

# Machine Learning Theory

## Least Squares with Misspecification and Non-Gaussian Noise

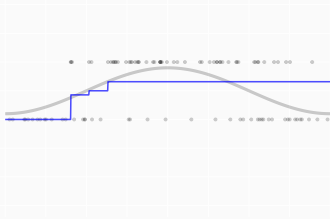
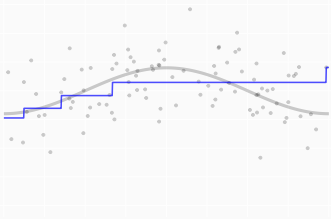
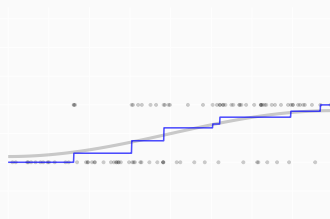
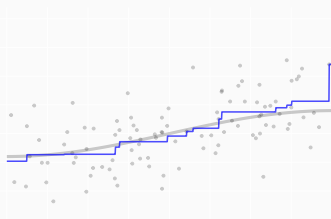
---

David A. Hirshberg

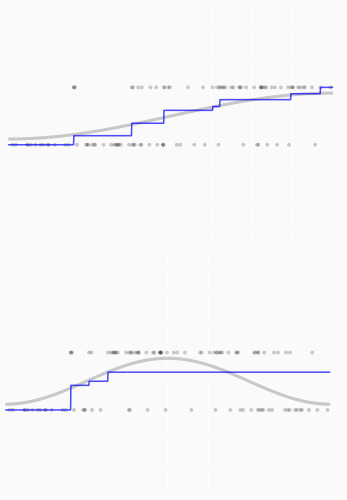
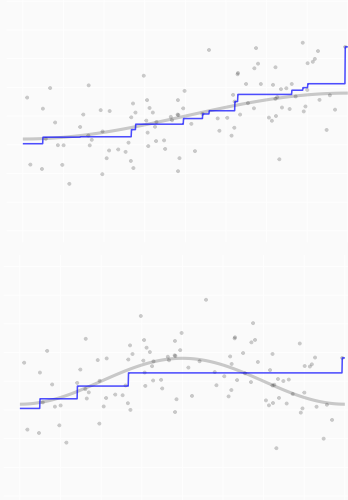
April 1, 2025

Emory University

# When Does Our Theory Apply?

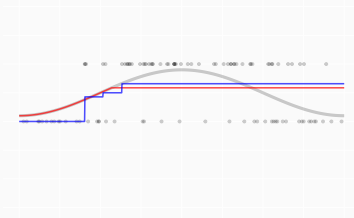
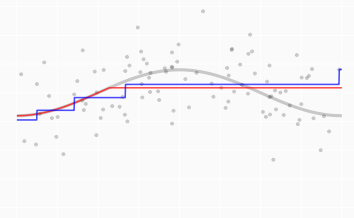
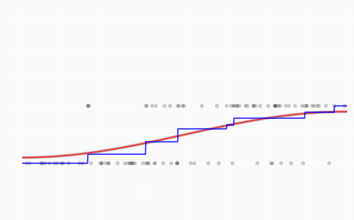
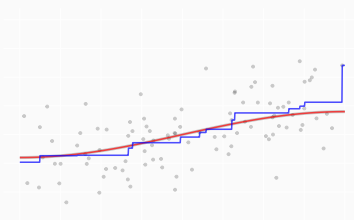


# When Does Our Theory Apply?

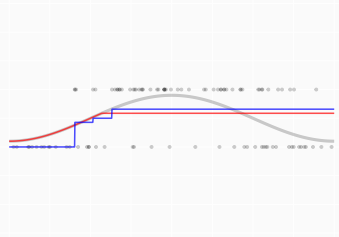
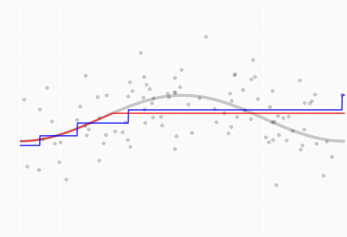
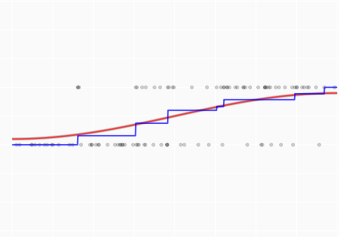
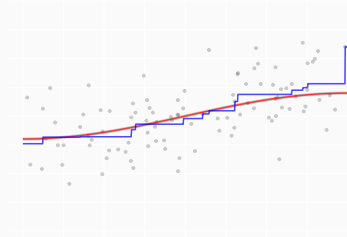


- The second column is out. We've assumed correct specification.
- The second row is out. We've assumed normality.

# Today, We Fix That



# Today, We Fix That

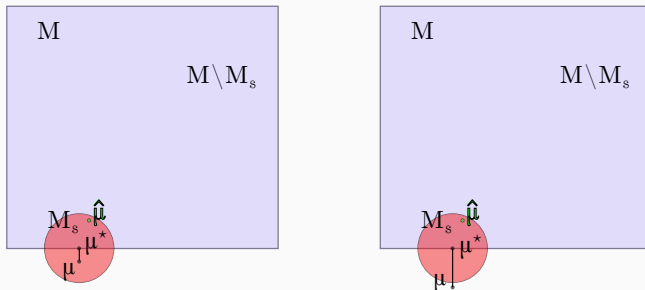


- With misspecification, we estimate the model's **best approximation** to  $\mu$ .
- Non-normality doesn't really matter much. We'll look at how it affects our bound.

## Misspecification

---

# What happens when $\mu$ isn't in the model?



- Our error in estimating  $\mu$  is bounded by a sum of two terms.
  - The critical radius  $s$ , i.e., the one satisfying  $s^2/2\sigma \geq w(\mathcal{M}_s^\circ) + s\sqrt{\frac{2M_n^2}{\delta n}}$ .
  - The distance from  $\mu$  to its best approximation in the model. Or really 3 times that.

We showed this in the model selection lab using the Cauchy-Schwarz inequality.

- In convex models, we can say more.  
Our error in estimating  $\mu^*$  does not depend on its distance to  $\mu$ .

# The Argument

For any  $\mu^* \in \mathcal{M}$ , we can expand our mean squared error difference as before.

$$\ell(m) - \ell(\mu^*) = \|m - \mu^*\|_{L_2(\mathbf{P}_n)}^2 - \frac{2}{n} \sum_{i=1}^n \varepsilon_i^* \{m(X_i) - \mu^*(X_i)\} \quad \text{for } \varepsilon_i^* = Y_i - \mu^*(X_i).$$

But our new ‘noise’  $\varepsilon_i^*$  doesn’t have mean zero. It’s our old noise  $\varepsilon_i$ , minus something.

$$\varepsilon_i^* = \underbrace{\{Y_i - \mu(X_i)\}}_{\varepsilon_i} - \underbrace{\{\mu^*(X_i) - \mu(X_i)\}}_{\text{something}}.$$

So we can think of our mean squared error difference as having three terms:

$$\begin{aligned} \ell(m) - \ell(\mu^*) &= \|m - \mu^*\|_{L_2(\mathbf{P}_n)}^2 && \text{squared distance, like before;} \\ &- \frac{2}{n} \sum_{i=1}^n \varepsilon_i \{m(X_i) - \mu^*(X_i)\} && \text{a mean zero term, like before;} \\ &+ \frac{2}{n} \sum_{i=1}^n \{\mu^*(X_i) - \mu(X_i)\} \{m(X_i) - \mu^*(X_i)\} && \text{and something else.} \end{aligned}$$

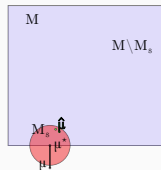
We can use our argument, ignoring the new term, if that term is always *non-negative*.

Why?



Why.

$$\begin{aligned}\ell(m) - \ell(\mu^*) &= \|m - \mu^*\|_{L_2(P_n)}^2 \\ &\quad - \frac{2}{n} \sum_{i=1}^n \varepsilon_i \{m(X_i) - \mu^*(X_i)\} \\ &\quad + \frac{2}{n} \sum_{i=1}^n \{\mu^*(X_i) - \mu(X_i)\} \{m(X_i) - \mu^*(X_i)\}\end{aligned}$$



We want to show that if distance from  $m$  to  $\mu^*$  is big enough, it wins.

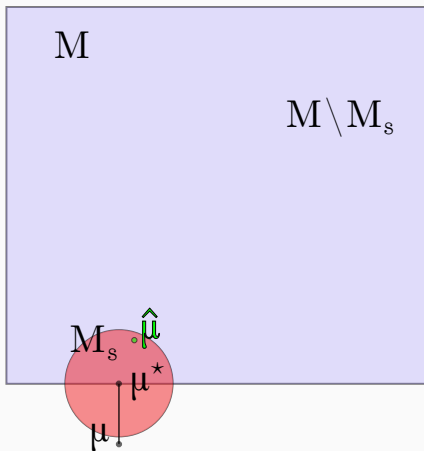
- In particular, it wins in the sense that the loss difference  $\ell(m) - \ell(\mu^*)$  is positive.
- That implies distance from  $\hat{\mu}$  to  $\mu^*$  is smaller, as distance doesn't win in that case.

