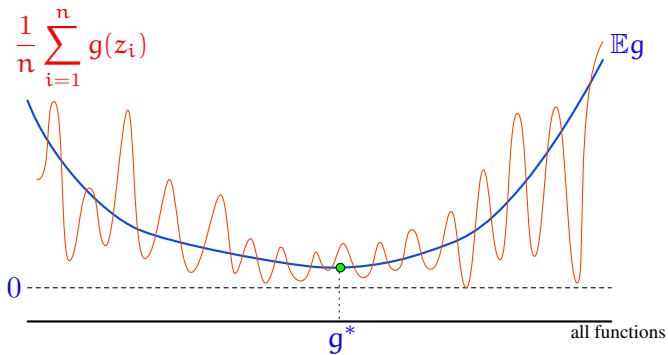
$$\sup_{g \in \mathcal{G}} \left\{ \mathbb{E}g(\mathbf{z}) - \frac{1}{n} \sum_{i=1}^n g(\mathbf{z}_i) \right\}$$

On the next slide, we visualize deviations $\mathbb{E}g(\mathbf{z}) - \frac{1}{n} \sum_{i=1}^n g(\mathbf{z}_i)$ for all possible functions g and discuss all the concepts introduced so far.

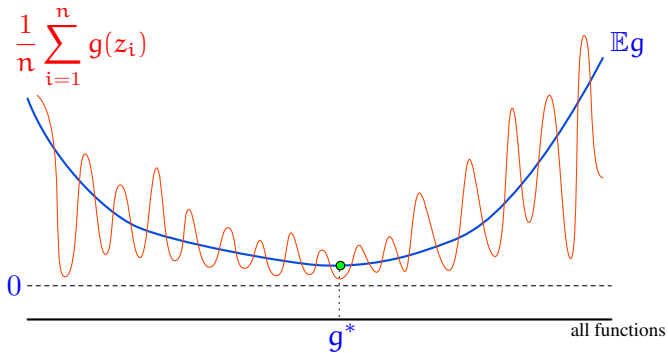
Empirical Process Viewpoint



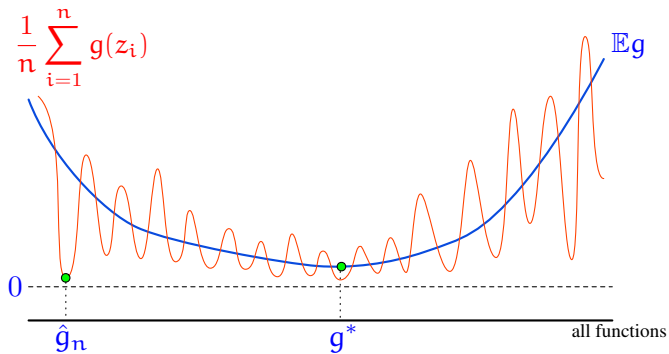
Empirical Process Viewpoint



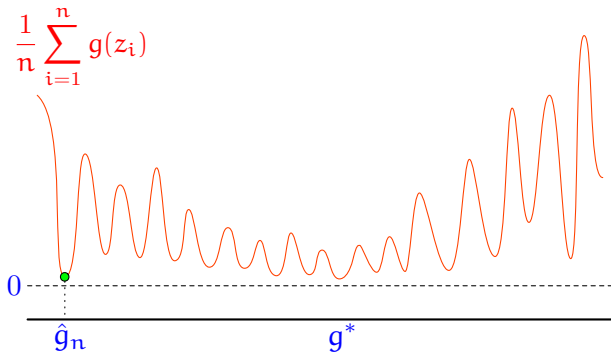
Empirical Process Viewpoint



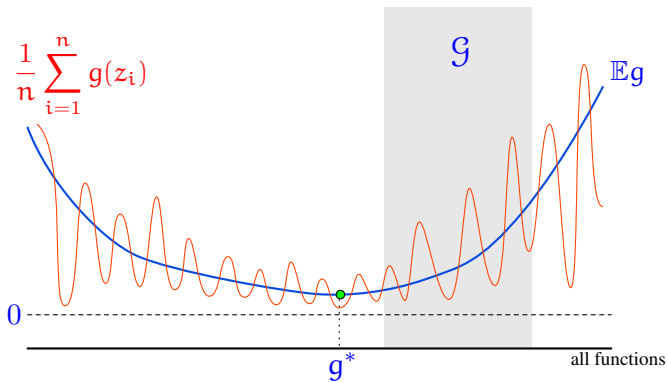
Empirical Process Viewpoint



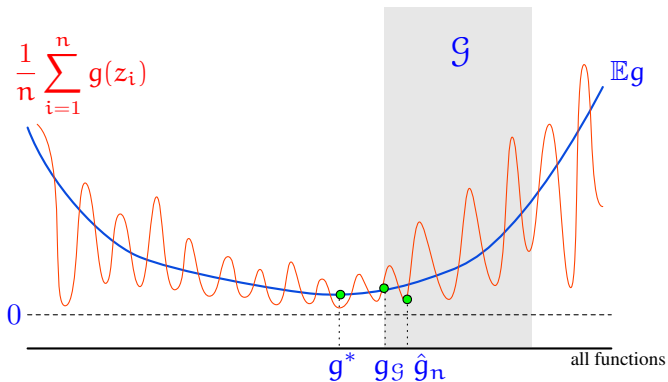
Empirical Process Viewpoint



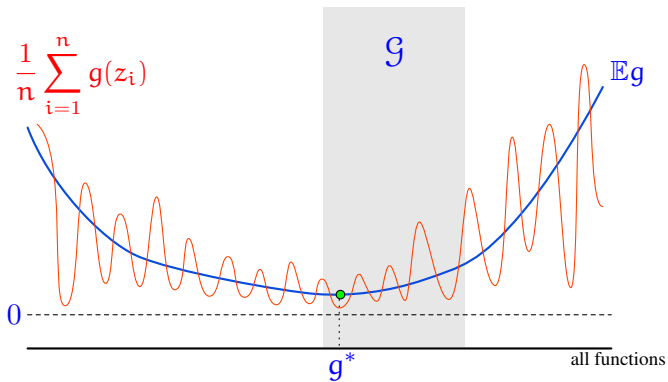
Empirical Process Viewpoint



Empirical Process Viewpoint



Empirical Process Viewpoint



Empirical Process Viewpoint

A *stochastic process* is a collection of random variables indexed by some set.

An *empirical process* is a stochastic process

$$\left\{ \mathbb{E}g(z) - \frac{1}{n} \sum_{i=1}^n g(z_i) \right\}_{g \in \mathcal{G}}$$

indexed by a function class \mathcal{G} .

Uniform Law of Large Numbers:

$$\sup_{g \in \mathcal{G}} \left| \mathbb{E}g - \frac{1}{n} \sum_{i=1}^n g(z_i) \right| \rightarrow 0$$

in probability.

Empirical Process Viewpoint

A *stochastic process* is a collection of random variables indexed by some set.

An *empirical process* is a stochastic process

$$\left\{ \mathbb{E}g(z) - \frac{1}{n} \sum_{i=1}^n g(z_i) \right\}_{g \in \mathcal{G}}$$

indexed by a function class \mathcal{G} .

Uniform Law of Large Numbers:

$$\sup_{g \in \mathcal{G}} \left| \mathbb{E}g - \frac{1}{n} \sum_{i=1}^n g(z_i) \right| \rightarrow 0$$

in probability.

Key question: How “big” can \mathcal{G} be for the supremum of the empirical process to still be manageable?

Union Bound (Boole's inequality)

Boole's inequality: for a finite or countable set of events,

$$\mathbb{P}(\cup_j A_j) \leq \sum_j \mathbb{P}(A_j)$$

Let $\mathcal{G} = \{g_1, \dots, g_N\}$. Then

$$\mathbb{P}\left(\exists g \in \mathcal{G} : \mathbb{E}g - \frac{1}{n} \sum_{i=1}^n g(z_i) > \epsilon\right) \leq \sum_{j=1}^N \mathbb{P}\left(\mathbb{E}g_j - \frac{1}{n} \sum_{i=1}^n g_j(z_i) > \epsilon\right)$$

Assuming $\mathbb{P}(a \leq g(z_i) \leq b) = 1$ for every $g \in \mathcal{G}$,

$$\mathbb{P}\left(\sup_{g \in \mathcal{G}} \left\{ \mathbb{E}g - \frac{1}{n} \sum_{i=1}^n g(z_i) \right\} > \epsilon\right) \leq N \exp\left(-\frac{2n\epsilon^2}{(b-a)^2}\right)$$

Finite Class

Alternatively, we set $\delta = N \exp\left(-\frac{2n\epsilon^2}{(b-a)^2}\right)$ and write

$$\mathbb{P}\left(\sup_{g \in \mathcal{G}} \left\{ \mathbb{E}g - \frac{1}{n} \sum_{i=1}^n g(z_i) \right\} > (b-a) \sqrt{\frac{\log(N) + \log(1/\delta)}{2n}}\right) \leq \delta$$

Another way to write it: with probability at least $1 - \delta$,

$$\sup_{g \in \mathcal{G}} \left\{ \mathbb{E}g - \frac{1}{n} \sum_{i=1}^n g(z_i) \right\} \leq (b-a) \sqrt{\frac{\log(N) + \log(1/\delta)}{2n}}$$

Hence, with probability at least $1 - \delta$, the ERM algorithm \hat{f}_n for a class \mathcal{F} of cardinality N satisfies

$$\mathbf{L}(\hat{f}_n) - \inf_{f \in \mathcal{F}} \mathbf{L}(f) \leq 2(b-a) \sqrt{\frac{\log(N) + \log(1/\delta)}{2n}}$$

assuming $a \leq \ell(f(x), y) \leq b$ for all $f \in \mathcal{F}$, $x \in \mathcal{X}$, $y \in \mathcal{Y}$.

The constant 2 is due to the $\mathbf{L}(f_{\mathcal{F}}) - \hat{\mathbf{L}}(f_{\mathcal{F}})$ term. This is a loose upper bound.

Once again...

