

553.740 Course Notes

September 14, 2023

Contents

1	Lecture 1	1
1.1	Administrativa	1
1.1.1	Git Repo	1
1.2	Course Overview	1
1.2.1	Introduction	1
1.2.2	The Standard Diagram	2
1.2.3	Themes	2
1.3	Zeroth Assignments	3
1.4	Intuition for Measure from Calculus	4
1.5	The Probability Measure	5
2	Lecture 2	8
2.1	Random Variables	8
2.2	Expectation	10
2.3	Marginalization	11
2.4	Conditional Probability	11
2.5	Independence	13
3	Lecture 3	15
3.1	Independence	15
3.2	Data	16
3.2.1	Law of Large Numbers (LLN)	16
3.3	Concentration Inequalities	17
3.4	Intuition in High Dimension	19
4	Lecture 4	20
5	Lecture 5	25

1 Lecture 1

1.1 Administrativa

1.1.1 Git Repo

Assignments, starter code, and data will be housed in a Git repository:

<https://github.com/schmidttgenstein/fa23-mli/> Please do not push anything to this repo!

1.2 Course Overview

1.2.1 Introduction

From the syllabus:

Machine Learning describes a mishmash of computational techniques for “finding patterns in data.” The scope of use, analytic tools, algorithms, and results are almost too numerous to meaningfully batch all such applications under a common appellation. Still, we try. This course focuses on *supervised* machine learning (sml) which roughly deals with using historical *labeled* data to construct a predictor which will correctly label future data.

Formally, we will be operating in space $\mathcal{X} \times \mathcal{Y}$, where \mathcal{X} denotes “input” space, and \mathcal{Y} “output.” While these notions are heuristic, they will frame the situation that we may easily sample data at will from \mathcal{X} , while sampling from \mathcal{Y} may be difficult or expensive, and often we would like to decide according to how we believe $x \in \mathcal{X}$ is associated with label $y \in \mathcal{Y}$.

In this course, you will learn how to formulate the supervised learning problem in mathematical terms, how to describe a measure of performance, restrict search space for constructing models for prediction, optimize performance measure in search space, and how to check for generalization. You will learn, also, how to implement some of these methods in code, from the ground up, as well as incorporating pre-built libraries (such as pytorch) for such tasks. Finally, you will learn how to articulate learning guarantees, and understand some of the limits of learning claims. While this course is primarily theory-centric, there will be no dearth of opportunity for employing concrete computational techniques.

1.2.2 The Standard Diagram

Describing our problem space in more detail, consider the following diagram

$$\begin{array}{ccc} \mathcal{X} \times \mathcal{Y} & \xrightarrow{\pi_{\mathcal{Y}}} & \mathcal{Y} \\ \downarrow \pi_{\mathcal{X}} & \nearrow \tilde{y} & \\ \mathcal{X} & & \end{array} \quad (1)$$

We will refer to this diagram often, and to do so give it the somewhat non-descriptive, but in our context wholly unambiguous, name ‘the standard diagram.’

Traditionally, $\pi_{\mathcal{X}} : \mathcal{X} \times \mathcal{Y} \rightarrow \mathcal{X}$, defined by mapping $(x, y) \mapsto x$ (read: $\pi_{\mathcal{X}}(x, y) := x$), is taken to be “easy, efficient, or cheap” to evaluate or sample while $\pi_{\mathcal{Y}} : \mathcal{X} \times \mathcal{Y} \rightarrow \mathcal{Y}$, defined by mapping $(x, y) \mapsto y$, is computationally expensive, expensive otherwise, difficult for other reasons, or altogether infeasible. The original space $\mathcal{X} \times \mathcal{Y}$ is itself inaccessible, except for some *labeled* data $S = \{(x_1, y_1), \dots, (x_m, y_m)\} \subset \mathcal{X} \times \mathcal{Y}$ which provides a proxy (and incomplete!) illustration of what $\mathcal{X} \times \mathcal{Y}$ looks like. The map $\tilde{y} : \mathcal{X} \dashrightarrow \mathcal{Y}$ is a critter we’d like to construct from data S so that both $\tilde{y}(x) \approx y$ for $(x, y) \in S$ and for $(x, y) \in \mathcal{X} \times \mathcal{Y}$. What “ \approx ” means, how to construct \tilde{y} , conditions on S which are needed to make this problem feasible, etc. are all aspects of the supervised machine learning problem which we will explore in this course.

As an example, suppose $\mathcal{X} = \mathbb{R}$ denotes credit score and $\mathcal{Y} = \{0, 1\}$ loan repayment (say ‘1’ corresponds to repayment of loan, ‘0’ to default). Then a *point* $(x, y) \in \mathbb{R}$ represents data corresponding to a loan whose account holder has credit score x and for which the loan was either paid in full ($y = 1$) or not ($y = 0$). The reason we say $\pi_{\mathcal{X}}$ is “easy” to sample is that you may ask any person what their credit score is (more realistically: as creditor, you would see this information *at the time of application*), while loan repayment information (the “label”) would not be observed until potentially many years later when the loan is finally repaid or defaults.

It is worth noting, and perhaps lingering upon the observation, that the “input-output” relation (x, y) is not necessarily functional, i.e. for two points $(x, y), (x', y') \in \mathcal{X} \times \mathcal{Y}, x = x'$ does not imply that $y = y'$. The stand-in for determinism is probability, i.e. we will presume that there is some joint probability measure $\mathbb{P}_{\mathcal{X} \times \mathcal{Y}}$ on $\mathcal{X} \times \mathcal{Y}$, and e.g. that $\mathbb{P}_{\mathcal{Y}|\mathcal{X}}(y = f(x)|x) \neq 1$.

1.2.3 Themes

Generalization Given model $\tilde{y} : \mathcal{X} \rightarrow \mathcal{Y}$, how well does \tilde{y} match the (finite) data we have $S \subseteq \mathcal{X} \times \mathcal{Y}$ —i.e. $\tilde{y}(x) \approx y$ for $(x, y) \in S$ —and the data we don’t have, $(x, y) \in \mathcal{X} \times \mathcal{Y}$?

Dimensionality Computation in high dimensions becomes harder, in part because computation is more expensive, and because there are more “corners” for data to hide in (which exacerbates the

computational problem). The geometry of high dimensionality will be a recurring theme; for now, we simply observe sources of dimensionality:

1. The data “set” itself $S = \{(x_1, y_1), \dots, (x_m, y_m)\} \subset \mathcal{X} \times \mathcal{Y}$. Properly speaking, this data will be presumed to be sampled $(x_i, y_i) \sim_{\text{iid}} \mathbb{P}_{\mathcal{X} \times \mathcal{Y}}$ which *means* that the data “set” is a *point* in the space $(\mathcal{X} \times \mathcal{Y})^m$. A reasonable question to ask, then, is: what is the measure $\mathbb{P}_{(\mathcal{X} \times \mathcal{Y})^m}$?

Independence tells us that it is $\prod_{j=1}^m \mathbb{P}_{\mathcal{X} \times \mathcal{Y}}$. More on this later.

2. Size of space itself. This could include high dimensionality of \mathcal{X} and/or \mathcal{Y} . Examples abound of high dimensional input data: numerous columned tabular data, imagery data, audio data, video data.

3. Parameter space for model \tilde{y} . In the case of linear regression, a model $\tilde{y}(x) = \sum_{j=0}^n a_j x^j$ may have arbitrarily large degree n . Or a fully connected neural network with many nodes and many layers. And so on. More generally, the space of functions $\mathcal{Y}^{\mathcal{X}} := \{\tilde{y} : \mathcal{X} \rightarrow \mathcal{Y}\}$ is even bigger.

Trade-offs Assumptions must be made and compromises allowed for in order to gain tractability in the learning problem. There is no universal solution (“no free lunch,” and as you may imagine, there’s a theorem for that) and formulating the setup to address one challenge may introduce other ones elsewhere (ML can sometimes feel like one giant game of whackamole).

The famed bias-variance trade-off is one example: a high complexity model may well represent the data S , which in one sense is good, but in another is bad if said model represents data *too well*, i.e. at the exclusion of modeling ‘from where the data comes.’

1.3 Zeroth Assignments

Worksheet Please bring hard copy of this worksheet with you to class on Wednesday. You may not answer a question about Independence with only equality or inequality between $f_{\mathcal{X} \times \mathcal{Y}}$ and $f_{\mathcal{X}} \cdot f_{\mathcal{Y}}$: please say something to indicate how you know such equality or inequality.

Programming Assignment You may find the first programming assignment under pa0 in git repo, and starter code in the git. I suggest you follow the tutorial at Real Python [real python](#)¹ which shows you how to spin up a logistic regression model using sklearn. Scikit-Learn (also known as sklearn) is an open source ML library for python, and contains functionality for constructing numerous models. This is perhaps the only time in this course you will be asked to use this library, and if you have a preferred alternative library, you are more than welcome to use it for this assignment.

The purpose of the assignment is threefold:

1. Gain initial exposure to the *structure* of machine learning code, including object oriented programming and the typical methods included.
2. Shake off any residual rust using Python.
3. To gain deeper appreciation for the *aim* of building model $\tilde{y} : \mathcal{X} \rightarrow \mathcal{Y}$ as in diagram (6), and metrics that illustrate success.

Recall the Standard Diagram

$$\begin{array}{ccc} \mathcal{X} \times \mathcal{Y} & \xrightarrow{\pi_{\mathcal{Y}}} & \mathcal{Y} \\ \downarrow \pi_{\mathcal{X}} & \nearrow \tilde{y} & \\ \mathcal{X} & & \end{array} \quad (2)$$

This diagram provides formalism for talking about “approximating” y with model $\tilde{y} : \mathcal{X} \rightarrow \mathcal{Y}$ when $y \neq y(x)$ is not necessarily functionally determined by x .

¹You may need to sign up in order to view this content, but it is not behind a paywall.

Consider a concrete example to illustrate the problem: suppose that $\mathcal{X} = \mathbb{R}$ denotes credit score and $\mathcal{Y} = \{0, 1\}$ denotes repayment on loan, 1 denotes full repayment and 0 denotes default. A data point $(x, y) \in \mathcal{X} \times \mathcal{Y}$ corresponds to a credit score-loan repayment pair, and in real life a loan account would have these attributes associated with it, i.e. the debtor would have some credit score and their loan will (eventually) be repaid or not. (Note that “eventuality” is what, in this case, makes $\pi_{\mathcal{Y}}$ hard or expensive to evaluate.) It is possible for two different loans, belonging to two different people, to agree on credit score but disagree on outcome $y \in \mathcal{Y}$. In fact, we will likely observe both outcomes $y = 0$ and $y = 1$ associated to *any* credit score. Presumably, there should be some relation between the relative *counts* of $\#y = 1$ and credit score; in other words, one might suppose that lower credit scores correspond to accounts which in actual fact get repaid less frequently than those with high credit scores. We need mathematical language to describe and work with this phenomenon. The language is probability.

Thus, we suppose that there is some joint probability measure $\mathbb{P}_{\mathcal{X} \times \mathcal{Y}}$ on $\mathcal{X} \times \mathcal{Y}$. One way of contextualizing supervised machine learning is as a study of probability on the standard diagram. In this lecture, we will review probability, give intuition for probability as measure, as well as notation $\mathbb{P}_{\mathcal{X}} = \int d\mathbb{P}_{\mathcal{X}}(x)$, and define expectation $\mathbb{E}(f)$ of a random variable $f : \mathcal{X} \rightarrow \mathbb{R}$ as a *Lebesgue* integral $\mathbb{E}(f) := \int_{\mathcal{X}} f(x) d\mathbb{P}_{\mathcal{X}}(x)$. We kick this lecture off by emphasizing that probability has nothing to do with randomness...yet. When we return to admitting randomness into our lexicon, it will be as a *result*. For the moment, we forget any association between probability and chance, stochasticity, randomness, or any other (for now) anathema word affiliated with the notion of uncertainty.

1.4 Intuition for Measure from Calculus

We start reviewing integration in calculus to preview notation for measure. One interpretation of the integral is as “area under the curve.” Another (what amounts to similar) is as measure. And one interpretation of dQ is as an infinitesimal Q element, whatever Q is. Another is as an indicator of what kind of measurement we are taking. We then express length ℓ as $\int d\ell$, area a as $\int da$, volume v as $\int dv$, measure m as $\int dm$, and so on. Note that the expression $\int dm$ is defined in terms of measure m , i.e. we suppose prior (at the very least conceptual) knowledge of the measure, irrespective of whether or not given a particular object to measure A , we are actually able to evaluate its measure $m(A)$.

Example: length We’ve learned that the length $\ell([a, b])$ of an interval $[a, b]$ is the difference of endpoints $b - a$. We can write this with an integral as $\ell([a, b]) = \int_a^b dx$ or more in line with notation we will use, as $\ell = \int_{[a, b]} d\ell(x)$. In this notation, the subscript is the object we are measuring, ℓ is the type of measure, and x is a dummy variable. Until we start integrating functions, we won’t need it, and we could have just written $\int_{[a, b]} d\ell$. The dummy variable was kept only to clearly delineate that $d\ell$ is not the same thing as dx : we are using calculus intuition only for loose guidance.

General Formula For measure m and object to be measured A , we write $m(A) = \int_A dm$. Right now, do not put too much stock in the—what in calculus would be thought of as infinitesimal— dm term: it is *defined* as a *phrase* with the integral \int ; neither in isolation makes sense, and the definition goes this way: $\int_A dm := m(A)$. Thus, you need to first know the measure. We’ll come back to this momentarily.

Example: Area

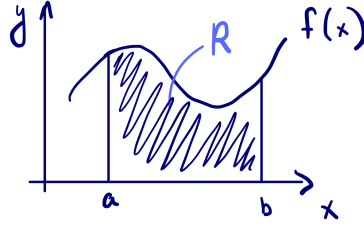


Figure 1: Area

In our notation $A(R) = \int_R dA$, where $A(R)$ is the area. If you want to import calculus intuition, when we interpret dA as a differential area element, we may express it likewise as the product $f(x)dx$, where $f(x)$ represents height and dx the differential width. As area is equal to height times width (or length) a differential area is a length element times a differential length element ie $A(R) = \int_a^b f(x)dx$.

