

Least Squares and Gaussian Width

April 15, 2025

1 Introduction

1.1 Review

In this week's lectures, we proved a bound on the error of the least squares estimator $\hat{\mu}$ in a convex model.

$$\hat{\mu} = \operatorname{argmin}_{m \in \mathcal{M}} \frac{1}{n} \sum_{i=1}^n \{Y_i - m(X_i)\}^2 \quad \text{where } \mathcal{M} \text{ is a convex set.} \quad (1)$$

To keep things simple, we focused on a stylized gaussian-noise model.

$$Y_i = \mu(X_i) + \varepsilon_i \quad \text{where } \varepsilon_i \stackrel{iid}{\sim} N(0, \sigma^2).$$

And we proved that the following high probability error bound in terms of the gaussian width of a centered neighborhood of μ .

$$\|\hat{\mu} - \mu\|_{L_2(\mathbb{P}_n)} < s \quad \text{w.p. } 1 - \delta \quad \text{if} \quad \frac{s^2}{2\sigma} \geq w(\mathcal{M}_s^\circ - \mu) + s \sqrt{\frac{2\{1 + 2\log(2n)\}}{\delta n}} \quad (2)$$

when $\mu \in \mathcal{M}$.

Here $\mathcal{M}_s^\circ = \{m - \mu : m \in \mathcal{M} \text{ and } \|m - \mu\| = s\}$. Furthermore, we showed that even if μ is not in the model, we have a bound like this on the distance between our estimator ($\hat{\mu}$) and our model's best approximation to the signal (μ_\star).

$$\|\hat{\mu} - \mu_\star\|_{L_2(\mathbb{P}_n)} < s \quad \text{w.p. } 1 - \delta \quad \text{if} \quad \frac{s^2}{2\sigma} \geq w(\mathcal{M}_s^\circ - \mu) + s \sqrt{\frac{2\{1 + 2\log(2n)\}}{\delta n}} \quad (3)$$

for $\mu_\star = \operatorname{argmin}_{m \in \mathcal{M}} \|m - \mu\|_{L_2(\mathbb{P}_n)}$.

Here $\mathcal{M}_s^\circ = \{m - \mu_\star : m \in \mathcal{M} \text{ and } \|m - \mu_\star\| = s\}$. This is a generalization of the previous bound. When $\mu \in \mathcal{M}$, $\mu_\star = \mu$ and (3) is equivalent to (2).

I also claimed (without proof) that this implies the following bound. The

advantage is that the inequality characterizing s is a bit simpler.

$$\begin{aligned} \|\hat{\mu} - \mu_\star\|_{L_2(\mathbb{P}_n)} &< s + 2\sigma \sqrt{\frac{2\{1 + 2\log(2n)\}}{\delta n}} \quad \text{w.p.} \quad 1 - \delta \quad \text{if} \quad \frac{s^2}{2\sigma} \geq w(\mathcal{M}_s - \mu) \\ \text{for} \quad \mu_\star &= \operatorname{argmin}_{m \in \mathcal{M}} \|m - \mu\|_{L_2(\mathbb{P}_n)}. \end{aligned} \tag{4}$$

Here $\mathcal{M}_s = \{m - \mu_\star : m \in \mathcal{M} \text{ and } \|m - \mu_\star\| \leq s\}$.

1.2 Assignment Summary

In this assignment, we'll take a few steps toward getting concrete, meaningful error bounds. Because it's unrealistic to expect real data to look exactly like signal plus gaussian noise, we'll derive a more meaningful version of (3) (and consequently (4)) in Section 3: a bound that holds when $\epsilon_1 \dots \epsilon_n$ are independent with mean zero, but don't have to be gaussian or even all have the same distribution. And because our error bound (3) is a little too abstract to make sense directly, we'll bound the gaussian width of neighborhoods $\mathcal{M}_s - \mu_\star$ in a few models and use the result to derive concrete model-specific error bounds. To prepare for all that, we'll start by proving a few properties of gaussian width. And, for good measure, we'll use them to derive our simplified error bound (4) from the one we proved in lecture (3).

2 Properties of Gaussian Width

In this section, we'll prove a few properties of gaussian width.

$$w(\mathcal{V}) = \mathbb{E} \max_{v \in \mathcal{V}} \langle g, v \rangle \quad \text{for} \quad g_i \stackrel{iid}{\sim} N(0, 1).$$

2.1 Basic Properties

- It's *increasing*. It's a maximum over the set \mathcal{V} , so it gets bigger if \mathcal{V} does.

$$w(\mathcal{V}) \leq w(\mathcal{V}^+) \quad \text{if} \quad \mathcal{V} \subseteq \mathcal{V}^+$$

- It's *homogeneous*. If we scale the vectors in \mathcal{V} , we scale its width.

$$w(\alpha \mathcal{V}) = \alpha w(\mathcal{V}) \quad \text{where} \quad \alpha \mathcal{V} := \{\alpha v : v \in \mathcal{V}\} \quad \text{for} \quad \alpha \geq 0.$$

- It's *translation invariant*. It doesn't care about how we center our vectors.

$$w(\mathcal{V} + x) = w(\mathcal{V}) \quad \text{for} \quad \mathcal{V} + x := \{v + x : v \in \mathcal{V}\}.$$

Exercise 1 Prove that gaussian width $w(\mathcal{V})$ has these three properties.

Solution 1 1. For any vector $g \in \mathbb{R}^n$, if $\mathcal{V} \subseteq \mathcal{V}^+$,

$$\max_{v \in \mathcal{V}} \langle g, v \rangle \leq \max_{v \in \mathcal{V}^+} \langle g, v \rangle.$$

Every value 'considered' in the first maximum is also considered in the second.