A take-away message is that the following two statements are worlds apart:

with probability at least $1 - \delta$, for any $g \in \mathcal{G}$, $\mathbb{E}g - \frac{1}{n} \sum_{i=1}^n g(z_i) \leq \epsilon$

vs

for any $g \in \mathcal{G}$, with probability at least $1 - \delta$, $\mathbb{E}g - \frac{1}{n} \sum_{i=1}^n g(z_i) \leq \epsilon$

The second statement follows from CLT, while the first statement is often difficult to obtain and only holds for some \mathcal{G} .

Outline

Introduction

Statistical Learning Theory

The Setting of SLT

Consistency, No Free Lunch Theorems, Bias-Variance Tradeoff

Tools from Probability, Empirical Processes

From Finite to Infinite Classes

Uniform Convergence, Symmetrization, and Rademacher Complexity

Large Margin Theory for Classification

Properties of Rademacher Complexity

Covering Numbers and Scale-Sensitive Dimensions

Faster Rates

Model Selection

Sequential Prediction / Online Learning

Motivation

Supervised Learning

Online Convex and Linear Optimization

Online-to-Batch Conversion, SVM optimization

Countable Class: Weighted Union Bound

Let \mathcal{G} be countable and fix a distribution w on \mathcal{G} such that $\sum_{g \in \mathcal{G}} w(g) \leq 1$.
For any $\delta > 0$, for any $g \in \mathcal{G}$

$$\mathbb{P} \left(\mathbb{E}g - \frac{1}{n} \sum_{i=1}^n g(z_i) \geq (b - a) \sqrt{\frac{\log 1/w(g) + \log(1/\delta)}{2n}} \right) \leq \delta \cdot w(g)$$

by Hoeffding's inequality (easy to verify!). By the Union Bound,

$$\mathbb{P} \left(\exists g \in \mathcal{G} : \mathbb{E}g - \frac{1}{n} \sum_{i=1}^n g(z_i) \geq (b - a) \sqrt{\frac{\log 1/w(g) + \log(1/\delta)}{2n}} \right) \leq \delta \sum_{g \in \mathcal{G}} w(g) \leq \delta$$

Therefore, with probability at least $1 - \delta$, for all $f \in \mathcal{F}$

$$\mathbf{L}(f) - \hat{\mathbf{L}}(f) \leq \underbrace{(b - a) \sqrt{\frac{\log 1/w(f) + \log(1/\delta)}{2n}}}_{\text{pen}_n(f)}$$

Countable Class: Weighted Union Bound

If \hat{f}_n is a regularized ERM,

$$\begin{aligned}\mathbf{L}(\hat{f}_n) - \mathbf{L}(f_{\mathcal{F}}) &\leq \{\mathbf{L}(\hat{f}_n) - \hat{\mathbf{L}}(\hat{f}_n) - \text{pen}_n(\hat{f}_n)\} \\ &\quad + \{\hat{\mathbf{L}}(\hat{f}_n) + \text{pen}_n(\hat{f}_n) - \hat{\mathbf{L}}(f_{\mathcal{F}}) - \text{pen}_n(f_{\mathcal{F}})\} \\ &\quad + \{\hat{\mathbf{L}}(f_{\mathcal{F}}) - \mathbf{L}(f_{\mathcal{F}})\} + \text{pen}_n(f_{\mathcal{F}}) \\ &\leq \sup_{f \in \mathcal{F}} \{\mathbf{L}(f) - \hat{\mathbf{L}}(f) - \text{pen}_n(f)\} + \{\hat{\mathbf{L}}(f_{\mathcal{F}}) - \mathbf{L}(f_{\mathcal{F}})\} + \text{pen}_n(f_{\mathcal{F}})\end{aligned}$$

So, (E) implies a bound on (B) when \hat{f}_n is regularized ERM.
From the weighted union bound for a countable class:

$$\begin{aligned}\mathbf{L}(\hat{f}_n) - \mathbf{L}(f_{\mathcal{F}}) &\leq \{\hat{\mathbf{L}}(f_{\mathcal{F}}) - \mathbf{L}(f_{\mathcal{F}})\} + \text{pen}_n(f_{\mathcal{F}}) \\ &\leq 2(b-a) \sqrt{\frac{\log 1/w(f_{\mathcal{F}}) + \log(1/\delta)}{2n}}\end{aligned}$$

Uncountable Class: Compression Bounds

Let us make the dependence of the algorithm \hat{f}_n on the training set $S = \{(x_1, y_1), \dots, (x_n, y_n)\}$ explicit: $\hat{f}_n = \hat{f}_n[S]$.

Suppose \mathcal{F} has the property that there exists a “compression function” C_k which selects from any dataset S of any size n a subset of k labeled examples $C_k(S) \subseteq S$ such that the algorithm can be written as

$$\hat{f}_n[S] = \hat{f}_k[C_k(S)]$$

Then,

$$\begin{aligned} L(\hat{f}_n) - \hat{L}(\hat{f}_n) &= \mathbb{E} \ell(\hat{f}_k[C_k(S)](x), y) - \frac{1}{n} \sum_{i=1}^n \ell(\hat{f}_k[C_k(S)](x_i), y_i) \\ &\leq \max_{I \subseteq \{1, \dots, n\}, |I| \leq k} \left\{ \mathbb{E} \ell(\hat{f}_k[S_I](x), y) - \frac{1}{n} \sum_{i=1}^n \ell(\hat{f}_k[S_I](x_i), y_i) \right\} \end{aligned}$$

Uncountable Class: Compression Bounds

Since $\hat{f}_k[S_I]$ only depends on k out of n points, the empirical average is “mostly out of sample”. Adding and subtracting

$$\frac{1}{n} \sum_{(x', y') \in W} \ell(\hat{f}_k[S_I](x'), y')$$

for an additional set of i.i.d. random variables $W = \{(x'_1, y'_1), \dots, (x'_k, y'_k)\}$ results in an upper bound

$$\max_{I \subset \{1, \dots, n\}, |I| \leq k} \left\{ \mathbb{E} \ell(\hat{f}_k[S_I](x), y) - \frac{1}{n} \sum_{(x, y) \in S \setminus S_I \cup W_{|I|}} \ell(\hat{f}_k[S_I](x), y) \right\} + \frac{(b-a)k}{n}$$

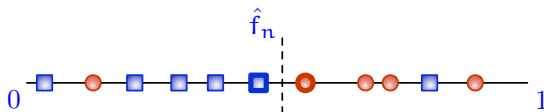
We appeal to the union bound over the $\binom{n}{k}$ possibilities, with a Hoeffding's bound for each. Then with probability at least $1 - \delta$,

$$\mathbf{L}(\hat{f}_n) - \inf_{f \in \mathcal{F}} \mathbf{L}(f) \leq 2(b-a) \sqrt{\frac{k \log(en/k) + \log(1/\delta)}{2n}} + \frac{(b-a)k}{n}$$

assuming $a \leq \ell(f(x), y) \leq b$ for all $f \in \mathcal{F}$, $x \in \mathcal{X}$, $y \in \mathcal{Y}$.

Example: Classification with Thresholds in 1D

- ▶ $\mathcal{X} = [0, 1]$, $\mathcal{Y} = \{0, 1\}$
- ▶ $\mathcal{F} = \{f_\theta : f_\theta(x) = \mathbf{I}_{\{x \geq \theta\}}, \theta \in [0, 1]\}$
- ▶ $\ell(f_\theta(x), y) = \mathbf{I}_{\{f_\theta(x) \neq y\}}$



For any set of data $(x_1, y_1), \dots, (x_n, y_n)$, the ERM solution \hat{f}_n has the property that the first occurrence x_l on the left of the threshold has label $y_l = 0$, while first occurrence x_r on the right – label $y_r = 1$.

Enough to take $k = 2$ and define $\hat{f}_n[S] = \hat{f}_2[(x_l, 0), (x_r, 1)]$.

Stability

Yet another way to limit the bias of $\ell(\hat{f}_n(x_i), y_i)$ as an estimate of $\mathbf{L}(\hat{f}_n)$ is through a notion of stability.

An algorithm \hat{f}_n is *stable* if a change (or removal) of a single data point does not change (in a certain mathematical sense) the function \hat{f}_n by much.

Of course, a dumb algorithm which outputs $\hat{f}_n = f_0$ without even looking at data is very stable and $\ell(\hat{f}_n(x_i), y_i)$ are independent random variables... But it is not a good algorithm! We would like to have an algorithm that both approximately minimizes the empirical error and is stable.

Turns out, certain types of regularization methods are stable. Example:

$$\hat{f}_n = \arg \min_{f \in \mathcal{F}} \frac{1}{n} \sum_{i=1}^n (f(x_i) - y_i)^2 + \lambda \|f\|_{\mathcal{K}}^2$$

where $\|\cdot\|$ is the norm induced by the kernel of a reproducing kernel Hilbert space (RKHS) \mathcal{F} .

Summary so far

We proved upper bounds on $\mathbf{L}(\hat{f}_n) - \mathbf{L}(f_{\mathcal{F}})$ for

- ▶ ERM over a finite class
- ▶ Regularized ERM over a countable class (weighted union bound)
- ▶ ERM over classes \mathcal{F} with the compression property
- ▶ ERM or Regularized ERM that are stable (only sketched it)

What about a more general situation? Is there a way to measure *complexity* of \mathcal{F} that tells us whether ERM will succeed?

Outline

Introduction

Statistical Learning Theory

The Setting of SLT

Consistency, No Free Lunch Theorems, Bias-Variance Tradeoff

Tools from Probability, Empirical Processes

From Finite to Infinite Classes

Uniform Convergence, Symmetrization, and Rademacher Complexity

Large Margin Theory for Classification

Properties of Rademacher Complexity

Covering Numbers and Scale-Sensitive Dimensions

Faster Rates

Model Selection

Sequential Prediction / Online Learning

Motivation

Supervised Learning

Online Convex and Linear Optimization

Online-to-Batch Conversion, SVM optimization

Uniform Convergence and Symmetrization

Let z'_1, \dots, z'_n be another set of n i.i.d. random variables from \mathcal{P} .

