

# Machine Learning Theory

## Sampling, Misspecification, and Non-Gaussian Noise

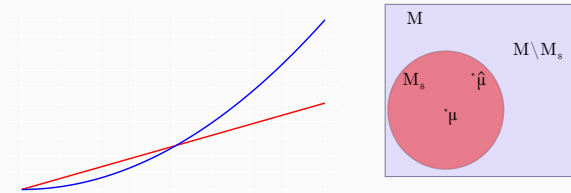
---

David A. Hirshberg

April 10, 2025

Emory University

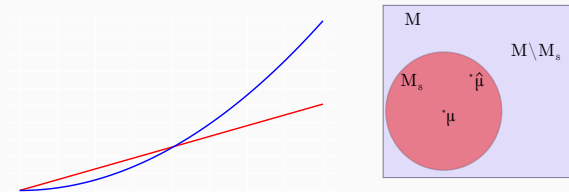
## Where We Left Off



What do we know about the error of this least squares estimator  $\hat{\mu}$ ?

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

# Where We Left Off



What do we know about the error of this least squares estimator  $\hat{\mu}$ ?

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

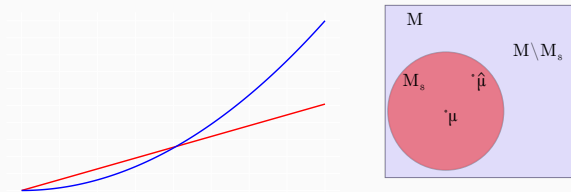
Here's what we proved in lecture.

$$\|\hat{\mu} - \mu^*\|_{L_2(P_n)} < s \quad \text{w.p. } 1 - \delta \quad \text{for} \quad \frac{s^2}{2\sigma} \geq \mathbf{w}(\mathcal{M}_s^\circ) + s \sqrt{\frac{2\Sigma_n}{\delta n}}$$

where  $\mathbf{w}(\mathcal{V}) = \mathbb{E} \max_{v \in \mathcal{V}} \langle g, v \rangle_{L_2(P_n)}$  and  $\Sigma_n = \sigma^2 \{1 + 2 \log(2n)\}$  for  $g_i \stackrel{iid}{\sim} N(0, 1)$

if  $Y_i = \mu(X_i) + \varepsilon_i$  for  $\varepsilon_i \stackrel{iid}{\sim} N(0, \sigma^2)$  for  $\mu \in \mathcal{M}$

# Where We Left Off



What do we know about the error of this least squares estimator  $\hat{\mu}$ ?

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

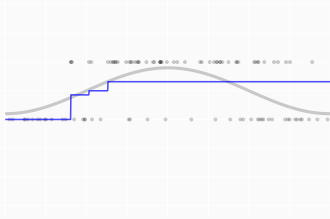
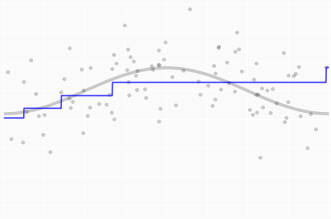
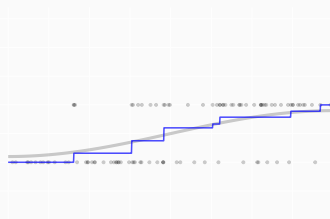
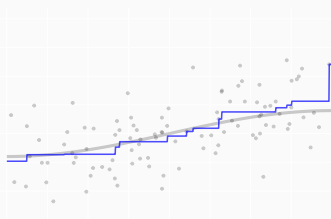
Here's a simplified version of you're proving for homework.

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

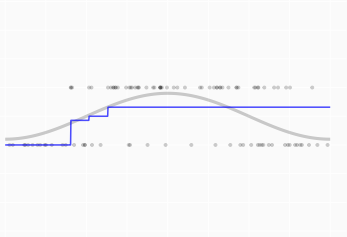
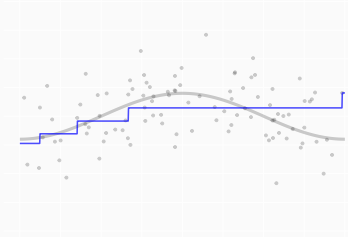
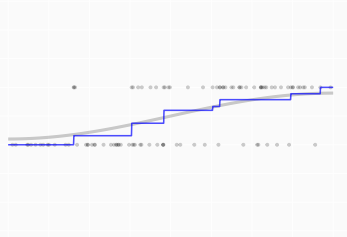
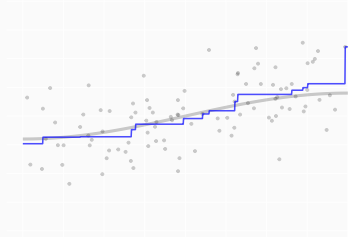
where  $\mathbf{w}(\mathcal{V}) = \mathbb{E} \max_{v \in \mathcal{V}} \langle g, v \rangle_{L_2(P_n)}$  and  $\Sigma_n = \sigma^2 \{1 + 2 \log(2n)\}$  for  $g_i \stackrel{iid}{\sim} N(0, 1)$