2. For any vector $g \in \mathbb{R}^n$,

$$\max_{v \in \alpha \mathcal{V}} \langle g, v \rangle = \max_{v \in \mathcal{V}} \alpha \langle g, v \rangle = \alpha \max_{v \in \mathcal{V}} \langle g, v \rangle.$$

The first equality follows from linearity of the inner product and the second because, if $\alpha \geq 0$, $\max_x \alpha f(x) = \alpha \max_x f(x)$ for any function f . The situation is reversed if $\alpha < 0$: we get α times the minimum of the function f .

3. For any vector $g \in \mathbb{R}^n$,

$$\max_{v \in \mathcal{V} + x} \langle g, v \rangle = \max_{v \in \mathcal{V}} \langle g, v \rangle + \langle g, x \rangle \quad \text{and therefore} \quad \mathbb{E} \max_{v \in \mathcal{V} + x} \langle g, v \rangle = \mathbb{E} \max_{v \in \mathcal{V}} \langle g, v \rangle + \mathbb{E} \langle g, x \rangle.$$

It's about linearity of expectations and inner products.

Notation. Often, because it's a bit more compact, we'll write $s^2 \geq 2\sigma w(\mathcal{M}_s^\circ)$ instead of $s^2 \geq 2\sigma w(\mathcal{M}_s^\circ - \mu_\star)$ in bounds like (3).

Exercise 2 Explain why it makes no difference whether we write $w(\mathcal{M}_s^\circ - \mu_\star)$ or $w(\mathcal{M}_s^\circ)$. A sentence should do.

Solution 2 Gaussian width is translation invariant.

2.2 Sublinearity

Exercise 3 Let \mathcal{M} be a convex set, μ_\star be a point in \mathcal{M} , and ρ be a seminorm defined on its elements. Prove that, for a neighborhood $\mathcal{M}_s = \{m - \mu_\star \in \mathcal{M} : \rho(m - \mu_\star) \leq s\}$ of μ_\star , $f(s) = w(\mathcal{M}_s - \mu_\star)$ is a sublinear function of s . That is, prove that $f(s)/s$, a function on the positive real numbers, is (non-necessarily-strictly) decreasing.

Tips.

1. $f(s)/s$ is decreasing if $f(s)/s \geq f(t)/t$ [or equivalently $f(s) \geq (s/t)f(t)$] whenever $s \leq t$. Is $(s/t)f(t)$ the gaussian width of some set? If so, what set? And how is it related to $\mathcal{M}_s - \mu_\star$? Use the properties of gaussian width you proved in Exercise 1.
2. It's important that \mathcal{M} is a convex set containing μ_\star . Why? If m is in \mathcal{M} , then so is $m_\lambda = \mu_\star + \lambda(m - \mu_\star)$ for any $\lambda \in [0, 1]$. Or equivalently, if $m - \mu_\star$ is in $\mathcal{M} - \mu_\star$, so is $m_t - \mu_\star = \lambda(m - \mu_\star)$.¹

Solution 3 We'll use the increasingness and homogeneity of gaussian width to show that $(s/t)f(t) \leq f(s)$ if $(s/t)(\mathcal{M}_t - \mu_\star) \subseteq \mathcal{M}_s - \mu_\star$.

$$(s/t)f(t) = w(\cdot)\{(s/t)(\mathcal{M}_t - \mu_\star)\} \leq w(\cdot)\{\mathcal{M}_s - \mu_\star\} = f(s) \text{ if } (s/t)(\mathcal{M}_t - \mu_\star) \subseteq \mathcal{M}_s - \mu_\star.$$

To conclude, we need to show that whenever $s \leq t$, the containment holds. That is, that every vector $v \in (s/t)(\mathcal{M}_t - \mu_\star)$ is also in $\mathcal{M}_s - \mu_\star$. And because $\mathcal{M}_s - \mu_\star$ is the set of vectors v for which (i) $v + \mu_\star \in \mathcal{M}$ and (ii) $\rho(v) \leq s$, that amounts to showing that every vector $v \in (s/t)(\mathcal{M}_t - \mu_\star)$ has these two properties.

To do this, recall that definitionally each vector $v \in (s/t)(\mathcal{M}_t - \mu_\star)$ has the form $v = (s/t)(u - \mu_\star)$ for some $u \in \mathcal{M}$ with $\rho(u - \mu_\star) \leq t$.

(i)

$$\text{Letting } \lambda = s/t \in [0, 1], \quad v + \mu_\star = \lambda(u - \mu_\star) + \mu_\star = \lambda u + (1 - \lambda)v,$$

i.e., $v + \mu_\star$ is a convex combination of two vectors in \mathcal{M} and therefore is in the model.

(ii)

$$\rho(v) = \rho\{\lambda(u - \mu_\star)\} = \lambda\rho(u - \mu_\star) \leq s \quad \text{because} \quad \lambda = s/t \text{ and } \rho(u - \mu_\star) \leq t.$$

Now you should have what you need to prove that (3) implies (4).

Exercise 4 Prove that if $s^2 \geq 2\sigma w(\mathcal{M}_s - \mu_\star)$, then $(s+x)^2 \geq 2\sigma w(\mathcal{M}_{s+x} - \mu_\star) + sx$ for any $x \geq 0$. Then briefly explain why this and (3) together imply (4). A sentence or two should be enough for this explanation.

