

Lecture 1: Proxy-Objectives and Decision-Aware Learning

Instructor: Vishal Gupta

guptavis@usc.edu

Disclaimer: *These notes may contain typos and errors. Please do not distribute without written permission of the Instructor.*

1. COURSE OVERVIEW

“Decision-aware” learning (also known as optimization-aware learning, end-to-end learning, task-based learning, decision-focused learning, and operational statistics) is an increasingly popular research area within data-driven optimization. Indeed, empirical studies of “decision-aware” algorithms often show they can substantially outperform traditional “predict-then-optimize” approaches to decision-making under uncertainty, but a complete theoretical understanding of their strengths and weaknesses is still developing.

This mini-course is meant as a quick & dirty introduction to the use of proxy-objectives for decision-aware learning, and analyzing those proxy-objectives via empirical process theory. Our focus is on how to use the theory to prove performance guarantees about various methods. The aim is to empower students to use these methods in their own research.

2. WARM-UP: SAMPLE AVERAGE APPROXIMATION (SAA)

Traditional Stochastic Optimization studies:

$$\mathbf{x}^* \in \operatorname{argmin}_{\mathbf{x} \in \mathcal{X}} \mathbb{E} [c(\mathbf{x}, \boldsymbol{\xi})], \quad (1-1)$$

where

- \mathbf{x} is our decision variable,
- \mathcal{X} is a *known* feasible region,
- $c(\cdot, \cdot)$ is a *known* cost function,
- and $\boldsymbol{\xi} \sim \mathbb{P}$ is an *exogenous* random variable.

Exercise: Describe some applications of stochastic optimization as described below. Be specific. What do each of \mathbf{x} , $\boldsymbol{\xi}$, $c(\cdot, \cdot)$ and \mathcal{X} model?

Importantly, in the data-driven setting, \mathbb{P} is *unknown*, so we cannot compute \mathbf{x}^* directly. Rather, we have access to data $\{\hat{\boldsymbol{\xi}}_1, \dots, \hat{\boldsymbol{\xi}}_n\}$ drawn i.i.d. from \mathbb{P} .

With these data, we might use the Sample Average Approximation (SAA) (also called Empirical Risk Minimization or ERM):

$$\mathbf{x}_n \in \operatorname{argmin}_{\mathbf{x} \in \mathcal{X}} \frac{1}{n} \sum_{j=1}^n c(\mathbf{x}, \hat{\boldsymbol{\xi}}_j). \quad (1-2)$$

How should we assess the quality of \mathbf{x}_n ? Since \mathbf{x}_n is feasible in Problem (1-1) by construction, we often just upper bound its out-of-sample performance $\mathbb{E} [c(\mathbf{x}_n, \boldsymbol{\xi}) \mid \hat{\boldsymbol{\xi}}_{1:n}]$.

IMPORTANT: *Out-of-sample performance is a random variable. (Why?) Hence, an upper bound is usually bounds its expectation or (better) its tail behavior.*

This upper bound is usually computed relative to a suitable oracle benchmark. For example, in a few sessions, we will prove the following theorem:

Theorem 1.1 (SAA for Discrete Feasible Regions). *Suppose $|\mathcal{X}| < \infty$, and $\max_{\mathbf{x} \in \mathcal{X}} |c(\mathbf{x}, \boldsymbol{\xi})| \leq c_{\max}$ almost surely. Then there exists a universal constant $C > 0$ such that for any $0 < \delta < 1$, with probability at least $1 - \delta$,*

$$0 \leq \mathbb{E} [c(\mathbf{x}_n, \boldsymbol{\xi}) \mid \hat{\boldsymbol{\xi}}_{1:n}] - \mathbb{E} [c(\mathbf{x}^*, \boldsymbol{\xi})] \leq C \cdot c_{\max} \sqrt{\frac{\log(2/\delta) \log(|\mathcal{X}|)}{n}}.$$

Here the “oracle benchmark” is the performance of \mathbf{x}^* , i.e., the performance we would have gotten had we known \mathbb{P} upfront. Thus, in words, the theorem says that with high probability, the out-of-sample performance is not much worse than the best possible performance, provided n is much larger than $\log(|\mathcal{X}|)$. Note that we can provide this bound, even though we *can't* evaluate \mathbf{x}^* !

3. THE PROXY-OBJECTIVE PERSPECTIVE

SAA is our first example of the use of a proxy-objective in decision-aware learning. Specifically,

- We would like to solve Problem (1-1) but we cannot because we don't know \mathbb{P} , and thus don't know the objective function.
- Instead, we solve a different optimization problem (i.e., Problem (1-2)) with an objective function built from the data and meant to approximate the unknown objective function.
- We provide a performance guarantee on the suboptimality of our proxy problem's solution in the original problem.

Many methods in data-driven optimization follow this loose structure and, consequently, can be analyzed using a small set of similar tools. We next give some other examples.

3.1. Parametric Probability Models. Suppose we believe the unknown \mathbb{P} follows some parametric model, e.g., that it is a multivariate normal distribution with known covariance equal to the identity.

Then, letting $\mathbb{P} = \mathbb{P}_{\boldsymbol{\theta}^*}$ for some parameters $\boldsymbol{\theta}^* \in \Theta$, we rewrite Problem (1-1) as

$$\mathbf{x}(\boldsymbol{\theta}^*) \in \operatorname{argmin}_{\mathbf{x} \in \mathcal{X}} \mathbb{E}^{\mathbb{P}_{\boldsymbol{\theta}^*}} [c(\mathbf{x}, \boldsymbol{\xi})]. \tag{1-3}$$

We might then use our favorite statistical learning procedure (e.g. Maximum Likelihood) to form an estimate $\hat{\boldsymbol{\theta}}_n = \hat{\boldsymbol{\theta}}_n(\hat{\boldsymbol{\xi}}_{1:n})$ from the data and compute a predict-then-optimize decision:

$$\mathbf{x}(\hat{\boldsymbol{\theta}}_n) \in \operatorname{argmin}_{\mathbf{x} \in \mathcal{X}} \mathbb{E}^{\mathbb{P}_{\hat{\boldsymbol{\theta}}_n}} [c(\mathbf{x}, \boldsymbol{\xi})] \tag{1-4}$$

In other words, we use Problem (1-4) as a proxy for Problem (1-3). We might then seek high probability bounds on

$$\mathbb{E} [c(\mathbf{x}(\hat{\boldsymbol{\theta}}_n), \boldsymbol{\xi}) \mid \hat{\boldsymbol{\xi}}_{1:n}] - \mathbb{E} [c(\mathbf{x}(\boldsymbol{\theta}^*), \boldsymbol{\xi})]$$

to establish a performance guarantee.

3.2. Policy Classes and Best-in-Class Policies. What happens if n is the same order as $\log(|\mathcal{X}|)$ in Theorem 1.1? Our performance guarantee is not informative since it's proportional to c_{\max} and $2c_{\max}$ is a trivial performance guarantee.

This observation isn't a weakness in our analysis.

Exercise: Construct an example of binary optimization of the form Problem (1-1) where the performance of SAA is (with high probability) far from the full-information optimum when the number of data points is equal to the number of binary variables.

In fact, in certain data settings, e.g., when data are scarce, *every* method has a large sub-optimality gap relative to the full-information optimal solution.

So what do we do in these cases? Since, full-information isn't achievable, it makes sense to consider a weaker oracle benchmark. Inspired by Theorem 1.1, one approach might be to identify a subset $\mathcal{X}_0 \subseteq \mathcal{X}$ of the feasible region, and limit attention to solutions in \mathcal{X}_0 .

Following this line of reasoning, we define

$$\mathbf{x}_0^* = \mathbf{x}_0^*(\boldsymbol{\xi}_{1:n}) \in \underset{\mathbf{x} \in \mathcal{X}_0}{\operatorname{argmin}} \mathbb{E} [c(\mathbf{x}, \boldsymbol{\xi})] \tag{1-5}$$

to be the best solution in \mathcal{X}_0 . The “best-in-class” solution cannot be computed because \mathbb{P} is unknown, but performs (path by path) every other policy with values in \mathcal{X}_0 .

Since we can't solve Problem (1-5), we proxy it by

$$\hat{\mathbf{x}}_0 = \hat{\mathbf{x}}_0(\hat{\boldsymbol{\xi}}_{1:n}) \in \underset{\mathbf{x} \in \mathcal{X}_0}{\operatorname{argmin}} \frac{1}{n} \sum_{j=1}^n c(\mathbf{x}, \hat{\boldsymbol{\xi}}_j), \tag{1-6}$$

i.e., the “corresponding” SAA solution. Theorem 1.1 then gives a performance guarantee: with probability at least $1 - \delta$,

$$0 \leq \mathbb{E} [c(\hat{\mathbf{x}}_0, \boldsymbol{\xi}) \mid \hat{\boldsymbol{\xi}}_{1:n}] - \mathbb{E} [c(\mathbf{x}_0^*, \boldsymbol{\xi})] \leq C \cdot c_{\max} \sqrt{\frac{\log(2/\delta) \log(|\mathcal{X}_0|)}{n}}. \tag{1-7}$$

Thus, if $\log(|\mathcal{X}_0|) \ll n$, we can learn a solution $\hat{\mathbf{x}}_0$ whose performance is comparable to the best-in-class solution \mathbf{x}_0^* .

Of course, there are no free lunches. The above result tells us nothing about the gap $\mathbb{E} [c(\mathbf{x}_0, \boldsymbol{\xi})] - \mathbb{E} [c(\mathbf{x}^*, \boldsymbol{\xi})]$. If $|\mathcal{X}_0|$ is small, we might expect this gap to be big. But, if $|\mathcal{X}_0|$ is big, then the performance guarantee Eq. (1-7) will be poor and our previous exercise suggests it may be hard to learn a best-in-class policy. This tension is somehow unavoidable.

In summary, in settings where the data are rich/plentiful enough, we should try to learn a policy with performance similar to the full-information solution \mathbf{x}^* . In settings where this is not possible, we should instead use our domain knowledge to restrict attention to a smaller set of policies, and

try to learn a policy with performance comparable to the best policy in that set. There are of course *many* ways to choose the “smaller set of policies,” and a good choice is typically application dependent.