Example: Volume

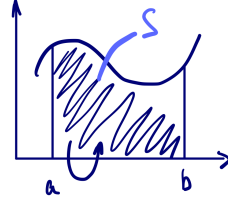


Figure 2: Volume

Using the above picture, $V(S) = \int_S dV$, where $V(S)$ is the volume. Therefore, by generalizing the above idea we can write $\mathbb{P}(A) = \int_A d\mathbb{P}$.

1.5 The Probability Measure

Now we formalize what we mean by probability measure. The first point to make is that it is a measure.

Definition 1.1. Let \mathcal{X} be a set. We define a *probability measure*

$$\mathbb{P}_{\mathcal{X}} : ([\text{Some}] \text{ Subsets of } \mathcal{X}) \rightarrow [0, 1]$$

on \mathcal{X} to be a map from (a subset of)² the power set of \mathcal{X} to the closed interval $[0, 1]$ satisfying the following two properties:

1. $\mathbb{P}_{\mathcal{X}}(\mathcal{X}) = 1$ (space is measure finite and normalized), and
2. $\mathbb{P}_{\mathcal{X}}\left(\bigsqcup_{j=1}^{\infty} A_j\right) = \sum_{j=1}^{\infty} \mathbb{P}_{\mathcal{X}}(A_j)$ where $\bigsqcup_{j=1}^{\infty} A_j$ denotes disjoint union, i.e. as a set $\bigsqcup_{j=1}^{\infty} A_j = \bigcup_{j=1}^{\infty} A_j$ and $A_i \cap A_j = \emptyset$ for $i \neq j$ (countable additivity).

²This subtle point is a technicality beyond the scope of this course, and quite frankly unnecessary for reasonably understanding measure. Measure measures subsets. But it's possible that it may not be able to measure *all* subsets. So the domain of the measure may not be *all* subsets. A lot of work goes into specification of the structure of the collection of subsets you can measure, and like topology is characterized by closure operations, e.g., if you can measure $\{A_j\}_{j \in \mathbb{N}}$ then you can measure its union $\bigcup_{j \in \mathbb{N}} A_j$. For the curious, you may look into sigma algebras for more detail.

Recall that the power set $2^{\mathcal{X}}$ of a set \mathcal{X} is defined to be the set of all subsets of \mathcal{X} , including \emptyset and \mathcal{X} itself. For example, when $\mathcal{X} = \{1, 2, 3\}$,

$$2^{\mathcal{X}} = \{\emptyset, \{1\}, \{2\}, \{3\}, \{1, 2\}, \{1, 3\}, \{2, 3\}, \mathcal{X}\}.$$

When $\mathcal{X} = \mathbb{R}$, the power set includes any reasonable combination (e.g. unions) of intervals (a, b) or $[a, b]$, but also many many more (see [Cantor Set](#) for a fun, but not particularly relevant, excursion into the wonders of measure theory).

The second condition is called ‘countable additivity’ (informally: a conservation of stuff principle) and represents the intuitive idea that if you slice and dice an object for measurement, measure each constituent piece without double counting, and add your results, you’ll end up with the same result as if you just measured the whole original unadulterated tamale.

You should check that countable additivity implies *finite* additivity $\mathbb{P}_{\mathcal{X}}\left(\sum_{j=1}^n A_j\right) = \sum_{j=1}^n \mathbb{P}_{\mathcal{X}}(A_j)$.

You will need the fact that $\mathbb{P}_{\mathcal{X}}(\emptyset) = 0$, which itself is implied by conditions 1. and 2. Indeed, $\mathcal{X} = \mathcal{X} \sqcup \bigsqcup_{j=2}^{\infty} \emptyset$. And the second fact is called the Union Bound: $\mathbb{P}(\cup_{j \in \mathbb{N}} A_j) \leq \sum_{j \in \mathbb{N}} \mathbb{P}(A_j)$.

Remark 1.1. We are being fairly blasé about the *domain* “some subsets of \mathcal{X} ” of $\mathbb{P}_{\mathcal{X}}$. Perhaps we shouldn’t be. Much of the architecture constructing measure depends on the fact that some structure must be placed on the set of subsets of \mathcal{X} which are suitable for measurement; in particular, that such a set comprises a so-called σ -algebra, is not necessarily (and in many cases in fact is not) the entire power set $2^{\mathcal{X}} = \{A \subset \mathcal{X}\}$, and so on. We pay lip-service to this nuance, but fret little over the possibility that we will accidentally run across both a measure $\mathbb{P}_{\mathcal{X}}$ and subset $A \subset \mathcal{X}$ for which $\mathbb{P}_{\mathcal{X}}(A)$ does not make sense (read: which $\mathbb{P}_{\mathcal{X}}$ is “incapable” of measuring). One must try very hard—you might find such a question on a measure theory qualifying exam—to come up with an example. Therefore you may reasonably suppose that any set you’d come across in real life is in fact measurable. Still, know *that* there is a potential problem: if you can construct a non-measurable set, then you can cut an apple into finitely many pieces and reassemble those finitely many pieces into *two* apples of the same size (see [Banach Tarski](#) for more information). In other words, weird things can happen with things that aren’t measurable.

Definition 1.2. We define a *probability space* to be a pair $(\mathcal{X}, \mathbb{P}_{\mathcal{X}})$ where \mathcal{X} is a set and

$$\mathbb{P}_{\mathcal{X}} : ([\text{Some/Many}] \text{ Subsets of } \mathcal{X}) \rightarrow [0, 1]$$

is a probability measure (c.f. definition 1.1).

In probabilistic terminology, measurable subsets $A \subset \mathcal{X}$ are often called ‘events’ and individual points $x \in \mathcal{X}$ called ‘outcomes.’ An outcome $x \in \mathcal{X}$ defines the *event* $\{x\} \subset \mathcal{X}$ with the single outcome $x \in \{x\}$.

Example 1 Let $(\mathcal{X} = [0, 1], \mathbb{P}_{\mathcal{X}}([a, b]) := b - a)$ for $0 \leq a \leq b \leq 1$. We define notation $\int_{[a, b]} d\mathbb{P}(x) = \mathbb{P}([a, b])$.

Example 2 Let $(\mathcal{X} = \mathbb{R}, \mathbb{P}_{\mathcal{X}}([a, b]) := \frac{1}{\sqrt{2\pi}} \int_a^b e^{-x^2/2})$, the so-called normal distribution with zero mean and unit variance. Observe that we use a Riemann integral to *compute* or give the rule for realizing the probability measure. The integrand $\frac{1}{\sqrt{2\pi}} e^{-x^2/2}$ of the Riemann integral is called a *probability density function*. The integral $\int_{[a, b]} d\mathbb{P}_{\mathcal{X}}(x)$, by contrast, is not a Riemann integral; it is *defined* by the measure $\mathbb{P}_{\mathcal{X}}([a, b])$. (When you ask: but what is the measure?, we gave the rule for how to calculate it!)

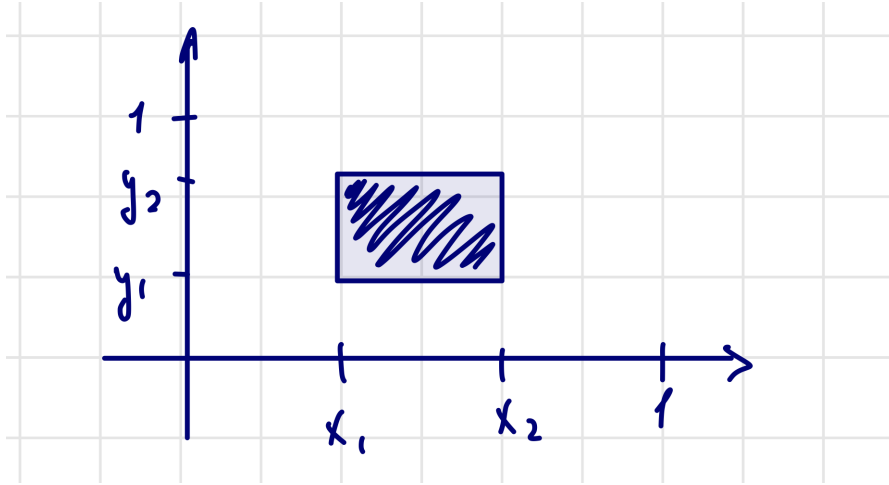


Figure 3: The area measure

Example 3 Let $\mathcal{Y} = [0, 1]$, $\mathbb{P}_{\mathcal{Y}}([x_1, x_2] \times [y_1, y_2]) = (x_2 - x_1) \cdot (y_2 - y_1)$. This example is a preliminary look into independence as $\mathbb{P}_{\mathcal{X}}([x_1, x_2]) = x_2 - x_1$, $\mathbb{P}_{\mathcal{Y}}([y_1, y_2]) = y_2 - y_1$ and $\mathbb{P}_{\mathcal{X}}([x_1, x_2] \times [y_1, y_2]) = \mathbb{P}_{\mathcal{X}}([x_1, x_2])\mathbb{P}_{\mathcal{Y}}([y_1, y_2])$. We will return to this example when we discuss independence, and will want to situate as a notion which is at home in high dimension.

Independence is a high dimensional phenomenon as it is clear that we can generalize the picture in 3 to n dimensions.

Example 4 : Bernouli Random Variable (flipping a fair or unfair coin): Let $\mathcal{X} = \{0, 1\}$ where $\mathbb{P}_{\mathcal{X}}(\{1\}) = p$ and $\mathbb{P}_{\mathcal{X}}(\{0\}) = 1 - p$. This is a probability space with two outcomes that clearly satisfies the axioms stated earlier since the events or subsets of the space \mathcal{X} , $\{0\}$ and $\{1\}$ are disjoint and the sum of their probabilities is 1:

$$\mathbb{P}_{\mathcal{X}}(1) + \mathbb{P}_{\mathcal{X}}(0) = p + (1 - p) = 1 = \mathbb{P}_{\mathcal{X}}(\{0\} \sqcup \{1\})$$

2 Lecture 2

2.1 Random Variables

Recall that a probability space $(\mathcal{X}, \mathbb{P}_{\mathcal{X}})$ consists of a set \mathcal{X} and a measure $\mathbb{P}_{\mathcal{X}}$ (see definition 1.2). We highlighted at the beginning of the lecture that probability fundamentally operates as a theory of measure, which in essence equate it to a theory of integration, a concept that primarily revolves around numbers (values). So far, we've said nothing about the nature of the set \mathcal{X} , and we don't need to. We do need is a way to attribute values to outcomes $x \in \mathcal{X}$.

Definition 2.1. Let $(\mathcal{X}, \mathbb{P}_{\mathcal{X}})$ be a probability space. We define a *random variable* to be a function $f : \mathcal{X} \rightarrow \mathbb{R}$, whose codomain is \mathbb{R} .

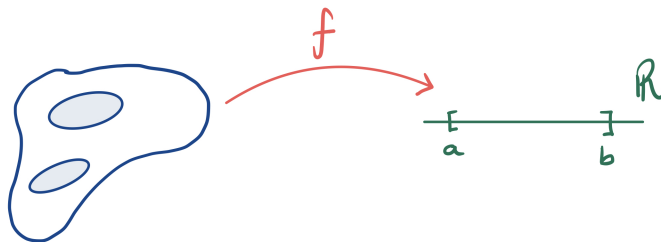


Figure 4: Random variable f

Such a function induces a measure $\mathbb{P}_{\mathbb{R}}$ on \mathbb{R} , defined by

$$\mathbb{P}_{\mathbb{R}}([a, b]) := \mathbb{P}_{\mathcal{X}}(f^{-1}([a, b])) = \mathbb{P}_{\mathcal{X}}(\{x \in \mathcal{X} : f(x) \in [a, b]\}). \quad (3)$$

One should check that this defines an honest probability measure on \mathbb{R} .

Remark 2.1. It is worth noting that certain conditions must be met by the map $f : \mathcal{X} \rightarrow \mathbb{R}$ in order to ensure that the induced measure $\mathbb{P}_{\mathbb{R}}$ is well-defined, i.e. it is a (probability) measure. In essence, we require a condition known as ‘measurability,’ which essentially means that f must be a *measurable* function (which really just means: f is such that the induced measure is a measure—this isn't circular!, it all comes down to saying, you cannot with impunity claim that any function whatsoever will induce a measure). Just as with the domain of $\mathbb{P}_{\mathcal{X}}$, where we suppose with little guilt that “all subsets” we encounter are measurable, we will also suppose that the functions we come across in practice are measurable. Again, there is nuance to be appreciated, but the supposition we make will very unlikely harm any of our day to day calculations.

(For the ultra-curious, the condition of measurability stipulates that any potentially measurable set in \mathbb{R} has preimage (by f^{-1}) which is measurable in \mathcal{X} .)

Remark 2.2. We noted that $f : \mathcal{X} \rightarrow \mathbb{R}$ induces a measure. In fact, this holds more generally for any (measurable) map $f : \mathcal{X} \rightarrow \mathcal{Y}$, $\mathbb{P}_{\mathcal{Y}}(B) := \mathbb{P}_{\mathcal{X}}(f^{-1}(B))$.

Definition 2.2. Let $f : \mathcal{X} \rightarrow \mathbb{R}$ be a random variable. We define *expectation* of f , denoted $\mathbb{E}(f)$, to be

$$\mathbb{E}(f) := \int_{\mathcal{X}} f(x) d\mathbb{P}_{\mathcal{X}}(x). \quad (4)$$

Equation (4) defines expectation, but we have some odd sort of critter we've not seen before on the right hand side. We must define it.