Let $\epsilon_1, \dots, \epsilon_n$ be i.i.d. Rademacher random variables:

$$\mathbb{P}(\epsilon_i = -1) = \mathbb{P}(\epsilon_i = +1) = 1/2$$

Let's get through a few manipulations:

$$\mathbb{E} \sup_{g \in \mathcal{G}} \left\{ \mathbb{E} g(z) - \frac{1}{n} \sum_{i=1}^n g(z_i) \right\} = \mathbb{E}_{z_{1:n}} \sup_{g \in \mathcal{G}} \left\{ \mathbb{E}_{z'_{1:n}} \left\{ \frac{1}{n} \sum_{i=1}^n g(z'_i) \right\} - \frac{1}{n} \sum_{i=1}^n g(z_i) \right\}$$

By Jensen's inequality, this is upper bounded by

$$\mathbb{E}_{z_{1:n}, z'_{1:n}} \sup_{g \in \mathcal{G}} \left\{ \frac{1}{n} \sum_{i=1}^n g(z'_i) - \frac{1}{n} \sum_{i=1}^n g(z_i) \right\}$$

which is equal to

$$\mathbb{E}_{\epsilon_{1:n}} \mathbb{E}_{z_{1:n}, z'_{1:n}} \sup_{g \in \mathcal{G}} \left\{ \frac{1}{n} \sum_{i=1}^n \epsilon_i (g(z'_i) - g(z_i)) \right\}$$

Uniform Convergence and Symmetrization

$$\begin{aligned} & \mathbb{E}_{\epsilon_{1:n}} \mathbb{E}_{z_{1:n}, z'_{1:n}} \sup_{g \in \mathcal{G}} \left\{ \frac{1}{n} \sum_{i=1}^n \epsilon_i (g(z'_i) - g(z_i)) \right\} \\ & \leq \mathbb{E} \sup_{g \in \mathcal{G}} \left\{ \frac{1}{n} \sum_{i=1}^n \epsilon_i g(z'_i) \right\} + \mathbb{E} \sup_{g \in \mathcal{G}} \left\{ \frac{1}{n} \sum_{i=1}^n -\epsilon_i g(z_i) \right\} \\ & = 2 \mathbb{E} \sup_{g \in \mathcal{G}} \left\{ \frac{1}{n} \sum_{i=1}^n \epsilon_i g(z_i) \right\} \end{aligned}$$

The *empirical Rademacher averages* of \mathcal{G} are defined as

$$\widehat{\mathcal{R}}_n(\mathcal{G}) = \mathbb{E} \left[\sup_{g \in \mathcal{G}} \left\{ \frac{1}{n} \sum_{i=1}^n \epsilon_i g(z_i) \right\} \mid z_1, \dots, z_n \right]$$

The *Rademacher average* (or *Rademacher complexity*) of \mathcal{G} is

$$\mathcal{R}_n(\mathcal{G}) = \mathbb{E}_{z_{1:n}} \widehat{\mathcal{R}}_n(\mathcal{G})$$

Classification: Loss Function Disappears

Let us focus on binary classification with indicator loss and let \mathcal{F} be a class of $\{0, 1\}$ -valued functions. We have

$$\ell(f(x), y) = \mathbf{I}_{\{f(x) \neq y\}} = (1 - 2y)f(x) + y$$

and thus

$$\begin{aligned}\widehat{\mathcal{R}}_n(\mathcal{G}) &= \mathbb{E} \left[\sup_{f \in \mathcal{F}} \left\{ \frac{1}{n} \sum_{i=1}^n \epsilon_i (f(x_i)(1 - 2y_i) + y_i) \right\} \middle| (x_1, y_1), \dots, (x_n, y_n) \right] \\ &= \mathbb{E} \left[\sup_{f \in \mathcal{F}} \left\{ \frac{1}{n} \sum_{i=1}^n \epsilon_i f(x_i) \right\} \middle| x_1, \dots, x_n \right] = \widehat{\mathcal{R}}_n(\mathcal{F})\end{aligned}$$

because, given y_1, \dots, y_n , the distribution of $\epsilon_i(1 - 2y_i)$ is the same as ϵ_i .

Vapnik-Chervonenkis Theory for Classification

We are now left examining

$$\mathbb{E} \left[\sup_{f \in \mathcal{F}} \left\{ \frac{1}{n} \sum_{i=1}^n \epsilon_i f(x_i) \right\} \mid x_1, \dots, x_n \right]$$

Given x_1, \dots, x_n , define the projection of \mathcal{F} onto sample:

$$\mathcal{F}|_{x_{1:n}} = \{(f(x_1), \dots, f(x_n)) \in \{0, 1\}^n : f \in \mathcal{F}\} \subseteq \{0, 1\}^n$$

Clearly, this is a finite set and

$$\widehat{\mathcal{R}}_n(\mathcal{F}) = \mathbb{E}_{\epsilon_{1:n}} \max_{v \in \mathcal{F}|_{x_{1:n}}} \frac{1}{n} \sum_{i=1}^n \epsilon_i v_i \leq \sqrt{\frac{2 \log \text{card}(\mathcal{F}|_{x_{1:n}})}{n}}$$

This is because a maximum of N (sub)Gaussian random variables $\sim \sqrt{\log N}$.

The bound is nontrivial as long as $\log \text{card}(\mathcal{F}|_{x_{1:n}}) = o(n)$.

Vapnik-Chervonenkis Theory for Classification

The *growth function* is defined as

$$\Pi_{\mathcal{F}}(\mathbf{n}) = \max \left\{ \text{card}(\mathcal{F}|_{x_1, \dots, x_n}) : x_1, \dots, x_n \in \mathcal{X} \right\}$$

The growth function measures *expressiveness* of \mathcal{F} . In particular, if \mathcal{F} can produce all possible signs (that is, $\Pi_{\mathcal{F}}(\mathbf{n}) = 2^n$), the bound becomes useless.

We say that \mathcal{F} *shatters* some set x_1, \dots, x_n if $\mathcal{F}|_{x^n} = \{0, 1\}^n$.

The *Vapnik-Chervonenkis (VC) dimension* of the class \mathcal{F} is defined as

$$\text{vc}(\mathcal{F}) = \max \left\{ d : \Pi_{\mathcal{F}}(t) = 2^t \right\}$$

Vapnik-Chervonenkis-Sauer-Shelah Lemma: If $d = \text{vc}(\mathcal{F}) < \infty$, then

$$\Pi_{\mathcal{F}}(\mathbf{n}) \leq \sum_{i=0}^d \binom{n}{i} \leq \left(\frac{en}{d} \right)^d$$

Vapnik-Chervonenkis Theory for Classification

Conclusion: for any \mathcal{F} with $\text{vc}(\mathcal{F}) < \infty$, the ERM algorithm satisfies

$$\mathbb{E} \left\{ \mathbf{L}(\hat{f}_n) - \inf_{f \in \mathcal{F}} \mathbf{L}(f) \right\} \leq 2 \sqrt{\frac{2d \log(en/d)}{n}}$$

While we proved the result in expectation, the same type of bound holds with high probability.

VC dimension is a combinatorial dimension of a binary-valued function class. Its finiteness is necessary and sufficient for learnability if we place no assumptions on the distribution \mathbb{P} .

Remark: the bound is similar to that obtained through compression. In fact, the exact relationship between compression and VC dimension is still an open question.

Vapnik-Chervonenkis Theory for Classification

Examples of VC classes:

- ▶ Half-spaces $\mathcal{F} = \{\mathbf{I}_{\{(w,x)+b \geq 0\}} : w \in \mathbb{R}^d, \|w\| = 1, b \in \mathbb{R}\}$ has $\text{vc}(\mathcal{F}) = d + 1$
- ▶ For a vector space \mathcal{H} of dimension d , VC dimension of $\mathcal{F} = \{\mathbf{I}_{\{h(x) \geq 0\}} : h \in \mathcal{H}\}$ is at most d
- ▶ The set of Euclidean balls $\mathcal{F} = \left\{ \mathbf{I}_{\{\sum_{i=1}^d \|x_i - a_i\|^2 \leq b\}} : a \in \mathbb{R}^d, b \in \mathbb{R} \right\}$ has VC dimension at most $d + 2$.
- ▶ Functions that can be computed using a finite number of arithmetic operations (see (*Goldberg and Jerrum, 1995*))

However: $\mathcal{F} = \{f_\alpha(x) = \mathbf{I}_{\{\sin(\alpha x) \geq 0\}} : \alpha \in \mathbb{R}\}$ has infinite VC dimension, so it is not correct to think of VC dimension as the number of parameters!

Vapnik-Chervonenkis Theory for Classification

Examples of VC classes:

- ▶ Half-spaces $\mathcal{F} = \{\mathbf{I}_{\{(w,x)+b \geq 0\}} : w \in \mathbb{R}^d, \|w\| = 1, b \in \mathbb{R}\}$ has $\text{vc}(\mathcal{F}) = d + 1$
- ▶ For a vector space \mathcal{H} of dimension d , VC dimension of $\mathcal{F} = \{\mathbf{I}_{\{h(x) \geq 0\}} : h \in \mathcal{H}\}$ is at most d
- ▶ The set of Euclidean balls $\mathcal{F} = \left\{ \mathbf{I}_{\{\sum_{i=1}^d \|x_i - a_i\|^2 \leq b\}} : a \in \mathbb{R}^d, b \in \mathbb{R} \right\}$ has VC dimension at most $d + 2$.
- ▶ Functions that can be computed using a finite number of arithmetic operations (see (*Goldberg and Jerrum, 1995*))

However: $\mathcal{F} = \{f_\alpha(x) = \mathbf{I}_{\{\sin(\alpha x) \geq 0\}} : \alpha \in \mathbb{R}\}$ has infinite VC dimension, so it is not correct to think of VC dimension as the number of parameters!

Unfortunately, the VC theory is unable to explain the good performance of neural networks and Support Vector Machines! This prompted the development of a margin-based theory.