For now, what I want to stress is that this example again follows our proxy-objective framework. We proxy Problem (1-5) by Problem (1-6) and prove a performance guarantee on the gap to best-in-class.

Remark 1.2 (Error from Policy Restriction). *After class, a smart student asked about bounding the gap $\mathbb{E}[c(\mathbf{x}_0^*, \boldsymbol{\xi})] - \mathbb{E}[c(\mathbf{x}^*, \boldsymbol{\xi})]$. In some very special situations, it is possible to bound this error, however, this is usually the exception, not the rule. In this course, we will focus on the cases where \mathcal{X}_0 is determined by domain knowledge and practical considerations, and focus simply on learning a best-in-class policy.*

3.3. Plug-In Policy Classes. Above, the set \mathcal{X}_0 does not depend on the data. (But, both the policies \mathbf{x}_0^* and $\hat{\mathbf{x}}_0$ do depend on the data. Why?) Identifying a good \mathcal{X}_0 so that $|\mathcal{X}_0|$ is small but still sufficiently rich so that $\mathbb{E}[c(\mathbf{x}_0, \boldsymbol{\xi})] - \mathbb{E}[c(\mathbf{x}^*, \boldsymbol{\xi})]$ is small from just prior domain knowledge might be tricky.

We next introduce a particular class of policies called “plug-in policies” that take values in a set $\mathcal{X}_0(\hat{\boldsymbol{\xi}}_{1:n})$ which does depend on the data. These policies have some nice practical features. For example, plug-in policies generalize the popular “estimate-then-optimize” approach to a decision-aware setting. Some authors prefer the term “smart predict-then-optimize” to describe best-in-class plug-in policies [EG22].

To that end, rewrite the full-information optimal solution to Problem (1-1) as

$$\mathbf{x}(\mathbb{P}) \in \operatorname{argmin}_{\mathbf{x} \in \mathcal{X}} \mathbb{E}_{\boldsymbol{\xi} \sim \mathbb{P}} [c(\mathbf{x}, \boldsymbol{\xi})],$$

to stress the dependence on \mathbb{P} . That is, $\mathbf{x}(\mathbb{P})$ is the optimal solution when the data follows distribution \mathbb{P} .

More generally, when $\boldsymbol{\xi} \sim \mathbb{Q}$, define the plug-in policy for \mathbb{Q} ,

$$\mathbf{x}(\mathbb{Q}) \in \operatorname{argmin}_{\mathbf{x} \in \mathcal{X}} \mathbb{E}_{\boldsymbol{\xi} \sim \mathbb{Q}} [c(\mathbf{x}, \boldsymbol{\xi})]. \tag{1-8}$$

These policies essentially “plug-in” \mathbb{Q} for the unknown \mathbb{P} and solve the optimization problem.

We can now define a class of plug-in policies. Consider a set of probability distributions (dependent on data) $\{\hat{\mathbb{P}}_\tau(\hat{\boldsymbol{\xi}}_{1:n}) : \tau \in \mathcal{T}\}$, where for each τ , $\hat{P}_\tau(\hat{\boldsymbol{\xi}}_{1:n})$ refers to a different way to fit a probability distribution to the observed data. For example, one τ might correspond to fitting a mixture of normal distributions to the data, while another might correspond to simply using the empirical distribution.

This set induces a set of (data-dependent) plug-in policies

$$\mathcal{X}_0(\hat{\boldsymbol{\xi}}_{1:n}) = \left\{ \mathbf{x}(\mathbb{Q}) : \mathbb{Q} \in \{\hat{\mathbb{P}}_\tau(\hat{\boldsymbol{\xi}}_{1:n}) : \tau \in \mathcal{T}\} \right\} = \left\{ \mathbf{x}(\mathbb{P}_\tau(\hat{\boldsymbol{\xi}}_{1:n})) : \tau \in \mathcal{T} \right\}.$$

Both the policy class $\mathcal{X}_0(\hat{\boldsymbol{\xi}}_{1:n})$ and each plug-in policy $\mathbf{x}(\mathbb{P}_\tau(\hat{\boldsymbol{\xi}}_{1:n})) \in \mathcal{X}_0(\hat{\boldsymbol{\xi}}_{1:n})$ depend on the data.

We define the (oracle) best-in-class decision $\mathbf{x}(\mathbb{P}_{\tau^{\text{OR}}}(\hat{\boldsymbol{\xi}}_{1:n}))$, where

$$\tau^{\text{OR}} \in \operatorname{argmin}_{\tau \in \mathcal{T}} \mathbb{E} \left[c(\mathbf{x}(\mathbb{P}_{\tau}(\hat{\boldsymbol{\xi}}_{1:n})), \boldsymbol{\xi}) \mid \hat{\boldsymbol{\xi}}_{1:n} \right]. \quad (1-9)$$

By construction, this oracle benchmark outperforms any other plug-in policy in $\mathcal{X}_0(\hat{\boldsymbol{\xi}}_{1:n})$. However, it is an “oracle” benchmark because one must know the distribution \mathbb{P} to compute τ^{OR} .

Since we can’t solve Problem (1-9) directly, we create a proxy. One natural proxy might be the leave-one-out objective:

$$\tau^{\text{LOO}} \in \operatorname{argmin}_{\tau \in \mathcal{T}} \sum_{j=1}^n c(\mathbf{x}(\hat{\mathbb{P}}_{\tau}(\hat{\boldsymbol{\xi}}_{-j}, \hat{\boldsymbol{\xi}}_j)), \hat{\boldsymbol{\xi}}_j), \quad (1-10)$$

where $\hat{\boldsymbol{\xi}}_{-j}$ is the data $\hat{\boldsymbol{\xi}}_{1:n}$ *excluding* data point j .

Thus, we have proxied Problem (1-9) by Problem (1-10). We might then seek to bound

$$\mathbb{E} \left[c(\mathbf{x}(\hat{\boldsymbol{\xi}}_{1:n}, \tau^{\text{LOO}}), \boldsymbol{\xi}) \right] - \mathbb{E} \left[c(\mathbf{x}(\hat{\boldsymbol{\xi}}_{1:n}, \tau^{\text{OR}}), \boldsymbol{\xi}) \right]$$

to provide a performance guarantee of our method relative to the best-in-class plug-in policy.

The notation in this example is a bit heavy. For now, things I want to stress:

- When considering “best-in-class” learning, we need to create a proxy for the best-in-class optimization problem. By contrast, in our example in Section 2 needed to create a proxy for the full-information problem. You should look to see how this difference also changes the kind of performance guarantee we seek.
- The optimization Problem (1-10) can be quite difficult unless \mathcal{T} and $\mathbf{x}(\hat{\boldsymbol{\xi}}_{1:n}, \tau)$ both have some nice structure.

Let’s consider a specific example of the above framework.

Example 1.3 (Decision-Aware Regularization). *Suppose $c(x, \xi) = (x - \xi)^+ + (\xi - x)^+$ is the newsvendor (a.k.a. pinball loss) function.*

We let $\mathcal{T} = \mathbb{R}_+$, and for each $\tau \in \mathcal{T}$, consider fitting an exponential distribution¹ to $\hat{\boldsymbol{\xi}}_{1:n}$ with regularized maximum likelihood. Recall, the negative log-likelihood of the data is given by

$$-\loglik(\lambda; \hat{\boldsymbol{\xi}}_{1:n}) = -n \log(\lambda) + \lambda \sum_j \hat{\xi}_j.$$

We will use a regularization proportional to $\log(\lambda)$ which essentially shrinks towards 0.

The regularized maximum likelihood estimate is then

$$\begin{aligned} \hat{\lambda}(\tau) &= \hat{\lambda}(\tau, \hat{\boldsymbol{\xi}}_{1:n}) \in \operatorname{argmin}_{\lambda \geq 0} -n \log(\lambda) + \lambda \sum_{j=1}^n \hat{\xi}_j + \tau \log(\lambda) \\ &= \frac{(n - \tau)^+}{\sum_j \hat{\xi}_j}. \end{aligned}$$

Hence, we take \hat{P}_{τ} to be an $\operatorname{Exp}\left(\frac{(n-\tau)^+}{\sum_j \hat{\xi}_j}\right)$ distribution.

¹Recall $\xi \sim \operatorname{Exp}(\lambda)$ means that $\mathbb{E}[\xi] = 1/\lambda$.

Notice then that $\mathbf{x}(\hat{P}_\tau(\hat{\xi}_{1:n})) = \frac{\log 2 \sum_j \hat{\xi}_j}{n-\tau}$ and, similarly, $\mathbf{x}(\hat{\xi}_{-j}, \tau) = \frac{\log 2 \sum_{i \neq j} \hat{\xi}_i}{n-\tau}$. Both of these are very easy to compute, because of the way we picked our policy class.

Since we don't know \mathbb{P} , we can't evaluate τ^{OR} . But, we can numerically find τ^{LOO} by

$$\tau^{\text{LOO}} \in \operatorname{argmin}_{0 \leq \tau \leq n} \sum_{i=1}^n \left[\frac{\log 2 \sum_{i \neq j} \hat{\xi}_i}{n-\tau} - \hat{\xi}_i \right]^+ + \left[\hat{\xi}_i - \frac{\log 2 \sum_{i \neq j} \hat{\xi}_i}{n-\tau} \right]^+.$$

This problem is a simple one dimensional optimization that we can (approximately) solve by enumeration.

Remark 1.4 (Benefits of Plug-in Policies). As mentioned, plug-in policy classes are only one type of policy class. There are many others. There are at least three important benefits of plug-in policies.

First, Eq. (1-8) is of the same form as Problem (1-1). Thus, if one already has a specialized algorithm for solving problems of the form Problem (1-1), we can leverage the same algorithm for computing our plug-in policies. This benefit is especially important in large-scale settings.

Second, tied to the previous point, the “structure” of the plug-in policy $\mathbf{x}(\mathbb{Q})$ is often similar to the “structure” of $\mathbf{x}(\mathbb{P})$ which can be important in practice. For example, when \mathcal{X} is a polyhedron and $c(\mathbf{x}, \xi) = \mathbf{x}^\top \xi$, both $\mathbf{x}(\mathbb{P})$ and $\mathbf{x}(\mathbb{Q})$ will be extreme points. Other, regularization or robust optimization approaches, might not enforce this structure. In network optimization problems, e.g., this distinction can be important, since we often need our solution to be an extreme point to be interpretable as a path or spanning tree.