Definition 2.3. Let $f : \mathcal{X} \rightarrow \mathbb{R}$ be a *simple* random variable, i.e. taking *finitely* many values a_1, \dots, a_k . Then we define the *Lebesgue* integral of f , denoted $\int_{\mathcal{X}} f(x) d\mathbb{P}_{\mathcal{X}}(x)$, to be

$$\int_{\mathcal{X}} f(x) d\mathbb{P}_{\mathcal{X}}(x) := \sum_{j=1}^k a_j \mathbb{P}_{\mathcal{X}}(f = a_j) = \sum_{j=1}^k a_j \mathbb{P}_{\mathcal{X}}(\{x \in \mathcal{X} : f(x) = a_j\}).$$

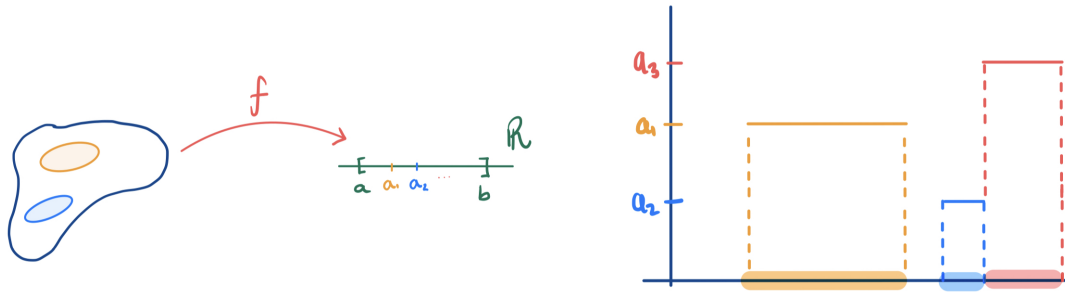


Figure 5: Simple random variable f

One may recognize this definition as the *expectation of a discrete random variable*. Indeed, that is exactly what it is. It is the Lebesgue integral! The Lebesgue integral extends to continuous random variables, but it requires some abstraction. For next time. In the meantime, note that for continuous random variable $f : \mathcal{X} \rightarrow \mathbb{R}$ (not necessarily taking finitely many values), we can *approximate* $\mathbb{E}(f)$ using the Lebesgue integral of a simple random variable (expectation of discrete r.v.), see fig. 6.

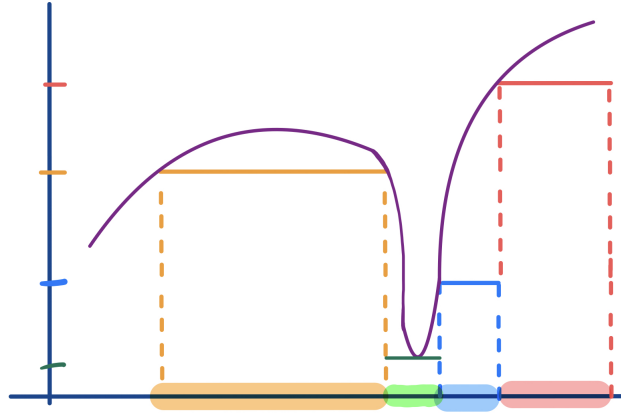


Figure 6: Approximation via simple functions

Remark 2.3. In the following, we will continue the definition of Lebesgue integral for continuous random variable. While there are theoretical reasons for insisting on the use of this abstraction, ours are more practically oriented. In fact, one may (very) often compute expectation *using* instead a Riemann integral. For example, the expectation of a mean 0 unit variance normally distributed random variable is $\frac{1}{\sqrt{2\pi}} \int_{-\infty}^{\infty} x e^{-x^2/2} dx$ (which of course is zero). In other words, if push comes to shove and you're asked to *compute* the expectation for a continuous random variable, very often you'll have a density function in your pocket and can just multiply it by the random variable, and integrate, business as usual. We introduce Lebesgue integration for two reasons:

1. To articulate the fact that the extension from discrete to continuous random variables may *seem* confusing, and that is due in part to the nauseating head-spinning move from Lebesgue to who knows (but usually Riemann) integration without even lip-service paid to the fact that expectation of discrete r.v.s itself is a non-trivial extension of our conceptual apparatus.
2. Ease of notation: $d\mathbb{P}_{\mathcal{X}}$ always makes sense and makes immediately obvious what our measure is. In other words, I don't always want to say: suppose that the density of a probability space exists, and anyway it might not and that doesn't matter!, for $\int_{\mathcal{A}} d\mathbb{P}_{\mathcal{X}}(x)$ makes sense regard-

less of whether we can integrate (Riemann-wise) as $\int \varphi(x)dx$ (for density $\varphi(x)$). Related: the notation $d\mathbb{P}_{\mathcal{X}}$ collapses the distinction between continuous and discrete random variables. The distinction, in my view, is convoluted and confusing, especially when we have mixed discrete-continuous nonsense going on (e.g. in the case of binary classification $\mathcal{X} = \mathbb{R}$, $\mathcal{Y} = \{0, 1\}$). Of course, successfully working with the collapse requires an comfort with the abstraction, and that may by itself be initially confusing as well.

Remark 2.4. Observe also that $\mathbb{E}(f)$ is somewhat uninformative, in a way that $\int_{\mathcal{X}} f(x)d\mathbb{P}_{\mathcal{X}}(x)$ is not. At the moment the added notational baggage of the latter may seem inconvenient, but once we start squinting at joint probability spaces $\mathcal{X} \times \mathcal{Y}$, turning them upside down, and so on, it will be imperative to be clear on how we are integrating. For example, we will see

$$\int_{\mathcal{X} \times \mathcal{Y}} f(x, y) d\mathbb{P}_{\mathcal{X} \times \mathcal{Y}}(x, y) = \int_{\mathcal{X}} \int_{\mathcal{Y}|\mathcal{X}} f(x, y) d\mathbb{P}_{\mathcal{Y}|\mathcal{X}}(y|x) d\mathbb{P}_{\mathcal{X}}(x).$$

This is sortof an extension on the notational point above. I will be relying on this notation aggressively, so it's crucial to anticipate eventually becoming comfortable with standard manipulations, and recalling (any time there is lingering confusion) that $d\mathbb{P}_{\mathcal{X}}$ is defined, not in isolation, but as a package $\int_A d\mathbb{P}_{\mathcal{X}} := \mathbb{P}_{\mathcal{X}}(A)$.

2.2 Expectation

Recall definition 2.2 for which $\mathbb{E}(f) = \int_{\mathcal{X}} f(x)d\mathbb{P}_{\mathcal{X}}(x)$ for a random variable $f : \mathcal{X} \rightarrow \mathbb{R}$ on the probability space $(\mathcal{X}, \mathbb{P}_{\mathcal{X}})$. While this definition lays out the concept of expectation, it leaves us pondering what defines the integral itself. We can draw insight from discrete (or simple) random variables like in definition 2.3.

Definition 2.4. Let $(\mathcal{X}, \mathbb{P}_{\mathcal{X}})$ be a probability space and $f : \mathcal{X} \rightarrow \mathbb{R}^{\geq 0}$ a nonnegative random variable. Then we define the *Lebesgue integral* of f as

$$\int_{\mathcal{X}} f(x) d\mathbb{P}_{\mathcal{X}}(x) := \sup \{ \mathbb{E}(\bar{f}) : 0 \leq \bar{f} \leq f \text{ is simple random variable} \}. \quad (5)$$

A visual illustration is provided in fig. 7.³

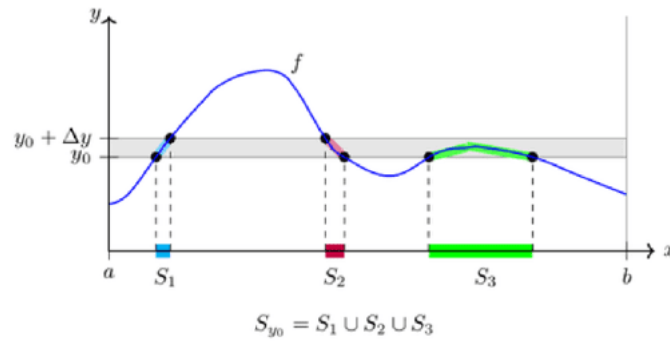


Figure 7: Visualizing Lebesgue Integration

We provide this definition because it is the definition given to us of expectation. The move to continuous random variables may seem abstract—indeed, we discretize the codomain \mathbb{R} instead of domain \mathcal{X} as we are used to doing with Riemann integration from calculus—but one may rest assured that in many instances expectation *may* be computed as a Riemann integral instead. In particular, when a *density* function exists, $\rho : \mathcal{X} \rightarrow \mathbb{R}$ such that $\mathbb{P}_{\mathcal{X}}([a, b]) = \int_a^b \rho(x)dx$, we may compute expectation using the following procedure:

³An Introduction to the Lebesgue Integral, Ikhlas Adi, 2017

1. multiply the random variable f by density ρ ; this will give us a function $f \cdot \rho : \mathcal{X} \rightarrow \mathbb{R}$, and
2. compute the Riemann integral $\int_{-\infty}^{\infty} f(x)\rho(x)dx$.

When the density does not exist, the definition of expectation still makes sense, and one cannot resort to this procedure for computation. Thus for the sake of understanding, one may choose to delve into the nuance of Lebesgue integration, or pretend, as we've been pretending about measurable sets and measurable functions, that "everything" can be computed as a Riemann integral.

Example 2.1. We once again point out that probability may be computed as expectation:

$$\mathbb{P}_{\mathcal{X}}(A) = 1 \cdot \mathbb{P}_{\mathcal{X}}(A) + 0 \cdot \mathbb{P}_{\mathcal{X}}(\mathcal{X} \setminus A) = \int_{\mathcal{X}} \mathbb{1}_{x \in A} d\mathbb{P}_{\mathcal{X}}(x) = \mathbb{E}(\mathbb{1}_{x \in A}).$$

Observe that this Lebesgue integral uses the definition from definition 2.3; we do not need to rely on any limiting procedure (as in e.g. definition 2.4).

For the sake of completeness, we define expectation for arbitrary r.v. $f : \mathcal{X} \rightarrow \mathbb{R}$.

Definition 2.5. Let $f : \mathcal{X} \rightarrow \mathbb{R}$ and suppose that $\mathbb{E}(f \cdot \mathbb{1}_{f \geq 0}) < \infty$ and $\mathbb{E}(-f \cdot \mathbb{1}_{f < 0}) < \infty$ (both expectations of non-negative random variables). Then we define

$$\mathbb{E}(f) := \mathbb{E}(f \cdot \mathbb{1}_{f \geq 0}) - \mathbb{E}(-f \cdot \mathbb{1}_{f < 0})$$

Remark 2.5. Expectation is already defined as $\mathbb{E}(f) = \int_{\mathcal{X}} f(x) d\mathbb{P}_{\mathcal{X}}(x)$. This definition is defining the right hand side.

Definition 2.6. Let $(\mathcal{X}, \mathbb{P}_{\mathcal{X}})$ and $(\mathcal{Y}, \mathbb{P}_{\mathcal{Y}})$ be probability spaces. A map $f : (\mathcal{X}, \mathbb{P}_{\mathcal{X}}) \rightarrow (\mathcal{Y}, \mathbb{P}_{\mathcal{Y}})$ is a map from \mathcal{X} to \mathcal{Y} that *respects the measure* i.e.

$$\mathbb{P}_{\mathcal{Y}}(B) = \mathbb{P}_{\mathcal{X}}(f^{-1}(B))$$

for all $B \subseteq \mathcal{Y}$.

2.3 Marginalization

We now turn to induced probability measures on the standard diagram

$$\begin{array}{ccc} \mathcal{X} \times \mathcal{Y} & \xrightarrow{\pi_{\mathcal{Y}}} & \mathcal{Y} \\ \downarrow \pi_{\mathcal{X}} & \nearrow \tilde{g} & \\ \mathcal{X} & & \end{array} \quad (6)$$

Let $(\mathcal{X} \times \mathcal{Y}, \mathbb{P}_{\mathcal{X} \times \mathcal{Y}})$ be a joint probability measure. The projection maps $\mathcal{X} \times \mathcal{Y} \xrightarrow{\pi_{\mathcal{X}}} \mathcal{X}$ and $\mathcal{X} \times \mathcal{Y} \xrightarrow{\pi_{\mathcal{Y}}} \mathcal{Y}$ induce probability measures on \mathcal{X} and \mathcal{Y} defined as:

$$\mathbb{P}_{\mathcal{X}}(A) := \mathbb{P}_{\mathcal{X} \times \mathcal{Y}}(\pi_{\mathcal{X}}^{-1}(A)) = \mathbb{P}_{\mathcal{X} \times \mathcal{Y}}(A \times \mathcal{Y}). \quad (7)$$

The second equality follows from the fact that $\pi_{\mathcal{X}}^{-1}(A) := \{(x, y) \in \mathcal{X} \times \mathcal{Y} : \pi_{\mathcal{X}}(x, y) \in A\}$. Of course, the marginalization for $\mathbb{P}_{\mathcal{Y}}$ is defined analogously.

Quite literally, marginalization is projection: it's a mechanism for putting a probability measure on the projected space assuming the existence of probability measure upstairs.

2.4 Conditional Probability

For a probability space $(\mathcal{X}, \mathbb{P}_{\mathcal{X}})$ you've likely seen the definition of conditional probability as $\mathbb{P}_{\mathcal{X}}(A|B) := \mathbb{P}_{\mathcal{X}}(A \cap B) / \mathbb{P}_{\mathcal{X}}(B)$ provided that the denominator is nonzero. It is easier to visualize this notion with joint probability. We continue with supposing that $(\mathcal{X} \times \mathcal{Y}, \mathbb{P}_{\mathcal{X} \times \mathcal{Y}})$ is a joint probability space. We've already defined marginal probability, so we can define the following conditional probability.