if  $Y_i = \mu(X_i) + \varepsilon_i$  for  $\varepsilon_i \stackrel{iid}{\sim} N(0, \sigma^2)$  for  $\mu \in \mathcal{M}$

# When Does This Bound Apply?

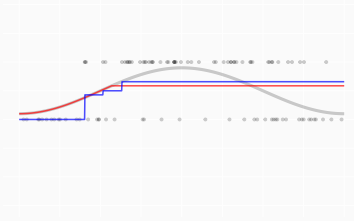
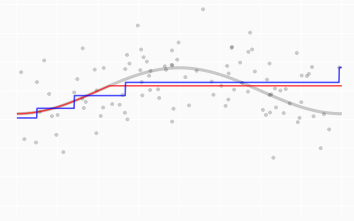
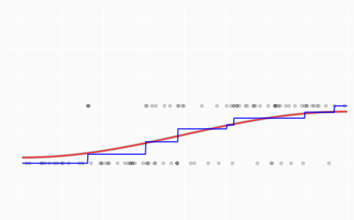
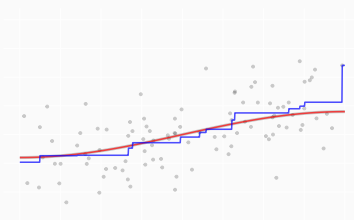


# When Does This Bound Apply?

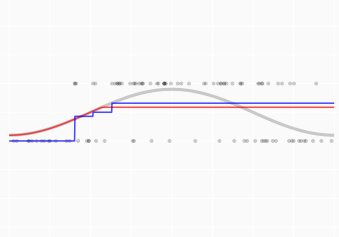
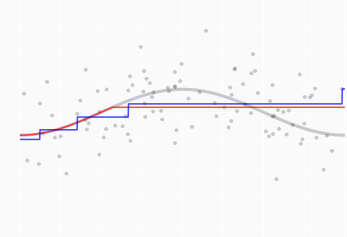
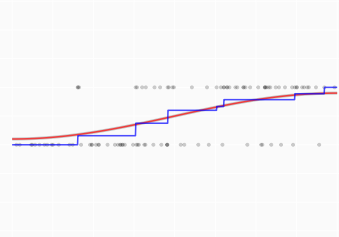
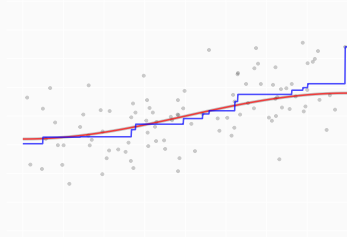


- The second column is out. We've assumed correct specification.
- The second row is out. We've assumed normality.

# Today, We Fix That



# Today, We Fix That



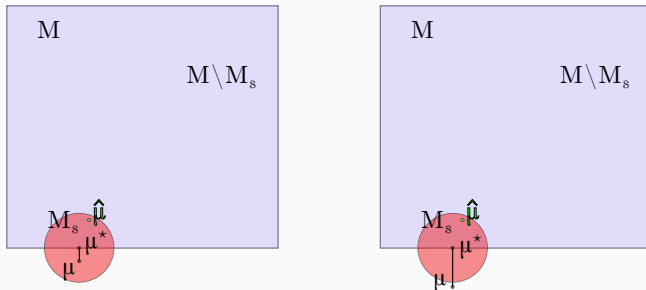
- With misspecification, we estimate the model's **best approximation** to  $\mu$ .
- Non-normality doesn't really matter much. We'll look at how it affects our bound.



## Misspecification

---

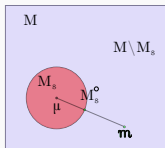
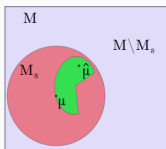
## What happens when $\mu$ isn't in the model?