Finally, plug-in policies generalize estimate-then-optimize procedures. We often are interested in the benefits of a decision-aware approach over estimate-then-optimize. By construction, our oracle benchmark is necessarily no worse than the predict-then-optimize policy, and so might offer a way to try and quantify the benefit of decision-aware approaches. Indeed, the if $\hat{\mathbb{P}}_\tau$ corresponds to a traditional statistical estimator for \mathbb{P} , then the gap $\mathbb{E} \left[c(\mathbf{x}(\hat{\mathbb{P}}_\tau), \xi) \right] - \mathbb{E} \left[c(\mathbf{x}(\hat{\mathbb{P}}_\tau^{\text{OR}}), \xi) \right]$ quantifies the potential benefits of decision-aware learning for the problem structure and class \mathcal{T} in an algorithm-independent way.

3.4. Other Examples. Many, many other algorithms in data-driven optimization can be seen through the lens of proxy-objectives including:

- The SPO+ Surrogate loss [EG22]
- The Stein Correction for Small-Data, Large-Scale Linear Optimization Problems [GR21]
- The Variance Gradient Correction [GHR22]
- The Predictive to Prescriptive Method [BK20]
- Policy Learning with Doubly Robust Estimates [DLL11]
- and more.

That said, other popular algorithms are not (easily) seen as proxy-objective methods, e.g. applying stochastic gradient descent directly to Problem (1-1). These lecture notes exclusively focus on proxy-objective methods.

4. OVERVIEW OF ANALYSIS TECHNIQUE

Our goal in viewing these varied algorithms through a common lens of proxy-objectives was to (hopefully) develop some common tools for their analysis. To that end, we next sketch the fundamental lemma that we will use throughout this course to prove performance guarantees for methods based on proxy-objectives.

To express the idea in a general setting, suppose we are interested in the target optimization problem

$$\mathbf{x}^* \in \operatorname{argmin}_{\mathbf{x} \in \mathcal{X}} f(\mathbf{x}),$$

but we instead solve the proxy objective problem

$$\hat{\mathbf{x}} \in \operatorname{argmin}_{\mathbf{x} \in \mathcal{X}} \hat{f}(\mathbf{x}).$$

We then bound the suboptimality of $\hat{\mathbf{x}}$ by the approximation quality of this proxy.

Lemma 1.5 (Bounding Suboptimality by Uniform Error). *We have*

$$0 \leq f(\hat{\mathbf{x}}) - f(\mathbf{x}^*) \leq 2 \sup_{\mathbf{x} \in \mathcal{X}} |f(\mathbf{x}) - \hat{f}(\mathbf{x})|.$$

Proof. The first inequality is clear by optimality of \mathbf{x}^* . For the second, note that

$$\begin{aligned} f(\hat{\mathbf{x}}) - f(\mathbf{x}^*) &= f(\hat{\mathbf{x}}) - \hat{f}(\hat{\mathbf{x}}) + \underbrace{\hat{f}(\hat{\mathbf{x}}) - \hat{f}(\mathbf{x}^*)}_{\leq 0 \text{ by optimality of } \hat{\mathbf{x}}} + \hat{f}(\mathbf{x}^*) - f(\mathbf{x}^*) \\ &\leq |f(\hat{\mathbf{x}}) - \hat{f}(\hat{\mathbf{x}})| + |f(\mathbf{x}^*) - \hat{f}(\mathbf{x}^*)| \\ &\leq 2 \sup_{\mathbf{x} \in \mathcal{X}} |f(\mathbf{x}) - \hat{f}(\mathbf{x})|. \end{aligned}$$

This concludes the proof. □

Thus, to provide a performance guarantee on a proxy-objective method, it suffices to show that the proxy is a good approximation to the true function everywhere over the domain. Notice the lemma makes no assumptions on the structure of \mathcal{X} or the convexity/continuity/regularity of f or \hat{f} .

The above lemma crucially uses the optimality of $\hat{\mathbf{x}}$ and \mathbf{x}^* . (Make sure you know where.) But it does not leverage much else about the structure of $\hat{\mathbf{x}}$ or \mathbf{x}^* , other than $\hat{\mathbf{x}}, \mathbf{x}^* \in \mathcal{X}$. As a consequence, in some data-driven settings, the bound can be quite bad.

Exercise: Draw a picture of $f(\cdot)$ and $\hat{f}(\cdot)$ where $\mathbf{x} = \hat{\mathbf{x}}$ but $\sup_{\mathbf{x} \in \mathcal{X}} |f(\mathbf{x}) - \hat{f}(\mathbf{x})| = \infty$.

Exercise: Look back at our plug-in policy class example (cf. Example 1.3). How would we apply the lemma? What are $f(\cdot)$, $\hat{f}(\cdot)$ and $\sup_{\mathbf{x} \in \mathcal{X}} |f(\mathbf{x}) - \hat{f}(\mathbf{x})|$ for that example?

REFERENCES

- [BK20] Dimitris Bertsimas and Nathan Kallus. “From predictive to prescriptive analytics”. In: *Management Science* 66.3 (2020), pp. 1025–1044.
- [DLL11] Miroslav Dudík, John Langford, and Lihong Li. “Doubly robust policy evaluation and learning”. In: *arXiv preprint arXiv:1103.4601* (2011).
- [EG22] Adam N Elmachtoub and Paul Grigas. “Smart “Predict, then Optimize””. In: *Management Science* 68.1 (2022), pp. 9–26.
- [GHR22] Vishal Gupta, Michael Huang, and Paat Rusmevichientong. “Debiasing in-sample policy performance for small-data, large-scale optimization”. In: *Operations Research* (2022). Forthcoming.
- [GR21] Vishal Gupta and Paat Rusmevichientong. “Small-data, large-scale linear optimization with uncertain objectives”. In: *Management Science* 67.1 (2021), pp. 220–241.

Lecture 2: Introduction to Concentration

Instructor: Vishal Gupta

guptavis@usc.edu

Disclaimer: *These notes may contain typos and errors. Please do not distribute without written permission of the Instructor.*

1. SOME NOTATION

To simplify many proofs, we will write $a \lesssim b$ to mean that there exists a universal constant C (not depending on any problem data) such that $a \leq Cb$.

So for example, when $d \geq 2$, we have that $\log(2d) \lesssim \log(d)$. (Why?) But, if we only knew that $d > 1$, we *cannot* say that $\log(2d) \lesssim \log(d)$. (Why?)

2. REVIEW OF THE CHERNOFF TECHNIQUE

Recall Markov's Inequality: For any *nonnegative* random variable X , $\mathbb{P}(X > t) \leq \mathbb{E}[X]/t$. Often we establish tail bounds on a random variable by applying Markov's inequality to a monotone transformation $f(X)$. The most common transformation is $f(t) = \exp(\theta t)$ for some (well-chosen) $\theta > 0$. This technique is usually called the Chernoff technique in textbooks. The Chernoff technique is a great general purpose approach that generally turns a bound on the mgf of a random variable $\mathbb{E}[\exp(\theta X)]$ into a tail bound.

Example 2.1 (Chernoff Bound for the Normal Distribution). *Suppose $X \sim \mathcal{N}(0, \sigma^2)$. Then, for any $\theta > 0$,*

$$\begin{aligned} \mathbb{P}(X > t) &= \mathbb{P}(\exp(\theta X) > \exp(\theta t)) \\ &\leq \exp(-\theta t) \mathbb{E}[\exp(\theta X)] \\ &= \exp\left(-\theta t + \frac{\theta^2 \sigma^2}{2}\right), \end{aligned}$$

where we've used the MGF of a normal distribution. We then optimize θ to obtain the best possible bound. Choosing $\theta = t/\sigma^2$ yields the bound

$$\mathbb{P}(X > t) \leq \exp\left(-\frac{t^2}{2\sigma^2}\right).$$

Said differently, for any $0 < \delta < 1$, with probability at least $1 - \delta$ we have

$$X \lesssim \sigma \sqrt{\log(2/\delta)}.$$

There are tighter bounds for the tails of a normal distribution, but this bound is fairly tight.

The Chernoff technique motivates us to consider classes of random variables that admit “nice” bounds on their mgf, since we can turn such bounds into tail bounds. We next introduce a particular class of random variables with “nice” mgfs.

3. SUBGAUSSIAN RANDOM VARIABLES

Definition 2.2. A mean-zero random variable X is subGaussian with variance proxy σ^2 if

$$\mathbb{E} \left[e^{\theta X} \right] \leq e^{\frac{\theta^2 \sigma^2}{2}} \quad \forall \theta \in \mathbb{R}.$$

As desired, this class of random variables admit a nice tail bound:

Lemma 2.3. If $X - \mathbb{E}[X]$ is subGaussian with variance proxy σ^2 , then for any $0 < \delta < 1$, we have with probability at least $1 - \delta$,

$$X - \mathbb{E}[X] \leq \sigma \sqrt{\log(2/\delta)}.$$

Proof. Left for the reader. (Should be immediate.) □

Example 2.4 (Normal Random Variables). If $X \sim \mathcal{N}(\mu, \sigma^2)$, then $X - \mathbb{E}[X]$ is subGaussian with variance proxy σ^2 . (Why?)

Example 2.5 (Rademacher Random Variables). Consider a Rademacher random variable η , i.e. $\mathbb{P}(\eta = 1) = \mathbb{P}(\eta = -1) = 1/2$. On HW 1, you proved that η is subGaussian with variance proxy 1.

Interestingly, any bounded variable is subGaussian.

Theorem 2.6 (Hoeffding’s Inequality). Suppose X is a random variable such that $a \leq X \leq b$ almost surely. Then $X - \mathbb{E}[X]$ is subGaussian with variance proxy σ^2 satisfying $\sigma \lesssim (b - a)$.

Proof. Our proof uses a technique called “symmetrization” to reduce the analysis of a complicated random to the analysis of a Rademacher variable.

Let η be a Rademacher random variable and let \bar{X} be an i.i.d. copy of X . Then, for any θ ,

$$\begin{aligned} \mathbb{E} [\exp(\theta(X - \mathbb{E}[X]))] &\leq \mathbb{E} [\exp(\theta(X - \bar{X}))] && \text{(Jensen’s Inequality)} \\ &= \mathbb{E} [\exp(\theta\eta |X - \bar{X}|)]. \end{aligned}$$

Notice crucially how the last line uses the fact that $X - \bar{X}$ is symmetric.