If this new term is non-negative, it helps distance win.

- If the MSE difference is positive when we ignore a non-negative term, then it's positive when we don't.

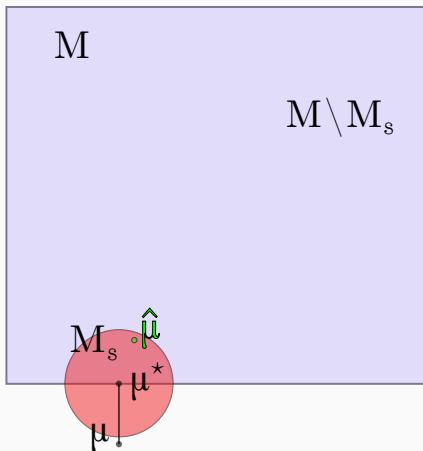
So we want to make sure this new term is non-negative. And we get to choose  $\mu^*$ .

## This sounds weird



- It sounds like we choose what our estimator converges to when we analyze it.
- Obviously we don't really get to do that. It's not really a choice—it's a guess.
- If  $\hat{\mu}$  converges to some curve  $\mu^*$ , then it can't converge to anything else.

## The right choice



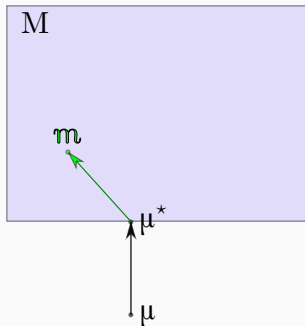
It's the best approximation to  $\mu$  in the model.

$$\mu^* = \operatorname{argmin}_{m \in \mathcal{M}} \|m - \mu\|_{L_2(P_n)}^2.$$

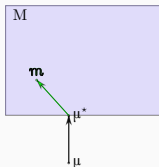
With this choice, the new term is always non-negative

$$\frac{2}{n} \sum_{i=1}^n \{\mu^*(X_i) - \mu(X_i)\} \{m(X_i) - \mu^*(X_i)\} = 2 \langle \mu^* - \mu, m - \mu^* \rangle_{L_2(\mathcal{P}_n)}$$

It's proportional to the dot product between two vectors:  $\mu \rightarrow \mu^*$  and  $\mu^* \rightarrow m$ .



When the model  $\mathcal{M}$  is convex, these vectors are always in the same direction.  
That is, this dot product is non-negative for all  $m \in \mathcal{M}$ .



**Claim.** For any convex set  $\mathcal{M}$  in an inner product space,<sup>1</sup>

$$\mu^* = \operatorname{argmin}_{m \in \mathcal{M}} \|m - \mu\| \quad \text{satisfies}$$

$$\langle \mu^* - \mu, m - \mu^* \rangle \geq 0 \quad \text{for all } m \in \mathcal{M}.$$

**Proof.** Let  $m_\lambda = \lambda(m - \mu^*) + \mu^*$ .

$$\begin{aligned} \|m_\lambda - \mu\|^2 &= \langle \lambda(m - \mu^*) + (\mu^* - \mu), \lambda(m - \mu^*) + (\mu^* - \mu) \rangle \\ &= \lambda^2 \|m - \mu^*\|^2 + \|\mu^* - \mu\|^2 + 2\lambda \langle m - \mu^*, \mu^* - \mu \rangle. \end{aligned}$$

Because  $m_\lambda \in \mathcal{M}$ , it follows that this is at least as large as  $\|\mu - \mu^*\|^2$ , so

$$0 \leq \lambda^2 \|m - \mu^*\|^2 + 2\lambda \langle m - \mu^*, \mu^* - \mu \rangle$$

and therefore, dividing by  $\lambda > 0$ , that

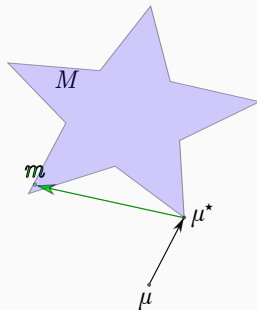
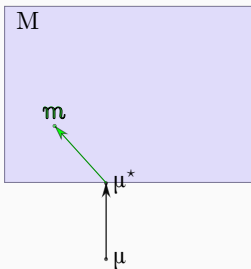
$$0 \leq \lambda \|m - \mu^*\|^2 + 2 \langle m - \mu^*, \mu^* - \mu \rangle.$$

Because this holds for arbitrarily small  $\lambda > 0$ , it must also hold for  $\lambda = 0$ .

<sup>1</sup>An inner product space is a vector space with a norm  $\|u\| = \sqrt{\langle u, u \rangle}$  induced by an inner product  $\langle u, v \rangle$ .

## That's not true for other choices

When  $\mu^* \in \mathcal{M}$  isn't the closest point to  $\mu$ ,  
these vectors can point in opposite directions.  
That is, this dot product can be negative for some  $m \in \mathcal{M}$ .



The same thing can happen *for the closest point* in a non-convex model.

When we use a convex model, the least squares estimator  $\hat{\mu}$  converges to the model's closest point to  $\mu$ .

- If  $\mu$  is in the model, that's  $\mu$ .
- Otherwise, it's something else.

We can bound our estimator's distance to that closest point  $\mu^\star$  just like we've been bounding distance to  $\mu$  when we assumed it was in the model.

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

for  $\mathcal{M}_s^\circ = \{m \in \mathcal{M} : \|m - \mu^\star\|_{L_2(\mathbb{P}_n)} = s\}$  and  $M_n = 1 + 2\log(2n)$ .

Let's get a feel for what that means by looking at some examples.

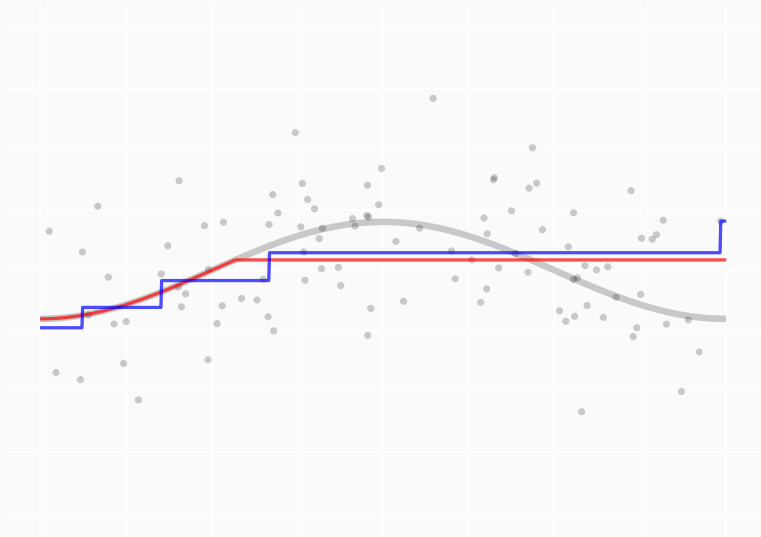


Figure 1: Increasing Curves.



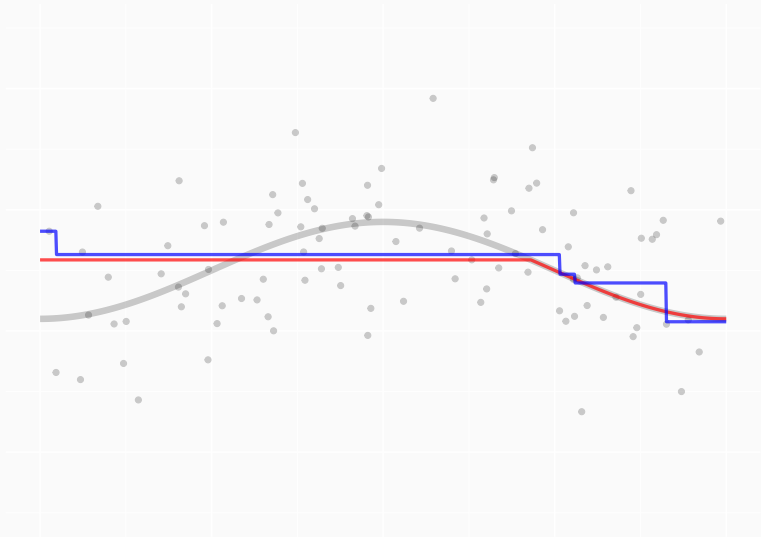


Figure 2: Decreasing Curves.