- Our error in estimating  $\mu$  is bounded by a sum of two terms.
  - The critical radius  $s$ , i.e., the one satisfying  $s^2/2\sigma \geq w(\mathcal{M}_s^\circ) + s\sqrt{\frac{2\Sigma n}{\delta n}}$ .
  - 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  $\mu$ .



$\hat{\mu}$  minimizes  $\ell(m) = \frac{1}{n} \sum_{i=1}^n \{Y_i - \mu(X_i)\}^2$  squared error loss  
among curves  $m$  in a convex set  $\mathcal{M}$ .

- If  $\mu$  is in the model, that tells us it's **one of the curves** with loss as small as  $\mu$ 's.

i.e.  $m = \hat{\mu}$  satisfies  $\ell(m) \leq \ell(\mu)$  if  $\mu \in \mathcal{M}$ .

- To prove  $\hat{\mu}$  is in the neighborhood  $\mathcal{M}_s$ , we show that ...

- ...none of **these curves** is in **the neighborhood's complement**  $M \setminus \mathcal{M}_s$ .

$\hat{\mu} \in \mathcal{M}_s$  if  $\ell(m) > \ell(\mu)$  for all  $m \in \mathcal{M} \setminus \mathcal{M}_s$ .

- i.e. we show the *loss difference* is strictly positive for curves in **the complement**.

- That's true if it's positive for curves on **the neighborhood's boundary**  $\mathcal{M}_s^\circ$ .

$\ell(m) - \ell(\mu) > 0$  for all  $m \in \mathcal{M} \setminus \mathcal{M}_s$  if  $\ell(m) > \ell(\mu)$  for all  $m \in \mathcal{M}_s^\circ$ .

- And that boils down to the neighborhood's *squared radius* exceeding ...

- ...twice its boundary's *maximal inner product* with noise  $\varepsilon = Y - m$ .

$$\ell(m) - \ell(\mu) = s^2 - \langle Y - \mu, m - \mu \rangle \geq s^2 - 2 \max_{m \in \mathcal{M}_s^\circ} \langle Y - \mu, m - \mu \rangle \quad \text{for all } m \in \mathcal{M}_s^\circ$$

- Then we do a little probability and get our error bound.

# The Argument with no if

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:

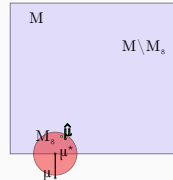
$$\begin{aligned} \ell(m) - \ell(\mu^*) &= \|m - \mu^*\|_{L_2(\mathbf{P}_n)}^2 && \text{squared distance, like before;} \\ &- \frac{2}{n} \sum_{i=1}^n \varepsilon_i \{m(X_i) - \mu^*(X_i)\} && \text{a mean zero term, like before;} \\ &+ \frac{2}{n} \sum_{i=1}^n \{\mu^*(X_i) - \mu(X_i)\} \{m(X_i) - \mu^*(X_i)\}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

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.

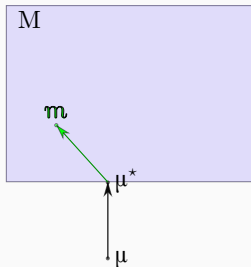
- If the loss difference is positive when we ignore a non-negative term ...
- ...then it's still positive when we don't.

$$\ell(m) - \ell(\mu^*) > 0 \quad \text{if} \quad \|m - \mu^*\|_{L_2(P_n)}^2 - \frac{2}{n} \sum_{i=1}^n \varepsilon_i \{m(X_i) - \mu^*(X_i)\} > 0 \quad \text{what we're used to}$$

$$\text{and} \quad \frac{2}{n} \sum_{i=1}^n \{\mu^*(X_i) - \mu(X_i)\} \{m(X_i) - \mu^*(X_i)\} \geq 0 \quad \text{new term}$$

This only works if the new term is non-negative. Can we choose  $\mu^* \in \mathcal{M}$  so it is?

## We can



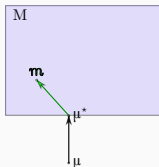
The new term is always non-negative when we compare to the *best approximation* to  $\mu$  in the model,

$$\mu^* = \operatorname{argmin}_{m \in \mathcal{M}} \|m - \mu\|_{L_2(P_n)}^2 \quad \text{satisfies} \quad \frac{2}{n} \sum_{i=1}^n \{\mu^*(X_i) - \mu(X_i)\} \{m(X_i) - \mu^*(X_i)\}$$

or in vector notation  $\frac{2}{n} \langle \mu^* - \mu, m - \mu^* \rangle_2 \geq 0 \quad \text{for all } m \in \mathcal{M}.$

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.
- They both point 'in' to the model. That means the dot product is non-negative.