Outline

Introduction

Statistical Learning Theory

The Setting of SLT

Consistency, No Free Lunch Theorems, Bias-Variance Tradeoff

Tools from Probability, Empirical Processes

From Finite to Infinite Classes

Uniform Convergence, Symmetrization, and Rademacher Complexity

Large Margin Theory for Classification

Properties of Rademacher Complexity

Covering Numbers and Scale-Sensitive Dimensions

Faster Rates

Model Selection

Sequential Prediction / Online Learning

Motivation

Supervised Learning

Online Convex and Linear Optimization

Online-to-Batch Conversion, SVM optimization

Classification with Real-Valued Functions

Many methods use

$$\mathbb{I}(\mathcal{F}) = \{\mathbf{I}_{\{f \geq 0\}} : f \in \mathcal{F}\}$$

for classification. The VC dimension can be very large, yet in practice the methods work well.

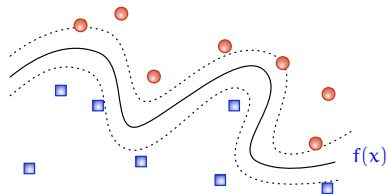
Example: $f(\mathbf{x}) = f_{\mathbf{w}}(\mathbf{x}) = \langle \mathbf{w}, \psi(\mathbf{x}) \rangle$ where ψ is a mapping to a high-dimensional feature space (see Kernel Methods). The VC dimension of the set is typically huge (equal to the dimensionality of $\psi(\mathbf{x})$) or infinite, yet the methods perform well!

Is there an explanation beyond VC theory?

Margins

Hard margin:

$$\exists f \in \mathcal{F} : \forall i, \quad y_i f(x_i) \geq \gamma$$



More generally, we hope to have

$$\exists f \in \mathcal{F} : \frac{\text{card}(\{i : y_i f(x_i) < \gamma\})}{n} \text{ is small}$$

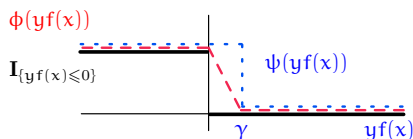
Surrogate Loss

Define

$$\phi(s) = \begin{cases} 1 & \text{if } s \leq 0 \\ 1 - s/\gamma & \text{if } 0 < s < \gamma \\ 0 & \text{if } s \geq \gamma \end{cases}$$

Then: $\mathbf{I}_{\{y \neq \text{sign}(f(x))\}} = \mathbf{I}_{\{yf(x) \leq 0\}} \leq \phi(yf(x)) \leq \psi(yf(x)) = \mathbf{I}_{\{yf(x) \leq \gamma\}}$

The function ϕ is an example of a *surrogate loss function*.



Let

$$\mathbf{L}_\phi(f) = \mathbb{E}\phi(yf(x)) \quad \text{and} \quad \hat{\mathbf{L}}_\phi(f) = \frac{1}{n} \sum_{i=1}^n \phi(y_i f(x_i))$$

Then

$$\mathbf{L}(f) \leq \mathbf{L}_\phi(f), \quad \hat{\mathbf{L}}_\phi(f) \leq \hat{\mathbf{L}}_\psi(f)$$

Surrogate Loss

Now consider uniform deviations for the surrogate loss:

$$\mathbb{E} \sup_{f \in \mathcal{F}} \{ \mathbf{L}_\phi(f) - \hat{\mathbf{L}}_\phi(f) \}$$

We had shown that this quantity is at most $2\mathcal{R}_n(\phi(\mathcal{F}))$ for

$$\phi(\mathcal{F}) = \{g(z) = \phi(yf(x)) : f \in \mathcal{F}\}$$

A useful property of Rademacher averages:

$$\mathcal{R}_n(\phi(\mathcal{F})) \leq L\mathcal{R}_n(\mathcal{F}) \quad \text{if } \phi \text{ is } L\text{-Lipschitz.}$$

Observe that in our example ϕ is $1/\gamma$ -Lipschitz. Hence,

$$\mathbb{E} \sup_{f \in \mathcal{F}} \{ \mathbf{L}_\phi(f) - \hat{\mathbf{L}}_\phi(f) \} \leq \frac{2}{\gamma} \mathcal{R}_n(\mathcal{F})$$

Margin Bound

Same result in high probability: with probability at least $1 - \delta$,

$$\sup_{f \in \mathcal{F}} \{L_\phi(f) - \hat{L}_\phi(f)\} \leq \frac{2}{\gamma} \mathcal{R}_n(\mathcal{F}) + \sqrt{\frac{\log(1/\delta)}{2n}}$$

With probability at least $1 - \delta$, for all $f \in \mathcal{F}$

$$L(f) \leq \hat{L}_\psi(f) + \frac{2}{\gamma} \mathcal{R}_n(\mathcal{F}) + \sqrt{\frac{\log(1/\delta)}{2n}}$$

If \hat{f}_n is minimizing margin loss

$$\hat{f}_n = \arg \min_{f \in \mathcal{F}} \frac{1}{n} \sum_{i=1}^n \phi(y_i f(x_i))$$

then with probability at least $1 - \delta$

$$L(\hat{f}_n) \leq \inf_{f \in \mathcal{F}} L_\psi(f) + \frac{4}{\gamma} \mathcal{R}_n(\mathcal{F}) + 2\sqrt{\frac{\log(1/\delta)}{2n}}$$

Note: ϕ assumes knowledge of γ , but this assumption can be removed.

Outline

Introduction

Statistical Learning Theory

The Setting of SLT

Consistency, No Free Lunch Theorems, Bias-Variance Tradeoff

Tools from Probability, Empirical Processes

From Finite to Infinite Classes

Uniform Convergence, Symmetrization, and Rademacher Complexity

Large Margin Theory for Classification

Properties of Rademacher Complexity

Covering Numbers and Scale-Sensitive Dimensions

Faster Rates

Model Selection

Sequential Prediction / Online Learning

Motivation

Supervised Learning

Online Convex and Linear Optimization

Online-to-Batch Conversion, SVM optimization

Useful Properties

1. If $\mathcal{F} \subseteq \mathcal{G}$, then $\widehat{\mathcal{R}}_n(\mathcal{F}) \leq \widehat{\mathcal{R}}_n(\mathcal{G})$
2. $\widehat{\mathcal{R}}_n(\mathcal{F}) = \widehat{\mathcal{R}}_n(\text{conv}(\mathcal{F}))$
3. For any $c \in \mathbb{R}$, $\widehat{\mathcal{R}}_n(c\mathcal{F}) = |c|\widehat{\mathcal{R}}_n(\mathcal{F})$
4. If $\phi : \mathbb{R} \mapsto \mathbb{R}$ is L -Lipschitz (that is, $\phi(\mathbf{a}) - \phi(\mathbf{b}) \leq L|\mathbf{a} - \mathbf{b}|$ for all $\mathbf{a}, \mathbf{b} \in \mathbb{R}$), then

$$\widehat{\mathcal{R}}_n(\phi \circ \mathcal{F}) \leq L\widehat{\mathcal{R}}_n(\mathcal{F})$$

Rademacher Complexity of Kernel Classes

- ▶ Feature map $\phi : \mathcal{X} \mapsto \ell_2$ and p.d. kernel $K(x_1, x_2) = \langle \phi(x_1), \phi(x_2) \rangle$
- ▶ The set $\mathcal{F}_B = \{f(x) = \langle w, \phi(x) \rangle : \|w\| \leq B\}$ is a ball in \mathcal{H}
- ▶ Reproducing property $f(x) = \langle f, K(x, \cdot) \rangle$

An easy calculation shows that empirical Rademacher averages are upper bounded as

$$\begin{aligned}\widehat{\mathcal{R}}_n(\mathcal{F}_B) &= \mathbb{E} \sup_{f \in \mathcal{F}_B} \frac{1}{n} \sum_{i=1}^n \epsilon_i f(x_i) = \mathbb{E} \sup_{f \in \mathcal{F}_B} \frac{1}{n} \sum_{i=1}^n \epsilon_i \langle f, K(x_i, \cdot) \rangle \\ &= \mathbb{E} \sup_{f \in \mathcal{F}_B} \left\langle f, \frac{1}{n} \sum_{i=1}^n \epsilon_i K(x_i, \cdot) \right\rangle = B \cdot \mathbb{E} \left\| \frac{1}{n} \sum_{i=1}^n \epsilon_i K(x_i, \cdot) \right\| \\ &= \frac{B}{n} \mathbb{E} \left(\sum_{i,j=1}^n \epsilon_i \epsilon_j \langle K(x_i, \cdot), K(x_j, \cdot) \rangle \right)^{-1/2} \\ &\leq \frac{B}{n} \left(\sum_{i=1}^n K(x_i, x_i) \right)^{-1/2}\end{aligned}$$

A data-independent bound of $O(B\kappa/\sqrt{n})$ can be obtained if $\sup_{x \in \mathcal{X}} K(x, x) \leq \kappa^2$. Then κ and B are the effective “dimensions”.

Other Examples

Using properties of Rademacher averages, we may establish guarantees for learning with neural networks, decision trees, and so on.

Powerful technique, typically requires only a few lines of algebra.

Occasionally, covering numbers and scale-sensitive dimensions can be easier to deal with.

Outline

Introduction

Statistical Learning Theory

The Setting of SLT

Consistency, No Free Lunch Theorems, Bias-Variance Tradeoff

Tools from Probability, Empirical Processes

From Finite to Infinite Classes

Uniform Convergence, Symmetrization, and Rademacher Complexity

Large Margin Theory for Classification

Properties of Rademacher Complexity

Covering Numbers and Scale-Sensitive Dimensions

Faster Rates

Model Selection

Sequential Prediction / Online Learning

Motivation

Supervised Learning

Online Convex and Linear Optimization

Online-to-Batch Conversion, SVM optimization

Real-Valued Functions: Covering Numbers

Consider

- ▶ a class \mathcal{F} of $[-1, 1]$ -valued functions
- ▶ let $\mathcal{Y} = [-1, 1]$, $\ell(f(x), y) = |f(x) - y|$

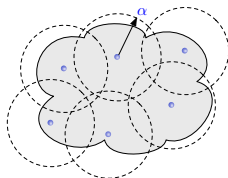
We have

$$\mathbb{E} \sup_{f \in \mathcal{F}} \mathbf{L}(f) - \hat{\mathbf{L}}(f) \leq 2 \mathbb{E}_{x_{1:n}} \widehat{\mathcal{R}}_n(\mathcal{F})$$

For real-valued functions the cardinality of $\mathcal{F}|_{x_{1:n}}$ is infinite. However, similar functions f and f' with

$$(f(x_1), \dots, f(x_n)) \approx (f'(x_1), \dots, f'(x_n))$$

should be treated as the same.