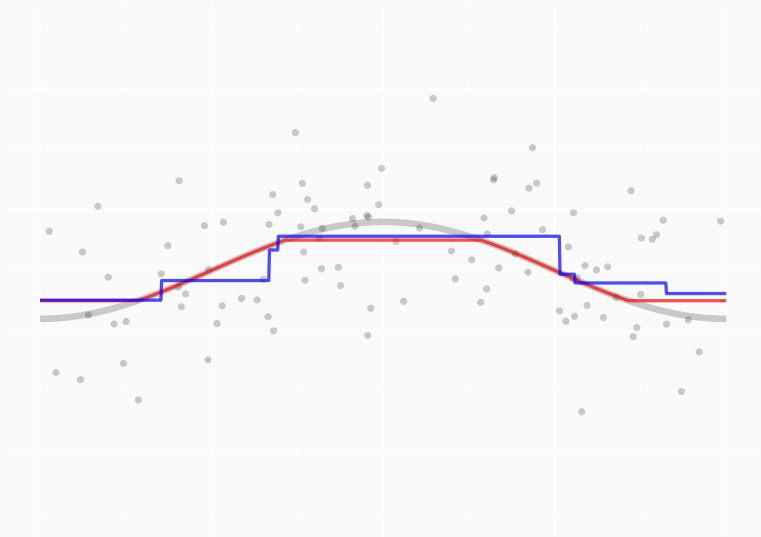


Figure 3: Bounded Variation Curves.  $\rho_{TV} \leq 1$

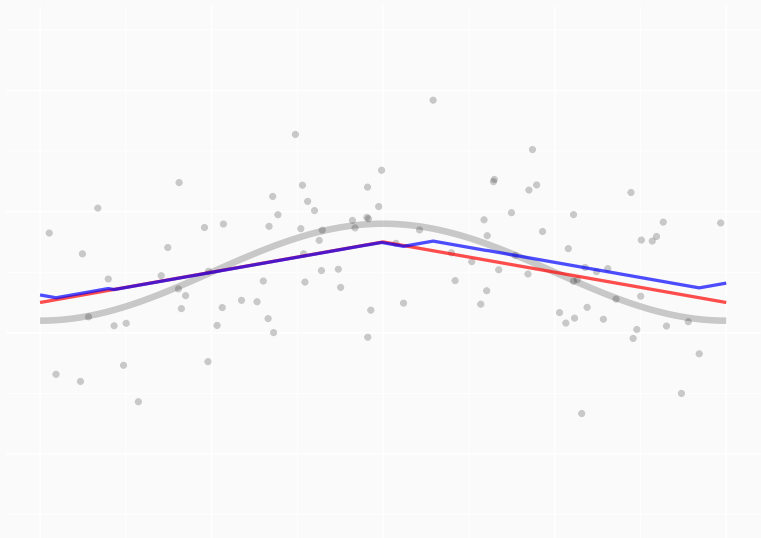


Figure 4: Lipschitz Curves.  $\rho_{\text{Lip}} \leq 1$

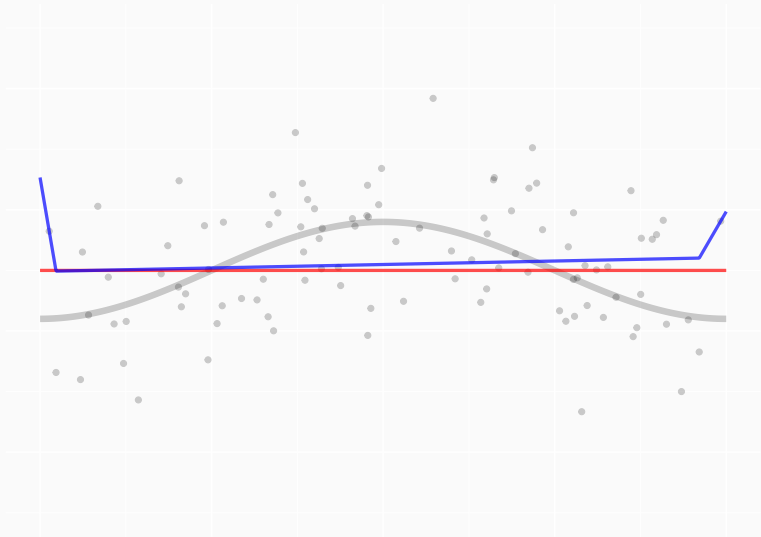
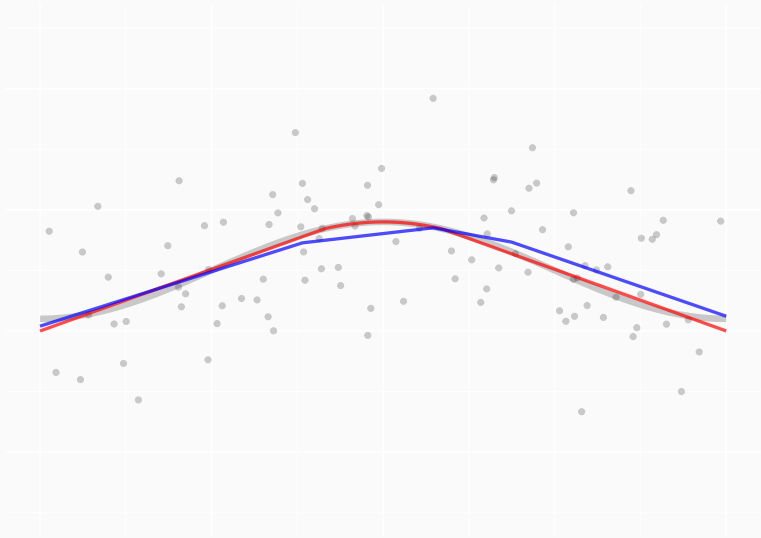


Figure 5: Convex Curves.



**Figure 6:** Concave Curves.

## Non-Gaussian Noise

---

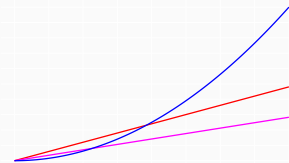
# Summary

$$\begin{aligned}\ell(m) - \ell(\mu^*) &= \|m - \mu^*\|_{L_2(\mathbb{P}_n)}^2 && \text{squared distance} \\ &- \frac{2}{n} \sum_{i=1}^n \varepsilon_i \{m(X_i) - \mu^*(X_i)\} && \text{a mean zero term} \\ &+ \frac{2}{n} \sum_{i=1}^n \{\mu^*(X_i) - \mu(X_i)\} \{m(X_i) - \mu^*(X_i)\} && \text{a non-negative term.}\end{aligned}$$

We can bound error using a corresponding *width*, no matter how noise is distributed.

$$\|\hat{\mu} - \mu^*\|_{L_2(\mathbb{P}_n)} < s + 2\sqrt{\frac{2M_n^2}{\delta n}} \quad \text{w.p. } 1 - \delta \quad \text{for } \frac{s^2}{2} \geq w_\epsilon(\mathcal{M}_s)$$

$$\text{where } w_\epsilon(\mathcal{V}) = \mathbb{E} \max_{v \in \mathcal{V}} \langle \epsilon, v \rangle_{L_2(\mathbb{P}_n)} \quad \text{and} \quad M_n^2 = \mathbb{E} \max_{i \in 1 \dots n} \varepsilon_i^2.$$



To take advantage of gaussian width calculations, we'll bound this width in terms of gaussian width like this:

$$w_\epsilon(\mathcal{M}_s) \leq \alpha w(\mathcal{M}_s)$$

for  $\alpha$  depending on  $\epsilon$  but not  $\mathcal{M}$  or  $s$ .

# Summary

$$\begin{aligned}\ell(m) - \ell(\mu^*) &= \|m - \mu^*\|_{L_2(\mathbb{P}_n)}^2 && \text{squared distance} \\ &- \frac{2}{n} \sum_{i=1}^n \varepsilon_i \{m(X_i) - \mu^*(X_i)\} && \text{a mean zero term} \\ &+ \frac{2}{n} \sum_{i=1}^n \{\mu^*(X_i) - \mu(X_i)\} \{m(X_i) - \mu^*(X_i)\} && \text{a non-negative term.}\end{aligned}$$

We can bound error using a corresponding *width*, no matter how noise is distributed.

$$\|\hat{\mu} - \mu^*\|_{L_2(\mathbb{P}_n)} < s + 2\sqrt{\frac{2M_n^2}{\delta n}} \quad \text{w.p. } 1 - \delta \quad \text{for} \quad \frac{s^2}{2} \geq \mathbf{w}_\epsilon(\mathcal{M}_s)$$

$$\text{where } \mathbf{w}_\epsilon(\mathcal{V}) = \mathbb{E} \max_{v \in \mathcal{V}} \langle \epsilon, v \rangle_{L_2(\mathbb{P}_n)} \quad \text{and} \quad M_n^2 = \mathbb{E} \max_{i \in 1 \dots n} \varepsilon_i^2.$$

At the heart of this comparison  $\mathbf{w}_\epsilon(\cdot) \leq \alpha \mathbf{w}(\cdot)$  are two ideas.

1. **Symmetrization.** We'll substitute for  $\epsilon_i$  a variant that's symmetric around zero.

$$\epsilon_i \rightarrow \epsilon_i - \epsilon'_i \quad \text{where} \quad \epsilon'_i \text{ is an independent copy of } \epsilon_i$$

This substitution *increases* width:  $\mathbf{w}_\epsilon(\cdot) \leq \mathbf{w}_{\epsilon - \epsilon'}(\cdot)$ .