**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

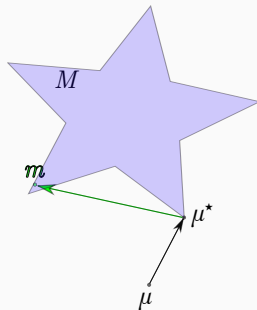
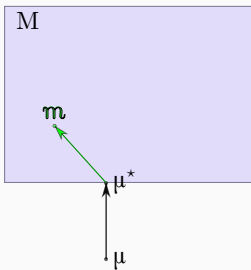
$$0 \leq \lambda \|m - \mu^*\|^2 + 2 \langle m - \mu^*, \mu^* - \mu \rangle.$$

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  $m \in \mathcal{M}$ .



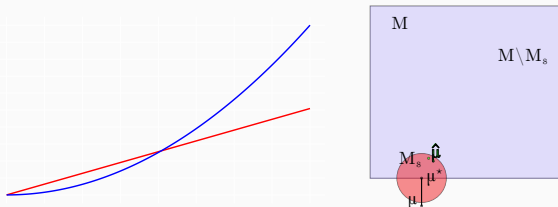
The same thing can happen *for the closest point* in a non-convex model.



# Summary

When we use a convex model, the least squares estimator  $\hat{\mu}$  converges to the model's closest point to  $\mu$ . This generalizes our result without misspecification.

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



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

$$\|\hat{\mu} - \mu^*\|_{L_2(P_n)} < s + 2\sqrt{\frac{2\Sigma_n}{\delta n}} \text{ w.p. } 1 - \delta \text{ if } s^2/2\sigma \geq w(\mathcal{M}_s)$$

for  $\mathcal{M}_s = \{m \in \mathcal{M} : \|m - \mu^*\|_{L_2(P_n)} \leq s\}$  and  $\Sigma_n = \sigma\{1 + 2\log(2n)\}$

if  $Y_i = \mu(X_i) + \varepsilon_i$  for  $\varepsilon_i \stackrel{iid}{\sim} N(0, \sigma^2)$  for some function  $\mu$ .

## Misspecification

---

### Examples

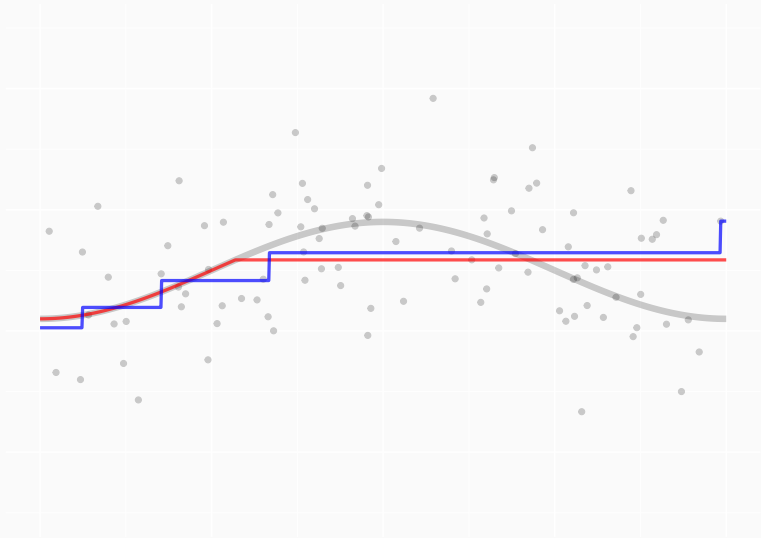


Figure 1: Increasing Curves ( $n = 100$ .)

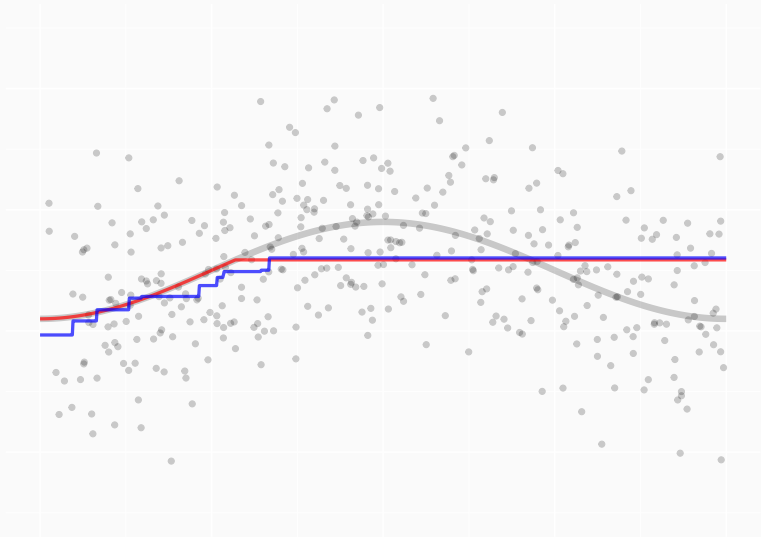


Figure 2: Increasing Curves ( $n = 400$ .)