Real-Valued Functions: Covering Numbers

Given $\alpha > 0$, suppose we can find $V \subset [-1, 1]^n$ of finite cardinality such that

$$\forall f, \exists v^f \in V, \text{ s.t. } \frac{1}{n} \sum_{i=1}^n |f(x_i) - v_i^f| \leq \alpha$$

Then

$$\begin{aligned} \widehat{\mathcal{R}}_n(\mathcal{F}) &= \mathbb{E}_{\epsilon_{1:n}} \sup_{f \in \mathcal{F}} \frac{1}{n} \sum_{i=1}^n \epsilon_i f(x_i) \\ &= \mathbb{E}_{\epsilon_{1:n}} \sup_{f \in \mathcal{F}} \frac{1}{n} \sum_{i=1}^n \epsilon_i (f(x_i) - v_i^f) + \mathbb{E}_{\epsilon_{1:n}} \sup_{f \in \mathcal{F}} \frac{1}{n} \sum_{i=1}^n \epsilon_i v_i^f \\ &\leq \alpha + \mathbb{E}_{\epsilon_{1:n}} \max_{v \in V} \frac{1}{n} \sum_{i=1}^n \epsilon_i v_i \end{aligned}$$

Now we are back to the set of finite cardinality:

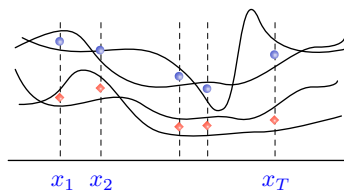
$$\widehat{\mathcal{R}}_n(\mathcal{F}) \leq \alpha + \sqrt{\frac{2 \log \text{card}(V)}{n}}$$

Real-Valued Functions: Covering Numbers

Such a set V is called an α -cover (or α -net). More precisely, a set V is an α -cover with respect to ℓ_p norm if

$$\forall f, \exists v^f \in V, \text{ s.t. } \frac{1}{n} \sum_{i=1}^n |f(x_i) - v_i^f|^p \leq \alpha^p$$

The size of the smallest α -cover is denoted by $\mathcal{N}_p(\mathcal{F}|_{x_{1:n}}, \alpha)$.



Above : Two sets of levels provide an α -cover for the four functions. Only the values of functions on x_1, \dots, x_T are relevant.

Real-Valued Functions: Covering Numbers

We have proved that for any x_1, \dots, x_n ,

$$\widehat{\mathcal{R}}_n(\mathcal{F}) \leq \inf_{\alpha \geq 0} \left\{ \alpha + \frac{1}{\sqrt{n}} \sqrt{2 \log \text{card}(\mathcal{N}_1(\mathcal{F}|_{x_{1:n}}, \alpha))} \right\}$$

A better bound (called Dudley entropy integral):

$$\widehat{\mathcal{R}}_n(\mathcal{F}) \leq \inf_{\alpha \geq 0} \left\{ 4\alpha + \frac{12}{\sqrt{n}} \int_{\alpha}^1 \sqrt{2 \log \text{card}(\mathcal{N}_2(\mathcal{F}|_{x_{1:n}}, \delta))} d\delta \right\}$$

Example: Nondecreasing functions.

Consider the set \mathcal{F} of nondecreasing functions $\mathbb{R} \mapsto [-1, 1]$.

While \mathcal{F} is a very large set, $\mathcal{F}|_{x_{1:n}}$ is not that large:

$$\mathcal{N}_1(\mathcal{F}|_{x_{1:n}}, \alpha) \leq \mathcal{N}_2(\mathcal{F}|_{x_{1:n}}, \alpha) \leq n^{2/\alpha}.$$

The first bound on the previous slide yields

$$\inf_{\alpha \geq 0} \left\{ \alpha + \frac{1}{\sqrt{\alpha n}} \sqrt{4 \log(n)} \right\} = \tilde{O}(n^{-1/3})$$

while the second bound (the Dudley entropy integral)

$$\inf_{\alpha \geq 0} \left\{ 4\alpha + \frac{12}{\sqrt{n}} \int_{\alpha}^1 \sqrt{4/\delta \log(n)} d\delta \right\} = \tilde{O}(n^{-1/2})$$

where the \tilde{O} notation hides logarithmic factors.

Scale-Sensitive Dimensions

We say that $\mathcal{F} \subseteq \mathbb{R}^{\mathcal{X}}$ α -shatters a set (x_1, \dots, x_T) if there exist $(y_1, \dots, y_T) \in \mathbb{R}^T$ (called a *witness to shattering*) with the following property:

$$\forall (b_1, \dots, b_T) \in \{0, 1\}^T, \exists f \in \mathcal{F} \text{ s.t.}$$

$$f(x_t) > y_t + \frac{\alpha}{2} \text{ if } b_t = 1 \quad \text{and} \quad f(x_t) < y_t - \frac{\alpha}{2} \text{ if } b_t = 0$$

The *fat-shattering dimension* of \mathcal{F} at scale α , denoted by $\text{fat}(\mathcal{F}, \alpha)$, is the size of the largest α -shattered set.

Scale-Sensitive Dimensions

We say that $\mathcal{F} \subseteq \mathbb{R}^{\mathcal{X}}$ α -shatters a set (x_1, \dots, x_T) if there exist $(y_1, \dots, y_T) \in \mathbb{R}^T$ (called a *witness to shattering*) with the following property:

$$\forall (b_1, \dots, b_T) \in \{0, 1\}^T, \exists f \in \mathcal{F} \text{ s.t.}$$

$$f(x_t) > y_t + \frac{\alpha}{2} \text{ if } b_t = 1 \quad \text{and} \quad f(x_t) < y_t - \frac{\alpha}{2} \text{ if } b_t = 0$$

The *fat-shattering dimension* of \mathcal{F} at scale α , denoted by $\text{fat}(\mathcal{F}, \alpha)$, is the size of the largest α -shattered set.

Wait, another measure of complexity of \mathcal{F} ? How is it related to covering numbers?

Scale-Sensitive Dimensions

We say that $\mathcal{F} \subseteq \mathbb{R}^{\mathcal{X}}$ α -shatters a set (x_1, \dots, x_T) if there exist $(y_1, \dots, y_T) \in \mathbb{R}^T$ (called a *witness to shattering*) with the following property:

$$\forall (b_1, \dots, b_T) \in \{0, 1\}^T, \exists f \in \mathcal{F} \quad \text{s.t.}$$

$$f(x_t) > y_t + \frac{\alpha}{2} \quad \text{if } b_t = 1 \quad \text{and} \quad f(x_t) < y_t - \frac{\alpha}{2} \quad \text{if } b_t = 0$$

The *fat-shattering dimension* of \mathcal{F} at scale α , denoted by $\text{fat}(\mathcal{F}, \alpha)$, is the size of the largest α -shattered set.

Wait, another measure of complexity of \mathcal{F} ? How is it related to covering numbers?

Theorem (Mendelson & Vershynin): For $\mathcal{F} \subseteq [-1, 1]^{\mathcal{X}}$ and any $0 < \alpha < 1$,

$$\mathcal{N}_2(\mathcal{F}|_{x_{1:n}}, \alpha) \leq \left(\frac{2}{\alpha}\right)^{K \cdot \text{fat}(\mathcal{F}, c\alpha)}$$

where K, c are positive absolute constants.

Quick Summary

We are after uniform deviations in order to understand performance of ERM. Rademacher averages is a nice measure with useful properties. They can be further upper bounded by covering numbers through the Dudley entropy integral. In turn, covering numbers can be controlled via the fat-shattering combinatorial dimension. Whew!

Outline

Introduction

Statistical Learning Theory

The Setting of SLT

Consistency, No Free Lunch Theorems, Bias-Variance Tradeoff

Tools from Probability, Empirical Processes

From Finite to Infinite Classes

Uniform Convergence, Symmetrization, and Rademacher Complexity

Large Margin Theory for Classification

Properties of Rademacher Complexity

Covering Numbers and Scale-Sensitive Dimensions

Faster Rates

Model Selection

Sequential Prediction / Online Learning

Motivation

Supervised Learning

Online Convex and Linear Optimization

Online-to-Batch Conversion, SVM optimization

Faster Rates

Are there situations when

$$\mathbb{E}L(\hat{f}_n) - \inf_{f \in \mathcal{F}} L(f)$$

approaches 0 faster than $O(1/\sqrt{n})$?

Yes! We can beat the Central Limit Theorem!

How is this possible??

Recall that the CLT tells us about convergence of average to the expectation for random variables with bounded second moment. What if this variance is small?