Tip. You can get a condition equivalent to the one you want to show by dividing both sides by $s + x$.

$$(s + x)^2 \geq 2\sigma w(\mathcal{M}_{s+x} - \mu_\star) + sx \quad \text{if and only if}$$

$$s + x \geq 2\sigma \frac{w(\mathcal{M}_{s+x} - \mu_\star)}{s + x} + \frac{s}{s + x}x.$$

Looking at this equivalent condition, compare the first term on the left side to the first term on the right and the second term on the left side to the second term on the right.

Solution 4 *To prove the claimed implication, we'll use the tip. Suppose s satisfies $s^2 \geq 2\sigma w(\mathcal{M}_s - \mu_\star)$ or equivalently $s \geq 2\sigma w(\mathcal{M}_s - \mu_\star)/s$.*

$$2\sigma \frac{w(\mathcal{M}_{s+x} - \mu_\star)}{s + x} \leq 2\sigma \frac{w(\mathcal{M}_s - \mu_\star)}{s} \leq s \quad \text{because width is sublinear and } s + x \geq s$$

$$\frac{s}{s + x}x \leq x \quad \text{because } \frac{s}{s + x} \leq 1 \text{ and } x \geq 0.$$

Summing both sides, we get the equivalent condition the tip suggests we show.

What remains is to show that the bound in (3) implies eq:error-bound-simplified. Suppose s satisfies $s^2 \geq 2\sigma w(\mathcal{M}_s - \mu_\star)$ as assumed in (4). Then, letting $x = 2\sigma \sqrt{\frac{2\{1+2\log(2n)\}}{\delta n}}$, it follows from what we've just proven that $(s + x)^2 \geq 2\sigma w(\mathcal{M}_{s+x} - \mu_\star) + sx$. That is, it follows that the radius $s + x$ that (4) claims is a bound on $\|\hat{\mu} - \mu_\star\|$ does, in fact, satisfy the condition required by (3) to actually be a bound on $\|\hat{\mu} - \mu_\star\|$.

3 A More Realistic Error Bound

Setting. In this section, we will consider the case that we observe pairs $(X_1, Y_1) \dots (X_n, Y_n)$ where $X_1 \dots X_n$ are deterministic and $Y_i = \mu(X_i) + \varepsilon_i$ for $\varepsilon_1 \dots \varepsilon_n$ that are independent, but not necessarily identically distributed, random variables with $E \varepsilon_i = 0$.

What's more realistic about this? It describes the kind of data we get with an actual usable sampling mechanism. If we draw pairs $(X_1, Y_1) \dots (X_n, Y_n)$ *uniformly at random with replacement* from a population $(x_1, y_1) \dots (x_m, y_m)$, then do our analysis *conditioning on* $X_1 \dots X_n$, this is the setting we find ourselves in. In that case, our signal is $\mu(x) = E[Y_i | X_i = x] = \frac{1}{m_x} \sum_{j: x_j = x} y_j$, the average outcome among people in the population with $x_j = x$. And the high probability error bounds we prove hold with *conditional probability* $1 - \delta$.¹

In this setting, we can prove the following error bound in terms of the random vector $\varepsilon \in \mathbb{R}^n$ with i th element ε_i .

$$\begin{aligned} \|\hat{\mu} - \mu_\star\|_{L_2(P_n)} &< s + 2\sqrt{\frac{2\Sigma}{\delta n}} \quad \text{w.p.} \quad 1 - \delta \quad \text{if} \quad \frac{s^2}{2} \geq w_\varepsilon(\mathcal{M}_s^\circ - \mu_\star) \\ \text{for} \quad \mu_\star &= \operatorname{argmin}_{m \in \mathcal{M}} \|m - \mu\|_{L_2(P_n)}, \\ w_\varepsilon(\mathcal{V}) &= E \max_{v \in \mathcal{V}} \langle \varepsilon, v \rangle_{L_2(P_n)}, \\ \text{and} \quad \Sigma &= E \max_{i \in 1 \dots n} \varepsilon_i^2. \end{aligned} \tag{5}$$

We'll start with a warm-up.

Exercise 5 Prove, by plugging in $\varepsilon_i \stackrel{iid}{\sim} N(0, \sigma^2)$, that this error bound (5) implies the bound we use in the gaussian case (2). A sentence or two should do.

You can use, without proof, the bound $E \max_{i \in 1 \dots n} \varepsilon_i^2 \leq \sigma^2 \{1 + 2 \log(2n)\}$ That's Lemma 11.3 of Boucheron, Lugosi, and Massart's Concentration Inequalities: A Nonasymptotic Theory of Independence.

A Correction. There's a little bit of a mixup about whether we're meant to be showing something like (4) or (3). **Option 1.** If we're meant to be showing something like (4), then the condition on s should be

$$\frac{s^2}{2} \geq w_\varepsilon(\mathcal{M}_s - \mu_\star)$$

referring to the neighborhood $\mathcal{M}_s - \mu_\star$ rather than its boundary $\mathcal{M}_s^\circ - \mu_\star$.

Option 2. If we're meant to be showing something like (3) or (2), then the condition on s should be

$$\frac{s^2}{2} \geq w_\varepsilon(\mathcal{M}_s^\circ - \mu_\star) + s\sqrt{\frac{2\Sigma}{\delta n}},$$

¹This implies they hold with unconditional probability $1 - \delta$ too. For any event A and conditioning set B , by the law of iterated expectations, the probability of any event A is the expected value of the conditional probability of A given B : $P(A) = E[1_A] = E[E[1_A | B]] = E[P(A | B)]$.

including a linear-in- s term as in those bounds. In my solution, I'll do the **Option 1** version.