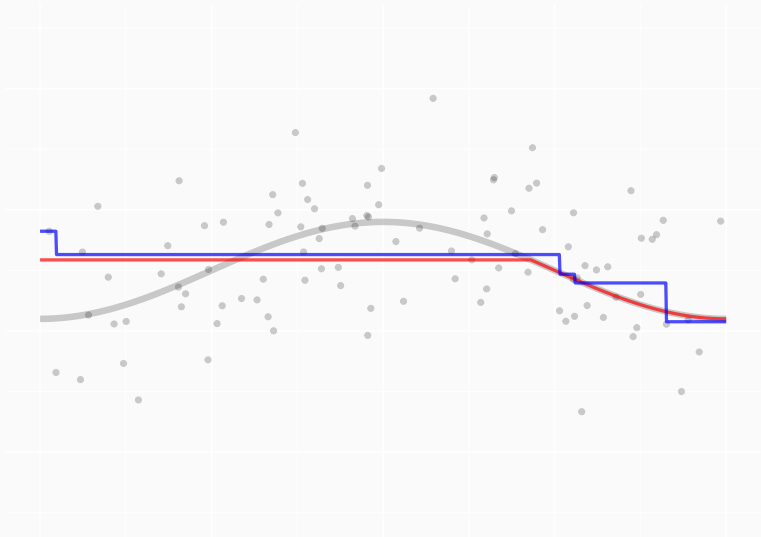


Figure 3: Decreasing Curves ( $n = 100$ .)

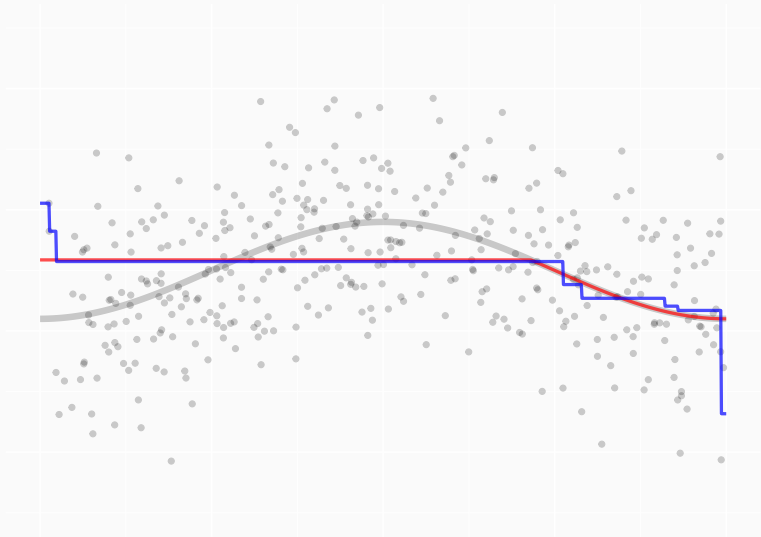


Figure 4: Decreasing Curves ( $n = 400$ .)

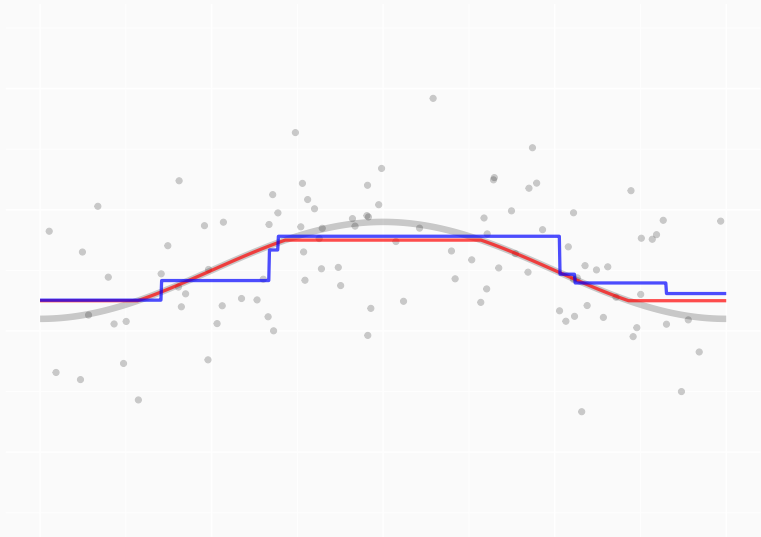


Figure 5: Bounded Variation Curves:  $\rho_{TV} \leq 1$  ( $n = 100$ .)