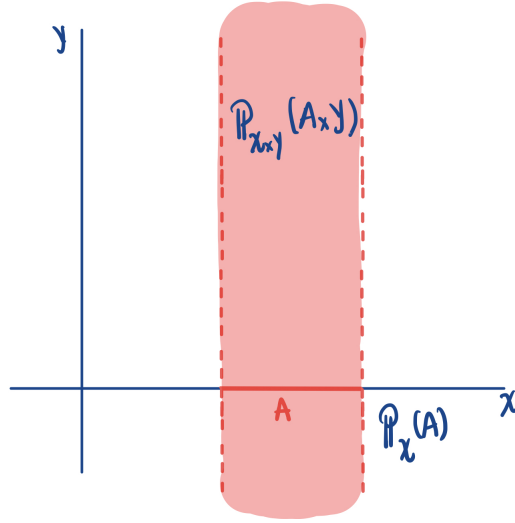


Figure 8

Definition 2.7. The conditional probability $\mathbb{P}_{\mathcal{Y}|\mathcal{X}}(B|A)$ is defined to be

$$\mathbb{P}_{\mathcal{Y}|\mathcal{X}}(B|A) := \frac{\mathbb{P}_{\mathcal{X} \times \mathcal{Y}}(A \times B)}{\mathbb{P}_{\mathcal{X}}(A)} \quad (8)$$

provided that the marginal $\mathbb{P}_{\mathcal{X}}(A) \neq 0$.

In general, you should be comfortable manipulating expressions with notation $d\mathbb{P}_{\mathcal{Y}|\mathcal{X}}$ for which we have implicit definition

$$\int_A \int_{B|A} d\mathbb{P}_{\mathcal{Y}|\mathcal{X}}(y|x) d\mathbb{P}_{\mathcal{X}}(x) := \int_{A \times B} d\mathbb{P}_{\mathcal{X} \times \mathcal{Y}}(x, y), \quad (9)$$

or for random variable $f : \mathcal{X} \times \mathcal{Y}$ we have

$$\int_{\mathcal{X}} \int_{\mathcal{Y}|\mathcal{X}} f(x, y) d\mathbb{P}_{\mathcal{Y}|\mathcal{X}}(y|x) d\mathbb{P}_{\mathcal{X}}(x) := \int_{\mathcal{X} \times \mathcal{Y}} f(x, y) d\mathbb{P}_{\mathcal{X} \times \mathcal{Y}}(x, y). \quad (10)$$

You may recognize the original definition in eq. (8) as secretly appearing here, since $d\mathbb{P}_{\mathcal{Y}|\mathcal{X}} d\mathbb{P}_{\mathcal{X}} = d\mathbb{P}_{\mathcal{X} \times \mathcal{Y}}$ may be “solved” for $d\mathbb{P}_{\mathcal{Y}|\mathcal{X}}$. But remember, in addition, that this ‘ $d\mathbb{P}$ ’ notation itself is defined as a package $\int_A d\mathbb{P} = \mathbb{P}(A)$.

Remark 2.6. There were many questions on interpreting the critter $A|B$. In pictures I drew the rectangle $A \times B$. It is fine to think of $A|B$ pictorially as $A \times B$, but it is important to concomitantly think of the ambient space in which $A|B$ lives: $A \times B$, $A|B$, $B|A$ are all “the same” rectangle, but $A \times B \subset \mathcal{X} \times \mathcal{Y}$, $A|B \subset \mathcal{X}|B$ (which you may visualize as the rectangle $\mathcal{X} \times B$ which rectangle has probability one, *not* $\mathbb{P}_{\mathcal{X} \times \mathcal{Y}}(\mathcal{X} \times B)$ which may be less than one). Similarly, $B|A$ is the rectangle $A \times B$, but instead of living in $\mathcal{X} \times \mathcal{Y}$ it lives in $\mathcal{Y}|A$ (or visually the rectangle $A \times \mathcal{Y}$ with unit probability).

In math, hand-waiving can sometimes be dangerous. On the other hand, it can sometimes allow us to think reasonably about, and operationalize our intuition of, notions whose formalism is “beyond the scope of this course,” all with the aim of performing computations and symbolic manipulations. If you are interested in more rigorously making sense of $d\mathbb{P}_{\mathcal{Y}|\mathcal{X}}$, I recommend diving into the measure theory. The level of rigor we need is: to fluidly manipulate integral expressions involving joint probability, relying on our multivariable intuition from calculus that we may decompose high dimensional integration as sequential one-dimensional integrals. What changes with probability, is that the (measure used for the) inner one-dimensional integral may depend on the outer value.

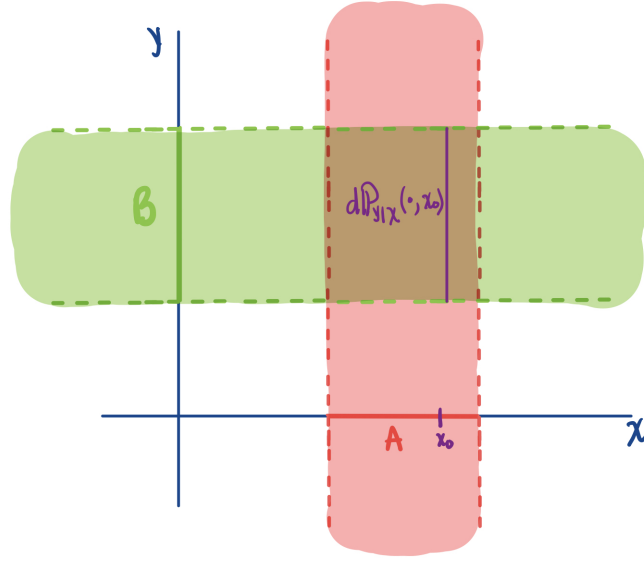


Figure 9

2.5 Independence

Forget, now, our reliance on joint probability space $(\mathcal{X} \times \mathcal{Y}, \mathbb{P}_{\mathcal{X} \times \mathcal{Y}})$, and suppose that we have two separate probability spaces $(\mathcal{X}, \mathbb{P}_{\mathcal{X}})$ and $(\mathcal{Y}, \mathbb{P}_{\mathcal{Y}})$. From these two spaces, one may reasonably ask if we can “go the other direction” and construct a (joint) probability measure $\mathbb{P}_{\mathcal{X} \times \mathcal{Y}}$ on $(\mathcal{X} \times \mathcal{Y})$. The answer is yes.

Before doing so, let us reason about the properties we would like this measure to have. The first thing we should expect is that projection $\pi_{\mathcal{X}} : \mathcal{X} \times \mathcal{Y} \rightarrow \mathcal{X}$ *preserves measure*, i.e. that

$$\mathbb{P}_{\mathcal{X} \times \mathcal{Y}}(A \times \mathcal{Y}) = \mathbb{P}_{\mathcal{X}}(\pi_{\mathcal{X}}(A \times \mathcal{Y})) = \mathbb{P}_{\mathcal{X}}(A).$$

Observe that in contrast to eq. (7), we do not define $\mathbb{P}_{\mathcal{X}}(A)$ in terms of $\mathbb{P}_{\mathcal{X} \times \mathcal{Y}}$. Instead, the definition goes the other way:

$$\mathbb{P}_{\mathcal{X} \times \mathcal{Y}}(A \times \mathcal{Y}) := \mathbb{P}_{\mathcal{X}}(A). \quad (11)$$

Of course, we would like the analogous equality to hold for $\mathcal{X} \times B$ with $\mathbb{P}_{\mathcal{Y}}(B)$.

Now this condition by itself only allows us to define sets of the form $A \times \mathcal{Y}$ or $\mathcal{X} \times B$. For general rectangle $A \times B \subset \mathcal{X} \times \mathcal{Y}$, define

$$\mathbb{P}_{\mathcal{X} \times \mathcal{Y}}(A \times B) := \mathbb{P}_{\mathcal{X}}(A) \cdot \mathbb{P}_{\mathcal{Y}}(B). \quad (12)$$

Observe, in connection with our first encounter of independence as $\mathbb{P}(A \cap B) = \mathbb{P}(A)\mathbb{P}(b)$, we can recover this relation by observing that

$$A \times B = (A \times \mathcal{Y}) \cap (\mathcal{X} \times B), \quad (13)$$

and compute

$$\begin{aligned} \mathbb{P}_{\mathcal{X} \times \mathcal{Y}}(A \times B) &= \mathbb{P}_{\mathcal{X}}(A)\mathbb{P}_{\mathcal{Y}}(B) \\ &= \mathbb{P}_{\mathcal{X} \times \mathcal{Y}}(A \times \mathcal{Y}) \cdot \mathbb{P}_{\mathcal{X} \times \mathcal{Y}}(\mathcal{X} \times B). \end{aligned} \quad (14)$$

The intuition we extract from this construction is that independence is a most natural way to construct measure on high dimensional space using measure from its lower dimensional components. In exactly the same way that we construct area from length or volume from area and length (or volume from length).

Definition 2.8. In general, we say that $(\mathcal{X}^m, \mathbb{P}_{\mathcal{X}^m})$ is *independent* if

$$\mathbb{P}_{\mathcal{X}^m} = (\mathbb{P}_{\mathcal{X}})^m$$

or more generally that $\left(\prod_{j=1}^m \mathcal{X}_j, \mathbb{P}_{\prod \mathcal{X}_j}\right)$ independent if

$$\mathbb{P}_{\prod \mathcal{X}_j} = \prod_{j=1}^m \mathbb{P}_{\mathcal{X}_j}.$$

Remark 2.7. One should verify that the (quasi-independence) conditions from the last problem of ws0 are a consequence of this definition (hint: marginalization).

3 Lecture 3

We continue discussing independence as a high-dimensional phenomenon, and introduce concentration.

We started with a quick review of conditional probability from last time, providing geometric interpretation for $A \times B$, $A|B$ and $B|A$ as “sets” in $\mathcal{X} \times \mathcal{Y}$ (remark 2.6).

3.1 Independence

Recall definition 2.8 which says that a joint probability measure $\mathbb{P}_{\mathcal{X} \times \mathcal{Y}}$ is independent if it decomposes as a product of marginals $\mathbb{P}_{\mathcal{X}} \cdot \mathbb{P}_{\mathcal{Y}}$. This extends to products of arbitrarily many factors.

To recognize independence, either of the following suffice:

1. conditional probability is independent of condition, or (what amounts to the same)
2. conditional probability is equal to the marginal

One may examine the independence question from worksheet for density

$$f(x, y) = cy^{-1/2}x\mathbb{1}_{y>0}\mathbb{1}_{x\geq 0}\mathbb{1}_{x^2+y^2\leq 1},$$

where the support of $f(x_0, \cdot)$ —i.e. the function $f_{x_0} : \mathcal{Y} \rightarrow \mathbb{R}$ defined by $f_{x_0}(y) = f(x_0, y)$ —is readily seen to differ according to the value of x_0 , namely $\left[0, \sqrt{1 - x_0^2}\right]$.

Ultimately, our purpose in going through this rigmarole is for computation. I would like you to feel comfortable doing symbolic manipulations e.g. of the form $\int_{\mathcal{X} \times \mathcal{Y}} = \int_{\mathcal{X}} \int_{\mathcal{Y}|\mathcal{X}}$. You might see this in other contexts as the “tower” property. If you are familiar with iterated and embedded expectations, then you need not refer to the geometric interpretation. Use whichever viewpoint you are most at home with when doing computations with joint probability spaces.

Example 3.1. Let us examine why we were destined to fail on the 0th programming assignment. We start with $(\mathcal{X} \times \mathcal{Y} = \mathbb{R}^k \times \{0, 1\}, \mathbb{P}_{\mathcal{X} \times \mathcal{Y}})$ independent, and suppose that we have “optimal” model $\tilde{y} : \mathcal{X} \rightarrow \mathcal{Y}$. We bracket, for the moment, what optimality here means. Computing accuracy, we have:

$$\begin{aligned} \mathbb{P}_{\mathcal{X} \times \mathcal{Y}}(\tilde{y} = y) &= \mathbb{E}(\mathbb{1}_{\tilde{y}(x)=y}) \\ &:= \int_{\mathcal{X} \times \mathcal{Y}} \mathbb{1}_{\tilde{y}(x)=y} d\mathbb{P}_{\mathcal{X} \times \mathcal{Y}}(x, y) \\ &= \int_{\mathcal{X}} \int_{\mathcal{Y}|\mathcal{X}} \mathbb{1}_{\tilde{y}(x)=y} d\mathbb{P}_{\mathcal{Y}|\mathcal{X}}(y|x) d\mathbb{P}_{\mathcal{X}}(x) \\ &= \int_{\mathcal{X}} \int_{\mathcal{Y}} \mathbb{1}_{\tilde{y}(x)=y} d\mathbb{P}_{\mathcal{Y}}(y) d\mathbb{P}_{\mathcal{X}}(x) \\ &= \int_{\mathcal{X}} \int_{\mathcal{Y}} \mathbb{1}_{\tilde{y}(x)=0}(1-p) + \mathbb{1}_{\tilde{y}(x)=1}p d\mathbb{P}_{\mathcal{X}}(x) \\ &= \int_{\mathcal{X}} \max\{p, 1-p\} d\mathbb{P}_{\mathcal{X}}(x) \\ &= \max\{p, 1-p\} \int_{\mathcal{X}} d\mathbb{P}_{\mathcal{X}}(x) \\ &= \max\{p, 1-p\} \end{aligned} \tag{15}$$

In the first equation, we recall that probability may be written as expectation of an indicator. This is useful, it allows us to do this computation. The second line is the definition of expectation (definition 2.2). The third line is our implicit definition of conditional probability eq. (10), the geometry of which you should think of as iterated integration in multivariable calculus. The fourth line is independence (condition #2 above). The fifth is expectation w.r.t. $\mathbb{P}_{\mathcal{Y}}$ of $\mathbb{1}_{\tilde{y}(x)=y}$. The sixth line is optimality of \tilde{y} ; we are trying to optimize accuracy, and as we’ve written it we get to a point where we can optimize *pointwise* (in x). In the seventh line, we pull out $\max\{p, 1-p\}$ since p is independent of x and in the final line we recall that $\mathbb{P}_{\mathcal{X}}(\mathcal{X}) = 1$.