Now consider conditioning on X, \bar{X} , and apply our bound on the mgf of a Rademacher random variable to obtain,

$$\mathbb{E} [\exp(\theta(X - \mathbb{E}[X]))] \leq \mathbb{E} \left[\exp \left(\theta^2 \frac{(X - \bar{X})^2}{2} \right) \right].$$

Finally, note that $a \leq X \leq b$ almost surely implies that $|X - \bar{X}| \leq (b - a)$. Hence, this last expectation is at most $\exp \left(\frac{\theta^2 (b-a)^2}{2} \right)$. In other words, $X - \mathbb{E}[X]$ is subGaussian, and its variance proxy is at most $(b - a)^2$. □

Remark 2.7. Our proof bounded the variance proxy up to a constant. This constant is not tight. See [BLM13].

Linear combinations of subGaussian random variables are subGaussian.

Theorem 2.8 (Linear Combinations of SubGaussian Random Variables). Suppose X_i for $i = 1, \dots, d < \infty$ are a set of mean-zero subGaussian random variables with variance proxies σ_i^2 . Then,

- (1) For any $a \in \mathbb{R}$, the r.v. aX_1 is subGaussian with variance proxy $a^2\sigma_1^2$.
(2) If the X_i are independent, then $\sum_{i=1}^d X_i$ is subGaussian with variance proxy $\sum_{i=1}^d \sigma_i^2$.

Proof. You'll prove parts ii) on your homework. Part i) is immediate from the definition. (Convince yourself of this.) \square

4. ORLICZ NORMS

There are many equivalent definitions of subGaussian random variables. One equivalence I find useful relates subGaussian random variables to random variables with finite Ψ -Orlicz norms.

Definition 2.9 (The Ψ -Orlicz Norm). *Let $\Psi(t) = \frac{1}{5} \exp(t^2)$. For any random variable X , the Ψ -Orlicz Norm of X , denoted $\|X\|_\Psi$, is given by*

$$\inf\{C > 0 : \mathbb{E}[\Psi(X/C)] \leq 1\}.$$

Importantly, $\|\cdot\|_\Psi$ is a *norm* on random variables, so it satisfies the “usual” properties of a norm (e.g., triangle-inequality). Orlicz norms can be defined with respect to other functions, but we will only focus on the case of $\Psi(\cdot)$. Some authors define $\Psi(\cdot)$ with a different constant than $\frac{1}{5}$. Note, X need not be mean-zero in this definition.

From Markov's inequality, we have that for any random variable X , with probability at least $1 - \delta$,

$$|X| \lesssim \|X\|_\Psi \sqrt{\log(2/\delta)}.$$

This observation motivates the following fact that you will prove on your homework:

Theorem 2.10 (Relating SubGaussian R.V.s and the $\|\cdot\|_\Psi$). *Suppose X is a mean-zero subGaussian random variable with variance proxy σ^2 . Then, $\|X\|_\Psi \lesssim \sigma$. Conversely, if X is mean-zero and $\|X\|_\Psi < \infty$, then X is subGaussian with variance proxy σ^2 such that $\sigma \lesssim \|X\|_\Psi$.*

One of the most important features of the Ψ -Orlicz Norms is that it behaves nicely under the max operation.

Theorem 2.11 (Orlicz Norm of the Max). *Suppose X_i for $i = 1, \dots, d < \infty$ are (not necessarily independent) random variables such that $\|X_i\|_\Psi < \infty$ for all i . Define $X_{\max} = \max_i X_i$. Then, for $d \geq 2$, $\|X_{\max}\|_\Psi \lesssim (\max_i \|X_i\|_\Psi) \sqrt{\log d}$.*

Proof. Fix some $C, K > 0$ to be specified later. From fundamental theorem of calculus,

$$\begin{aligned} \Psi(X_{\max}/C) &= \Psi(1) + \int_1^{X_{\max}/C} \Psi'(t) dt \\ &= \frac{e}{5} + \int_1^\infty \mathbb{I}\left\{t \leq \frac{X_{\max}}{C}\right\} \Psi'(t) dt \end{aligned}$$

Now note that because $\Psi(\cdot)$ is increasing,

$$t \leq \frac{X_{\max}}{C} \implies \Psi(Kt) \leq \Psi\left(\frac{KX_{\max}}{C}\right) \implies \frac{\Psi\left(\frac{KX_{\max}}{C}\right)}{\Psi(Kt)} \geq 1.$$

2-3

So, increasing the integrand, we have

$$\begin{aligned}\Psi(X_{\max}/C) &\leq \frac{e}{5} + \int_1^\infty \Psi\left(\frac{KX_{\max}}{C}\right) \frac{\Psi'(t)}{\Psi(Kt)} dt \\ &\leq \frac{e}{5} + \sum_{i=1}^d \int_1^\infty \Psi\left(\frac{KX_i}{C}\right) \frac{\Psi'(t)}{\Psi(Kt)} dt,\end{aligned}$$

by bounding the max by the sum. Now let $C = K \max_i \|X_i\|_\Psi$ and take expectations of both sides. We can pass the expectation inside the integral by the monotone convergence theorem. We thus have,

$$\begin{aligned}\mathbb{E}[\Psi(X_{\max}/C)] &\leq \frac{e}{5} + d \int_1^\infty \frac{\Psi'(t)}{\Psi(Kt)} \\ &= \frac{e}{5} + 2d \int_1^\infty t e^{-(K^2-1)t^2} dt \\ &= \frac{e}{5} + \frac{d}{K^2-1} e^{-(K^2-1)}.\end{aligned}$$

Taking $K = \sqrt{2 + \log d}$ makes this last expression at most 1. This shows $\|X_{\max}\|_\Psi \leq \sqrt{2 + \log d} \sim \sqrt{\log d}$. \square

Next week, we will show how to use the above tools to bound prove a performance guarantee for SAA. For now, let's see a quick toy example.

Example 2.12. Suppose $X_1 \sim \mathcal{N}(10, \sigma^2)$ and $X_2 \sim \text{Unif}[-3\sigma, 3\sigma]$, not necessarily independent with $\sigma \geq 1$. Define $Y = |X_1| + 2X_2$. We will prove a tail bound for Y .

Our strategy will be to bound $\|Y\|_\Psi$ and then use Theorem 2.10. Then, by triangle inequality,

$$\|Y\|_\Psi \leq \| |X_1| \|_\Psi + 2\|X_2\|_\Psi.$$

By Hoeffding's inequality, X_2 is subGaussian with variance proxy $\lesssim \sigma^2$, so Theorem 2.10 shows $\|X_2\|_\Psi \lesssim \sigma$.

Similarly, if we let $Z \sim \mathcal{N}(0, \sigma^2)$, then $|X_1| \leq 10 + |Z|$. By Theorem 2.10 $\| |Z| \|_\Psi = \|Z\|_\Psi \lesssim \sigma$. Hence $\| |X_1| \|_\Psi \lesssim 10 + \sigma \lesssim \sigma$.

Combining shows $\|Y\|_\Psi \sim \sigma$, so that with probability at least $1 - \delta$,

$$|Y| \lesssim \sigma \sqrt{\log(1/\delta)}.$$

5. (TIME PERMITTING): McDIARMID'S INEQUALITY

McDiarmid's Inequality is one of the most fundamental results in concentration. It is one way to formalize an important intuition in concentration: a function of many independent random variables that doesn't depend too strongly on any one variable will concentrate at its expectation.

Before stating McDiarmid's inequality, we need to define a function with bounded differences.

Definition 2.13 (Bounded Differences). Suppose $f(\mathbf{x}_1, \dots, \mathbf{x}_n)$ is a function of n random variables where $\mathbf{x}_i \in \mathcal{X}_i$ for each i . We say $f(\cdot)$ satisfies the bounded differences condition with respect to

constants c_1, \dots, c_n if for each $i = 1, \dots, n$

$$\sup_{\bar{x} \in \mathcal{X}_i} |f(\mathbf{x}_1, \dots, \mathbf{x}_n) - f(\mathbf{x}_1, \dots, \mathbf{x}_{i-1}, \bar{x}, \mathbf{x}_{i+1}, \dots, n)| \leq c_i.$$

Theorem 2.14 (McDiarmid's Inequality). *Suppose $X_i \in \mathcal{X}_i$ for $i = 1, \dots, n$ are independent random variables (not necessarily identically distributed) and consider a function $f(X_1, \dots, X_n)$. Suppose $f(\cdot)$ satisfies the bounded difference inequality with respect to a constants c_i for $i = 1, \dots, n$. Then, $f(\mathbf{X}_1, \dots, \mathbf{X}_n) - \mathbb{E}[f(\mathbf{X}_1, \dots, \mathbf{X}_n)]$ is subGaussian with variance proxy $\lesssim \|\mathbf{c}\|_2^2$.*

In particular, for any $0 < \delta < 1$, with probability at least $1 - \delta$

$$|f(\mathbf{X}_1, \dots, \mathbf{X}_n) - \mathbb{E}[f(\mathbf{X}_1, \dots, \mathbf{X}_n)]| \lesssim \|\mathbf{c}\|_2 \sqrt{\log(1/\delta)}.$$

Proof. We provide an elementary proof based on the Doob martingale. Define

$$Z_i = \mathbb{E}[f(\mathbf{X}_1, \dots, \mathbf{X}_n) \mid \mathbf{X}_{1:i}] - \mathbb{E}[f(\mathbf{X}_1, \dots, X_n) \mid \mathbf{X}_{1:(i-1)}].$$

Then, $f(\mathbf{X}_1, \dots, \mathbf{X}_n) - \mathbb{E}[f(\mathbf{X}_1, \dots, \mathbf{X}_n)] = \sum_{i=1}^n Z_i$. (Check this!). Moreover, Z_i is only a function of $\mathbf{X}_{1:i}$.

Now let $\bar{\mathbf{X}}_i$ be an i.i.d copy of \mathbf{X}_i . Then, because of the independence,

$$Z_i = \mathbb{E}[f(\mathbf{X}_1, \dots, \mathbf{X}_n) - f(\mathbf{X}_1, \dots, \mathbf{X}_{i-1}, \bar{\mathbf{X}}_i, \mathbf{X}_{i+1}, \dots, \mathbf{X}_n) \mid \mathbf{X}_{1:i}],$$

and by the bounded differences condition, $|Z_i| \leq c_i$.