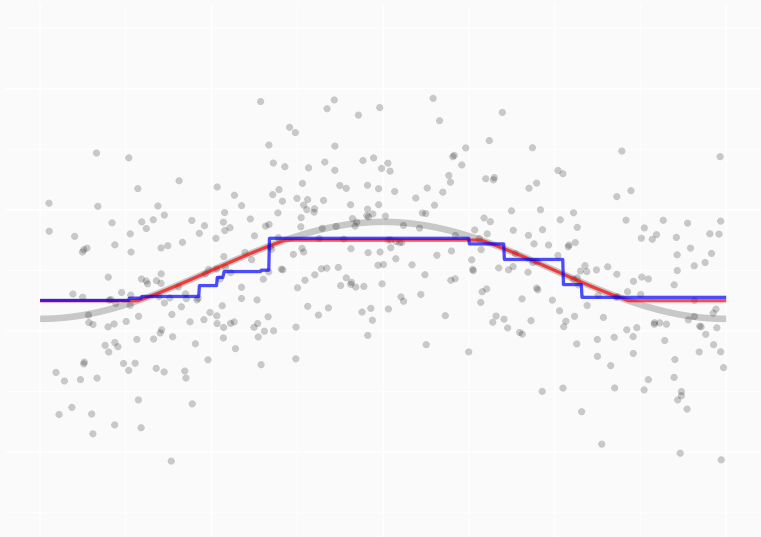


Figure 6: Bounded Variation Curves:  $\rho_{TV} \leq 1$ . ( $n = 400$ .)



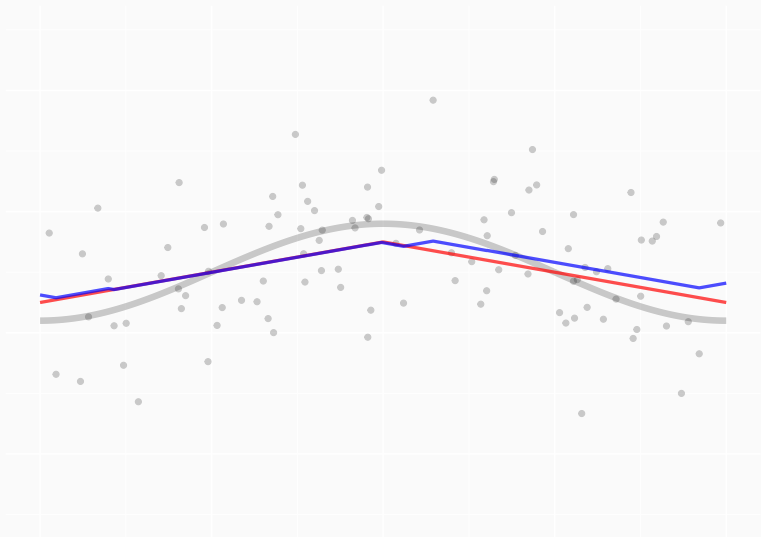


Figure 7: Lipschitz Curves:  $\rho_{\text{Lip}} \leq 1$  ( $n = 100$ .)