If you look at the scores for your model, you should have had something pretty close to p for almost all x .

Next we interpret independence for data.

3.2 Data

Referring to the standard diagram (6), we hope and expect that most instances where supervised machine learning comes in to play, the measure $\mathbb{P}_{\mathcal{X} \times \mathcal{Y}}$ on $\mathcal{X} \times \mathcal{Y}$ is not independent. Sometimes, as in the 0th programming assignment, you'll run into an evil data set, but such are not the norm. Instead, independence is a phenomenon we'd like to associate with *data*.

Definition 3.1. We say that a sequence of (labeled) points $\mathcal{S} = ((x_1, y_1), \dots, (x_m, y_m)) \sim_{\text{iid}} \mathbb{P}_{\mathcal{X} \times \mathcal{Y}}$ is (sampled) *independent and identically distributed* if $\mathcal{S} \in (\mathcal{X} \times \mathcal{Y})^m$ is a point in joint probability space $(\mathcal{X} \times \mathcal{Y})^m$ with independent measure $\mathbb{P}_{(\mathcal{X} \times \mathcal{Y})^m} = (\mathbb{P}_{\mathcal{X} \times \mathcal{Y}})^m$.

The following picture is for the intuition in the high-dimensional space.

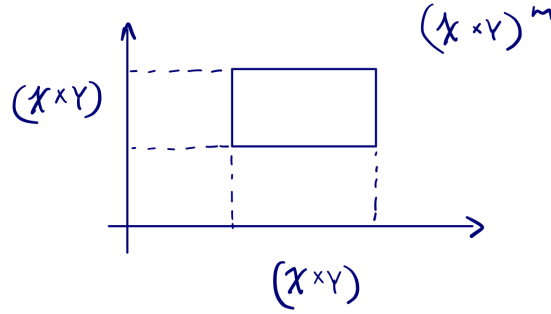


Figure 10

Remark 3.1. While seemingly pedantic, the nuance is important: a data set in ML is properly speaking a sequence. While we might not care about order, we absolutely do care about repetition.

Talk of sampling points suggests talk of randomness. You are allowed to select any point $x \in \mathcal{X}$ that you fancy. How you sample should in some sense be affiliated with your measure $\mathbb{P}_{\mathcal{X}}$. One way of handwaiving this affiliation is to say that given any set (event) $A \subset \mathcal{X}$, the probability that the point $x \in \mathcal{X}$ you picked happens to (also) be in A is $\mathbb{P}_{\mathcal{X}}(A)$. Properly speaking, we are saying in English that $\mathbb{E}(\mathbb{1}_{x \in A}) = \mathbb{P}_{\mathcal{X}}(A)$. Since we generally cannot evaluate $\mathbb{P}_{\mathcal{X}}$ directly, we would like a more data-centric way to talk about probability.

3.2.1 Law of Large Numbers (LLN)

Theorem 3.1. Let $(\mathcal{X}, \mathbb{P}_{\mathcal{X}})$ be a probability space and for $m \in \mathbb{N}$, $(\mathcal{X}^m, \mathbb{P}_{\mathcal{X}^m})$ independent (meaning: $\mathbb{P}_{\mathcal{X}^m} = \mathbb{P}_{\mathcal{X}}^m$), and suppose that both $\mu_{\mathcal{X}} := \mathbb{E}(x) = \int_{\mathcal{X}} x d\mathbb{P}_{\mathcal{X}}(x) < \infty$ and $\sigma_{\mathcal{X}}^2 := \mathbb{E}((x - \mu_{\mathcal{X}})^2) < \infty$. Define empirical means $s_m : \mathcal{X}^m \rightarrow \mathbb{R}$ by $(x_1, \dots, x_m) \mapsto \frac{1}{m}(x_1 + \dots + x_m)$. Then for $\varepsilon > 0$,

$$\lim_{m \rightarrow \infty} \mathbb{P}_{\mathcal{X}^m}(|s_m - \mu_{\mathcal{X}}| > \varepsilon) = 0.$$

Proof. WLOG suppose $\mu_{\mathcal{X}} = 0$. By Chebyshev

$$\begin{aligned} \mathbb{P}_{\mathcal{X}^m}(|s_m - \mu_{\mathcal{X}}| > \varepsilon) &\leq \frac{\mathbb{E}((s_m - \mu_{\mathcal{X}})(s_m - \mu_{\mathcal{X}}))}{\varepsilon^2} \\ &= \frac{\mathbb{E}\left(\frac{1}{m^2} \sum_{i=1}^m \sum_{j=1}^m x_i x_j - \frac{2}{m} \sum_{j=1}^m x_j \mu_{\mathcal{X}} + \mu_{\mathcal{X}}^2\right)}{\varepsilon^2} \\ &= \frac{\frac{1}{m^2} \sum_{i,j=1}^m \mathbb{E}(x_i x_j) - \frac{2}{m} \mu_{\mathcal{X}} \sum_{j=1}^m \mathbb{E}(x_j) + \mu_{\mathcal{X}}^2}{\varepsilon^2} = \frac{\frac{1}{m^2} \sum_{i \neq j} \mathbb{E}(x_i) \mathbb{E}(x_j) + \frac{1}{m^2} m \mathbb{E}(x^2)}{\varepsilon^2} \xrightarrow{m \rightarrow \infty} 0 \end{aligned}$$

□

Remark 3.2. To recall, the statement of LLN measures the set of tuples $(x_1, \dots, x_m) \in \mathcal{X}^m$ that are at least ε -far from $s_m^{-1}(\mu_{\mathcal{X}})$.

We can use the Law of Large numbers concretely with data: to say that $(x_1, \dots, x_n) \sim_{\text{iid}} \mathbb{P}_{\mathcal{X}}$ (for the moment, we're not talking about labeled data) means, among other things, that we can expect the law of large numbers to hold for probabilities: i.e. for $A \subset \mathcal{X}$,

$$\frac{1}{m} \sum_{j=1}^m \mathbb{1}_{x_j \in A} \xrightarrow{m \rightarrow \infty} \mathbb{E}(\mathbb{1}_{x \in A}) = \mathbb{P}_{\mathcal{X}}(A).$$

Note that convergence in this expression is “in probability,” i.e. fix $\varepsilon > 0$. Then for any $\delta > 0$, one can specify $M_\delta > 0$ so that

$$\mathbb{P}_{\mathcal{X}^m} \left(\left| \frac{1}{m} \sum_{j=1}^m \mathbb{1}_{x_j \in A} - \mathbb{P}_{\mathcal{X}}(A) \right| > \varepsilon \right) < \delta$$

whenever $m > M_\delta$.

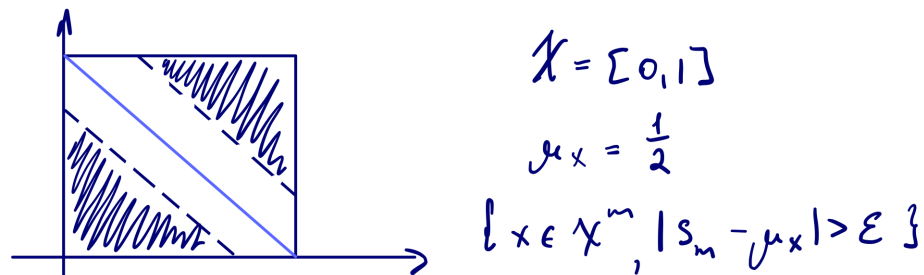


Figure 11

Remark 3.3. I didn't say this in class, but to check in simulation, generate a data sequence $(x_1, \dots, x_m) \sim_{\text{iid}} \mathbb{P}_{\mathcal{X}}$, and check/evaluate

$$\mathbb{1}_{|s_m(x_1, \dots, x_m) - \mathbb{P}_{\mathcal{X}}(A)| > \varepsilon}.$$

The answer will be zero or one. Most of the time it should be zero. If you do this many times, and take the empirical average of *these* results, the answer should be less than δ .

3.3 Concentration Inequalities

Concentration inequalities are sort of a heartbeat in ML. They are omnipresent and make a lot of important results work.

Proposition 3.1 (Markov). $(\mathcal{X} \subset \mathbb{R}^{\geq 0}, \mathbb{P}_{\mathcal{X}})$ a non-negative probability space, then

$$\mathbb{P}_{\mathcal{X}}(x \geq t) \leq \mathbb{E}(x)/t.$$

Proof. Compute: $\mathbb{E}(x) := \int_{\mathcal{X}} x d\mathbb{P}_{\mathcal{X}}(x) = \int_{[0, t)} x d\mathbb{P}_{\mathcal{X}}(x) + \int_{[t, \infty)} x d\mathbb{P}_{\mathcal{X}}(x)$, the last equality by linearity of integration. Since $x \geq 0$, this expression is bounded above by

$$\int_{[t, \infty)} x d\mathbb{P}_{\mathcal{X}}(x) \leq \int_{[t, \infty)} t d\mathbb{P}_{\mathcal{X}}(x),$$

where the latter follows from the fact that $x \geq t$ on $x \in [t, \infty)$. Taking t out and simplifying the integral as probability proves the inequality. \square

Proposition 3.2 (Chebyshev). $(\mathcal{X} \subset \mathbb{R}, \mathbb{P}_{\mathcal{X}})$ a probability space with finite variance $\sigma_{\mathcal{X}}^2$. Then for $\varepsilon > 0$,

$$\mathbb{P}_{\mathcal{X}}(|x - \mu_{\mathcal{X}}| > \varepsilon) \leq \sigma_{\mathcal{X}}^2 / \varepsilon^2.$$

Proof. Markov applied to $\{|x - \mu_{\mathcal{X}}| > t\} = \{(x - \mu_{\mathcal{X}})^2 > t^2\}$ \square

Proposition 3.3 (Hoeffding). $(\mathcal{X} \subset [0, 1], \mathbb{P}_{\mathcal{X}})$ a probability space, then

$$\mathbb{P}_{\mathcal{X}^m} \left(\left| \frac{1}{m} \sum_{j=1}^m x_j - \mu_{\mathcal{X}} \right| > \varepsilon \right) \leq 2e^{-2m\varepsilon^2}.$$

Hoeffding may appear a bit bizarre at first. What you should try to latch onto is that the right hand side looks like something you're already familiar with, the probability density function for normal distribution. While Chebyshev gives a quasi-distributionless way to get a handle on tail probabilities, Hoeffding is particularly nice because it gives a bound as a quadratically decreasing exponential. Moreover, it generalizes when $\mathcal{X} = [a, b]$.

Theorem 3.2 (Glivenko-Cantelli). $(\mathcal{X} = \mathbb{R}, \mathbb{P}_{\mathcal{X}})$ a probability space, for $t \in \mathbb{R}$, define $F(t) := \mathbb{P}_{\mathcal{X}}(x \leq t)$ the cdf and $F_m(t) : \mathcal{X}^m \rightarrow \mathbb{R}$ by

$$(x_1, \dots, x_m) \mapsto \frac{1}{m} \sum_{j=1}^m \mathbb{1}_{x_j \leq t}$$

the *empirical cdf*. Then

$$\mathbb{P}_{\mathcal{X}^m} \left(\sup_{t \in \mathbb{R}} |F_m(t) - F(t)| > \varepsilon \right) \leq 8(m+1)e^{-\frac{m\varepsilon^2}{32}}. \quad (16)$$

The statement of this theorem is rather striking: no matter what $\mathbb{P}_{\mathcal{X}}$, you can stipulate conditions for ensuring precision specification (ε) is satisfies with arbitrarily high confidence ($1 - \delta$) provided m is "sufficiently large." And again... totally independent of $\mathbb{P}_{\mathcal{X}}$. Sometimes, one must pause to marvel at the beauty of mathematics. The picture illustrating the idea of the theorem is below

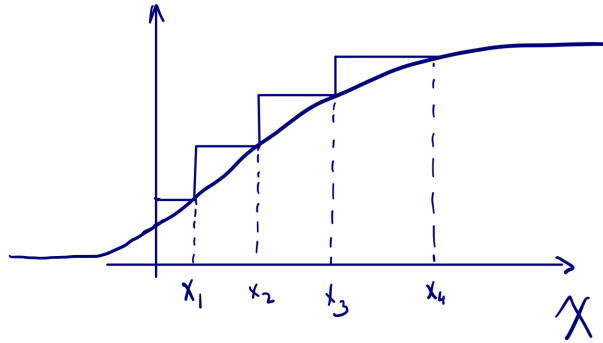


Figure 12

3.4 Intuition in High Dimension

Probability does weird things in high dimension. I made a tenuous sounding claim that a high dimensional gaussian has density concentrated at the origin (tracks our low-dimensional intuition) but probability concentrated *away* from it (what?!).

Let's say this a bit more formally. Suppose that $(\mathcal{X} = \mathbb{R}, \mathbb{P}_{\mathcal{X}})$ is probability space with $\mathbb{P}_{\mathcal{X}}$ normal, zero mean, unit variance, and $(\mathcal{X}^m, \mathbb{P}_{\mathcal{X}^m})$ independent. For $\varepsilon > 0$,

$$\lim_{m \rightarrow \infty} \mathbb{P}_{\mathcal{X}^m} \left(\left| \|\mathbf{x}\|^2 - m \right| > \varepsilon \right) = \lim_{m \rightarrow \infty} \mathbb{P}_{\mathcal{X}^m} \left(\left\{ \mathbf{x} \in \mathcal{X}^m : \left| \|\mathbf{x}\|^2 - m \right| > \varepsilon \right\} \right) = 0.$$

This says that in high dimension, a gaussian concentrates around the sphere of radius \sqrt{m} . In particular, the *solid* sphere of radius (strictly) less than \sqrt{m} is practically empty!