Solution 5 *I'll be a little verbose here.*

First, observe that the 'crossing-point condition' characterizing s in (5) and (4) are equivalent. Our noise vector ε is a scaled version of the standard normal vector g used in the definition of gaussian width: $\varepsilon = \sigma g$, as multiplying a standard normal g_i by σ give you a gaussian with mean zero and variance σ^2 . It follows that $w_\varepsilon(\mathcal{V}) = \sigma w(\mathcal{V})$ for all sets \mathcal{V} .

$$w_\varepsilon(\mathcal{V}) = \mathbb{E} \max_{v \in \mathcal{V}} \langle \sigma g, v \rangle_2 = \sigma \max_{v \in \mathcal{V}} \mathbb{E} \langle g, v \rangle_2 = \sigma w(\mathcal{V}).$$

Thus, the condition that s satisfies $s^2/2 \geq 2w_\varepsilon(\mathcal{M}_s - \mu_\star)$ [here in (5)] is equivalent to the condition that s satisfies $s^2/2 \geq 2\sigma w(\mathcal{M}_s - \mu_\star)$ [in (4)].

Second, given s satisfying this condition, the bound on $\|\hat{\mu} - \mu_\star\|$ given above in (4) is valid (i.e. exceeds $\|\hat{\mu} - \mu_\star\|$) whenever the one given here in (5) is valid, as the former bound $s + 2\sqrt{\frac{2\sigma^2\{1+2\log(2n)\}}{\delta n}}$ is, in light of the bound on Σ from the Concentration Inequalities book, larger than the latter bound $s + 2\sqrt{\frac{2\Sigma}{\delta n}}$.

It's time to prove the new bound. Virtually everything you'll need can be borrowed from lecture, but I'm going to ask you to write a complete proof, copying out the parts you need to. This is meant as encouragement to review the proofs from lecture and understand how they work.

Exercise 6 *Write out a complete proof of the new error bound (5). You don't need to include a proof of the Efron-Stein inequality or anything you've proven above in Section 2, but everything else should be included.*

Solution 6 *This one I'm omitting. My proof is on the slides.*

To make this bound useful, we'll need to bound $w_\varepsilon(\mathcal{M}_s^\circ - \mu_\star)$ for some models \mathcal{M} . Today, we'll do that for the gaussian case $\varepsilon_i \stackrel{iid}{\sim} N(0, \sigma^2)$. Later in the semester, we'll prove bounds of the form $w_\varepsilon(\mathcal{V}) \leq \alpha w(\mathcal{V})$ that hold for every set \mathcal{V} with a constant α that depends on the distribution of ε . This'll let us use our gaussian width bounds together with (5) to get concrete, realistic error bounds.

4 Gaussian Width Calculations

In this section, we're going to be talking about two sets of linear functions of K -dimensional covariates, i.e., functions of the form $m(x) = x^T \beta$ for $x \in \mathbb{R}^K$. The first, the kind of linear model we talk about in classes like QTM220, will be the set of all of these. Here's how we write it, both as a set of functions and as the set of vectors $[m(X_1), m(X_2), \dots, m(X_n)] \in \mathbb{R}^n$ that we get by evaluating it at our observations. We'll let X be the $K \times n$ matrix with columns $X_1 \dots X_n$.

$$\begin{aligned} \mathcal{M} &= \{m(x) = x^T \beta : \beta \in \mathbb{R}^K\} && \text{as a set of functions} \\ &= \{X^T \beta : \beta \in \mathbb{R}^K\} && \text{as a set of vectors} \end{aligned} \quad (6)$$

The second, the set of linear functions we work with when we use *the lasso*, is the subset of these with coefficients satisfying a one-norm bound $\|\beta\|_1 \leq B$.

$$\begin{aligned} \mathcal{M} &= \{m(x) = x^T \beta : \beta \in \mathbb{R}^K \text{ and } \|\beta\|_1 \leq B\} && \text{as a set of functions} \\ &= \{X^T \beta : \beta \in \mathbb{R}^K \text{ and } \|\beta\|_1 \leq B\} && \text{as a set of vectors} \end{aligned} \quad (7)$$

Throughout, we'll focus on the gaussian noise case: $\varepsilon_1 \dots \varepsilon_n \stackrel{\text{iid}}{\sim} \mathbf{N}(\mathbf{0}, \sigma^2)$.

4.1 The Linear Model

Exercise 7 Find an upper bound on the gaussian width $w(\mathcal{M}_s - \mu_\star)$ of a centered neighborhood $\mathcal{M}_s - \mu_\star$ in the linear model (6). Then give a bound on the error $\|\hat{\mu} - \mu_\star\|_{L_2(\mathbb{P}_n)}$ of the least squares estimator in this model that holds with probability $1 - \delta$.

Tips.

1. It'll be convenient to work with the Euclidean norm $\|\cdot\|_2$ and inner product $\langle \cdot, \cdot \rangle_2$ instead of the sample two norm and inner product. Rewrite the $s^2 \geq \max \dots$ condition in (4) in terms of these using the scaling-up relationships

$$\|\cdot\|_2 = \sqrt{n} \|\cdot\|_{L_2(\mathbb{P}_n)} \quad \text{and} \quad \langle \cdot, \cdot \rangle_2 = n \langle \cdot, \cdot \rangle_{L_2(\mathbb{P}_n)}.$$

$\mathcal{M}_s - \mu_\star = \{X^T v : v \in \mathbb{R}^K \text{ and } \|X^T v\|_2 \leq s\sqrt{n}\}$. Why?