Faster Rates: Classification

Consider the problem of binary classification with the indicator loss and a class \mathcal{F} of $\{0, 1\}$ -valued functions. For any $f \in \mathcal{F}$,

$$\frac{1}{n} \sum_{i=1}^n \ell(f(x_i), y_i)$$

is an average of n Bernoulli random variables with bias $p = \mathbb{E}\ell(f(x), y)$.
Exact expression for the binomial tails:

$$\mathbb{P}(\mathbf{L}(f) - \hat{\mathbf{L}}(f) > \epsilon) = \sum_{i=0}^{\lfloor n(p-\epsilon) \rfloor} \binom{n}{i} p^i (1-p)^{n-i}$$

Further upper bounds:

$$\exp\left\{-\frac{n\epsilon^2}{2p(1-p) + 2\epsilon/3}\right\} \quad \text{Bernstein}$$

$$\exp\{-2n\epsilon^2\} \quad \text{Hoeffding}$$

Faster Rates: Classification

Inverting

$$\exp\left\{-\frac{n\epsilon^2}{2p(1-p) + 2\epsilon/3}\right\} \leq \exp\left\{-\frac{n\epsilon^2}{2p + 2\epsilon/3}\right\} =: \delta$$

yields that for any $f \in \mathcal{F}$, with probability at least $1 - \delta$

$$\mathbf{L}(f) \leq \hat{\mathbf{L}}(f) + \sqrt{\frac{2\mathbf{L}(f) \log(1/\delta)}{n}} + \frac{2\log(1/\delta)}{3n}$$

For non-negative numbers A, B, C

$$A \leq B + C\sqrt{A} \quad \text{implies} \quad A \leq B + C^2 + \sqrt{BC}$$

Therefore for any $f \in \mathcal{F}$, with probability at least $1 - \delta$,

$$\mathbf{L}(f) \leq \hat{\mathbf{L}}(f) + \sqrt{\frac{2\hat{\mathbf{L}}(f) \log(1/\delta)}{n}} + \frac{4\log(1/\delta)}{n}$$

Faster Rates: Classification

By the Union Bound, for \mathcal{F} with finite $N = \text{card}(\mathcal{F})$, with probability at least $1 - \delta$,

$$\forall f \in \mathcal{F}: \quad \mathbf{L}(f) \leq \hat{\mathbf{L}}(f) + \sqrt{\frac{2\hat{\mathbf{L}}(f) \log(N/\delta)}{n}} + \frac{4\log(N/\delta)}{n}$$

For an empirical minimizer \hat{f}_n , with probability at least $1 - \delta$, a zero empirical loss $\hat{\mathbf{L}}(\hat{f}_n) = 0$ implies

$$\mathbf{L}(\hat{f}_n) \leq \frac{4\log(N/\delta)}{n}$$

This happens, for instance, in the so-called *noiseless case*: $\mathbf{L}(f_{\mathcal{F}}) = 0$. Indeed, then $\hat{\mathbf{L}}(f_{\mathcal{F}}) = 0$ and thus $\hat{\mathbf{L}}(\hat{f}_n) = 0$.

Summary: Minimax Viewpoint

Value of a game where we choose an algorithm, Nature chooses a distribution $\mathbf{P} \in \mathcal{P}$, and our payoff is the expected loss of our algorithm relative to the best in \mathcal{F} :

$$\mathcal{V}^{\text{iid}}(\mathcal{F}, \mathcal{P}, n) = \inf_{\hat{f}_n} \sup_{\mathbf{P} \in \mathcal{P}} \left\{ \mathbf{L}(\hat{f}_n) - \inf_{f \in \mathcal{F}} \mathbf{L}(f) \right\}$$

If we make no assumption on the distribution \mathbf{P} , then \mathcal{P} is the set of all distributions. Many of the results we obtained in this lecture are for this distribution-free case. However, one may view margin-based results and the above fast rates for the noiseless case as studying $\mathcal{V}^{\text{iid}}(\mathcal{F}, \mathcal{P}, n)$ when \mathcal{P} is “nicer”.

Outline

Introduction

Statistical Learning Theory

The Setting of SLT

Consistency, No Free Lunch Theorems, Bias-Variance Tradeoff

Tools from Probability, Empirical Processes

From Finite to Infinite Classes

Uniform Convergence, Symmetrization, and Rademacher Complexity

Large Margin Theory for Classification

Properties of Rademacher Complexity

Covering Numbers and Scale-Sensitive Dimensions

Faster Rates

Model Selection

Sequential Prediction / Online Learning

Motivation

Supervised Learning

Online Convex and Linear Optimization

Online-to-Batch Conversion, SVM optimization

Model Selection

For a given class \mathcal{F} , we have proved statements of the type

$$\mathbb{P} \left(\sup_{f \in \mathcal{F}} \{ \mathbf{L}(f) - \hat{\mathbf{L}}(f) \} \geq \phi(\delta, n, \mathcal{F}) \right) < \delta$$

Now, take a countable nested sieve of models

$$\mathcal{F}_1 \subseteq \mathcal{F}_2 \subseteq \dots$$

such that $\mathcal{H} = \cup_{i=1}^{\infty} \mathcal{F}_i$ is a very large set that will surely capture the Bayes function.

For a function $f \in \mathcal{H}$, let $k(f)$ be the smallest index of \mathcal{F}_k that contains f . Let us write $\phi_n(\delta, i)$ for $\phi(\delta, n, \mathcal{F}_i)$.

Let us put a distribution $w(i)$ on the models, with $\sum_{i=1}^{\infty} w(i) = 1$. Then for every i ,

$$\mathbb{P} \left(\sup_{f \in \mathcal{F}_i} \{ \mathbf{L}(f) - \hat{\mathbf{L}}(f) \} \geq \phi_n(\delta w(i), i) \right) < \delta \cdot w(i)$$

simply by replacing δ with $\delta w(i)$.

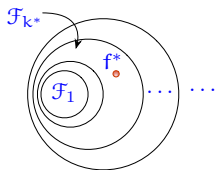
Now, taking a union bound:

$$\mathbb{P} \left(\sup_{f \in \mathcal{H}} \{ \mathbf{L}(f) - \hat{\mathbf{L}}(f) \} \geq \phi_n(\delta w(k(f)), k(f)) \right) < \sum_i \delta w(i) \leq \delta$$

Consider the penalized method

$$\begin{aligned} \hat{f}_n &= \arg \min_{f \in \mathcal{H}} \{ \hat{\mathbf{L}}(f) + \phi_n(\delta w(k(f)), k(f)) \} \\ &= \arg \min_{i, f \in \mathcal{F}_i} \{ \hat{\mathbf{L}}(f) + \phi_n(\delta w(i), i) \} \end{aligned}$$

This balances fit to data and the complexity of the model. Of course, this is exactly a regularized ERM form analyzed earlier.



Let $k^* = k(f^*)$ be the (smallest) model \mathcal{F}_i that contains the optimal function.

Exactly as on the slide “Countable Class: Weighted Union Bound”,

$$\begin{aligned}\mathbf{L}(\hat{f}_n) - \mathbf{L}(f^*) &\leq \{\mathbf{L}(\hat{f}_n) - \hat{\mathbf{L}}(\hat{f}_n) - \text{pen}_n(\hat{f}_n)\} \\ &\quad + \{\hat{\mathbf{L}}(\hat{f}_n) + \text{pen}_n(\hat{f}_n) - \hat{\mathbf{L}}(f_{\mathcal{F}}) - \text{pen}_n(f_{\mathcal{F}})\} \\ &\quad + \{\hat{\mathbf{L}}(f_{\mathcal{F}}) - \mathbf{L}(f_{\mathcal{F}})\} + \text{pen}_n(f_{\mathcal{F}}) \\ &\leq \hat{\mathbf{L}}(f^*) - \mathbf{L}(f^*) + \text{pen}_n(f^*) \\ &= \hat{\mathbf{L}}(f^*) - \mathbf{L}(f^*) + \phi_n(\delta w(k^*), k^*)\end{aligned}$$

The first part of this bound is $O_P(1/\sqrt{n})$ by the CLT, just as before.

If the dependence of ϕ on $1/\delta$ is logarithmic, then taking $w(i) = 2^{-i}$ simply implies an additional additive i^* , a penalty for not knowing the model in advance.

Conclusion: given uniform deviation bounds for a single class \mathcal{F} , as developed earlier, we can perform model selection by penalizing model complexity!

Outline

Introduction

Statistical Learning Theory

The Setting of SLT

Consistency, No Free Lunch Theorems, Bias-Variance Tradeoff

Tools from Probability, Empirical Processes

From Finite to Infinite Classes

Uniform Convergence, Symmetrization, and Rademacher Complexity

Large Margin Theory for Classification

Properties of Rademacher Complexity

Covering Numbers and Scale-Sensitive Dimensions

Faster Rates

Model Selection

Sequential Prediction / Online Learning

Motivation

Supervised Learning

Online Convex and Linear Optimization

Online-to-Batch Conversion, SVM optimization

Outline

Introduction

Statistical Learning Theory

The Setting of SLT

Consistency, No Free Lunch Theorems, Bias-Variance Tradeoff

Tools from Probability, Empirical Processes

From Finite to Infinite Classes

Uniform Convergence, Symmetrization, and Rademacher Complexity

Large Margin Theory for Classification

Properties of Rademacher Complexity

Covering Numbers and Scale-Sensitive Dimensions

Faster Rates

Model Selection

Sequential Prediction / Online Learning

Motivation

Supervised Learning

Online Convex and Linear Optimization

Online-to-Batch Conversion, SVM optimization

Looking back: Statistical Learning

- ▶ future looks like the past
- ▶ modeled as i.i.d. data
- ▶ evaluated on a random sample from the same distribution
- ▶ developed various measures of complexity of \mathcal{F}

Example #1: Bit Prediction

Predict a binary sequence $y_1, y_2, \dots \in \{0, 1\}$, which is revealed one by one. At step t , make a prediction z_t of the t -th bit, then y_t is revealed.

Let $c_t = \mathbf{I}_{\{z_t=y_t\}}$. Goal: make $\bar{c}_n = \frac{1}{n} \sum_{t=1}^n c_t$ large.

Suppose we are told that the sequence presented is Bernoulli with an unknown bias p . How should we choose predictions?

Example #1: Bit Prediction

Of course, we should do majority vote over the past outcomes

$$z_t = \mathbf{I}_{\{\bar{y}_{t-1} \geq 1/2\}}$$

where $\bar{y}_{t-1} = \frac{1}{t-1} \sum_{s=1}^{t-1} y_s$. This algorithm guarantees $\bar{c}_t \rightarrow \max\{p, 1-p\}$ and

$$\liminf_{t \rightarrow \infty} (\bar{c}_t - \max\{\bar{z}_t, 1 - \bar{z}_t\}) \geq 0 \quad \text{almost surely} \quad (*)$$

Example #1: Bit Prediction

Of course, we should do majority vote over the past outcomes

$$z_t = \mathbf{I}_{\{\bar{y}_{t-1} \geq 1/2\}}$$

where $\bar{y}_{t-1} = \frac{1}{t-1} \sum_{s=1}^{t-1} y_s$. This algorithm guarantees $\bar{c}_t \rightarrow \max\{p, 1-p\}$ and

$$\liminf_{t \rightarrow \infty} (\bar{c}_t - \max\{\bar{z}_t, 1 - \bar{z}_t\}) \geq 0 \quad \text{almost surely} \quad (*)$$

Claim: there is an algorithm that ensures $(*)$ for an **arbitrary sequence**.
Any idea how to do it?