2. **Contraction.** We'll substitute a gaussian vector for our symmetrized noise  $\epsilon - \epsilon'$ .  
We can bound the impact of this substitution in a model-invariant way.

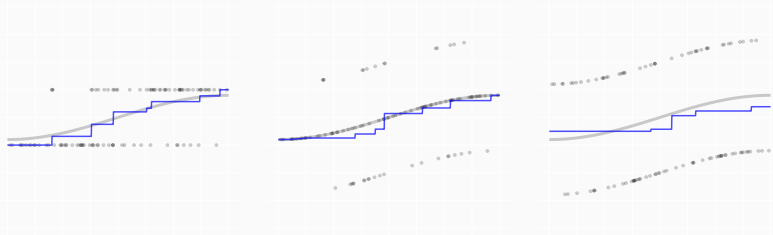
$$\mathbf{w}_{\epsilon - \epsilon'}(\cdot) \leq M'_n \times \mathbf{w}(\cdot) \quad \text{for} \quad M'_n = \mathbb{E} \max_{i \in 1 \dots n} |\epsilon_i - \epsilon'_i| \leq 2 \mathbb{E} \max_{i \in 1 \dots n} |\epsilon_i|.$$



## Non-Gaussian Noise

---

Example: Probabilistic Classification



**Figure 7:** classification noise  $\rightarrow$  symmetrized classification noise  $\rightarrow$  random-sign noise

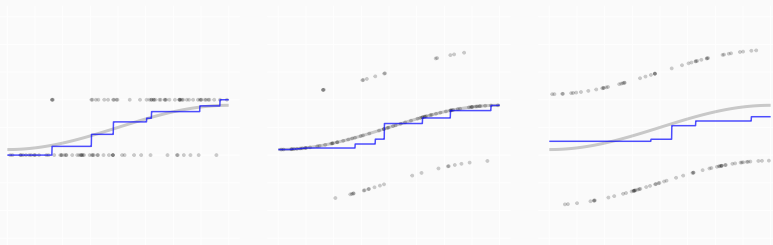
Suppose we have independent *binary observations*.

$$Y_i = \begin{cases} 1 & \text{with conditional probability } \mu(X_i) \\ 0 & \text{otherwise} \end{cases}$$
$$= \mu(X_i) + \varepsilon_i \quad \text{for} \quad \varepsilon_i = \begin{cases} 1 - \mu(X_i) & \text{with conditional probability } \mu(X_i) \\ -\mu(X_i) & \text{with conditional probability } 1 - \mu(X_i) \end{cases}.$$

Note that this *classification noise*  $\varepsilon_i$  has conditional mean zero.

$$\mathbb{E}[\varepsilon_i \mid X_i] = \mu(X_i)\{1 - \mu(X_i)\} + \{1 - \mu(X_i)\}\{-\mu(X_i)\} = 0.$$

# The Setting



**Figure 7:** classification noise  $\rightarrow$  symmetrized classification noise  $\rightarrow$  random-sign noise

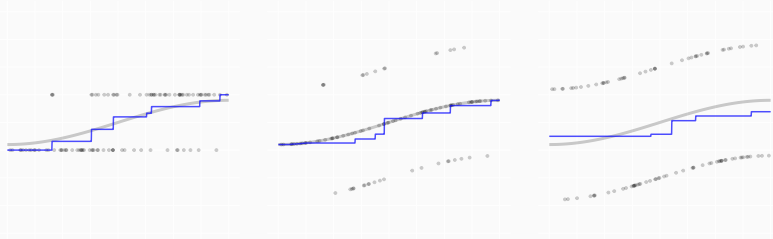
What we need to bound is *classification-noise width*

$$w_{\epsilon}(\mathcal{V}) = \frac{1}{n} \mathbb{E} \max_{v \in \mathcal{V}} \sum_{i=1}^n \epsilon_i v_i.$$

We'll show it's no bigger than a version with *symmetrized noise*.

$$\epsilon_i - \epsilon'_i = \begin{cases} +1 & \text{when } \epsilon_i = 1 - \mu(X_i), \epsilon'_i = \mu(X_i) \\ -1 & \text{when } \epsilon_i = \mu(X_i), \epsilon'_i = 1 - \mu(X_i) \\ 0 & \text{when } \epsilon_i = \epsilon'_i \end{cases}$$

# The Setting



**Figure 7:** classification noise  $\rightarrow$  symmetrized classification noise  $\rightarrow$  random-sign noise

And we'll show that *this* is no bigger than a version with *random sign noise*

$$w_{\epsilon}(\mathcal{V}) \leq w_{\epsilon-\epsilon'}(\mathcal{V}) \leq w_s(\mathcal{V}) \quad \text{where} \quad s_i = \pm 1 \text{ w.p. } 1/2.$$

The trick will be multiplying the symmetrized noise by a random sign.

It's already symmetric, so that doesn't change its distribution.

$$\epsilon_i - \epsilon'_i \stackrel{\text{dist}}{=} s_i(\epsilon_i - \epsilon'_i)$$

Then we'll *contract out* the symmetrized noise, leaving the random sign. You'll see.

## Step 1

We bound our maximum in terms of one involving symmetric noise.

We'll work with an *independent copy*  $\varepsilon'$  of our noise vector  $\varepsilon$ .

$$\begin{aligned} \mathbb{E}_{\varepsilon} \max_{v \in \mathcal{V}} \sum_{i=1}^n \varepsilon_i v_i &\stackrel{(a)}{=} \mathbb{E}_{\varepsilon} \max_{v \in \mathcal{V}} \sum_{i=1}^n (\varepsilon_i - \mathbb{E}_{\varepsilon'} \varepsilon'_i) v_i \\ &\stackrel{(b)}{=} \mathbb{E}_{\varepsilon} \max_{v \in \mathcal{V}} \mathbb{E}_{\varepsilon'} \sum_{i=1}^n (\varepsilon_i - \varepsilon'_i) v_i \\ &\stackrel{(c)}{\leq} \mathbb{E}_{\varepsilon} \mathbb{E}_{\varepsilon'} \max_{v \in \mathcal{V}} \sum_{i=1}^n (\varepsilon_i - \varepsilon'_i) v_i. \end{aligned}$$

Why do these steps work?

(a)  $\mathbb{E}_{\varepsilon'} \varepsilon'_i = 0$ .

## Step 1

We bound our maximum in terms of one involving symmetric noise.

We'll work with an *independent copy*  $\varepsilon'$  of our noise vector  $\varepsilon$ .

$$\begin{aligned} \mathbb{E}_{\varepsilon} \max_{v \in \mathcal{V}} \sum_{i=1}^n \varepsilon_i v_i &\stackrel{(a)}{=} \mathbb{E}_{\varepsilon} \max_{v \in \mathcal{V}} \sum_{i=1}^n (\varepsilon_i - \mathbb{E}_{\varepsilon'} \varepsilon'_i) v_i \\ &\stackrel{(b)}{=} \mathbb{E}_{\varepsilon} \max_{v \in \mathcal{V}} \mathbb{E}_{\varepsilon'} \sum_{i=1}^n (\varepsilon_i - \varepsilon'_i) v_i \\ &\stackrel{(c)}{\leq} \mathbb{E}_{\varepsilon} \mathbb{E}_{\varepsilon'} \max_{v \in \mathcal{V}} \sum_{i=1}^n (\varepsilon_i - \varepsilon'_i) v_i. \end{aligned}$$

Why do these steps work?

(a)  $\mathbb{E}_{\varepsilon'} \varepsilon'_i = 0$ .

(b) Expectation is linear.

## Step 1

We bound our maximum in terms of one involving symmetric noise.

We'll work with an *independent copy*  $\varepsilon'$  of our noise vector  $\varepsilon$ .

$$\begin{aligned} \mathbb{E}_{\varepsilon} \max_{v \in \mathcal{V}} \sum_{i=1}^n \varepsilon_i v_i &\stackrel{(a)}{=} \mathbb{E}_{\varepsilon} \max_{v \in \mathcal{V}} \sum_{i=1}^n (\varepsilon_i - \mathbb{E}_{\varepsilon'} \varepsilon'_i) v_i \\ &\stackrel{(b)}{=} \mathbb{E}_{\varepsilon} \max_{v \in \mathcal{V}} \mathbb{E}_{\varepsilon'} \sum_{i=1}^n (\varepsilon_i - \varepsilon'_i) v_i \\ &\stackrel{(c)}{\leq} \mathbb{E}_{\varepsilon} \mathbb{E}_{\varepsilon'} \max_{v \in \mathcal{V}} \sum_{i=1}^n (\varepsilon_i - \varepsilon'_i) v_i. \end{aligned}$$

Why do these steps work?

- (a)  $\mathbb{E}_{\varepsilon'} \varepsilon'_i = 0$ .
- (b) Expectation is linear.
- (c) Maximizing the average gives us something smaller than averaging the maxima.

## Step 1

We bound our maximum in terms of one involving symmetric noise.