2. You're going to want to use the Cauchy-Schwarz bound, but the bound $\langle \varepsilon, X^T v \rangle_2 \leq \|\varepsilon\|_2 \|X^T v\|_2 \leq \|\varepsilon\|_2 s\sqrt{n}$ isn't going to be good enough. The bound we get this way is the same one we got in lecture for the completely general model. Take a look at Appendix A. What is $\mathbb{E}\|\varepsilon^u\|_2$ where $\varepsilon^u = \sum_{j=1}^K \langle u_j, \varepsilon \rangle_2 u_j$ is the projection of ε onto the span of the columns of X ?
3. For any random variable Z including $Z = \|\varepsilon^u\|_2$,

$$\mathbb{E} Z^2 = (\mathbb{E} Z)^2 + \text{Var}(Z) \quad \text{and therefore} \quad (\mathbb{E} Z)^2 \leq \mathbb{E} Z^2.$$

Solution 7 We'll start by addressing the 'Why?' part of Tip 1, i.e. by showing that the set of vectors in $\mathcal{M}_s - \mu_\star$ is

$$\mathcal{M}_s - \mu_\star = \{X^T v : v \in \mathbb{R}^K \text{ and } \|X^T v\|_2 \leq s\sqrt{n}\}.$$

1. Because $\mu_\star \in \mathcal{M}$, we know that, in vector terms, it's $X^T \beta_\star$ for some $\beta_\star \in \mathbb{R}^K$. And it follows that $\mathcal{M} - \mu_\star = \{X^T(\beta - \beta_\star) : \beta \in \mathbb{R}^K\}$.
2. Of these, the vectors in our neighborhood satisfy

$$s^2 \geq \|m - \mu_\star\|_{L_2(\mathbb{P}_n)}^2 = \frac{1}{n} \sum_{i=1}^n (X_i^T \beta - X_i^T \beta_\star)^2 = \frac{1}{n} \sum_{i=1}^n \{X^T(\beta - \beta_\star)\}_i^2 = \frac{1}{n} \|X^T(\beta - \beta_\star)\|_2^2.$$

$$\text{Rearranging, we get } \sqrt{ns} \geq \|X^T(\beta - \beta_\star)\|_2.$$

Making the substitution $v = \beta - \beta_\star$ gives the claimed characterization.

The Width Calculation. Now let's characterize the gaussian width of this neighborhood in vector terms.

$$\begin{aligned} w(\mathcal{M}_s - \mu_\star) &= \mathbb{E} \max_{h \in \mathcal{M}_s - \mu_\star} \frac{1}{n} \langle \varepsilon, h \rangle_2 \quad \text{for } \varepsilon_i \stackrel{iid}{\sim} N(0, 1) \\ &= \frac{1}{n} \mathbb{E} \max_{\substack{h = X\beta \\ \|X^T v\|_2 \leq \sqrt{ns}}} \langle \varepsilon, X^T v \rangle_2 \end{aligned}$$

Here's where the projection stuff from Appendix A comes in. Let's work with the decomposition $\varepsilon = \varepsilon_\parallel + \varepsilon_\perp$ where ε_\parallel is the orthogonal projection of ε onto the image of X^T . This satisfies $\langle \varepsilon_\parallel, X^T v \rangle_2 = \langle \varepsilon, X^T v \rangle_2$ for all v , so

$$\begin{aligned} \max_{\substack{h = X\beta \\ \|X^T v\|_2 \leq \sqrt{ns}}} \langle \varepsilon, X^T v \rangle_2 &= \max_{\substack{h = X\beta \\ \|X^T v\|_2 \leq \sqrt{ns}}} \langle \varepsilon_\parallel, X^T v \rangle_2 \\ &\leq \max_{\substack{h = X\beta \\ \|X^T v\|_2 \leq \sqrt{ns}}} \|\varepsilon_\parallel\|_2 \|X^T v\|_2 \quad \text{via Cauchy-Schwarz} \\ &\leq \|\varepsilon_\parallel\|_2 \times \sqrt{ns} \quad \text{as a result of our neighborhood constraint.} \end{aligned}$$

Taking the expectation of this, we get a bound on gaussian width.

$$w(\mathcal{M}_s - \mu_\star) \leq \frac{1}{n} \mathbb{E} \|\varepsilon_\parallel\|_2 \times \sqrt{ns} = \frac{s}{\sqrt{n}} \mathbb{E} \|\varepsilon_\parallel\|_2.$$

What's $\mathbb{E} \|\varepsilon_\parallel\|_2$? We'll show it's less than \sqrt{K} . To do this, we'll use the bound $\{\mathbb{E}[Z]\}^2 \leq \mathbb{E}[Z^2]$ for $Z = \|\varepsilon_\parallel\|_2$ and the explicit formula $\varepsilon_\parallel = \sum_{j=1}^K \langle \varepsilon, u_j \rangle_2 u_j$

in terms of an orthonormal basis $u_1 \dots u_K$ for the image of X^T . Here goes.