Here's a concrete computation which sortof illustrates the point.

Example 3.2. Consider $(\mathcal{X} = [0, 1], \mathbb{P}_{\mathcal{X}})$ with uniform measure $\mathbb{P}_{\mathcal{X}}([a, b]) = (b - a)\mathbb{1}_{0 \leq a \leq b \leq 1}$. Suppose that $A \subset \mathcal{X}$ with $1 > \mathbb{P}_{\mathcal{X}}(A) \geq 1 - \varepsilon$. The "hypercube" $A^m \subset \mathcal{X}^m$ in high dimension has measure

$$\mathbb{P}_{\mathcal{X}^m}(A^m) = (\mathbb{P}_{\mathcal{X}}(A))^m = (1 - \varepsilon)^m,$$

the first equality by independence. For $\varepsilon > 0$, $\lim_{m \rightarrow \infty} (1 - \varepsilon)^m = 0$.

4 Lecture 4

In this lecture we prove Glivenko-Cantelli inequality. Given data samples $(x_1, \dots, x_m) \in \mathcal{X}^m$, we want to approximate the cdf $F(t) := \mathbb{P}_{\mathcal{X}}(\mathcal{X} \leq t)$ by its empirical average defined by

$$F_m(t) := \frac{1}{m} \sum_{j=1}^m \mathbb{1}_{x_j \leq t} \quad (17)$$

Properly speaking $F_m : \mathcal{X}^m \rightarrow [0, 1]^{\mathcal{X}}$ defined by mapping $\bar{x} = (x_1, \dots, x_m) \mapsto F_m^{\bar{x}} : \mathcal{X} \rightarrow [0, 1]$, the latter of which is defined by (17). Typically we drop notational dependence on \bar{x} .

Glivenko-Cantelli tells us that $\sup_{t \in \mathbb{R}} |F_m(t) - F(t)| \rightarrow 0$ with high probability. The rigorous statement:

Theorem 4.1 (Glivenko-Cantelli). Let $(\mathcal{X}, \mathbb{P}_{\mathcal{X}})$ be a random variable. Then

$$\mathbb{P}_{\mathcal{X}^m} \left(\sup_{t \in \mathbb{R}} |F_m(t) - F(t)| > \epsilon \right) \leq 8(m+1)e^{-\frac{m \cdot \epsilon^2}{32}} \quad (18)$$

for m sufficiently large.

The proof will be presented in a sequence of steps, which we initially enumerate and subsequently prove.

Outline of the proof

There are five main steps. For the first three, we fix $t \in \mathbb{R}$.

1. First we bound the (probability of) separation between (true) cdf $F(t)$ and empirical cdf $F_m(t)$ by a probability of separation between two *separate* empirical cdfs $F_m(t)$ and $F'_m(t)$. Notice that the probability statement in eq. (18) is w.r.t. the data sample $(x_1, \dots, x_m) \in \mathcal{X}^m$; in this step we consider (/transport our attention to) data samples $(x_1, \dots, x_m, x'_1, \dots, x'_m) \in \mathcal{X}^{2m}$. **This step is crucial, as we'll see in step , for turning $\sup_{t \in \mathbb{R}} = \bigcup_{t \in \mathbb{R}}$ into a finite union**

2. Having introduced a second (in-distribution identical) data sample $(x'_1, \dots, x'_m) \in \mathcal{X}^m$, we

observe that separation condition $\left| \frac{1}{m} \sum_{j=1}^m \mathbb{1}_{x_j \leq t} - \mathbb{1}_{x'_j \leq t} \right| > \epsilon'$ is symmetric in $x, x' \in \mathcal{X}^m$ ($\epsilon' > 0$ is not necessarily equal to ϵ). We thus, without affecting the probability statement, introduce (symmetric) *Rademacher* variables $s \in \mathcal{S} := \{-1, 1\}$ for which $\mathbb{P}_{\mathcal{S}}(s = 1) = \frac{1}{2}$.

With the Rademacher variables, we consider separation condition

$$\left| \left(\frac{1}{m} \sum_{j=1}^m s_j (\mathbb{1}_{x_j \leq t} - \mathbb{1}_{x'_j \leq t}) \right) \right| > \epsilon'$$

and using triangle inequality isolate each sample into its own separation condition

$$\left| \frac{1}{m} \sum_{j=1}^m s_j \mathbb{1}_{x_j \leq t} \right| > \epsilon'',$$

for some other $\epsilon'' > 0$. We will see that this step collapses the probability calculation on \mathcal{X}^{2m} to one on \mathcal{X}^m .

3. Next we apply law of total expectation, condition on \mathcal{X}^m , fixing sample $(x_1, \dots, x_m) \in \mathcal{X}^m$, and compute $\int_{\mathcal{X}^m \times \mathcal{S}^m} = \int_{\mathcal{X}^m} \int_{\mathcal{S}^m | \mathcal{X}^m}$. The inner expectation is

$$\int_{\mathcal{S}^m | \mathcal{X}^m} \mathbb{1}_{\left| \frac{1}{m} \sum_{j=1}^m s_j \mathbb{1}_{x_j \leq t} \right| > \epsilon''} d\mathbb{P}_{\mathcal{S}^m | \mathcal{X}^m}(\bar{s} | \bar{x}) = \mathbb{P}_{\mathcal{S}^m | \mathcal{X}^m} \left(\left| \frac{1}{m} \sum_{j=1}^m s_j \mathbb{1}_{x_j \leq t} \right| > \epsilon'' \mid x_1, \dots, x_m \right).$$

4. Supremizing over $t \in \mathbb{R}$, we observe that

$$\mathbb{P}_{\mathcal{S}^m | \mathcal{X}^m} \left(\sup_{t \in \mathbb{R}} \left| \frac{1}{m} \sum_{j=1}^m s_j \mathbb{1}_{x_j \leq t} \right| > \varepsilon'' |\bar{x}| \right) = \mathbb{P}_{\mathcal{S}^m | \mathcal{X}^m} \left(\bigcup_{t \in \{t_0, \dots, t_m\}} \left| \frac{1}{m} \sum_{j=1}^m s_j \mathbb{1}_{x_j \leq t} \right| > \varepsilon'' |\bar{x}| \right).$$

Generally, $\mathbb{P}(\sup_{t \in \mathbb{R}}) = \mathbb{P}\left(\bigcup_{t \in \mathbb{R}}\right)$; in this case, however, conditioned (i.e. *fixing*) on x_1, \dots, x_m , the values of $\sum_{j=1}^m \mathbb{1}_{x_j \leq t}$ change only at $t = x_1, \dots, x_m$. Therefore, this supremum turns out to be a finite union.

5. Finally we apply Hoeffding's inequality to the inner integral to get a bound in terms of decreasing exponential and then apply the outer integral $\int_{\mathcal{X}^m}$, which recovers a probability we seek to bound.

You should go through the details of the argument at least once, get a feel for the techniques used, and make yourself at home with the outline. The first step is an obnoxiously detailed calculation, but step 2, e.g., contains a kernel of a concept we'll see again when we cover PAC learnability and complexity of hypothesis classes.

Lemma 4.1.

$$\mathbb{P}_{\mathcal{X}^m}(|F_m(t) - F(t)| > \epsilon) \leq 2\mathbb{P}_{\mathcal{X}^{2m}}(|F'_m(t) - F_m(t)| > \frac{\epsilon}{2})$$

where

$$F'_m = \frac{1}{m} \sum_{j=1}^m \mathbb{1}_{x'_j \leq t}$$

for sample $\bar{x}' = (x'_1, \dots, x'_m) \in \mathcal{X}^m$.

Proof. Fix $t \in \mathbb{R}$, $\epsilon > 0$, and suppose that $\epsilon < |F(t) - F_m(t)|$. Inserting $0 = -F'_m(t) + F'_m(t)$ into the right-hand side and applying the triangle inequality gives

$$\epsilon < |F(t) - F_m(t)| = |F(t) - F'_m(t) + F'_m(t) - F_m(t)| \leq |F(t) - F'_m(t)| + |F_m(t) - F'_m(t)|. \quad (19)$$

Supposing (19), $|F(t) - F'_m(t)| \leq \frac{\epsilon}{2}$ implies that $|F_m(t) - F'_m(t)| > \frac{\epsilon}{2}$ (why?). Translating conditions into an indicator function, this implication is equivalent to inequality

$$\mathbb{1}_{|F(t) - F_m(t)| > \epsilon} \cdot \mathbb{1}_{|F(t) - F'_m(t)| \leq \frac{\epsilon}{2}} \leq \mathbb{1}_{|F_m(t) - F'_m(t)| > \frac{\epsilon}{2}}.$$

This is true for all $x_1, \dots, x_m, x'_1, \dots, x'_m$, so, taking the expectation of both sides with respect to $d\mathbb{P}_{\mathcal{X}^{2m}}(x_1, \dots, x_m, x'_1, \dots, x'_m)$ preserves the inequality:

$$\mathbb{E}(\mathbb{1}_{\epsilon < |F(t) - F_m(t)|} \cdot \mathbb{1}_{\frac{\epsilon}{2} > |F(t) - F'_m(t)|}) \leq \mathbb{E}(\mathbb{1}_{\frac{\epsilon}{2} < |F_m(t) - F'_m(t)|}).$$

Recalling that expectation of an indicator is probability, we obtain

$$\mathbb{P}_{\mathcal{X}^m}(|F(t) - F_m(t)| > \epsilon) \mathbb{P}_{\mathcal{X}^m}(|F(t) - F'_m(t)| \leq \frac{\epsilon}{2}) \leq \mathbb{P}_{\mathcal{X}^{2m}}(|F_m(t) - F'_m(t)| > \frac{\epsilon}{2}).$$

Because Bernoulli random variable $\mathbb{1}_{x \leq t}$ has expectation $F(t)$ with finite variance, we can apply Chebyshev's inequality to bound the second term in the above product

$$\mathbb{P}_{\mathcal{X}^m}(|F(t) - F'_m(t)| \leq \frac{\epsilon}{2}) = 1 - \mathbb{P}(|F(t) - F'_m(t)| \geq \frac{\epsilon}{2}) \geq \frac{\text{Var}(\mathbb{1}_{x \leq t})}{m(\epsilon/2)^2} \geq \frac{1}{2}$$

for m large enough. □

$$\mathbb{P}(|F(t) - F_m(t)| < \frac{\epsilon}{2}) \geq 1 - \frac{\text{Var}(F_m(t))}{(\epsilon/2)^2} \geq 1 - \frac{1/4m}{\epsilon^2/4} = 1 - \frac{1}{m\epsilon^2}$$

$$\mathbb{P}(|F(t) - F_m(t)| < \frac{\epsilon}{2}) \geq 1 - \frac{1}{m\epsilon^2} \geq 1 - \frac{1}{2} = \frac{1}{2}$$

One more lemma:

Lemma 4.2. (Symmetrization) Let $(\mathcal{X} \subset \mathbb{R}, \mathbb{P}_{\mathcal{X}})$ be a random variable, and $(\mathcal{S} = \{-1, 1\}, \mathbb{P}_{\mathcal{S}}(1) = \frac{1}{2})$ symmetric. Then random variable $p : \mathcal{X} \times \mathcal{S} \rightarrow \mathcal{X}$ defined by mapping $(x, s) \mapsto s \cdot x$ is symmetric mean zero, i.e.

$$\mathbb{P}_{\mathcal{X}}(sx > t) = \mathbb{P}_{\mathcal{X} \times \mathcal{S}}(sx < -t) \quad (20)$$

When $(\mathcal{X}, \mathbb{P}_{\mathcal{X}})$ is symmetric, both of these tail probabilities equal the original tail probability $\mathbb{P}_{\mathcal{X}}(x > t)$.

Proof. Let's consider the embeddings $\iota_+ : \mathcal{X} \hookrightarrow \mathcal{X} \times \mathcal{S}$, $\iota_- : \mathcal{X} \hookrightarrow \mathcal{X} \times \mathcal{S}$ defined by

$$\iota_+(x) := (x, 1) \text{ and } \iota_-(x) := (x, -1)$$

We want to induce via \mathcal{X} a measure on $\mathcal{X} \times \mathcal{S}$. If we naively try to define $\mathbb{P}_{\mathcal{X} \times \mathcal{S}}([a, b] \times \{s\}) := \mathbb{P}_{\mathcal{X}}(\iota_+^{-1}([a, b] \times \{s\}))$, then for $s = -1$ we will have

$$\mathbb{P}_{\mathcal{X} \times \mathcal{S}}([a, b] \times \{-1\}) = \mathbb{P}_{\mathcal{X}}(\iota_+^{-1}([a, b] \times \{-1\})) = \mathbb{P}_{\mathcal{X}}(\emptyset) = 0,$$

Instead, ι_+, ι_- will induce a conditional measure $\mathbb{P}_{\mathcal{X}|\mathcal{S}}$ on $\mathcal{X}|s$, and then we extend to $\mathcal{X} \times \mathcal{S}$ by the law of total probability. Specifically,

$$\mathbb{P}_{\mathcal{X} \times \mathcal{S}}([a, b] \times \{s\}) := \int_{\mathcal{S}} \mathbb{P}_{\mathcal{X}|\mathcal{S}}([a, b]|s) d\mathbb{P}_{\mathcal{S}}(s) = \frac{1}{2} \mathbb{P}_{\mathcal{X}|\mathcal{S}}([a, b]|1) + \frac{1}{2} \mathbb{P}_{\mathcal{X}|\mathcal{S}}([a, b]|-1)$$

In particular,

$$\mathbb{P}_{\mathcal{X}|\mathcal{S}}(x > t|s = 1) := \mathbb{P}_{\mathcal{X}}(x > t) =: \mathbb{P}_{\mathcal{X}|\mathcal{S}}(x > t|s = -1) \quad (21)$$

then

$$\begin{aligned} \mathbb{P}_{\mathcal{X}}(p > t) &:= \mathbb{P}_{\mathcal{X} \times \mathcal{S}}(x \cdot s > t) \\ &= \sum_{s \in \mathcal{S}} \mathbb{P}_{\mathcal{S}}(s) \mathbb{P}_{\mathcal{X}|\mathcal{S}}(x > t|s) \\ &= \frac{1}{2} \mathbb{P}_{\mathcal{X}|\mathcal{S}}(x > t|s = 1) + \frac{1}{2} \mathbb{P}_{\mathcal{X}|\mathcal{S}}(-x > t|s = -1) \\ &= \frac{1}{2} \mathbb{P}_{\mathcal{X}}(x > t) + \frac{1}{2} \mathbb{P}_{\mathcal{X}}(x < -t) \\ &= \frac{1}{2} \mathbb{P}_{\mathcal{X}|\mathcal{S}}(x > t|s = -1) + \frac{1}{2} \mathbb{P}_{\mathcal{X}|\mathcal{S}}(-x > t|s = 1) \\ &= \mathbb{P}_{\mathcal{X}}(p < -t) \end{aligned} \quad (22) \quad (23)$$

In the second step (22), we apply the law of total probability, and as a result of (21), we arrive at (23). \square

We may apply this result to symmetric random variable $\mathbb{1}_{x \leq t} - \mathbb{1}_{x' \leq t}$ on \mathcal{X}^2 with relevant modification:

$$\mathbb{P}_{\mathcal{X}^2}(|\mathbb{1}_{x \leq t} - \mathbb{1}_{x' \leq t}| > \epsilon) = \mathbb{P}_{\mathcal{X}^2 \times \mathcal{S}}(|s(\mathbb{1}_{x \leq t} - \mathbb{1}_{x' \leq t})| > \epsilon), \quad (24)$$

and extend application *mutatis mutandis* to $\mathcal{X}^{2m} \times \mathcal{S}^m$ for separation of empirical cdfs.