By the tower rule,

$$\mathbb{E}\left[\exp\left(\theta \sum_{i=1}^n Z_i\right)\right] = \mathbb{E}\left[\exp\left(\theta \sum_{i=1}^{n-1} Z_i\right) \mathbb{E}\left[\exp(\theta Z_n) \mid \mathbf{X}_{1:(n-1)}\right]\right].$$

We can then apply Hoeffding's inequality (conditionally) to conclude there exists a universal constant C such that

$$\mathbb{E}\left[\exp\left(\theta \sum_{i=1}^{n-1} Z_i\right) \mathbb{E}\left[\exp(\theta Z_n) \mid \mathbf{X}_{1:(n-1)}\right]\right] \leq e^{C\theta^2 c_i^2} \mathbb{E}\left[\exp\left(\theta \sum_{i=1}^{n-1} Z_i\right)\right].$$

We apply this process repeatedly and conclude that

$$\mathbb{E}\left[\exp(\theta(f(\mathbf{X}_{1:n}) - \mathbb{E}[f(\mathbf{X}_{1:n})]))\right] \leq \exp(C\theta^2 \|\mathbf{c}\|_2^2).$$

This completes the proof. □

6. SOME REMARKS

The theory of subGaussian random variables is incredibly rich. I highly encourage you to read the "Equivalent Characterizations of SubGaussian Random Variables" in [Wai19]. While I don't often use this theorem directly, it's very good for developing intuition about when something might subGaussian.

Orlicz norms are a handy tool (in my opinion) for manipulating potentially non-independent random variables using reasoning based on norms. While we focus on the $\|\cdot\|_\Psi$, Orlicz norms for other functions are often necessary when the tails don't have subGaussian behavior. Other common

Orlicz Norms used in the literature correspond to the functions $\Psi(t) = \exp(t)$ or $\Psi(t) = t^p$ for some $p \geq 1$. We won't use these other norms in this class.

Not all random variables are subGaussian. A much richer class of random variables are those with finite moment generating function, and such random variables are described subExponential or subGamma. [BLM13] treats subGamma random variables fairly well, and shows why you sometimes need to consider them. Essentially, if you start ending up with tail bounds that depend on $\log(1/\delta)$ instead of $\sqrt{1/\log(\delta)}$, you probably need to work with subGamma random variables.

There are *many* generalizations of McDiarmid's inequality. [BLM13] has some good ones, but there are still new ones being published in modern conferences (e.g. [MP21]). One of the ones I've found most useful in my research is [Com15]. The proof is particularly simple and beautiful (and short!). [Wai19] gives some excellent examples of how to use McDiarmid's inequality.

REFERENCES

- [BLM13] Stéphane Boucheron, Gábor Lugosi, and Pascal Massart. *Concentration inequalities: A nonasymptotic theory of independence*. Oxford university press, 2013.
- [Com15] Richard Combes. “An extension of McDiarmid’s inequality”. In: *arXiv preprint arXiv:1511.05240* (2015).
- [MP21] Andreas Maurer and Massimiliano Pontil. “Concentration inequalities under sub-Gaussian and sub-exponential conditions”. In: *Advances in Neural Information Processing Systems* 34 (2021), pp. 7588–7597.
- [Wai19] Martin J Wainwright. *High-Dimensional Statistics: A Non-Asymptotic Viewpoint*. Vol. 48. Cambridge University Press, 2019.

Lecture 3: Towards Performance Bounds: Packing

Instructor: Vishal Gupta

guptavis@usc.edu

Disclaimer: These notes may contain typos and errors. Please do not distribute without written permission of the Instructor.

1. PERFORMANCE BOUND: SAA FOR FINITE FEASIBLE REGIONS

Using the tools we've developed so far, we prove a performance bound on SAA. Recall, our target optimization is

$$\mathbf{x}^* \in \operatorname{argmin}_{\mathbf{x} \in \mathcal{X}} \mathbb{E} [c(\mathbf{x}, \boldsymbol{\xi})]$$

where $\boldsymbol{\xi} \sim \mathbb{P}$, and \mathbb{P} is unknown. We have access to data $\hat{\boldsymbol{\xi}}_j \sim \mathbb{P}$ i.i.d. for $j = 1, \dots, n$. We consider the SAA policy

$$\hat{\mathbf{x}}_n \in \operatorname{argmin}_{\mathbf{x} \in \mathcal{X}} \frac{1}{n} \sum_{j=1}^n c(\mathbf{x}, \hat{\boldsymbol{\xi}}_j).$$

We prove the following theorem from our first lecture.

Theorem 3.1 (SAA for Finite Feasible Regions). *Suppose $|\mathcal{X}| < \infty$ and that there exists $c_{\max} < \infty$ such that $\max_{\mathbf{x} \in \mathcal{X}} |c(\mathbf{x}, \boldsymbol{\xi})| \leq c_{\max}$ almost surely. Then, for any $0 < \delta < 1$, with probability at least $1 - \delta$,*

$$0 \leq \mathbb{E} [c(\hat{\mathbf{x}}_n, \boldsymbol{\xi}) \mid \hat{\boldsymbol{\xi}}_{1:n}] - \mathbb{E} [c(\mathbf{x}^*, \boldsymbol{\xi})] \lesssim c_{\max} \sqrt{\frac{\log(|\mathcal{X}|) \log(2/\delta)}{n}}.$$

Proof of Theorem 3.1. By Lemma “Bounding Suboptimality by Uniform Error” from Lecture 1, it suffices to bound

$$\max_{\mathbf{x} \in \mathcal{X}} \left| \frac{1}{n} \sum_{j=1}^n c(\mathbf{x}, \hat{\boldsymbol{\xi}}_j) - \mathbb{E} [c(\mathbf{x}, \boldsymbol{\xi})] \right| = \max_{\mathbf{x} \in \mathcal{X}} \left| \frac{1}{n} \sum_{j=1}^n \left(c(\mathbf{x}, \hat{\boldsymbol{\xi}}_j) - \mathbb{E} [c(\mathbf{x}, \hat{\boldsymbol{\xi}}_j)] \right) \right|.$$

Fix some $\mathbf{x} \in \mathcal{X}$. Note that for each j , $|c(\mathbf{x}, \hat{\boldsymbol{\xi}}_j)| \lesssim c_{\max}$. Hence, by Hoeffding's inequality, $c(\mathbf{x}, \hat{\boldsymbol{\xi}}_j) - \mathbb{E} [c(\mathbf{x}, \hat{\boldsymbol{\xi}}_j)]$ is subGaussian with variance proxy $\lesssim c_{\max}^2$. Since the j are independent, it follows that $\frac{1}{n} \sum_{j=1}^n c(\mathbf{x}, \hat{\boldsymbol{\xi}}_j) - \mathbb{E} [c(\mathbf{x}, \hat{\boldsymbol{\xi}}_j)]$ is subGaussian with variance proxy $\lesssim c_{\max}^2/n$. (Why?)

By the relation between subGaussian variables and the Ψ -Orlicz norm,

$$\left\| \frac{1}{n} \sum_{j=1}^n c(\mathbf{x}, \hat{\boldsymbol{\xi}}_j) - \mathbb{E} [c(\mathbf{x}, \hat{\boldsymbol{\xi}}_j)] \right\|_{\Psi} \lesssim c_{\max} / \sqrt{n}.$$

From the maxima of Orlicz Norms, it follows that

$$\left\| \max_{\mathbf{x} \in \mathcal{X}} \left| \frac{1}{n} \sum_{j=1}^n \left(c(\mathbf{x}, \hat{\boldsymbol{\xi}}_j) - \mathbb{E} \left[c(\mathbf{x}, \hat{\boldsymbol{\xi}}_j) \right] \right) \right| \right\|_{\Psi} \lesssim c_{\max} \sqrt{\frac{\log(|\mathcal{X}|)}{n}}.$$

The result then follows from the tail bound of a r.v. with finite Ψ -norm. \square

It's worth stressing how short and simple the above proof is. You should compare it to the proof in [KSH02] which establishes (qualitatively) similar results.

1.1. An Alternate Proof Based on Symmetrization. Before proceeding, we present a different proof of the above theorem. This alternate version more closely mirrors traditional proofs in the literature and paves the way for more advanced results in the course.

Before establishing the proof, we establish the following auxiliary lemma that is of independent interest.

Lemma 3.2 (Symmetrization Lemma). *For any convex, increasing function $\Phi(\cdot)$,*

$$\mathbb{E} \left[\Phi \left(\max_{\mathbf{x} \in \mathcal{X}} \left| \frac{1}{n} \sum_{j=1}^n \left(c(\mathbf{x}, \hat{\boldsymbol{\xi}}_j) - \mathbb{E} \left[c(\mathbf{x}, \hat{\boldsymbol{\xi}}_j) \right] \right) \right| \right) \right] \leq \mathbb{E} \left[\Phi \left(2 \max_{\mathbf{c} \in \mathcal{F}} \left| \frac{1}{n} \sum_{j=1}^n \eta_j c_j \right| \right) \right], \quad (3-11)$$

where η_j for $j = 1, \dots, n$ are independent Rademacher random variables and

$$\mathcal{F} = \mathcal{F}(\hat{\boldsymbol{\xi}}_{1:n}) = \{ (c(\mathbf{x}, \hat{\boldsymbol{\xi}}_1), \dots, c(\mathbf{x}, \hat{\boldsymbol{\xi}}_n))^{\top} : \mathbf{x} \in \mathcal{X} \} \subseteq \mathbb{R}^n.$$