We'll work with an *independent copy*  $\varepsilon'$  of our noise vector  $\varepsilon$ .

$$\begin{aligned} \mathbb{E}_{\varepsilon} \max_{v \in \mathcal{V}} \sum_{i=1}^n \varepsilon_i v_i &\stackrel{(a)}{=} \mathbb{E}_{\varepsilon} \max_{v \in \mathcal{V}} \sum_{i=1}^n (\varepsilon_i - \mathbb{E}_{\varepsilon'} \varepsilon'_i) v_i \\ &\stackrel{(b)}{=} \mathbb{E}_{\varepsilon} \max_{v \in \mathcal{V}} \mathbb{E}_{\varepsilon'} \sum_{i=1}^n (\varepsilon_i - \varepsilon'_i) v_i \\ &\stackrel{(c)}{\leq} \mathbb{E}_{\varepsilon} \mathbb{E}_{\varepsilon'} \max_{v \in \mathcal{V}} \sum_{i=1}^n (\varepsilon_i - \varepsilon'_i) v_i. \end{aligned}$$

Why do these steps work?

- (a)  $\mathbb{E}_{\varepsilon'} \varepsilon'_i = 0$ .
- (b) Expectation is linear.
- (c) Maximizing the average gives us something smaller than averaging the maxima.
  - In (c), we choose the maximizing  $v \in \mathcal{V}$  for each  $\varepsilon'$ .



## Step 1

We bound our maximum in terms of one involving symmetric noise.

We'll work with an *independent copy*  $\varepsilon'$  of our noise vector  $\varepsilon$ .

$$\begin{aligned} \mathbb{E}_{\varepsilon} \max_{v \in \mathcal{V}} \sum_{i=1}^n \varepsilon_i v_i &\stackrel{(a)}{=} \mathbb{E}_{\varepsilon} \max_{v \in \mathcal{V}} \sum_{i=1}^n (\varepsilon_i - \mathbb{E}_{\varepsilon'} \varepsilon'_i) v_i \\ &\stackrel{(b)}{=} \mathbb{E}_{\varepsilon} \max_{v \in \mathcal{V}} \mathbb{E}_{\varepsilon'} \sum_{i=1}^n (\varepsilon_i - \varepsilon'_i) v_i \\ &\stackrel{(c)}{\leq} \mathbb{E}_{\varepsilon} \mathbb{E}_{\varepsilon'} \max_{v \in \mathcal{V}} \sum_{i=1}^n (\varepsilon_i - \varepsilon'_i) v_i. \end{aligned}$$

Why do these steps work?

- (a)  $\mathbb{E}_{\varepsilon'} \varepsilon'_i = 0$ .
- (b) Expectation is linear.
- (c) Maximizing the average gives us something smaller than averaging the maxima.
  - In (c), we choose the maximizing  $v \in \mathcal{V}$  for each  $\varepsilon'$ .
  - If we wanted to choose the same one each time, like we do in (b), we could.

We introduce independent random signs  $s_i = \pm 1$  w.p.  $1/2$ , changing nothing.

$$\mathbb{E}_\varepsilon \mathbb{E}_{\varepsilon'} \max_{v \in \mathcal{V}} \sum_{i=1}^n (\varepsilon_i - \varepsilon'_i) v_i = \mathbb{E}_s \mathbb{E}_\varepsilon \mathbb{E}_{\varepsilon'} \max_{v \in \mathcal{V}} \sum_{i=1}^n s_i (\varepsilon_i - \varepsilon'_i) v_i.$$

Why does this change nothing?

We introduce independent random signs  $s_i = \pm 1$  w.p.  $1/2$ , changing nothing.

$$\mathbb{E}_\varepsilon \mathbb{E}_{\varepsilon'} \max_{v \in \mathcal{V}} \sum_{i=1}^n (\varepsilon_i - \varepsilon'_i) v_i = \mathbb{E}_s \mathbb{E}_\varepsilon \mathbb{E}_{\varepsilon'} \max_{v \in \mathcal{V}} \sum_{i=1}^n s_i (\varepsilon_i - \varepsilon'_i) v_i.$$

Why does this change nothing?

- Because the inner mean  $(\mathbb{E}_\varepsilon \mathbb{E}_{\varepsilon'})$  doesn't depend on the signs  $s_i$ .
- That's because  $\varepsilon_i$  and  $\varepsilon'_i$  have the same distribution.
- And this implies  $(\varepsilon_i - \varepsilon'_i)$  and  $(\varepsilon'_i - \varepsilon) = -(\varepsilon_i - \varepsilon'_i)$  do, too.

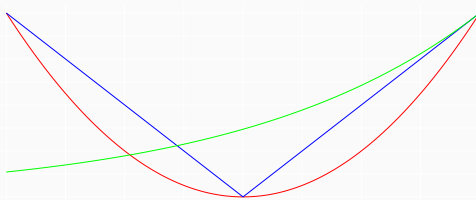
## Step 3

We swap the order of our averages and think about the inner average as a *function* of our vector of symmetric noise.

$$\begin{aligned} \mathbb{E}_s \mathbb{E}_\varepsilon \mathbb{E}_{\varepsilon'} \max_{v \in \mathcal{V}} \sum_{i=1}^n s_i (\varepsilon_i - \varepsilon'_i) v_i &= \mathbb{E}_\varepsilon \mathbb{E}_{\varepsilon'} \mathbb{E}_s \max_{v \in \mathcal{V}} \sum_{i=1}^n s_i (\varepsilon_i - \varepsilon'_i) v_i \\ &= \mathbb{E}_\varepsilon \mathbb{E}_{\varepsilon'} f(\varepsilon - \varepsilon') \quad \text{for} \quad f(u) = \mathbb{E}_s \max_{v \in \mathcal{V}} \sum_{i=1}^n s_i u_i v_i. \end{aligned}$$

This function  $f$  is convex.

What does that mean? These, for example, are all convex.



$$f\{(1-\lambda)a + \lambda b\} \leq (1-\lambda)f(a) + \lambda f(b) \quad \text{for} \quad \lambda \in [0, 1]. \quad \text{That's Convexity}$$

## Step 3

We swap the order of our averages and think about the inner average as a *function* of our vector of symmetric noise.

$$\begin{aligned} \mathbb{E}_s \mathbb{E}_\varepsilon \mathbb{E}_{\varepsilon'} \max_{v \in \mathcal{V}} \sum_{i=1}^n s_i (\varepsilon_i - \varepsilon'_i) v_i &= \mathbb{E}_\varepsilon \mathbb{E}_{\varepsilon'} \mathbb{E}_s \max_{v \in \mathcal{V}} \sum_{i=1}^n s_i (\varepsilon_i - \varepsilon'_i) v_i \\ &= \mathbb{E}_\varepsilon \mathbb{E}_{\varepsilon'} f(\varepsilon - \varepsilon') \quad \text{for} \quad f(u) = \mathbb{E}_s \max_{v \in \mathcal{V}} \sum_{i=1}^n s_i u_i v_i. \end{aligned}$$

This function  $f$  is convex.

How do we know? Maximizing each term is better than maximizing their sum.

$$\begin{aligned} f\{(1-\lambda)a + \lambda b\} &= \mathbb{E}_s \max_{v \in \mathcal{V}} \left\{ (1-\lambda) \sum_{i=1}^n s_i a_i v_i + \lambda \sum_{i=1}^n s_i b_i v_i \right\} \\ &\leq \mathbb{E}_s \left\{ \max_{v \in \mathcal{V}} (1-\lambda) \sum_{i=1}^n s_i a_i v_i + \max_{v \in \mathcal{V}} \lambda \sum_{i=1}^n s_i b_i v_i \right\} \\ &= (1-\lambda) \mathbb{E}_s \max_{v \in \mathcal{V}} \sum_{i=1}^n s_i a_i v_i + \lambda \mathbb{E}_s \max_{v \in \mathcal{V}} \sum_{i=1}^n s_i b_i v_i \\ &= (1-\lambda)f(a) + \lambda f(b). \end{aligned}$$

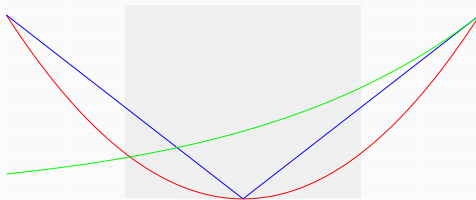
## Step 3

We swap the order of our averages and think about the inner average as a *function* of our vector of symmetric noise.

$$\begin{aligned} \mathbb{E}_s \mathbb{E}_\varepsilon \mathbb{E}_{\varepsilon'} \max_{v \in \mathcal{V}} \sum_{i=1}^n s_i (\varepsilon_i - \varepsilon'_i) v_i &= \mathbb{E}_\varepsilon \mathbb{E}_{\varepsilon'} \mathbb{E}_s \max_{v \in \mathcal{V}} \sum_{i=1}^n s_i (\varepsilon_i - \varepsilon'_i) v_i \\ &= \mathbb{E}_\varepsilon \mathbb{E}_{\varepsilon'} f(\varepsilon - \varepsilon') \quad \text{for} \quad f(u) = \mathbb{E}_s \max_{v \in \mathcal{V}} \sum_{i=1}^n s_i u_i v_i. \end{aligned}$$

This function  $f$  is convex.