Proof of Glivenko-Cantelli. In lemma 4.1, we showed that

$$\mathbb{P}_{\mathcal{X}^m}(|F(t) - F_m(t)| > \varepsilon) \leq 2\mathbb{P}_{\mathcal{X}^{2m}}(|F_m(t) - F'_m(t)| > \frac{\varepsilon}{2}).$$

Rewriting the right hand side, we wish to bound the probability of event

$$\left| \frac{1}{m} \sum_{j=1}^m \mathbb{1}_{x_j \leq t} - \mathbb{1}_{x'_j \leq t} \right| > \frac{\varepsilon}{2}.$$

Applying (the appropriate generalization of) eq. (24), we observe that

$$\mathbb{P}_{\mathcal{X}^{2m}}(|F_m(t) - F'_m(t)| > \frac{\varepsilon}{2}) = \mathbb{P}_{\mathcal{X}^{2m} \times \mathcal{S}^m} \left(\left| \frac{1}{m} \sum_{j=1}^m s_j (\mathbb{1}_{x_j \leq t} - \mathbb{1}_{x'_j \leq t}) \right| > \frac{\varepsilon}{2} \right).$$

Applying the triangle inequality to the term inside absolute value, we bound

$$\left| \frac{1}{m} \sum_{j=1}^m s_j (\mathbb{1}_{x_j \leq t} - \mathbb{1}_{x'_j \leq t}) \right| \leq \left| \frac{1}{m} \sum_{j=1}^m s_j (\mathbb{1}_{x_j \leq t}) \right| + \left| \frac{1}{m} \sum_{j=1}^m s_j (\mathbb{1}_{x'_j \leq t}) \right|.$$

Therefore, a bound of

$$\left| \frac{1}{m} \sum_{j=1}^m s_j (\mathbb{1}_{x_j \leq t} - \mathbb{1}_{x'_j \leq t}) \right| > \frac{\varepsilon}{2}$$

on the left hand side implies a bound of

$$\left| \frac{1}{m} \sum_{j=1}^m s_j (\mathbb{1}_{x_j \leq t}) \right| > \frac{\varepsilon}{4} \text{ or } \left| \frac{1}{m} \sum_{j=1}^m s_j (\mathbb{1}_{x'_j \leq t}) \right| > \frac{\varepsilon}{4}$$

on the right hand side.

Implication of conditions translates to inclusion of sets (events) which translates to a bound on probability:

$$\begin{aligned} & \mathbb{P}_{\mathcal{X}^{2m} \times \mathcal{S}^m} \left(\left| \frac{1}{m} \sum_{j=1}^m s_j (\mathbb{1}_{x_j \leq t} - \mathbb{1}_{x'_j \leq t}) \right| > \frac{\varepsilon}{2} \right) \\ & \leq \mathbb{P}_{\mathcal{X}^{2m} \times \mathcal{S}^m} \left(\left\{ \left| \frac{1}{m} \sum_{j=1}^m s_j (\mathbb{1}_{x_j \leq t}) \right| > \frac{\varepsilon}{4} \right\} \cup \left\{ \left| \frac{1}{m} \sum_{j=1}^m s_j (\mathbb{1}_{x'_j \leq t}) \right| > \frac{\varepsilon}{4} \right\} \right) \\ & \leq 2\mathbb{P}_{\mathcal{X}^{2m} \times \mathcal{S}^m} \left(\left| \frac{1}{m} \sum_{j=1}^m s_j (\mathbb{1}_{x_j \leq t}) \right| > \frac{\varepsilon}{4} \right) \\ & = 2\mathbb{P}_{\mathcal{X}^m \times \mathcal{S}^m} \left(\left| \frac{1}{m} \sum_{j=1}^m s_j (\mathbb{1}_{x_j \leq t}) \right| > \frac{\varepsilon}{4} \right). \end{aligned} \tag{25}$$

The third line follows from the union bound, together with the condition in x_j s defining event in \mathcal{X}^{2m} is identical—bracketing factors—to condition in x'_j s defining the analogous event, and therefore their measures are the same. The final line follows by marginalizing out the residual (extra) factors of \mathcal{X}^m .

Next, we decompose the last probability on right hand side of eq. (25) using law of total probability ($\mathbb{P}_{\mathcal{X} \times \mathcal{S}} = \mathbb{P}_{\mathcal{X}} \mathbb{P}_{\mathcal{S}|\mathcal{X}}$)/expectation $\left(\int_{\mathcal{X} \times \mathcal{S}} = \int_{\mathcal{X}} \int_{\mathcal{S}|\mathcal{X}} \right)$:

$$\mathbb{P}_{\mathcal{X}^m \times \mathcal{S}^m} \left(\left| \frac{1}{m} \sum_{j=1}^m s_j (\mathbb{1}_{x_j \leq t}) \right| > \frac{\varepsilon}{4} \right) = \int_{\mathcal{X}^m} \mathbb{P}_{\mathcal{S}^m|\mathcal{X}^m} \left(\left| \frac{1}{m} \sum_{j=1}^m s_j (\mathbb{1}_{x_j \leq t}) \right| > \frac{\varepsilon}{4} \middle| \bar{x} \right) d\mathbb{P}_{\mathcal{X}^m}(\bar{x}). \tag{26}$$

Up to this point $t \in \mathbb{R}$ has been arbitrary, fixed. But the statement of Glivenko-Cantelli considers supremum over $t \in \mathbb{R}$. The bound, therefore, over original event $\sup_{t \in \mathbb{R}} |F(t) - F_m(t)| > \varepsilon$ is in terms of

$$\int_{\mathcal{X}^m} \mathbb{P}_{\mathcal{S}^m | \mathcal{X}^m} \left(\sup_{t \in \mathbb{R}} \left| \frac{1}{m} \sum_{j=1}^m s_j (\mathbb{1}_{x_j \leq t}) \right| > \frac{\varepsilon}{4} |\bar{x}| \right) d\mathbb{P}_{\mathcal{X}^m}(\bar{x}).$$

Isolating the inner probability

$$\mathbb{P}_{\mathcal{S}^m | \mathcal{X}^m} \left(\sup_{t \in \mathbb{R}} \left| \frac{1}{m} \sum_{j=1}^m s_j (\mathbb{1}_{x_j \leq t}) \right| > \frac{\varepsilon}{4} |\bar{x}| \right),$$

we note that \bar{x} is fixed in the conditional, so the sequence $(\mathbb{1}_{x_1 \leq t}, \dots, \mathbb{1}_{x_m \leq t}) \in \{0, 1\}^m$ takes at most $m + 1$ values as t ranges over \mathbb{R} . Therefore the event (in $\mathcal{S}^m | \mathcal{X}^m$)

$$\left\{ \left| \sup_{t \in \mathbb{R}} \frac{1}{m} \sum_{j=1}^m s_j (\mathbb{1}_{x_j \leq t}) \right| > \frac{\varepsilon}{4} |\bar{x}| \right\} \subseteq \bigcup_{t \in \{t_0, \dots, t_m\}} \left\{ \left| \frac{1}{m} \sum_{j=1}^m s_j (\mathbb{1}_{x_j \leq t}) \right| |\bar{x}| > \frac{\varepsilon}{4} \right\},$$

to which we may apply the union bound in probability:⁴

$$\begin{aligned} \mathbb{P}_{\mathcal{S}^m | \mathcal{X}^m} \left(\left| \sup_{t \in \mathbb{R}} \frac{1}{m} \sum_{j=1}^m s_j (\mathbb{1}_{x_j \leq t}) \right| > \frac{\varepsilon}{4} |\bar{x}| \right) &\leq \sum_{t \in \{t_0, \dots, t_m\}} \mathbb{P}_{\mathcal{S}^m | \mathcal{X}^m} \left(\left| \frac{1}{m} \sum_{j=1}^m s_j (\mathbb{1}_{x_j \leq t}) \right| > \frac{\varepsilon}{4} |\bar{x}| \right) \\ &= (m + 1) \mathbb{P}_{\mathcal{S}^m | \mathcal{X}^m} \left(\left| \frac{1}{m} \sum_{j=1}^m s_j (\mathbb{1}_{x_j \leq t}) \right| > \frac{\varepsilon}{4} |\bar{x}| \right) \\ &\leq 2(m + 1) e^{-\frac{m(\varepsilon/4)^2}{2}} \\ &= 2(m + 1) e^{-\frac{m\varepsilon^2}{32}} \end{aligned}$$

In the third line, we cite Hoeffding's bound. When we substitute this back and calculate the outer integral in equation 26, and conclude, piecing everything together, that

$$\mathbb{P}_{\mathcal{X}^m} (|F(t) - F_m(t)| > \varepsilon) \leq 8(m + 1) e^{-\frac{m\varepsilon^2}{32}} \int_{\mathcal{X}^m} d\mathbb{P}_{\mathcal{X}^m}(x) = 8(m + 1) e^{-\frac{m\varepsilon^2}{32}}$$

as desired. □

⁴We reiterate that s_j are the only variables in this expression; we are conditioning on \mathcal{X}^m and therefore holding $x \in \mathcal{X}^m$ fixed; as an event we may write e.g. the left hand side explicitly as $\left\{ \bar{s} \in \mathcal{S}^m : \left| \sup_{t \in \mathbb{R}} \frac{1}{m} \sum_{j=1}^m \sigma_j (\mathbb{1}_{x_j \leq t}) \right| \right\}$

5 Lecture 5

Our setting: we have joint probability space $(\mathcal{X} \times \mathcal{Y}, \mathbb{P}_{\mathcal{X} \times \mathcal{Y}})$ and recall the standard diagram eq. (6). Our goal is to construct *model* $\tilde{y} : \mathcal{X} \rightarrow \mathcal{Y}$ for which $\tilde{y}(x) \approx y$ for “most” pairs $(x, y) \in \mathcal{X} \times \mathcal{Y}$. Probability makes sense (/ precisifies) ‘most.’

Once we have our hands on a model $\tilde{y} : \mathcal{X} \rightarrow \mathcal{Y}$, we would like some way to *measure* (not in the measure-theoretic sense!) how well \tilde{y} “does” on arbitrary labeled data point $(x, y) \in \mathcal{X} \times \mathcal{Y}$, which we capture by saying there is some random variable $l_{\tilde{y}} : \mathcal{X} \times \mathcal{Y} \rightarrow \mathbb{R}$, which we call the *cost function*, mapping $(x, y) \mapsto l_{\tilde{y}}(x, y) \in \mathbb{R}$. Often or usually this random variable will be nonnegative and we’d like it to be as small as possible on as many points as possible. Said differently, we want $\mathbb{E}(l_{\tilde{y}})$ to be small.

What is the “variable?” The model! So you can think of cost l as inducing a *map*

$$l_{(\cdot)} : \mathcal{Y}^{\mathcal{X}} \rightarrow \mathbb{R}^{\mathcal{X} \times \mathcal{Y}},$$

$\tilde{y} \rightarrow l_{\tilde{y}}$, where $\{\tilde{y} : \mathcal{X} \rightarrow \mathcal{Y}\}$ and $\{l : \mathcal{X} \times \mathcal{Y} \rightarrow \mathbb{R}\}$ are the cost maps from the set of models to the set of random variables, returning a random variable $l_{\tilde{y}}$ for each specified model $\tilde{y} \in \mathcal{X} \rightarrow \mathcal{Y}$. And thus stated we cast *the* problem of supervised machine learning (sml) as finding an optimal model $y^* \in \arg \min_{\tilde{y} : \mathcal{X} \rightarrow \mathcal{Y}} \mathbb{E}(l_{\tilde{y}})$. Note that our objective is stated “in expectation.” Crudely you can think of sml as curve fitting, and I have no problem with this plebian perspective as long as you distinguish fitting the data in your hands from the data “not at your immediate disposal.” Supervised machine learning deals not with fitting data as such but fitting the source, whence data comes. It amounts to fitting the measure! Hence how ml provides a concrete setting for understanding probability: you really cannot talk about ml without it.

Now, let’s consider a binary classification: $\mathcal{Y} = \{0, 1\}$.