$$\begin{aligned}
\{\mathbb{E}\|\varepsilon_{\parallel}\|_2\}^2 &\leq \mathbb{E}\|\varepsilon_{\parallel}\|_2^2 \\
&= \mathbb{E}\left\|\sum_{j=1}^K \langle \varepsilon, u_j \rangle_2 u_j\right\|_2^2 && \text{using our formula for } \varepsilon_{\parallel} \\
&= \mathbb{E}\left\langle \sum_{j=1}^K \langle \varepsilon, u_j \rangle_2 u_j, \sum_{k=1}^K \langle \varepsilon, u_k \rangle_2 u_k \right\rangle_2 && \text{using the definition } \|v\|_2^2 = \langle v, v \rangle_2 \\
&= \mathbb{E} \sum_{j=1}^K \sum_{k=1}^K \langle \varepsilon, u_j \rangle_2 \langle \varepsilon, u_k \rangle_2 \langle u_j, u_k \rangle_2 && \text{because inner products are linear} \\
&= \mathbb{E} \sum_{j=1}^K \langle \varepsilon, u_j \rangle_2^2 && \text{because our basis is orthonormal} \\
&= \sum_{j=1}^K \mathbb{E} \langle \varepsilon, u_j \rangle_2^2 && \text{because expectation is linear} \\
&= \sum_{j=1}^K \underbrace{\mathbb{E} \varepsilon_1^2}_{=1} = K && \text{because } \varepsilon \text{ is spherically symmetric.}
\end{aligned}$$

That's it. We've proven the bound $w(\mathcal{M}_s - \mu_{\star}) \leq s\sqrt{K/n}$.

The Error Bound. Substituting the width bound into our ‘crossing point condition’, we find that

$$s = 2\sigma\sqrt{\frac{K}{n}} \quad \text{satisfies} \quad \frac{s^2}{2\sigma} \geq \sqrt{\frac{K}{n}} \geq w(\mathcal{M}_s - \mu_{\star}).$$

Substituting this radius into our abstract error bound (4), we get a concrete error bound that applies when we do least squares regression with gaussian noise with variance σ^2 .

$$\|\hat{\mu} - \mu_{\star}\|_{L_2(\mathbb{P}_n)} \leq 2\sigma\sqrt{\frac{K}{n}} + 2\sigma\sqrt{\frac{2\{1 + 2\log(2n)\}}{\delta n}} \quad \text{with probability } 1 - \delta.$$

4.2 The Lasso

Exercise 8 Find an upper bound of the gaussian width $w(\mathcal{M})$ of the model (7) used in the lasso. Then use it to give a bound on the error $\|\hat{\mu} - \mu_\star\|_{L_2(\mathbb{P}_n)}$ of the least squares estimator in this model that holds with probability $1 - \delta$.

Tips.

1. How can we bound the dot product $\langle \varepsilon, X^T \beta \rangle_2 = \langle X \varepsilon, \beta \rangle_2$ when $\|\beta\|_1 \leq B$? Look over the Inner Product Spaces Homework and Appendix B.
2. You can get away with using the bound $w(\mathcal{M}_s - \mu) \leq w(\mathcal{M})$ when calculating your error bound. It turns out we can't do much better than this. The reason is, in essence, that this model is so 'pointy' that unless s is very small, it contains very few functions with $\|m - \mu\|_{L_2(\mathbb{P}_n)} > s$ anyway. Section 7.5 of High Dimensional Probability explains this nicely.²

Solution 8

The Width Calculation. As suggested, we'll use the radius-independent bound

$$w(\mathcal{M}_s - \mu_\star) \leq w(\mathcal{M} - \mu_\star) = w(\mathcal{M}).$$

The last identity follows from translation invariance of the gaussian width. Writing the gaussian width of our model in vector terms, observing that

$$\langle \varepsilon, X^T \beta \rangle_2 = (X^T \beta)^T \varepsilon = \beta^T (X \varepsilon) = \langle \beta, X \varepsilon \rangle_2$$

and using Hölder's inequality, we get a bound in terms of the maximal absolute value of K gaussian random variables.

$$\begin{aligned} w(\mathcal{M}) &= \frac{1}{n} \mathbb{E} \max_{\substack{h=X\beta \\ \|\beta\|_1 \leq 1}} \langle \varepsilon, X^T \beta \rangle_2 \\ &= \frac{1}{n} \mathbb{E} \max_{\substack{\beta \in \mathbb{R}^K \\ \|\beta\|_1 \leq 1}} \|\beta\|_1 \|X \varepsilon\|_\infty \\ &\leq \frac{1}{n} \mathbb{E} \|X \varepsilon\|_\infty = \mathbb{E} \left[\max_{j \in 1 \dots K} |Z_j| \right] \quad \text{for } Z_j = X_{\cdot j}^T \varepsilon. \end{aligned}$$

Interpretation. What we're calling $X_{\cdot j}$ isn't the same thing as X_j —it's not any one individual's covariate vector. It's the vector of length n containing the j th component of each individual's covariate vector X_i . And $\frac{1}{n} Z_j$ is a weighted average of those n values where each individual gets a gaussian weight ε_i .

Here's where the stuff from Appendix B comes in. If we have any K gaussian random variables $Z_1 \dots Z_K$ with variance less than or equal to V , then

$$\mathbb{E} \max_{j \in 1 \dots K} Z_j \leq 2\sqrt{2V \log(K)}.$$

That doesn't quite do what we want because we want the maximal absolute value, but there's a simple trick to make it work: the absolute value of Z_1 is the maximum of Z_1 and $-Z_1$, so the maximal absolute value of $Z_1 \dots Z_K$ is the maximum of $Z_1 \dots Z_K$ and $-Z_1 \dots Z_K$ — $2K$ gaussian random variables with variance less than or equal to V .

$$\mathbb{E} \max_{j \in 1 \dots K} |Z_j| \leq 2\sqrt{2V \log(2K)}.$$

What we need now is a bound σ^2 on the variance of each $Z_i = X_j^T \varepsilon$. For this, there's another trick. If we have a dot product $u^T v$, $(u^T v)^2 = u^T v v^T u$.