Why does this matter? The max of a convex function over a cube occurs at a corner.



What cube?

The vector of symmetric noise,  $\varepsilon - \varepsilon'$ , is in the *unit cube*  $[-1, 1]^n$ .

$$\varepsilon_i - \varepsilon'_i = \begin{cases} 0 & \text{when } \varepsilon_i = \varepsilon'_i \\ +1 & \text{when } \varepsilon_i = 1 - \mu(X_i), \varepsilon'_i = \mu(X_i) \\ -1 & \text{when } \varepsilon_i = \mu(X_i), \varepsilon'_i = 1 - \mu(X_i). \end{cases}$$

The average over this random vector is bounded by the maximum over the cube it's in.

$$\begin{aligned} \mathbb{E}_{\varepsilon} \mathbb{E}_{\varepsilon'} \mathbb{E}_s \max_{v \in \mathcal{V}} \sum_{i=1}^n s_i (\varepsilon_i - \varepsilon'_i) v_i &\leq \max_{u \in [-1, 1]^n} \mathbb{E}_s \max_{v \in \mathcal{V}} \sum_{i=1}^n s_i u_i v_i \\ &= \max_{u \in [-1, 1]^n} f(u) \quad \text{max over the cube} \\ &= \max_{u \in \{-1, 1\}^n} f(u) \quad \text{max over its corners} \end{aligned}$$

We characterize this maximum over corners. Remember what  $f$  is.

$$\begin{aligned}\max_{u \in \{-1,1\}^n} f(u) &= \max_{u \in \{-1,1\}^n} \mathbb{E}_s \max_{v \in \mathcal{V}} \sum_{i=1}^n s_i u_i v_i \\ &= \mathbb{E}_s \max_{v \in \mathcal{V}} \sum_{i=1}^n s_i v_i.\end{aligned}$$

Why?

Hint. What's the distribution of  $s_i$ ? And  $s_i u_i$  for  $u_i \in \{-1,1\}$ ?



We characterize this maximum over corners. Remember what  $f$  is.

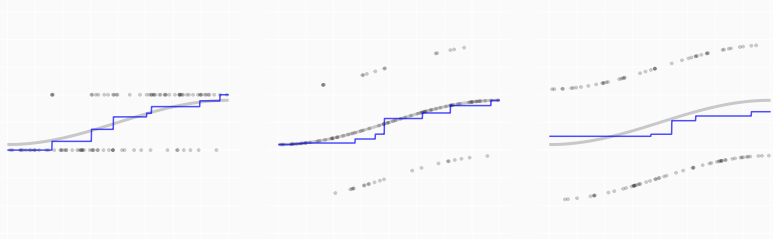
$$\begin{aligned}\max_{u \in \{-1, 1\}^n} f(u) &= \max_{u \in \{-1, 1\}^n} \mathbb{E}_s \max_{v \in \mathcal{V}} \sum_{i=1}^n s_i u_i v_i \\ &= \mathbb{E}_s \max_{v \in \mathcal{V}} \sum_{i=1}^n s_i v_i.\end{aligned}$$

Why?

Hint. What's the distribution of  $s_i$ ? And  $s_i u_i$  for  $u_i \in \{-1, 1\}$ ?

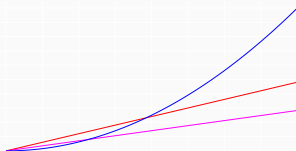
- For  $u_i \in \{-1, 1\}$ , the distributions of  $u_i$  and  $s_i u_i$  are the same.
- So the distribution of the sum, and its maximum, are the same at every corner  $u$ .
- Including the vector of all ones  $u = (1, 1, \dots, 1)$ .

# Summary



classification noise width  $\leq$  symmetrized classification noise width  $\leq$  random sign width  
 This means probabilistic classification is *easier* than regression with random sign noise. Or, at least, that we get a better bound.

$$\frac{s^2}{2} \geq w_s(\mathcal{M}_s) \quad \text{and} \quad w_s(\mathcal{M}_s) \geq w_\varepsilon(\mathcal{M}_s) \quad \Rightarrow \quad \frac{s^2}{2} \geq w_\varepsilon(\mathcal{M}_s)$$



People call random sign width, or something like it, *Rademacher Complexity*.

$$\text{Rademacher Complexity}(\mathcal{V}) = \mathbb{E} \max_{v \in \mathcal{V}} \langle s, v \rangle_{L_2(\mathbf{P}_n)} \quad \text{for i.i.d. } s_i = \begin{cases} +1 & \text{w.p. } 1/2 \\ -1 & \text{w.p. } 1/2 \end{cases}$$

or maybe  $= \mathbb{E} \max_{v \in \mathcal{V}} |\langle s, v \rangle_{L_2(\mathbf{P}_n)}|$

- This second definition is the same if  $\mathcal{V}$  is symmetric, i.e.  $v \in \mathcal{V} \implies -v \in \mathcal{V}$ .
- Otherwise, it can be a little bigger.
  - At most  $2\times$  bigger. Prove it!
  - Use the bound  $\max a, b \leq a + b$  and the symmetry of  $s$ 's distribution.

## Non-Gaussian Noise

---

### The General Case

# Symmetrization and Contraction: Examples

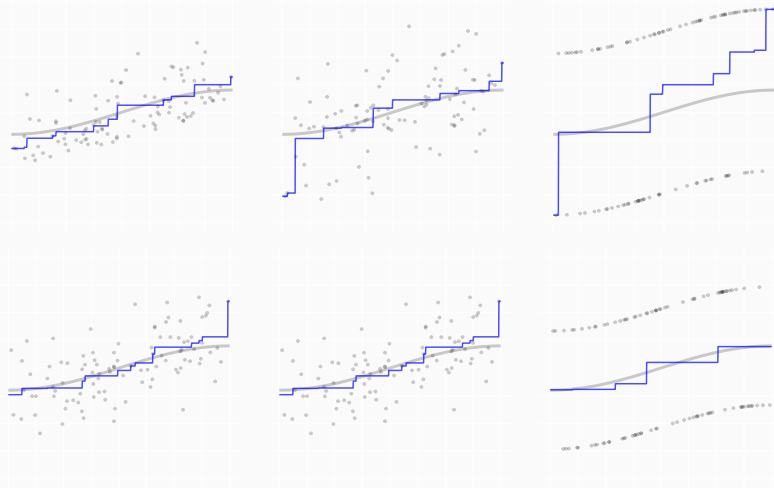


Figure 8: real noise  $\rightarrow$  symmetrized noise  $\rightarrow$  scaled sign noise

$$w_{\varepsilon}(\mathcal{V}) \leq w_{s(\varepsilon - \varepsilon')}(\mathcal{V}) \leq 2 w_{s\varepsilon}(\mathcal{V})$$

$$\begin{aligned} \mathbb{E} \max_{v \in \mathcal{V}} \sum_{i=1}^n \varepsilon_i v_i &= \mathbb{E} \max_{v \in \mathcal{V}} \sum_{i=1}^n (\varepsilon_i - \mathbb{E} \varepsilon'_i) v_i \\ &\stackrel{(a)}{\leq} \mathbb{E} \mathbb{E}' \max_{v \in \mathcal{V}} \sum_{i=1}^n (\varepsilon_i - \varepsilon'_i) v_i \\ &= \mathbb{E}_s \mathbb{E} \mathbb{E}' \max_{v \in \mathcal{V}} \sum_{i=1}^n s_i (\varepsilon_i - \varepsilon'_i) v_i \\ &\stackrel{(b)}{\leq} \mathbb{E}_s \mathbb{E} \max_{v \in \mathcal{V}} \sum_{i=1}^n s_i \varepsilon_i + \mathbb{E}_s \mathbb{E}' \max_{v \in \mathcal{V}} \sum_{i=1}^n s_i \varepsilon'_i v_i \\ &= 2 \mathbb{E}_s \mathbb{E} \max_{v \in \mathcal{V}} \sum_{i=1}^n \varepsilon_i s_i v_i. \end{aligned}$$

(a) Replacing  $\varepsilon_i$  with  $s_i(\varepsilon_i - \varepsilon'_i)$  is 'free'.

- We stopped here in our example because  $\varepsilon_i - \varepsilon'_i$  was easy to bound.
- Generally, we take an extra step to express things in terms of  $\varepsilon_i$  again.