Example #1: Bit Prediction

Of course, we should do majority vote over the past outcomes

$$z_t = \mathbf{I}_{\{\bar{y}_{t-1} \geq 1/2\}}$$

where $\bar{y}_{t-1} = \frac{1}{t-1} \sum_{s=1}^{t-1} y_s$. This algorithm guarantees $\bar{c}_t \rightarrow \max\{p, 1-p\}$ and

$$\liminf_{t \rightarrow \infty} (\bar{c}_t - \max\{\bar{z}_t, 1 - \bar{z}_t\}) \geq 0 \quad \text{almost surely} \quad (*)$$

Claim: there is an algorithm that ensures $(*)$ for an **arbitrary sequence**. Any idea how to do it?

Another way to formulate $(*)$: number of mistakes should be not much more than made by the best of the two “experts”, one predicting “1” all the time, the other constantly predicting “0”.

Example #1: Bit Prediction

Of course, we should do majority vote over the past outcomes

$$z_t = \mathbf{I}_{\{\bar{y}_{t-1} \geq 1/2\}}$$

where $\bar{y}_{t-1} = \frac{1}{t-1} \sum_{s=1}^{t-1} y_s$. This algorithm guarantees $\bar{c}_t \rightarrow \max\{p, 1-p\}$ and

$$\liminf_{t \rightarrow \infty} (\bar{c}_t - \max\{\bar{z}_t, 1 - \bar{z}_t\}) \geq 0 \quad \text{almost surely} \quad (*)$$

Claim: there is an algorithm that ensures $(*)$ for an **arbitrary sequence**. Any idea how to do it?

Another way to formulate $(*)$: number of mistakes should be not much more than made by the best of the two “experts”, one predicting “1” all the time, the other constantly predicting “0”.

Note the difference: estimating a hypothesized model vs competing against a reference set. We had seen this distinction in the previous lecture.

Example #2: Email Spam Detection



We are tasked with developing a spam detection program that needs to be adaptive to malicious attacks.

- ▶ x_1, \dots, x_n are email messages, revealed one-by-one
- ▶ upon observing the message x_t , the learner (spam detector) needs to decide whether it is spam or not spam ($\hat{y}_t \in \{0, 1\}$)
- ▶ the actual label $y_t \in \{0, 1\}$ is revealed (e.g. by the user)

Do it seem plausible that $(x_1, y_1), \dots, (x_n, y_n)$ are i.i.d. from some distribution P ?

Probably not... In fact, the sequence might even be *adversarially* chosen. In fact, spammers *adapt* and try to improve their strategies.

Outline

Introduction

Statistical Learning Theory

The Setting of SLT

Consistency, No Free Lunch Theorems, Bias-Variance Tradeoff

Tools from Probability, Empirical Processes

From Finite to Infinite Classes

Uniform Convergence, Symmetrization, and Rademacher Complexity

Large Margin Theory for Classification

Properties of Rademacher Complexity

Covering Numbers and Scale-Sensitive Dimensions

Faster Rates

Model Selection

Sequential Prediction / Online Learning

Motivation

Supervised Learning

Online Convex and Linear Optimization

Online-to-Batch Conversion, SVM optimization

Online Learning (Supervised)

- ▶ No assumption that there is a single distribution \mathcal{P}
- ▶ Data not given all at once, but rather in the online fashion
- ▶ As before, \mathcal{X} is the space of inputs, \mathcal{Y} the space of outputs
- ▶ Loss function $\ell(\mathbf{y}_1, \mathbf{y}_2)$

Online protocol (*supervised learning*):

For $t = 1, \dots, n$
Observe \mathbf{x}_t , predict $\hat{\mathbf{y}}_t$, observe \mathbf{y}_t

Goal: keep **regret** small:

$$\text{Reg}_n = \frac{1}{n} \sum_{t=1}^n \ell(\hat{\mathbf{y}}_t, \mathbf{y}_t) - \inf_{f \in \mathcal{F}} \frac{1}{n} \sum_{t=1}^n \ell(f(\mathbf{x}_t), \mathbf{y}_t)$$

A bound on Reg_n should hold for **any** sequence $(\mathbf{x}_1, \mathbf{y}_1), \dots, (\mathbf{x}_n, \mathbf{y}_n)$!

Pros/Cons of Online Learning

The good:

- ▶ An upper bound on regret implies good performance relative to the set \mathcal{F} *no matter how adversarial the sequence is.*
- ▶ Online methods are typically computationally attractive as they process one data point at a time. Used when data sets are huge.
- ▶ Interesting research connections to Game Theory, Information Theory, Statistics, Computer Science.

The bad:

- ▶ A regret bound implies good performance only if one of the elements of \mathcal{F} has good performance (just as in Statistical Learning). However, for non-iid sequences a single $f \in \mathcal{F}$ might not be good at all! To alleviate this problem, the comparator set \mathcal{F} can be made into a set of more complex strategies.
- ▶ There might be some (non-i.i.d.) structure of sequences that we are not exploiting (this is an interesting area of research!)

Setting Up the Minimax Value

First, it turns out that \hat{y}_t has to be a *randomized* prediction: we need to decide on a distribution $q_t \in \Delta(\mathcal{Y})$ and then draw \hat{y}_t from q_t .

The minimax best that both the learner and the adversary (or, Nature) can do is

$$\mathcal{V}(\mathcal{F}, n) = \left\langle \left\langle \sup_{x_t \in \mathcal{X}} \inf_{q_t \in \Delta} \sup_{y_t \in \mathcal{Y}} \mathbb{E}_{y_t \sim q_t} \right\rangle \right\rangle_{t=1}^n \left\{ \frac{1}{n} \sum_{t=1}^n \ell(\hat{y}_t, y_t) - \inf_{f \in \mathcal{F}} \frac{1}{n} \sum_{t=1}^n \ell(f(x_t), y_t) \right\}$$

This is an awkward and long expression, so no need to be worried. All you need to know right now is:

- ▶ An upper bound on $\mathcal{V}(\mathcal{F}, n)$ guarantees existence of a strategy (learning algorithm) that will suffer at most that much regret.
- ▶ A lower bound on $\mathcal{V}(\mathcal{F}, n)$ means the adversary can inflict at least that much damage, no matter what the learning algorithm does.

It is interesting to study $\mathcal{V}(\mathcal{F}, n)$! It turns out, many of the tools we used in Statistical Learning can be extended to study Online Learning!

Sequential Rademacher Complexity

A (complete binary) \mathcal{X} -valued tree \mathbf{x} of depth n is a collection of functions $\mathbf{x}_1, \dots, \mathbf{x}_n$ such that $\mathbf{x}_i : \{\pm 1\}^{i-1} \mapsto \mathcal{X}$ and \mathbf{x}_1 is a constant function.

A sequence $\epsilon = (\epsilon_1, \dots, \epsilon_n)$ defines a path in \mathbf{x} :

$$\mathbf{x}_1, \mathbf{x}_2(\epsilon_1), \mathbf{x}_3(\epsilon_1, \epsilon_2), \dots, \mathbf{x}_n(\epsilon_1, \dots, \epsilon_{n-1})$$

Define *sequential Rademacher complexity* as

$$\mathcal{R}_n^{\text{seq}}(\mathcal{F}, n) = \sup_{\mathbf{x}} \mathbb{E}_{\epsilon_{1:n}} \sup_{f \in \mathcal{F}} \left\{ \frac{1}{n} \sum_{t=1}^n \epsilon_t f(\mathbf{x}_t(\epsilon_{1:t-1})) \right\}$$

where the supremum is over all \mathcal{X} -valued trees of depth n .

Theorem

Let $\mathcal{Y} = \{0, 1\}$ and \mathcal{F} is a class of binary-valued functions. Let ℓ be the indicator loss. Then

$$\mathcal{V}(\mathcal{F}, n) \leq 2\mathcal{R}_n^{\text{seq}}(\mathcal{F}, n)$$

Finite Class

Suppose \mathcal{F} is finite, $N = \text{card}(\mathcal{F})$. Then for any tree \mathbf{x} ,

$$\mathbb{E}_{\epsilon_{1:n}} \sup_{f \in \mathcal{F}} \left\{ \frac{1}{n} \sum_{t=1}^n \epsilon_t f(\mathbf{x}_t(\epsilon_{1:t-1})) \right\} \leq \sqrt{\frac{2 \log N}{n}}$$

because, again, this is a maximum of N (sub)Gaussian Random variables!

Hence,

$$\mathcal{V}(\mathcal{F}, n) \leq 2\sqrt{\frac{2 \log N}{n}}$$

This bound is basically the same as that for Statistical Learning with a finite number of functions!

Therefore, there must exist an algorithm for predicting \hat{y}_t given \mathbf{x}_t such that regret scales as $O\left(\sqrt{\frac{\log N}{n}}\right)$. What is it?

Exponential Weights, or the Experts Algorithm

We think of each element $\{f_1, \dots, f_N\} = \mathcal{F}$ as an expert who gives a prediction $f_i(x_t)$ given side information x_t . We keep distribution w_t over experts, according to their performance.

Let $w_1 = (1/N, \dots, 1/N)$, $\eta = \sqrt{(8 \log N)/T}$.

To predict at round t , observe x_t , pick $i_t \sim w_t$ and set $\hat{y}_t = f_{i_t}(x_t)$.

Update

$$w_{t+1}(i) \propto w_t(i) \exp \{-\eta \mathbf{I}_{\{f_i(x_t) \neq y_t\}}\}$$

Claim: for any sequence $(x_1, y_1), \dots, (x_n, y_n)$, with probability at least $1 - \delta$

$$\frac{1}{n} \sum_{t=1}^n \mathbf{I}_{\{\hat{y}_t \neq y_t\}} - \inf_{f \in \mathcal{F}} \frac{1}{n} \sum_{t=1}^n \mathbf{I}_{\{f(x_t) \neq y_t\}} \leq \sqrt{\frac{\log N}{2n}} + \sqrt{\frac{\log(1/\delta)}{2n}}$$

Useful Properties of Sequential Rademacher Complexity

Sequential Rademacher complexity enjoys the same nice properties as its iid cousin, except for the Lipschitz contraction (4). At the moment we can only prove

$$\mathcal{R}_n^{\text{seq}}(\phi \circ \mathcal{F}) \leq L \mathcal{R}_n^{\text{seq}}(\mathcal{F}) \times O(\log^{3/2} n)$$

It is an open question whether this logarithmic factor can be removed...