Taking $u = X_{\cdot j}$ and $v = \varepsilon$, we get

$$\begin{aligned} \text{Var}[X_{\cdot j}^T \varepsilon] &= \mathbb{E} [(X_{\cdot j}^T \varepsilon)^2] \\ &= \mathbb{E} [X_{\cdot j}^T \varepsilon \varepsilon^T X_{\cdot j}] && \text{using this trick} \\ &= X_{\cdot j}^T \Sigma X_{\cdot j} && \text{where } \Sigma = \mathbb{E} [\varepsilon \varepsilon^T] \text{ is the covariance matrix of } \varepsilon \\ &= X_{\cdot j}^T I X_{\cdot j} = \|X_{\cdot j}\|_2^2 && \text{because } \varepsilon \text{'s covariance matrix is the identity.} \end{aligned}$$

How do we know that Σ is the identity? Let's do the calculation to show its elements are one on the diagonal and zero elsewhere.

$$\Sigma_{ij} = \mathbb{E} \varepsilon_i \varepsilon_j = \begin{cases} \mathbb{E} \varepsilon_i^2 = 1 & \text{if } i = j \\ \mathbb{E} \varepsilon_i \varepsilon_j = 0 & \text{if } i \neq j \end{cases}$$

Let's put it all together.

$$\begin{aligned} w(\mathcal{M}) &\leq \frac{1}{n} \times 2\sqrt{2V \log(2K)} \quad \text{for } V = \max_j \|X_{\cdot j}\|_2^2 = n \times \max_j \|X_{\cdot j}\|_{L_2(\mathbb{P}_n)}^2 \\ &= \frac{2B \sqrt{\frac{2 \log(2K)}{n}}}{\text{for}} B = \max_j \|X_{\cdot j}\|_{L_2(\mathbb{P}_n)} \end{aligned}$$

Here you can think of B as the typical magnitude of the largest component of our individual's covariate vectors X_i . In particular, if the elements of X are in $[-1, 1]$, $B \leq 1$.

The Error Bound.

Proceeding as in the previous problem, we find a radius s that satisfies our crossing point condition. The crossing point picture is admittedly a bit simplistic here because our bound on a neighborhood's width doesn't depend on its radius—we're looking at the point where the function $f(s) = s^2$ crosses a horizontal line—but it works.

$$s^2 \geq \frac{2B \sqrt{\frac{2 \log(K)}{n}}}{\geq} w(\mathcal{M}_s - \mu_\star) \quad \text{is satisfied for } s = \sqrt{2\sqrt{2}B} \sqrt[4]{\frac{\log(K)}{n}}.$$

And it follows from (4) that, when we use the lasso when we have gaussian noise with variance σ^2 , we get this error bound.

$$\|\hat{\mu} - \mu_\star\|_{L_2(P_n)} \leq 2\sigma\sqrt{2\sqrt{2}B}\sqrt[4]{\frac{\log(2K)}{n}} + 2\sigma\sqrt{\frac{2\{1 + 2\log(2n)\}}{\delta n}} \quad \text{with probability } 1 - \delta.$$

Interpretation. This is, for what it's worth, called the slow rate or assumptionless analysis of the lasso. The essential lesson is that if we're willing to predict Y_i using some absolutely convex combinations of the components of our covariate vectors X_i , then it doesn't really matter how many components there are. That's what it means for our bound to be proportional to $\sqrt[4]{\log(2K)}$, which is so slow-growing that it goes from roughly 1 to 3 as we increase K from 1 to the mass of the earth in grams.

What we don't like about this bound is the way it varies with sample size n : a fourth-root rate. There's a fancier argument we can use when X is a very special matrix and most components of β_\star are zero that leads to a what is called a sparsity-dependent fast rate bound, which can be better.

A Projections

It's often useful to decompose a vector into relevant and irrelevant parts. For example, if we're interested in an inner product $\langle u, Av \rangle$, it's helpful to decompose u as a sum $u_{\parallel} + u_{\perp}$ where $\langle u_{\perp}, Av \rangle = 0$ for all v . This is particularly nice if we're going to use a Cauchy-Schwarz bound, as we can get a better bound by first getting rid of the irrelevant part u_{\perp} .

$$\langle u, Av \rangle = \langle u_{\parallel}, Av \rangle + \langle u_{\perp}, Av \rangle = \langle u_{\parallel}, Av \rangle \leq \|u_{\parallel}\| \|Av\|.$$

The best way to do this, in the sense that $\|u_{\parallel}\|$ is smallest, is to take u_{\parallel} to be the *orthogonal projection* onto the *image* of A —the image of A is the set of all vectors we can write as matrix-vector projects Av . To do that, we'll want an orthonormal basis for the image of A , i.e., a set of vectors u_1, u_2, \dots with the property that $\langle u_i, u_j \rangle$ is one if $i = j$ and zero otherwise. To get a basis like this, we can run any set of vectors that spans the image of A , e.g. the columns of A , through the Gram-Schmidt Process. Then we write u_{\parallel} as a linear combination of these vectors, $u_{\parallel} = \sum_k u_k \langle u_k, u \rangle$. To check that $\langle u_{\parallel}, Av \rangle = \langle u, Av \rangle$ for all v , observe that because u_1, u_2, \dots is a basis for the image of A , we can express Av as a linear combination $\sum_k \alpha_k u_k$ of these basis vectors. And we can calculate $\langle u, Av \rangle$ and $\langle u_{\parallel}, Av \rangle$ and compare. They're the same.

$$\begin{aligned} \langle u, Av \rangle &= \left\langle u, \sum_k \alpha_k u_k \right\rangle = \sum_k \alpha_k \langle u, u_k \rangle \\ \langle u_{\parallel}, Av \rangle &= \left\langle \sum_j u_j \langle u_j, u \rangle, \sum_k \alpha_k u_k \right\rangle = \sum_j \sum_k \alpha_k \langle u_j, u_k \rangle \langle u_j, u \rangle = \sum_k \alpha_k \langle u_k, u \rangle \end{aligned}$$

When we simplified the double sum above, we observed that terms with $j \neq k$ were zero because $\langle u_j, u_k \rangle = 0$ and that $\langle u_j, u_k \rangle = 1$ in terms with $j = k$.