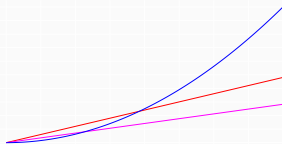
(b) Replacing  $\varepsilon_i$  with  $s_i \varepsilon_i$  increases width by at most  $2 \times$ .

$$w_\eta(\mathcal{V}) = w_{s\eta}(\mathcal{V}) \leq \mathbb{E} \|\eta\|_\infty w_\eta(\mathcal{V}) \quad \text{if } \eta \stackrel{\text{dist}}{=} -\eta.$$

$$\begin{aligned} \mathbb{E}_s \mathbb{E}_\eta \max_{v \in \mathcal{V}} \sum_{i=1}^n \eta_i s_i v_i &\leq \mathbb{E}_\eta \max_{\substack{u \in \mathbb{R}^n \\ \|u_i\| \leq \|\eta\|_\infty}} \mathbb{E}_s \max_{v \in \mathcal{V}} \sum_{i=1}^n u_i s_i v_i \\ &= \mathbb{E}_\eta \|\eta\|_\infty \max_{u \in [-1, 1]^n} \mathbb{E}_s \max_{v \in \mathcal{V}} \sum_{i=1}^n u_i s_i v_i \\ &= \mathbb{E}_\eta \|\eta\|_\infty \times \max_{u \in [-1, 1]^n} \mathbb{E}_s \max_{v \in \mathcal{V}} \sum_{i=1}^n u_i s_i v_i \\ &= \mathbb{E}_\eta \|\eta\|_\infty \times \mathbb{E}_s \max_{v \in \mathcal{V}} \sum_{i=1}^n s_i v_i \end{aligned}$$

- We can 'contract out' any **symmetrically distributed** noise vector  $\eta$  by ...
  1. multiplying in independent random signs  $s_i$ . Symmetry  $\implies s_i \eta_i \stackrel{\text{dist}}{=} \eta_i$ .
  2. maximizing over a cube containing  $\eta$ .
- We just have to use a big enough cube.
  - In our example,  $\eta = \varepsilon - \varepsilon'$  was in the unit cube  $[-1, 1]^n$  deterministically.
  - Generally, we maximize over a random cube  $[-\|\eta\|_\infty, \|\eta\|_\infty]^n$ .
  - And we can pull out the cube's radius  $\|\eta\|_\infty$  as a multiplicative factor.

# Implications for Regression



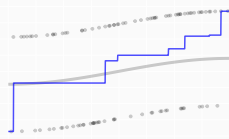
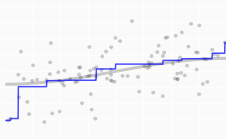
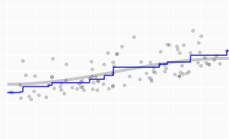
$$w_{\varepsilon}(\mathcal{V}) \leq \mathbb{E} \|\varepsilon_i - \varepsilon'_i\|_{\infty} w_s(\mathcal{V}) \leq 2 \mathbb{E} \|\varepsilon_i\|_{\infty} w_s(\mathcal{V})$$

Regression with arbitrary independent noise, i.e.

$$Y_i = \mu(X_i) + \varepsilon_i \quad \text{where} \quad \varepsilon_1 \dots \varepsilon_n \text{ are independent,}$$

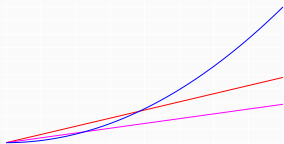
is no harder than with scaled-up random sign noise, i.e.

$$Y_i = \mu(X_i) + Ms_i \quad \text{for} \quad M = \mathbb{E} \|\varepsilon_i - \varepsilon'_i\|_{\infty} \quad \text{and} \quad s_i = \begin{cases} +1 & \text{w.p. } 1/2 \\ -1 & \text{w.p. } 1/2 \end{cases}.$$





# The Symmetric Case



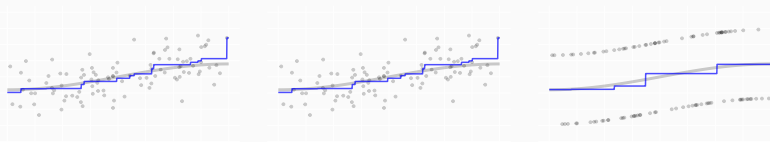
$$w_\varepsilon(\mathcal{V}) \leq \mathbb{E}\|\varepsilon_i\|_\infty w_s(\mathcal{V})$$

Regression with arbitrary independent *symmetric* noise, i.e.

$$Y_i = \mu(X_i) + \varepsilon_i \quad \text{where} \quad \varepsilon_1 \dots \varepsilon_n \text{ are independent with } \varepsilon_i \stackrel{\text{dist}}{=} -\varepsilon_i,$$

is no harder than with scaled-up random sign noise, i.e.

$$Y_i = \mu(X_i) + Ms_i \quad \text{for}^2 \quad M = \mathbb{E}\|\varepsilon_i\|_\infty \quad \text{and} \quad s_i = \begin{cases} +1 & \text{w.p. } 1/2 \\ -1 & \text{w.p. } 1/2 \end{cases}.$$



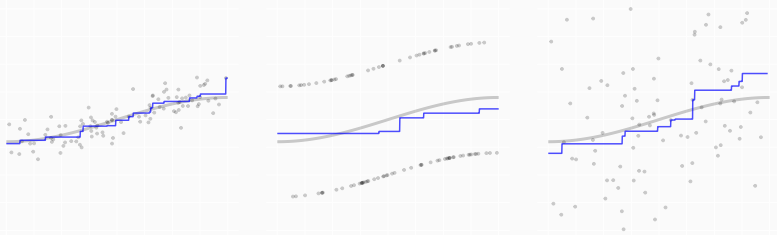
**Figure 9:** real noise  $\rightarrow$  symmetrized noise  $\rightarrow$  scaled sign noise

<sup>2</sup> $M = \mathbb{E}\|\varepsilon_i\|_\infty \leq 2\sigma\sqrt{2\log(2n)}$  for  $\varepsilon_i \sim N(0, \sigma^2)$ . See Appendix B of the Gaussian Width Homework.

## Non-Gaussian Noise

---

Comparison to the Gaussian Case



- So far, we've bounded arbitrary-noise width in terms of random-sign width.
- But often, it's easier to understand gaussian width. That's good enough.<sup>3</sup>

$$\frac{1}{2\sqrt{\log(2n)}} w_g(\mathcal{V}) \leq w_s(\mathcal{V}) \leq \sqrt{\frac{\pi}{2}} w_g(\mathcal{V})$$

$\approx .2 \text{ for } n=100$   $\approx 1.25$

- We just saw it can't be **that much bigger** than random-sign width.
- And we can show it's **at least 4/5 as big**.

$$\mathbb{E} \max_{v \in \mathcal{V}} \sum_{i=1}^n g_i v_i = \mathbb{E}_s \mathbb{E}_g \max_{v \in \mathcal{V}} \sum_{i=1}^n |g_i| s_i v_i \geq \mathbb{E}_s \max_{v \in \mathcal{V}} \sum_{i=1}^n \mathbb{E}_g |g_i| s_i v_i = \sqrt{\frac{2}{\pi}}$$

<sup>3</sup>We can show  $.125 w_g(\mathcal{V}) \leq w_s(\mathcal{V}) \leq 1.25 w_g(\mathcal{V})$  for  $n \leq 10$  trillion by bounding  $\mathbb{E} \|g\|_\infty$  more carefully.

# Comparison in Steps

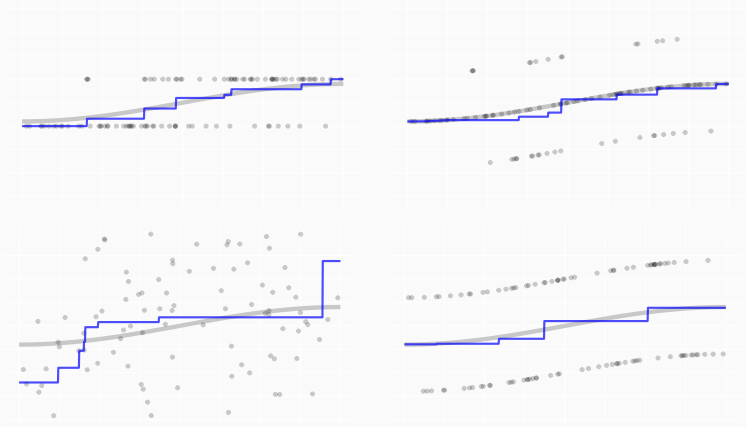


Figure 10: real noise  $\rightarrow$  symmetrized noise  $\downarrow$  scaled sign noise  $\leftarrow$  scaled gaussian noise

$$w_{\varepsilon}(\mathcal{V}) \leq w_{\varepsilon - \varepsilon'}(\mathcal{V}) \leq \underset{\leq 2 \mathbb{E} \|\varepsilon\|_{\infty}}{\mathbb{E} \|\varepsilon - \varepsilon'\|_{\infty}} \quad w_s(\mathcal{V}) \leq \sqrt{\frac{\pi}{2}} \mathbb{E} \|\varepsilon - \varepsilon'\|_{\infty} \quad w_g(\mathcal{V}) \leq \sqrt{2\pi} \approx 2.5 \times \mathbb{E} \|\varepsilon\|_{\infty}$$

# Implications for Regression

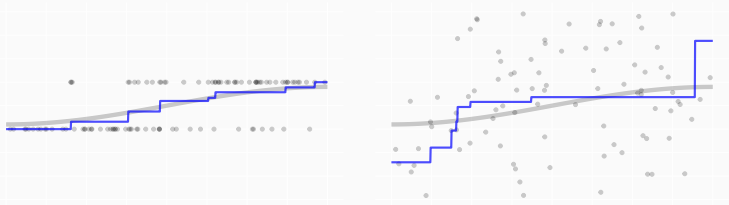


Figure 11: real noise  $\rightarrow$  scaled gaussian noise

For any noise vector  $\varepsilon$  with independent components  $\varepsilon_i$ ,

$$w_{\varepsilon}(\mathcal{V}) \leq 2 \mathbb{E} \|\varepsilon\|_{\infty} \cdot w_s(\mathcal{V}) \leq \sqrt{2\pi} \mathbb{E} \|\varepsilon\|_{\infty} \cdot w_g(\mathcal{V}).$$

- We can bound the width  $w_{\varepsilon}$  in terms of
  1. random-sign width
  2. the maximum absolute value of  $\varepsilon$ 's components.
- And we can bound random-sign width in terms of gaussian width.