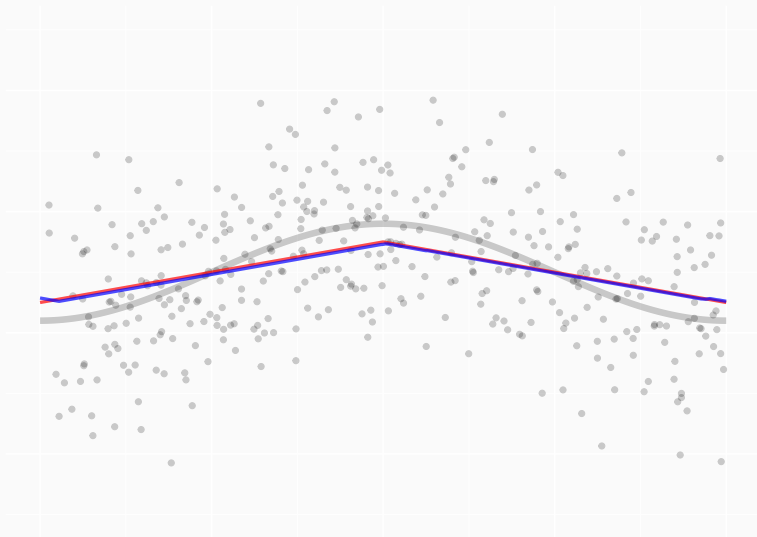


Figure 8: Lipschitz Curves:  $\rho_{\text{Lip}} \leq 1$  ( $n = 400$ .)

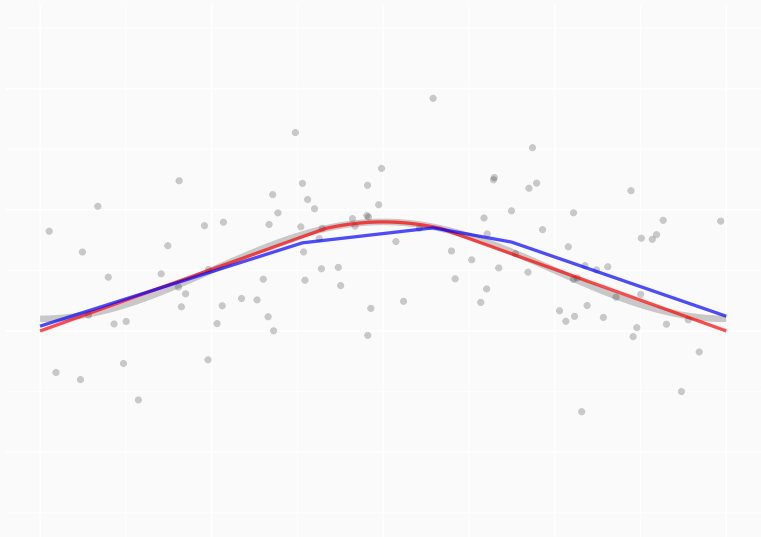


Figure 9: Concave Curves ( $n = 100$ .)

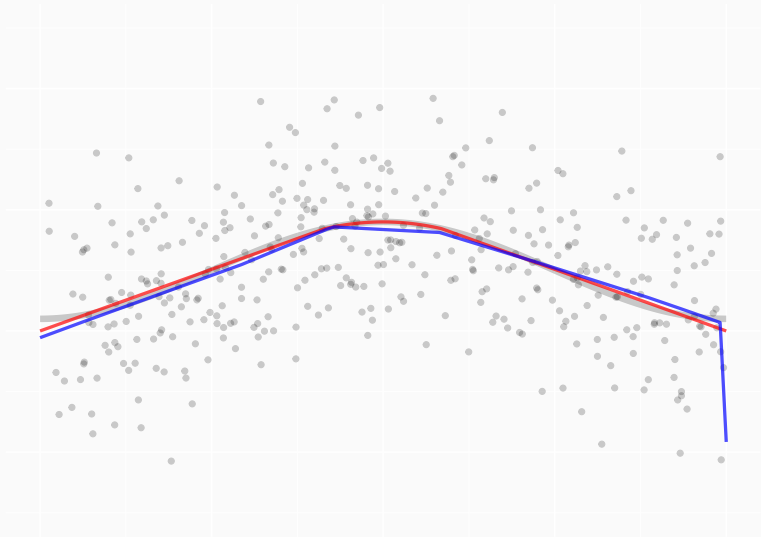


Figure 10: Concave Curves ( $n = 400$ .)

## Non-Gaussian Noise

---

# Starting Point

$$\ell(m) - \ell(\mu^*) = \|m - \mu^*\|_{L_2(P_n)}^2$$

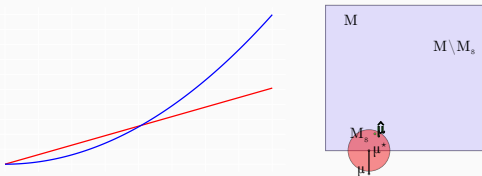
squared distance

$$- \frac{2}{n} \sum_{i=1}^n \varepsilon_i \{m(X_i) - \mu^*(X_i)\}$$

a mean zero term

$$+ \frac{2}{n} \sum_{i=1}^n \{\mu^*(X_i) - \mu(X_i)\} \{m(X_i) - \mu^*(X_i)\}$$

a non-negative term.