I'll leave it to you to convince yourself that this is the best we can do, i.e., that there is no vector \tilde{u}_{\parallel} satisfying $\langle \tilde{u}_{\parallel}, Av \rangle = \langle u, Av \rangle$ for all v with $\|\tilde{u}_{\parallel}\| < \|u_{\parallel}\|$.

This all works with any inner product $\langle u, v \rangle$ and associated norm $\|v\| = \sqrt{\langle v, v \rangle}$. In this homework, we'll use it to talk about the dot product between gaussian vectors and vectors of the form Av . Note that if A is a $m \times n$ matrix, then our basis u_1, u_2, \dots contains at most $\min(m, n)$ vectors.

B Bounding Expectations by Integrating Tail Bounds

Going from bounds on tail probabilities to bounds on expectations is basically just a matter of integration, as $E Z = \int_0^\infty P(Z > z) dz$ for any positive random variable Z . Here's a proof.

$$E Z = E \int_0^Z 1 dz = E \int_0^\infty 1(Z > z) dz = \int_0^\infty E 1(Z > z) dz = \int_0^\infty P(Z > z) dz.$$

More generally, $E Z \leq \int_0^\infty P(Z > z) dz$ for any random variable Z . To show that, we can use the formula for nonnegative random variables on $Z_+ = \max\{Z, 0\}$ and observe that (i) $E Z \leq E Z_+$ and (ii) $P(Z_+ > z) = P(Z > z)$ for $z \geq 0$.

Let's use it to derive a bound on the expected value of the maximum $M_K = \max_{j \in 1 \dots K} Z_j$ of K mean-zero normals $Z_j \sim N(0, \sigma^2)$ using a tail bound from a previous lecture. We'll use the tail bound $P(M_K \geq z) \leq K e^{-z^2/2\sigma^2}$ that we substituted $z = 2\sigma\sqrt{\log(K)}$ into to get the bound $P(M_K \geq 2\sigma\sqrt{\log(K)}) \leq 1/K$ in our lecture on least squares in finite models.² Here's the bound.

$$E M_K \leq 2\sigma\sqrt{2\log(K)}.$$

To prove it, we'll break the integral $E M_K = \int_0^\infty P(M_K > z) dz$ into a sum of two integrals, one up to z_0 and one from there to ∞ . We'll bound the first using the simple observation that probabilities are less than one. And we'll bound the second using the fact that $z/z_0 \geq 0$ on the domain of integration $[z_0, \infty)$ and the identity $(d/dz)e^{-z^2/2} = -ze^{-z^2/2}$.

$$\begin{aligned} E M_K &= \int_0^{z_0} P(M_K \geq z) + \int_{z_0}^\infty P(M_K \geq z) \\ &\leq \int_0^{z_0} 1 + \int_{z_0}^\infty \frac{z}{z_0} K e^{-z^2/2\sigma^2} \\ &\leq z_0 + \frac{\sigma^2 K}{z_0} \int_{z_0}^\infty \frac{z}{\sigma^2} e^{-z^2/2\sigma^2} \\ &= z_0 - \frac{\sigma^2 K}{z_0} \int_{z_0}^\infty \frac{d}{dz} e^{-z^2/2\sigma^2} \\ &= z_0 - \frac{\sigma^2 K}{z_0} e^{-z^2/2\sigma^2} \Big|_{z_0}^\infty \\ &= z_0 + \frac{\sigma^2 K e^{-z_0^2/2\sigma^2}}{z_0}. \end{aligned}$$

Taking $z_0 = \sigma\sqrt{2\log(K)}$, $e^{-z_0^2/2\sigma^2} = e^{-\log(K)} = 1/K$. And we get this bound.

$$E M_K \leq \sigma\sqrt{2\log(K)} + \frac{\sigma^2 K \times 1/K}{\sigma\sqrt{2\log(K)}} = \sigma\left(\sqrt{2\log(K)} + \frac{1}{\sqrt{2\log(K)}}\right).$$

²More generally, this tail bound holds for any random variables $Z_1 \dots Z_K$ satisfying the bound $P(Z_j \geq z) \leq \frac{1}{2\pi\sigma^2} e^{-z^2/2\sigma^2}$, like a gaussian with mean zero and variance *less than or equal to* σ^2 would.

The bound $\mathbb{E} Z_K \leq 2\sigma\sqrt{2\log(K)}$ is a simplified version of this. Let's focus on the case that $K \geq 2$. Because $2\log(K) \geq 2\log(2) > 1$ for $K \geq 2$, the second term in curly brackets is smaller than the first, so their sum is bounded by twice the first. Thus, $\mathbb{E} Z_K \leq 2\sigma\sqrt{2\log(K)}$ as claimed.

This bound applies for $K = 1$ as well, as in that case $M_K = Z_1$ and $\mathbb{E} M_K = \mathbb{E} Z_1 = 0$.