This means we don't have to bound a million different kinds of widths for each model.  
We can bound random-sign width or gaussian width. Whichever is easier.

## Width Comparisons imply Radius Comparisons

**Claim.** If  $w_\varepsilon \leq \alpha w_\eta$  for  $\alpha \geq 1$ , then the critical radius using noise  $\varepsilon$  is at most  $\alpha$  times the critical radius using noise  $\eta$ , i.e.

$$\frac{(\alpha s)^2}{2} \geq w_\varepsilon(\mathcal{M}_{\alpha s}) \quad \text{if} \quad \frac{s^2}{2} \geq w_\eta(\mathcal{M}_s) \quad \text{and} \quad w_\varepsilon \leq \alpha w_\eta \quad \text{for} \quad \alpha \geq 1.$$

**Proof.** If  $s^2/2 \geq w_\eta(\mathcal{M}_s)$ , then

## Width Comparisons imply Radius Comparisons

**Claim.** If  $w_\varepsilon \leq \alpha w_\eta$  for  $\alpha \geq 1$ , then the critical radius using noise  $\varepsilon$  is at most  $\alpha$  times the critical radius using noise  $\eta$ , i.e.

$$\frac{(\alpha s)^2}{2} \geq w_\varepsilon(\mathcal{M}_{\alpha s}) \quad \text{if} \quad \frac{s^2}{2} \geq w_\eta(\mathcal{M}_s) \quad \text{and} \quad w_\varepsilon \leq \alpha w_\eta \quad \text{for} \quad \alpha \geq 1.$$

**Proof.** If  $s^2/2 \geq w_\eta(\mathcal{M}_s)$ , then

|   |  |
|---|--|
| $\alpha s/2 \geq \alpha w_\eta(\mathcal{M}_s)/s$        | multiplying both sides by $\alpha/s$                   |
| $\geq \alpha w_\eta(\mathcal{M}_{\alpha s})/(\alpha s)$ | using sublinearity of $f(s) = w_\eta(\mathcal{M}_s)$   |
| $\geq w_\varepsilon(\mathcal{M}_{\alpha s})/(\alpha s)$ | using our premise $\alpha w_\eta \geq w_\varepsilon$ . |

Multiplying both sides by  $\alpha s$ , we get our claim.

Background: Convex Functions Are  
Maximized At Extreme Points

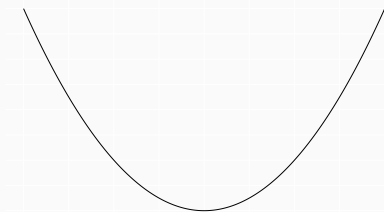
---



## Definition

A function  $f$  is convex if *secants* lie above the curve.

$$f\{(1 - \lambda)a + \lambda b\} \leq (1 - \lambda)f(a) + \lambda f(b) \quad \text{for } \lambda \in [0, 1]$$



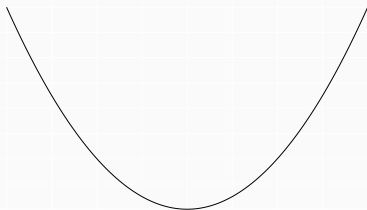
We can give this a *probabilistic interpretation* for a random variable  $Z_\lambda$ .

$$f(E Z_\lambda) \leq E f(Z_\lambda) \quad \text{where } Z_\lambda =$$

## Definition

A function  $f$  is convex if *secants* lie above the curve.

$$f\{(1 - \lambda)a + \lambda b\} \leq (1 - \lambda)f(a) + \lambda f(b) \quad \text{for } \lambda \in [0, 1]$$



We can give this a *probabilistic interpretation* for a random variable  $Z_\lambda$ .

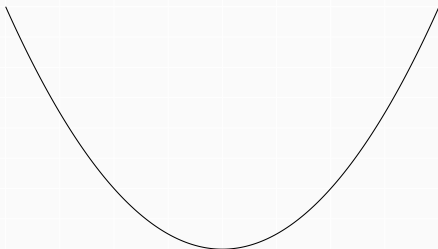
$$f(\mathbb{E} Z_\lambda) \leq \mathbb{E} f(Z_\lambda) \quad \text{where} \quad Z_\lambda = \begin{cases} a & \text{w.p. } 1 - \lambda \\ b & \text{w.p. } \lambda \end{cases}$$

# Jensen's Inequality

In fact, this is true all random variables  $Z$ .  
If  $f$  is convex, its mean value exceeds its value at the mean.

$$f(\mathbb{E} Z) \leq \mathbb{E} f(Z)$$

That's called Jensen's Inequality.



You can prove it for discrete random variables via induction.

# Jensen's Inequality Proof

## Base case.

It's true for random variables taking on 2 values.

$$f(\lambda_1 z_1 + \lambda_2 z_2) \leq \lambda_1 f(z_1) + \lambda_2 f(z_2) \quad \text{if} \quad \lambda_1, \lambda_2 \geq 0 \quad \text{satisfy} \quad \lambda_1 + \lambda_2 = 1$$

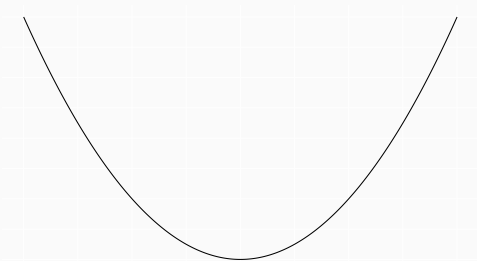
## Inductive Step.

We'll show that if it's true for random variables taking on  $n - 1$  values, then it's also true for ones taking on  $n$  values.

$$\begin{aligned} f\left\{\sum_{i=1}^n \lambda_i z_i\right\} &= f\left\{(1 - \lambda_n)\left(\sum_{i=1}^{n-1} \frac{\lambda_i}{1 - \lambda_n} z_i\right) + \lambda_n z_n\right\} \\ &\leq (1 - \lambda_n) f\left(\sum_{i=1}^{n-1} \frac{\lambda_i}{1 - \lambda_n} z_i\right) + \lambda_n f(z_n) \\ &\leq (1 - \lambda_n) \sum_{i=1}^{n-1} \frac{\lambda_i}{1 - \lambda_n} f(z_i) + \lambda_n f(z_n) \\ &= \sum_{i=1}^{n-1} \lambda_i f(z_i) + \lambda_n f(z_n) \end{aligned}$$

# Maxima of Convex Functions

Convex functions have no local maxima.



That means the maximum of a convex function over an interval occurs at an endpoint.

**Proof.**

$$\max_{x \in [a,b]} f(x) = \max_{\lambda \in [0,1]} f\{(1-\lambda)a + \lambda b\} \leq \max_{\lambda \in [0,1]} (1-\lambda)f(a) + \lambda f(b) = \max\{f(a), f(b)\}$$

This is essentially true in higher dimensions as well.

We just need the right generalizations of *interval* and its *endpoints*.

# Convex Polytopes

The natural generalizations a *convex polytope* and its *extreme points*.

## Definitions.

A **convex polytope** is the set of all weighted averages of some set of vectors  $u_1 \dots u_K$ .

$$\mathcal{U} = \left\{ \sum_i \lambda_i u_i : \lambda \in \Lambda \right\} \quad \text{where} \quad \Lambda = \left\{ \lambda : \lambda_i \geq 0 \text{ for all } i \text{ and } \sum_i \lambda_i = 1 \right\}$$

Its **extreme points** are the subset of these vectors that are not redundant.  
That is, they're the ones we cannot write as weighted averages of the others.

## Examples.

- A triangle is the set of weighted averages of its three vertices, its extreme points.
- A square is the set of weighted averages of its four vertices, its extreme points.
- A cube in  $\mathbb{R}^n$  is the set of weighted averages of its  $2^n$  vertices, its extreme points.

# Maxima of Convex Functions over Polytopes

The maximum of a convex function over a convex polytope occurs at an extreme point.

**Proof.**

It's more-or-less the same as the one-dimensional case. We use Jensen's inequality.

$$\max_{u \in \mathcal{U}} f(u) = \max_{\lambda \in \Lambda} f\left(\sum_i \lambda_i u_i\right) \leq \max_{\lambda \in \Lambda} \sum_i \lambda_i f(u_i) \leq \max_i f(u_i)$$