Proof of Lemma. Let $\bar{\boldsymbol{\xi}}_j$ for $j = 1, \dots, n$ be a second i.i.d. sample from \mathbb{P} and fix the function $\Phi(\cdot)$. Then,

$$\begin{aligned} & \mathbb{E} \left[\Phi \left(\max_{\mathbf{x} \in \mathcal{X}} \left| \frac{1}{n} \sum_{j=1}^n \left(c(\mathbf{x}, \hat{\boldsymbol{\xi}}_j) - \mathbb{E} \left[c(\mathbf{x}, \hat{\boldsymbol{\xi}}_j) \right] \right) \right| \right) \right] \\ &= \mathbb{E} \left[\Phi \left(\max_{\mathbf{x} \in \mathcal{X}} \left| \frac{1}{n} \sum_{j=1}^n \left(c(\mathbf{x}, \hat{\boldsymbol{\xi}}_j) - \mathbb{E} \left[c(\mathbf{x}, \bar{\boldsymbol{\xi}}_j) \mid \hat{\boldsymbol{\xi}}_{1:n} \right] \right) \right| \right) \right] \\ &\leq \mathbb{E} \left[\Phi \left(\max_{\mathbf{x} \in \mathcal{X}} \left| \frac{1}{n} \sum_{j=1}^n \left(c(\mathbf{x}, \hat{\boldsymbol{\xi}}_j) - c(\mathbf{x}, \bar{\boldsymbol{\xi}}_j) \right) \right| \right) \right] \quad (\text{Jensen's Inequality}). \end{aligned}$$

$$\mathbb{E} \left[\Phi \left(\max_{\mathbf{x} \in \mathcal{X}} \left| \frac{1}{n} \sum_{j=1}^n \left(c(\mathbf{x}, \hat{\boldsymbol{\xi}}_j) - c(\mathbf{x}, \bar{\boldsymbol{\xi}}_j) \right) \right| \right) \right] \leq \mathbb{E} \left[\Phi \left(\max_{\mathbf{x} \in \mathcal{X}} \left| \frac{1}{n} \sum_{j=1}^n \eta_j \left(c(\mathbf{x}, \hat{\boldsymbol{\xi}}_j) - c(\mathbf{x}, \bar{\boldsymbol{\xi}}_j) \right) \right| \right) \right]$$

Now, focus on the argument of $\Phi(\cdot)$. Observe

$$\begin{aligned} \max_{\mathbf{x} \in \mathcal{X}} \left| \frac{1}{n} \sum_{j=1}^n \eta_j \left(c(\mathbf{x}, \hat{\boldsymbol{\xi}}_j) - c(\mathbf{x}, \bar{\boldsymbol{\xi}}_j) \right) \right| &\leq \max_{\mathbf{x} \in \mathcal{X}} \left| \frac{1}{n} \sum_{j=1}^n \eta_j(\mathbf{x}, \hat{\boldsymbol{\xi}}_j) \right| + \max_{\mathbf{x} \in \mathcal{X}} \left| \frac{1}{n} \sum_{j=1}^n \eta_j(\mathbf{x}, \bar{\boldsymbol{\xi}}_j) \right| \\ &= \frac{1}{2} \left(2 \max_{\mathbf{x} \in \mathcal{X}} \left| \frac{1}{n} \sum_{j=1}^n \eta_j(\mathbf{x}, \hat{\boldsymbol{\xi}}_j) \right| + 2 \max_{\mathbf{x} \in \mathcal{X}} \left| \frac{1}{n} \sum_{j=1}^n \eta_j(\mathbf{x}, \bar{\boldsymbol{\xi}}_j) \right| \right) \end{aligned}$$

Hence, increasing the argument of $\Phi(\cdot)$ and using convexity shows that

$$\begin{aligned} &\mathbb{E} \left[\Phi \left(\max_{\mathbf{x} \in \mathcal{X}} \left| \frac{1}{n} \sum_{j=1}^n \eta_j \left(c(\mathbf{x}, \hat{\boldsymbol{\xi}}_j) - c(\mathbf{x}, \bar{\boldsymbol{\xi}}_j) \right) \right| \right) \right] \\ &\leq \frac{1}{2} \mathbb{E} \left[\Phi \left(2 \max_{\mathbf{x} \in \mathcal{X}} \left| \frac{1}{n} \sum_{j=1}^n \eta_j c(\mathbf{x}, \hat{\boldsymbol{\xi}}_j) \right| \right) \right] + \frac{1}{2} \mathbb{E} \left[\Phi \left(2 \max_{\mathbf{x} \in \mathcal{X}} \left| \frac{1}{n} \sum_{j=1}^n \eta_j c(\mathbf{x}, \bar{\boldsymbol{\xi}}_j) \right| \right) \right] \\ &= \mathbb{E} \left[\Phi \left(2 \max_{\mathbf{x} \in \mathcal{X}} \left| \frac{1}{n} \sum_{j=1}^n \eta_j c(\mathbf{x}, \hat{\boldsymbol{\xi}}_j) \right| \right) \right]. \end{aligned}$$

This completes the proof. \square

Lemma 3.2 is used regularly in the literature when analyzing sums of independent random variables. You should study the set \mathcal{F} defined in the lemma. This set is random (it depends on the data) and also depends on the particular cost function and feasible region.

IMPORTANT: Our lemma did not require that $|\mathcal{X}| < \infty!$

Equipped with the Lemma 3.2, we provide our alternate proof of Theorem 3.1.

Alternate Proof of Theorem 3.1. We again focus on bounding the Ψ Orlicz norm of

$$\max_{\mathbf{x} \in \mathcal{X}} \left| \frac{1}{n} \sum_{j=1}^n c(\mathbf{x}, \hat{\boldsymbol{\xi}}_j) - \mathbb{E} \left[c(\mathbf{x}, \hat{\boldsymbol{\xi}}_j) \right] \right|.$$

Motivated by Eq. (3-11), we'll first study

$$\max_{\mathbf{c} \in \mathcal{F}(\hat{\boldsymbol{\xi}}_{1:n})} \left| \frac{1}{n} \mathbf{c}^\top \boldsymbol{\eta} \right|.$$

There are two sources of randomness here: $\hat{\boldsymbol{\xi}}_{1:n}$ and the Rademacher variables $\eta_{1:n}$. We first argue conditionally, conditioning on $\hat{\boldsymbol{\xi}}_{1:n}$.

Specifically, for any fixed \mathbf{c} , $\frac{1}{n} \boldsymbol{\eta}^\top \mathbf{c}$ is subGaussian and $\|\frac{1}{n} \boldsymbol{\eta}^\top \mathbf{c}\|_\Psi \lesssim \frac{1}{n} \|\mathbf{c}\|_2$. (Why?)

Hence, from our result on the Orlicz norm of the max, it follows that the *conditional* Orlicz norm satisfies

$$\left\| \max_{\mathbf{c} \in \mathcal{F}} \left| \frac{1}{n} \sum_{j=1}^n \eta_j c_j \right| \right\|_{\Psi|\hat{\boldsymbol{\xi}}_{1:n}} \lesssim \frac{\sqrt{\log |\mathcal{F}|}}{n} \cdot \max_{\mathbf{c} \in \mathcal{F}} \|\mathbf{c}\|_2 \lesssim c_{\max} \sqrt{\frac{\log |\mathcal{F}|}{n}}.$$

Said differently, there exists a universal constant K such that

$$\left\| \max_{\mathbf{c} \in \mathcal{F}} \left| \frac{1}{n} \sum_{j=1}^n \eta_j c_j \right| \right\|_{\Psi|\hat{\boldsymbol{\xi}}_{1:n}} = K c_{\max} \sqrt{\frac{\log |\mathcal{F}|}{n}}.$$

For convenience, let $J = J(\hat{\boldsymbol{\xi}}_{1:n}) = K c_{\max} \sqrt{\frac{\log |\mathcal{F}|}{n}}$. Then, taking an expectation of the definition of the Orlicz Norm, we've thus far shown

$$\mathbb{E} \left[\exp \left(\frac{\max_{\mathbf{c} \in \mathcal{F}} \left(\frac{1}{n} \sum_{j=1}^n \eta_j c_j \right)^2}{J(\hat{\boldsymbol{\xi}}_{1:n})^2} \right) \right] \leq 5. \quad (3-12)$$

It would be great at this point if we could apply Lemma 3.2 with the convex function $t \mapsto \exp\left(\frac{t^2}{2J^2}\right)$ and then use Eq. (3-12) to finish the proof. Unfortunately, we can't do this because $J = J(\hat{\boldsymbol{\xi}}_{1:n})$ is random. However, since

$$|\mathcal{F}| \leq |\mathcal{X}|, \quad (3-13)$$

J is upper bounded by $K c_{\max} \sqrt{\frac{\log |\mathcal{X}|}{n}}$. Hence, we can instead take

$$\Phi(t) \equiv \exp \left(\frac{t^2}{2K^2 c_{\max}^2 \frac{\log |\mathcal{X}|}{n}} \right).$$

Then, applying Lemma 3.2 with this convex, increasing function, and using Eq. (3-12) shows

$$\mathbb{E} \left[\Phi \left(\max_{\mathbf{x} \in \mathcal{X}} \left| \frac{1}{n} \sum_{j=1}^n \left(c(\mathbf{x}, \hat{\boldsymbol{\xi}}_j) - \mathbb{E} [c(\mathbf{x}, \hat{\boldsymbol{\xi}}_j)] \right) \right| \right) \right] \leq 5.$$

Applying Markov's Inequality completes the proof. \square

The above proof might look a lot more complicated, but really the only "hard" part is establishing Lemma 3.2.

Philosophically, the advantage of the above proof is we replaced the analysis of a (potentially complicated) random variable $c(\mathbf{x}, \hat{\boldsymbol{\xi}}_j)$ with analysis of i) Rademacher random variables η_j and ii) the set $\mathcal{F}(\hat{\boldsymbol{\xi}}_{1:n})$. The Rademacher random variables are simple to analyze, and in most arguments of this form, understanding $\mathcal{F}(\hat{\boldsymbol{\xi}}_{1:n})$ comes down to geometry instead of probability.

Practically, the above proof really relies on the size of the set \mathcal{F} , not the size of the set \mathcal{X} . In some cases, it's easier to analyze one set rather than the other. On your homework, you'll see an example where $|\mathcal{X}| = \infty$, but $|\mathcal{F}| < \infty$, so in fact you *have* to analyze things with the second proof.

The downside of the alternate proof is that Eq. (3-11) *heavily* leverages the fact that we're looking at a sum of independent random variables. By contrast, our first proof just required us to manipulate tail bounds and might be adapted more easily to other settings.