Theory for Online Learning

There is now a theory with combinatorial parameters, covering numbers, and even a recipe for developing online algorithms!

Many of the relevant concepts (e.g. sequential Rademacher complexity) are generalizations of the i.i.d. analogues to the case of dependent data.

Coupled with the online-to-batch conversion we introduce in a few slides, there is now an interesting possibility of developing new computationally attractive algorithms for statistical learning. One such example will be presented.

Theory for Online Learning

| Statistical Learning | Online Learning |
|--|--|
| i.i.d. data | arbitrary sequences |
| tuples of data | binary trees |
| Rademacher averages | sequential Rademacher complexity |
| covering / packing numbers | tree cover |
| Dudley entropy integral | analogous result with tree cover |
| VC dimension | Littlestone's dimension |
| Scale-sensitive dimension | analogue for trees |
| Vapnik-Chervonenkis-Sauer-Shelah Lemma | analogous combinatorial result for trees |
| ERM and regularized ERM | many interesting algorithms |

Outline

Introduction

Statistical Learning Theory

The Setting of SLT

Consistency, No Free Lunch Theorems, Bias-Variance Tradeoff

Tools from Probability, Empirical Processes

From Finite to Infinite Classes

Uniform Convergence, Symmetrization, and Rademacher Complexity

Large Margin Theory for Classification

Properties of Rademacher Complexity

Covering Numbers and Scale-Sensitive Dimensions

Faster Rates

Model Selection

Sequential Prediction / Online Learning

Motivation

Supervised Learning

Online Convex and Linear Optimization

Online-to-Batch Conversion, SVM optimization

Online Convex and Linear Optimization

For many problems, $\ell(f, (x, y))$ is convex in f and \mathcal{F} is a convex set. Let us simply write $\ell(f, z)$, where the move z need not be of the form (x, y) .

- ▷ e.g. square loss $\ell(f, (x, y)) = (\langle f, x \rangle - y)^2$ for linear regression.
- ▷ e.g. hinge loss $\ell(f, (x, y)) = \max\{0, 1 - y \langle f, x \rangle\}$, a surrogate loss for classification.

We may then use optimization algorithms for updating our hypothesis after seeing each additional data point.

Online Convex and Linear Optimization

Online protocol (*Online Convex Optimization*):

For $t = 1, \dots, n$
Predict $f_t \in \mathcal{F}$, observe z_t

Goal: keep **regret** small:

$$\text{Reg}_n = \frac{1}{n} \sum_{t=1}^n \ell(f_t, z_t) - \inf_{f \in \mathcal{F}} \frac{1}{n} \sum_{t=1}^n \ell(f, z_t)$$

Online Linear Optimization is a particular case when $\ell(f, z) = \langle f, z \rangle$.

Gradient Descent

At time $t = 1, \dots, n$, predict $f_t \in \mathcal{F}$, observe z_t , update

$$f'_{t+1} = f_t - \eta \nabla \ell(f_t, z_t)$$

and project f'_{t+1} to the set \mathcal{F} , yielding f_{t+1} .

- ▶ η is a learning rate (step size)
- ▶ gradient is with respect to the first coordinate

This simple algorithm guarantees that for any $f \in \mathcal{F}$

$$\begin{aligned} \frac{1}{n} \sum_{t=1}^n \ell(f_t, z_t) - \frac{1}{n} \sum_{t=1}^n \ell(f, z_t) &\leq \frac{1}{n} \sum_{t=1}^n \langle f_t, \nabla \ell(f_t, z_t) \rangle - \frac{1}{n} \sum_{t=1}^n \langle f, \nabla \ell(f_t, z_t) \rangle \\ &\leq O(n^{-1/2}) \end{aligned}$$

as long as $\|\nabla \ell(f_t, z_t)\| \leq c$ for some constant c , for all t , and \mathcal{F} has a bounded diameter.

Gradient Descent for Strongly Convex Functions

Assume that for any z , $\ell(\cdot, z)$ is strongly convex in the first argument. That is, $\ell(f, z) - \frac{1}{2}\|f\|^2$ is a convex function.

The same gradient descent algorithm with a different step size η guarantees that for any $f^* \in \mathcal{F}$

$$\frac{1}{n} \sum_{t=1}^n \ell(f_t, z_t) - \frac{1}{n} \sum_{t=1}^n \ell(f, z_t) \leq O\left(\frac{\log(n)}{n}\right),$$

a faster rate.

Outline

Introduction

Statistical Learning Theory

The Setting of SLT

Consistency, No Free Lunch Theorems, Bias-Variance Tradeoff

Tools from Probability, Empirical Processes

From Finite to Infinite Classes

Uniform Convergence, Symmetrization, and Rademacher Complexity

Large Margin Theory for Classification

Properties of Rademacher Complexity

Covering Numbers and Scale-Sensitive Dimensions

Faster Rates

Model Selection

Sequential Prediction / Online Learning

Motivation

Supervised Learning

Online Convex and Linear Optimization

Online-to-Batch Conversion, SVM optimization

How to use regret bounds for i.i.d. data

Suppose we have a regret bound

$$\frac{1}{n} \sum_{t=1}^n \ell(f_t, z_t) - \inf_{f \in \mathcal{F}} \frac{1}{n} \sum_{t=1}^n \ell(f, z_t) \leq R_n$$

that holds for all sequences z_1, \dots, z_n , for some $R_n \rightarrow 0$.

Assume z_1, \dots, z_n are i.i.d. with distribution P . Run the regret minimization algorithm on these data and let $\bar{f} = \frac{1}{n} \sum_{t=1}^n f_t$. Then

$$\mathbb{E}_{z, z_1, \dots, z_n} \ell(\bar{f}, z) \leq \mathbb{E} \left\{ \frac{1}{n} \sum_{t=1}^n \ell(f_t, z) \right\} = \mathbb{E} \left\{ \frac{1}{n} \sum_{t=1}^n \ell(f_t, z_t) \right\}$$

where the last step holds because f_t only depends on z_1, \dots, z_{t-1} . Also,

$$\mathbb{E} \left\{ \inf_{f \in \mathcal{F}} \frac{1}{n} \sum_{t=1}^n \ell(f, z_t) \right\} \leq \inf_{f \in \mathcal{F}} \mathbb{E} \left\{ \frac{1}{n} \sum_{t=1}^n \ell(f, z_t) \right\} = \mathbb{E}_z \ell(f_{\mathcal{F}}, z)$$

Combining,

$$\mathbb{E} L(\bar{f}) - \inf_{f \in \mathcal{F}} L(f) \leq R_n$$

How to use regret bounds for i.i.d. data

This gives an alternative way of proving bounds on

$$\mathbb{E}L(\hat{f}_n) - \inf_{f \in \mathcal{F}} L(f)$$

by using $\hat{f}_n = \bar{f}$, the average of the trajectory of an online learning algorithm.

Next, we present an interesting application of this idea.

Pegasos

Support Vector Machine is a fancy name for the algorithm

$$\hat{f}_n = \arg \min_{f \in \mathbb{R}^d} \frac{1}{m} \sum_{i=1}^m \max\{0, 1 - y_i \langle f, x_i \rangle\} + \frac{\lambda}{2} \|f\|^2$$

in the linear case.

The objective can be “kernelized” for representing linear separators in higher-dimensional feature space. The hinge loss is convex in f .

Write

$$\ell(f, z) = \max\{0, 1 - y \langle f, x \rangle\} + \frac{\lambda}{2} \|f\|^2$$

for $z = (x, y)$. Then the objective of SVM can be written as

$$\min_f \mathbb{E} \ell(f, z)$$

The expectation is with respect to the *empirical distribution* $\frac{1}{m} \sum_{i=1}^m \delta_{(x_i, y_i)}$.

Then an i.i.d. sample z_1, \dots, z_n from the empirical distribution is simply a draw with replacement from the dataset $\{(x_1, y_1), \dots, (x_m, y_m)\}$.

Pegasos

A gradient descent $f_{t+1} = f_t - \eta \nabla \ell(f_t, z_t)$ with

$$\nabla \ell(f_t, z_t) = -y_t x_t \mathbf{I}_{\{y_t \langle f_t, x_t \rangle < 1\}} + \lambda f_t$$

then gives a guarantee

$$\mathbb{E} \ell(\bar{f}, z) - \inf_{f \in \mathcal{F}} \mathbb{E} \ell(f, z) \leq R_n$$

Since $\ell(f, z)$ is λ -strongly convex, the rate $R_n = O(\log(n)/n)$.

Pegasos (Shalev-Shwartz et al, 2010)

For $t = 1, \dots, n$

Choose a random example (x_{i_t}, y_{i_t}) from the dataset. Set $\eta = 1/(\lambda t)$

If $y_{i_t} \langle f_t, x_{i_t} \rangle < 1$, update $f_{t+1} = (1 - \eta_t \lambda) f_t + \eta_t x_{i_t} y_{i_t}$

else, update $f_{t+1} = (1 - \eta_t \lambda) f_t$

The algorithm and analysis are due to (S. Shalev-Shwartz, Singer, Srebro, Cotter, 2010)

Pegasos

We conclude that $\bar{f} = \frac{1}{n} \sum_{t=1}^n f_t$ computed using the gradient descent algorithm is an $\tilde{O}(n^{-1})$ -approximate minimizer of the SVM objective after n steps.

This gives an $O(d/(\lambda\epsilon))$ time to converge to an ϵ -minimizer. Very fast SVM solver, attractive for large datasets!

Summary

Key points for both statistical and online learning:

- ▶ obtained performance guarantees with minimal assumptions
- ▶ prior knowledge is captured by the comparator term
- ▶ understanding the inherent complexity of the comparator set
- ▶ key techniques: empirical processes for iid and non-iid data
- ▶ interesting relationships between statistical and online learning
- ▶ computation and statistics – a basis of machine learning