Example of the cost function:

1. $l_{\tilde{y}}(x, y) = \mathbb{1}_{\tilde{y}(x) \neq y}$
2. $l_{\tilde{y}}(x, y) = (\tilde{y}(x) - y)^2$

To solve binary classification problem, where $\mathcal{Y} = \mathbb{R}$ and $\mathbb{P}_{\mathcal{Y}}(0, 1) = 0$ ie $\mathbb{P}_{\mathcal{Y}}(\{0, 1\}) = 1$:

1. Construct score function $\tilde{y} : \mathcal{X} \rightarrow [0, 1]$
2. Obtain classification threshold ie set $t \in \mathbb{R} : \tilde{y} : \mathcal{X} \rightarrow \{0, 1\} : x \mapsto \tilde{y}(x) = \mathbb{1}_{\tilde{y}(x) \geq t}$

The objective of our optimization problem is the following:

$$\min_{\tilde{y} : \mathcal{X} \rightarrow \mathcal{Y}} \mathbb{E}(l_{\tilde{y}}) = \int_{\mathcal{X} \times \mathcal{Y}} l_{\tilde{y}}(x, y) d\mathbb{P}_{\mathcal{X} \times \mathcal{Y}}(x, y)$$

We are going to consider cost function $l_{\tilde{y}}(x, y) = (\tilde{y}(x) - y)^2$ and find $y^* = \arg \min_{\tilde{y} : \mathcal{X} \rightarrow \mathcal{Y}} \mathbb{E}(l_{\tilde{y}})$

$$\mathbb{E}(l_{\tilde{y}}) = \int_{\mathcal{X} \times \mathcal{Y}} (\tilde{y}(x) - y)^2 d\mathbb{P}_{\mathcal{X} \times \mathcal{Y}}(x, y) = \int_{\mathcal{X}} \int_{\mathcal{Y}|\mathcal{X}} (\tilde{y}(x) - y)^2 d\mathbb{P}_{\mathcal{Y}|\mathcal{X}}(y|x) d\mathbb{P}_{\mathcal{X}}(x)$$

We need to minimize globally by minimizing pointwise in \mathcal{X} ie we need to define $\tilde{y}(x) \in \mathcal{X} \forall x \in \mathcal{X}$.

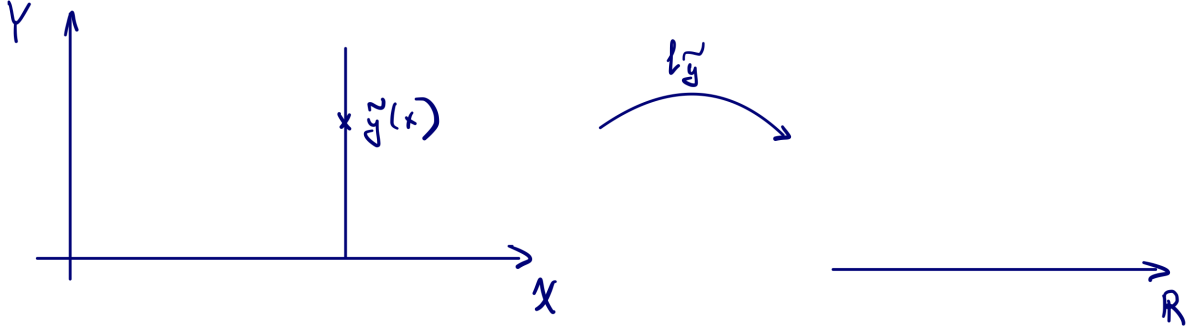


Figure 13

$$\begin{aligned} \min_{\tilde{y}(x) \in \mathcal{Y}} \int_{\mathcal{Y}|\mathcal{X}} (\tilde{y}(x) - y)^2 d\mathbb{P}_{\mathcal{Y}|\mathcal{X}}(y|x) &= \mathbb{P}_{\mathcal{Y}|\mathcal{X}}(y = 0|x)(\tilde{y}(x))^2 + \mathbb{P}_{\mathcal{Y}|\mathcal{X}}(y = 1|x)(\tilde{y}(x) - 1)^2 \\ &= (1 - p(x))\tilde{y}(x)^2 + p(x)(\tilde{y}(x) - 1)^2 \end{aligned}$$

To find an optimal $y^*(x) \in \mathcal{Y}$ we need to differentiate w.r.t. $\tilde{y}(x)$ and set that derivative to 0.

$$\begin{aligned} \frac{d((1-p(x))\tilde{y}(x)^2 + p(x)(\tilde{y}(x)-1)^2)}{d(\tilde{y}(x))} &= (1-p(x))2\tilde{y}(x) + p(x)2(\tilde{y}(x)-1) = 0 \\ 2\tilde{y}(x) - 2p(x)\tilde{y}(x) + 2p(x)\tilde{y}(x) - 2p(x) &= 0 \\ \tilde{y}(x) - p(x) &= 0 \\ \tilde{y}(x) &= p(x) \end{aligned}$$

Therefore,

$$y^* : \mathcal{X} \rightarrow \mathcal{Y} : y^*(x) = \mathbb{P}_{\mathcal{Y}|\mathcal{X}}(y = 1|x) = \mathbb{E}(\mathcal{Y}|\mathcal{X}) \text{ is optimal.}$$

In general,

$$y^*(x) = \int_{\mathcal{Y}|\mathcal{X}} y d\mathbb{P}_{\mathcal{Y}|\mathcal{X}}(y|x)$$

But we have a problem: we typically don't know true measure $\mathbb{P}_{\mathcal{X} \times \mathcal{Y}}$ (or $\mathbb{P}_{\mathcal{Y}|\mathcal{X}}$).

Now, consider $\mathbb{P}_{\mathcal{X}|\mathcal{Y}}(\tilde{y}(x)|y)$, where $\tilde{y} : \mathcal{X} \rightarrow \mathcal{Y}$ induces a measure $\mathbb{P}_{\mathcal{Y}^2}$ s.t. $\mathbb{P}_{\mathcal{Y}^2}(\tilde{y} \in [a, b], y) = \mathbb{P}_{\mathcal{X} \times \mathcal{Y}}(\tilde{y}^{-1}([a, b]), y)$. Denote $F_0(t) = \mathbb{P}_{\mathcal{X}|\mathcal{Y}}(\tilde{y}(x) \leq t|y = 0)$ and $F_1(t) = \mathbb{P}_{\mathcal{X}|\mathcal{Y}}(\tilde{y}(x) \leq t|y = 1)$.

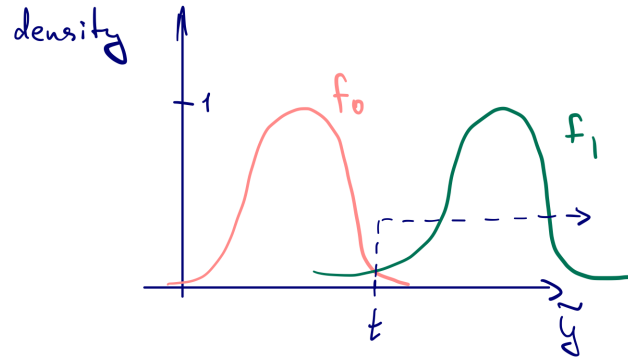


Figure 14

We consider various metrics for evaluating score $\tilde{y} \in [0, 1]$:

1. ks score $F_y(t)P_{\tilde{y}|y}(\tilde{y} \leq t|y)$ and set $ks := \sup_{t \in \mathbb{R}} |F_0(t) - F_1(t)|$
Ideally, classifier \tilde{y} separates class data perfectly.

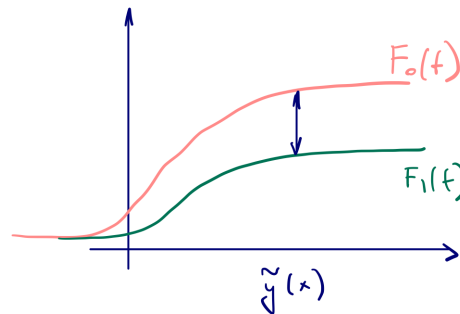


Figure 15

2. Consider the score itself. We obtain a classifier by thresholding ie $\tilde{y}_t(x) := \mathbb{1}_{(\tilde{y}(x)) \geq t}$. Then,

$$\int_t^\infty f_y(s) ds = \begin{cases} \mathbb{P}_{\mathcal{X}|Y}(\tilde{y}_t = 1 | y = 1) = \text{tpr}(t), \\ \mathbb{P}_{\mathcal{X}|Y}(\tilde{y}_t = 1 | y = 0) = \text{fpr}(t) \end{cases} \quad (27)$$

We want tpr to be as big as possible and fpr to be as small as possible.

Key point: true/false positive rate parametrized by the score range ie tpr/fpr: $\tilde{y} \rightarrow [0, 1]$.

We also have a notation for precision: $\mathbb{P}_{Y|\mathcal{X}}(y = 1 | \tilde{y}_t = 1)$.

3. Parametrization on (27) induces ROC/AUC curves.

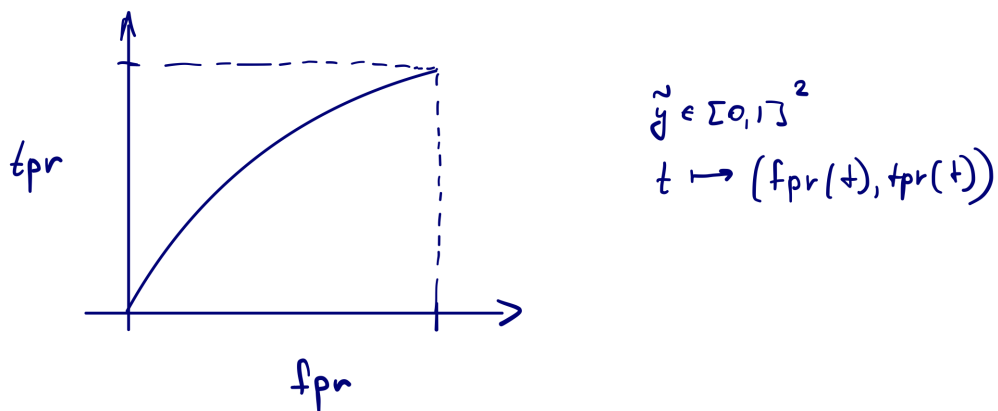


Figure 16

Gold standard: $\exists t^* : (fpr(t^*), tpr(t^*)) = (0, 1)$.

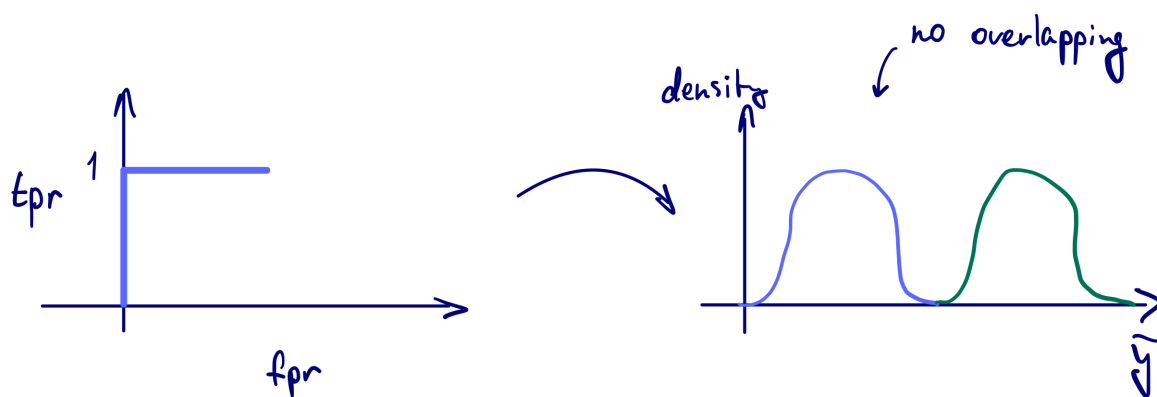


Figure 17

4. Typically think of \tilde{y} as a likelihood score ie likelihood of belonging to a certain class. To make it rigorous:

$$y^*(x) = \mathbb{P}_{\mathcal{Y}|\mathcal{X}}(y = 1|x) =: p(x) \quad (28)$$

But it's not necessarily true (was true in pa0):

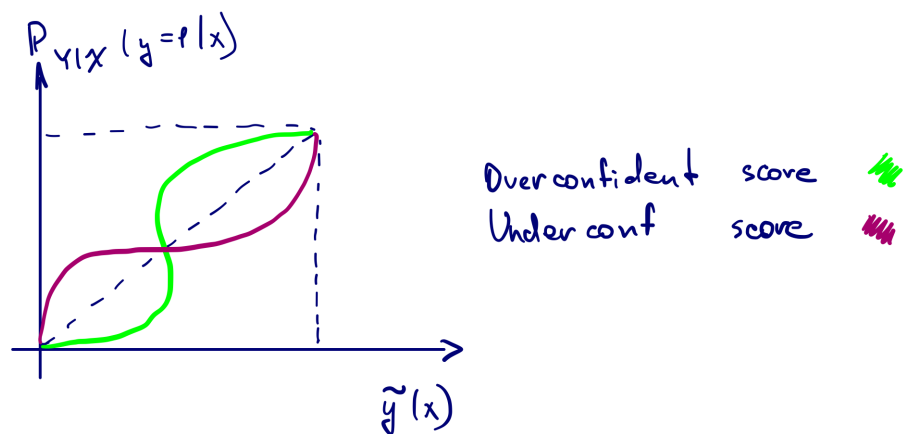


Figure 18

If (28) is satisfied we say that \tilde{y} is calibrated.

5. Equal error rate is defined in ws3. Finally, following from that, we have

$$\mathbb{P}_{Y|X}(y=1|x) = \frac{\mathbb{P}_Y(y=1)f_1(x)}{\mathbb{P}_Y(y=0)f_0(x) + \mathbb{P}_Y(y=1)f_1(x)}$$