2. BEYOND FINITE FEASIBLE REGIONS: PACKING IN METRIC SPACES

A lot of the theory of uniform laws of large numbers involves generalizing the machinery of the above proofs to the case where \mathcal{X} (or \mathcal{F}) is not finite by approximating it by a finite subset. To this end, we introduce the notion of a packing and covering in Euclidean spaces.

Definition 3.3 (Packing Number). *Given a set $\mathcal{F} \subseteq \mathbb{R}^n$ and a norm $\|\cdot\|$ on \mathbb{R}^n , the ϵ packing number of \mathcal{F} (denoted $D(\epsilon, \mathcal{F}, \|\cdot\|)$) is the maximum number of points in \mathcal{F} such that every pair of points is at least ϵ far apart with respect to $\|\cdot\|$.*

Definition 3.4 (Covering Number). *Given a set $\mathcal{F} \subseteq \mathbb{R}^n$ and a norm $\|\cdot\|$ on \mathbb{R}^n , the ϵ -covering number of \mathcal{F} (denoted $N(\epsilon, \mathcal{F}, \|\cdot\|)$) is the smallest number of closed balls of radius ϵ whose union contains \mathcal{F} .*

Exercise: Draw a picture in your notes to understand the definitions.

It's possible to generalize both definitions to general metric spaces, but we won't need it. In this course, we'll mostly be concerned with packings, but the two numbers are closely related. Indeed, you should prove to yourself that for any set \mathcal{F} ,

$$N(\epsilon, \mathcal{F}, \|\cdot\|) \leq D(\epsilon, \mathcal{F}, \|\cdot\|) \leq N(\epsilon/2, \mathcal{F}, \|\cdot\|).$$

Example 3.5 (Packing an Interval). *What is the ϵ packing number of $[0, 1]$ with respect to ℓ_2 norm? Are you sure?*

Exercise: Let \mathcal{F} be a compact subset of \mathbb{R}^n . Draw a plot where the x -axis is ϵ and the y -axis is $D(\epsilon, \mathcal{F}, \|\cdot\|_2)$.

In our proofs, we will mostly use the ℓ_2 packing number. You'll prove some useful facts about packing numbers on your homework. For now, one very important fact is the following:

Theorem 3.6 (Packing of Euclidean Balls). *The ϵ packing number of a ball of radius $R > \epsilon$ in \mathbb{R}^d with respect to ℓ_2 satisfies*

$$D(\epsilon, B(\mathbf{0}, R), \|\cdot\|_2) \leq \left(\frac{3R}{\epsilon}\right)^d.$$

Proof. Consider some packing of $B(\mathbf{0}, R)$ and denote the points as f_1, \dots, f_m where $m = D(\epsilon, B(\mathbf{0}, R), \|\cdot\|_2)$. Now consider balls of radius $\epsilon/2$ centered at each f_i . By construction, these balls are disjoint. Their total volume is $Km2^{-d}\epsilon^d$ for some constant K .

Now each f_i lies within a distance R of f_1 because the f_i are in $B(\mathbf{0}, R)$. Hence, the ball $B(f_i, \epsilon/2)$ lies within the larger ball $B(f_1, 3R/2)$ for each i . This larger ball has volume $KR^d 1.5^d$. Comparing completes the proof. \square

In class, we'll leverage packing numbers to extend our proof of SAA to non-finite feasible regions with small packing numbers.

REFERENCES

- [KSH02] Anton J Kleywegt, Alexander Shapiro, and Tito Homem-de-Mello. “The sample average approximation method for stochastic discrete optimization”. In: *SIAM Journal on Optimization* 12.2 (2002), pp. 479–502.

Lecture 4: More Uniform Laws of Large Numbers

Instructor: Vishal Gupta

guptavis@usc.edu

Disclaimer: *These notes may contain typos and errors. Please do not distribute without written permission of the Instructor.*

1. REVIEW

Let's recap: So far, we've shown how analyzing a data-driven method based on a proxy-objective reduces to studying the uniform approximation error between the true objective and proxy-objective. In the case of SAA, this amounted to studying the uniform error between a sample average and a true expectation over the feasible region.

In this lecture, we're going to push this idea further to develop stronger bounds on this uniform error. These will translate into other performance guarantees for SAA.

2. SAA WITH BOUNDED FEASIBLE REGIONS AND LIPSCHITZ COST FUNCTIONS

As a first example, we'll study our usual stochastic optimization

$$\mathbf{x}^* \in \operatorname{argmin}_{\mathbf{x} \in \mathcal{X}} \mathbb{E} [c(\mathbf{x}, \boldsymbol{\xi})]$$

under the assumptions that $\mathbf{x} \mapsto c(\mathbf{x}, \boldsymbol{\xi})$ is L -Lipschitz almost surely, and $\mathcal{X} \subseteq \mathbb{R}^d$ is a bounded feasible region.

Definition 4.1. *A set $\mathcal{X} \subseteq \mathbb{R}^d$ is bounded if $R \equiv \frac{1}{2} \sup_{\mathbf{x}, \bar{\mathbf{x}} \in \mathcal{X}} \|\mathbf{x} - \bar{\mathbf{x}}\| < \infty$. In that case, we say R is radius of \mathcal{X} .*

Recall from our earlier discussion on packing, $\log D(\epsilon, \mathcal{X}) \leq d \log(3R/\epsilon)$. Finally, let

$$\hat{\mathbf{x}}_n \in \operatorname{argmin}_{\mathbf{x} \in \mathcal{X}} \frac{1}{n} \sum_{j=1}^n c(\mathbf{x}, \hat{\boldsymbol{\xi}}_j)$$

be the SAA solution.

We then have the following theorem:

Theorem 4.2 (SAA for Bounded Feasible Regions, Lipschitz Cost Functions). *Suppose \mathcal{X} is bounded with radius $R < \infty$, and $\mathbf{x} \mapsto c(\mathbf{x}, \boldsymbol{\xi})$ is L -Lipschitz almost surely. Further, assume $\max_{\mathbf{x} \in \mathcal{X}} |c(\mathbf{x}, \boldsymbol{\xi})| \leq c_{\max}$ almost surely and $n \geq 2$.*

Then,

$$\left\| \max_{\mathbf{x} \in \mathcal{X}} \left| \frac{1}{n} \sum_{j=1}^n c(\mathbf{x}, \hat{\boldsymbol{\xi}}_j) - \mathbb{E} [c(\mathbf{x}, \hat{\boldsymbol{\xi}}_j)] \right| \right\|_{\Psi} \lesssim (L + c_{\max} \sqrt{\log(3R+2)}) \sqrt{\frac{d \log n}{n}}.$$

4-1

Consequently, for any $0 < \delta < 1$, with probability at least $1 - \delta$,

$$0 \leq \mathbb{E} \left[c(\hat{\mathbf{x}}_n, \boldsymbol{\xi}) \mid \hat{\boldsymbol{\xi}}_{1:n} \right] - \mathbb{E} [c(\mathbf{x}^*, \boldsymbol{\xi})] \leq (L + c_{\max} \sqrt{\log(3R + 2)}) \sqrt{\frac{d \log n}{n}} \cdot \sqrt{\log(2/\delta)}.$$

Proof. By the Symmetrization Lemma, it will suffice to bound

$$\frac{1}{n} \left\| \max_{\mathbf{x} \in \mathcal{X}} \mathbf{c}(\mathbf{x})^\top \boldsymbol{\eta} \right\|_{\Psi}$$

where $\mathbf{c}(\mathbf{x}) \equiv (c(\mathbf{x}, \hat{\boldsymbol{\xi}}_1), \dots, c(\mathbf{x}, \hat{\boldsymbol{\xi}}_n))^\top \in \mathbb{R}^n$, and $\boldsymbol{\eta} \in \mathbb{R}^n$ is a vector of i.i.d. Rademacher random variables. (Convince yourself this is true.)

Now fix some $0 < \epsilon < 1$ (to be specified later), and let $\mathcal{X}_0 = \{\mathbf{x}_1, \dots, \mathbf{x}_m\}$ be an ϵ -packing of \mathcal{X} . By our earlier remarks,

$$\log(m) \leq d \log(3R/\epsilon) \lesssim d \log(e + R) \log(2/\epsilon).$$

Then, for any $\mathbf{x} \in \mathcal{X}$, there exists a $\mathbf{x}_k \in \mathcal{X}_0$ such that $\|\mathbf{x} - \mathbf{x}_k\| \leq \epsilon$. (Why?) For such a k , we have that

$$\begin{aligned} \|\mathbf{c}(\mathbf{x}) - \mathbf{c}(\mathbf{x}_k)\|^2 &= \sum_{j=1}^n (c(\mathbf{x}, \hat{\boldsymbol{\xi}}_j) - c(\mathbf{x}_k, \hat{\boldsymbol{\xi}}_j))^2 \\ &\leq \sum_{j=1}^n L^2 \|\mathbf{x} - \mathbf{x}_k\|^2 \\ &\leq L^2 \epsilon^2 n. \end{aligned}$$

This implies $\|\mathbf{c}(\mathbf{x}) - \mathbf{c}(\mathbf{x}_k)\|_2 \leq L\epsilon\sqrt{n}$.

Now write,

$$\begin{aligned} \mathbf{c}(\mathbf{x})^\top \boldsymbol{\eta} &= (\mathbf{c}(\mathbf{x}) - \mathbf{c}(\mathbf{x}_k))^\top \boldsymbol{\eta} + \mathbf{c}(\mathbf{x}_k)^\top \boldsymbol{\eta} \\ &\leq L\epsilon\sqrt{n}\|\boldsymbol{\eta}\| + \mathbf{c}(\mathbf{x}_k)^\top \boldsymbol{\eta} \quad (\text{Cauchy-Schwarz}) \\ &\leq L\epsilon\sqrt{n}\|\boldsymbol{\eta}\| + \max_{i=1, \dots, m} \mathbf{c}(\mathbf{x}_i)^\top \boldsymbol{\eta}. \end{aligned}$$

Notice, $\|\boldsymbol{\eta}\|_2 = \sqrt{n}$ almost surely. Hence, if we take the maximum over $\mathbf{x} \in \mathcal{X}$ of both sides we have that

$$\max_{\mathbf{x} \in \mathcal{X}} \mathbf{c}(\mathbf{x})^\top \boldsymbol{\eta} \leq L\epsilon n + \max_{i=1, \dots, m} \mathbf{c}(\mathbf{x}_i)^\top \boldsymbol{\eta}.$$

Now take the conditional Ψ Orlicz norm (conditional on $\hat{\boldsymbol{\xi}}_{1:n}$) of both sides. For each i , $\|\mathbf{c}(\mathbf{x}_i)^\top \boldsymbol{\eta}\|_{\Psi|\hat{\boldsymbol{\xi}}_{1:n}} \lesssim \|\mathbf{c}(\mathbf{x}_i)\|_2 \lesssim c_{\max}\sqrt{n}$. Hence, using our result for the Orlicz Norm of the max, we have that

$$\left\| \max_{\mathbf{x} \in \mathcal{X}} \mathbf{c}(\mathbf{x})^\top \boldsymbol{\eta} \right\|_{\Psi|\hat{\boldsymbol{\xi}}_{1:n}} \lesssim L\epsilon n + c_{\max} \sqrt{n \log m} \lesssim L\epsilon n + c_{\max} \sqrt{nd \log(3R) \log(2/\epsilon)}.$$

We are now free to choose $\epsilon > 0$ to optimize the bound. Optimizing it exactly seems tricky, so instead we pick a (suboptimal) choice of $\epsilon = \frac{1}{\sqrt{n}}$. This yields,

$$\begin{aligned} \left\| \max_{\mathbf{x} \in \mathcal{X}} \mathbf{c}(\mathbf{x})^\top \boldsymbol{\eta} \right\|_{\Psi|\hat{\boldsymbol{\xi}}_{1:n}} &\lesssim L\sqrt{n} + c_{\max} \sqrt{nd \log(e+R) \log(2\sqrt{n})} \\ &\lesssim (L + c_{\max} \sqrt{\log(e+R)}) \sqrt{\frac{d \log n}{n}}. \end{aligned}$$

By writing out the definition of the conditional Ψ -norm above and taking an expectation with respect to $\hat{\boldsymbol{\xi}}_{1:n}$, we find the unconditional Ψ -norm above is also bounded by the same quantity.

Applying the Symmetrization Lemma then shows

$$\left\| \max_{\mathbf{x} \in \mathcal{X}} \left| \frac{1}{n} \sum_{j=1}^n c(\mathbf{x}, \hat{\boldsymbol{\xi}}_j) \right| \right\|_{\Psi} \leq (L + c_{\max} \sqrt{\log(e+R)}) \sqrt{\frac{d \log n}{n}}.$$

This proves the first statement. The final statement follows from Markov's Inequality. \square

Some remarks: It's worth noting that we can view the approximating finite set \mathcal{X}_0 either as a packing of \mathcal{X} , or we can view the induced set

$$\mathcal{F}_0 = \left\{ \left(c(\mathbf{x}, \hat{\boldsymbol{\xi}}_1), \dots, c(\mathbf{x}, \hat{\boldsymbol{\xi}}_n) \right)^\top : \mathbf{x} \in \mathcal{X}_0 \right\}$$

as a $L\epsilon\sqrt{n}$ -covering of $\mathcal{F}(\hat{\boldsymbol{\xi}}_{1:n})$. This second perspective is perhaps more illuminating for the results below.

Secondly, the assumptions on \mathcal{X} weren't strictly necessary. What we really needed was that the packing numbers of \mathcal{X} didn't grow too rapidly as $\epsilon \rightarrow 0$. In the above case, they looked roughly like ϵ^{-d} , and this was sufficient so long as $d \ll n$. We could make other assumptions to induce "small" packing numbers for \mathcal{X} , or, equivalently, small packing numbers for \mathcal{F} .

3. MORE ADVANCED RESULTS

As is hopefully clear by now, many performance guarantees come down to bounding the maximal deviation of a quantity like

$$Z = Z(\hat{\boldsymbol{\xi}}_{1:n}) = \sup_{t \in \mathcal{T}} \left| \sum_{j=1}^n f_j(t, \hat{\boldsymbol{\xi}}_j) - \mathbb{E} \left[f_j(t, \hat{\boldsymbol{\xi}}_j) \right] \right| \quad (4-14)$$

for some functions f_j and independent random variables $\hat{\boldsymbol{\xi}}_j$. Below we quote some stronger results from [Pol90] that bound the tails of Z when the $\hat{\boldsymbol{\xi}}_j$ are independent, but not necessarily identically distributed. We won't prove these results – they rely on some ideas we haven't talked about like chaining – but they're useful to know. When you write a paper, you should use these results from [Pol90].

To state the results, first, recall that symmetrization encouraged us to think about the (random) set

$$\mathcal{F} = \mathcal{F}(\hat{\boldsymbol{\xi}}_{1:n}) = \left\{ \left(f_1(t, \hat{\boldsymbol{\xi}}_1), \dots, f_n(t, \hat{\boldsymbol{\xi}}_n) \right)^\top : t \in \mathcal{T} \right\}.$$

We define the entropy integral of \mathcal{F} by

$$J = J(\hat{\xi}_{1:n}) \equiv 9 \int_0^\delta \sqrt{\log D(\epsilon, \mathcal{F})} d\epsilon \quad \text{where } \delta = \delta(\hat{\xi}_{1:n}) = \sup_{\mathbf{f} \in \mathcal{F}} \|\mathbf{f}\|_2.$$

On first reading, entropy integrals are intimidating. Remember, though, that this integral is not defined probabilistically. For a fixed $\hat{\xi}_{1:n}$, it's just an ordinary integral from calculus. Intuitively, packing numbers measure the complexity of \mathcal{F} at scale ϵ . The entropy integral measures complexity at varying scales, and, if the packing numbers “blow up” too fast as $\epsilon \rightarrow 0$, then this integral diverges.

Using a standard change of variables from calculus,

$$J = 9\delta \int_0^1 \sqrt{\log D(\delta\epsilon, \mathcal{F})} d\epsilon.$$

Perhaps surprisingly, although this integral depends on the random set \mathcal{F} and random “radius” δ , we can often upper bound it by a deterministic constant, so the only randomness in J comes from the leading δ term.

We give one such example below:

Lemma 4.3 (Entropy Integral for Euclidean Classes). *Suppose that \mathcal{F} is Euclidean, i.e., there exists constants A and W such that with probability 1,*

$$D(\epsilon\delta, \mathcal{F}) \leq A\epsilon^{-W} \quad 0 < \epsilon < 1.$$

Then,

$$J \lesssim \delta \cdot \frac{\log A + W/2}{\sqrt{\log A}}.$$

Proof. You'll prove this on your homework. (It's just an exercise in evaluating/bounding an integral!) If you're really struggling, look at Theorem A.2 in [GR21], but try your best before looking it up. \square

Regardless of whether \mathcal{F} is Euclidean, we get nice bounds on Z in terms of bounds on J .

Theorem 4.4 (Bounds on Suprema of Empirical Process). *Fix any $0 < \beta < 1$.*

- (1) *(When J bounded by a constant) Suppose there exists a deterministic constant K such that $J \leq K$ almost surely. Then, with probability at least $1 - \beta$,*

$$Z \lesssim K \sqrt{\log(2/\beta)}.$$

- (2) *(When J is SubGaussian) Suppose that $\|J\|_\Psi < \infty$. Then, with probability at least $1 - \beta$,*

$$Z \lesssim \|J\|_\Psi \log(2/\beta).$$

- (3) *(When J has finite p^{th} moment) Suppose that $\|J\|_{\mathcal{L}^p} \equiv \mathbb{E}[J^p]^{1/p} < \infty$. Then, with probability at least $1 - \beta$,*

$$Z \lesssim \sqrt{p}\beta^{-1/p} \|J\|_{\mathcal{L}^p}.$$

Of course, when \mathcal{F} is Euclidean we can simplify the results above appropriately. In particular, in that case the bounds above really depend on the behavior of δ , i.e, how big do the largest elements of the process get? I'll leave the details to you...

Let's step back: What have we done? The above theorem is a general purpose tool. It's important to us because we can often reduce the analysis of a performance guarantee for an optimization method to studying the tails of a random variable like Z . and then we can use the above theorem to get a bound.

You'll follow this recipe on your homework to reprove a stronger performance guarantee under the setting of Theorem 4.2. If you do it right, you'll find the result is only really stronger by a logarithmic factor.... And this insight is important (to me). Namely, if you use some of the strongest methods available in empirical process, you often only get small improvements over the more naive (by-hand) approaches. And it's often more straightforward to adapt these "by-hand" approaches when your problem doesn't *exactly* match the setting of the theorems in a textbook. So, don't be afraid to do things by hand.

4. OPTIONAL REMARKS

The proof of the above theorem in [Pol90] is very readable (especially if you ignore the section on pseudo-dimension which isn't critical). However, it's somewhat different from the "standard" proofs in other books that generally restrict attention to the case where δ is bounded by a constant. I prefer the extra generality above (since it comes in useful sometimes).

That said, the above theorem isn't really best possible. To get some intuition, let's think about the case where \mathcal{F} is Euclidean. Then, in that case, the tails of Z are basically determined by δ , i.e. the maximal value of the process. This is in some ways analogous to Hoeffding's Inequality: the tails of a bounded random variable are bounded by the maximal size of the support. But we know that for random variables with small variance, Hoeffding's inequality can be quite loose. In the same way, for certain processes with small "variance", the the above theorem can be quite loose. You can make this heuristic argument rigorous by considering the special case when \mathcal{F} is a singleton in Theorem 4.4 and comparing the given result to Bernstein's inequality.

[BLM13] present tighter bounds on the random variable Z that leverage "variance" information to help close the gap. The tightest known bounds are due to Talagrand (Talagrand's Inequality for the Suprema of the Empirical Process) but the proof is pretty dense, leverages a lot of special features, and (to my knowledge) requires $\hat{\xi}_{1:n}$ to be i.i.d, not just independent.

REFERENCES

- [BLM13] Stéphane Boucheron, Gábor Lugosi, and Pascal Massart. *Concentration inequalities: A nonasymptotic theory of independence*. Oxford university press, 2013.
- [GR21] Vishal Gupta and Paat Rusmevichientong. “Small-data, large-scale linear optimization with uncertain objectives”. In: *Management Science* 67.1 (2021), pp. 220–241.
- [Pol90] David Pollard. “Empirical Processes: Theory and Applications”. In: *Ims*. 1990.