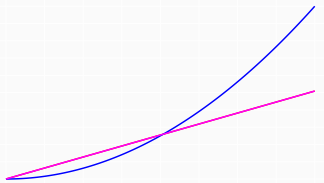


We can bound error using a corresponding *width*, no matter how noise is distributed.

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

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

This bound depends on the model  $\mathcal{M}$  and the distribution of the noise  $\varepsilon$  in a complex, entangled way: through the width  $w_\varepsilon(\mathcal{M}_s)$ .



To disentangle the impact of the model and noise distribution, we'll bound this width in terms of gaussian width.

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

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

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

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

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

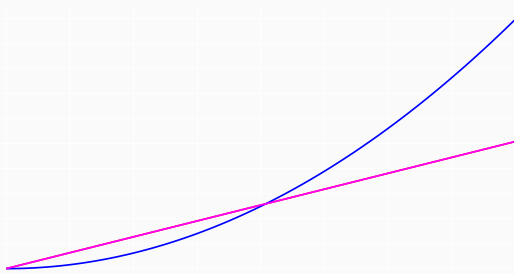
This substitution *increases* width:  $w_{\epsilon}(\cdot) \leq w_{\epsilon - \epsilon'}(\cdot)$ .

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.

$$w_{\epsilon - \epsilon'}(\cdot) \leq \sqrt{2\pi} M_n \times w(\cdot) \quad \text{for} \quad M_n = \mathbb{E} \max_{i \in 1 \dots n} |\epsilon_i|$$

This lets us re-use our gaussian width calculations to analyze regression with any noise distribution.

# A Simple Consequence: Width Comparison implies Radius Comparison



- If you have a width comparison  $w_\epsilon \leq \alpha w_\eta$  for some  $\alpha \geq 1$ .
- This implies a radius comparison  $s_\epsilon \leq \alpha s_\eta$  for all convex models  $\mathcal{M}$ .

$$s_\epsilon = \alpha s_\eta \quad \text{satisfies} \quad \frac{s_\epsilon^2}{2} \geq w_\epsilon(\mathcal{M}_{s_\epsilon}) \quad \text{if} \quad \frac{s_\eta^2}{2} \geq w_\eta(\mathcal{M}_{s_\eta}) \quad \text{for convex } \mathcal{M}$$

and  $w_\epsilon \leq \alpha w_\eta$  for  $\alpha \geq 1$ .

- *Interpretation.*  
The noise  $\epsilon$  makes regression at most ' $\alpha$  times harder' than the noise  $\eta$ .
- *This is simplistic and 'lossy'.*  
For most models, our width comparison implies a better radius comparison.



## Proof: Width Comparisons imply Radius Comparisons

**Claim.** If  $w_\varepsilon \leq \alpha w_\eta$  for  $\alpha \geq 1$ , then for any convex model  $\mathcal{M}$ , 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

## Proof: Width Comparisons imply Radius Comparisons

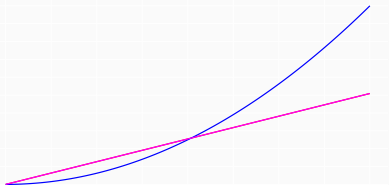
**Claim.** If  $w_\varepsilon \leq \alpha w_\eta$  for  $\alpha \geq 1$ , then for any convex model  $\mathcal{M}$ , the critical radius using noise  $\varepsilon$  is at most  $\alpha$  times the critical radius using noise  $\eta$ , i.e.

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

**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.



**Where we are.** We have a bound that depends on the model  $\mathcal{M}$  and the distribution of the noise  $\epsilon$  in a complex and entangled way.

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

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

**Where we're going.** We'll derive a bound that depends on the model  $\mathcal{M}$  and the distribution of the noise  $\epsilon$  in simpler and disentangled way.

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

$$\text{where } w(\mathcal{V}) = \mathbb{E} \max_{v \in \mathcal{V}} \langle g, v \rangle_{L_2(\mathcal{P}_n)} \quad \text{and} \quad M_n = \mathbb{E} \max_{i \in 1 \dots n} |\epsilon_i|.$$

**Better yet.** We can simplify it into a bound that depends on only one measure of noise complexity.

$$\|\hat{\mu} - \mu^*\|_{L_2(\mathcal{P}_n)} \leq \sqrt{2\pi\Sigma_n} \left( s + \sqrt{\frac{2}{\delta n}} \right) \quad \text{because} \quad M_n \leq \sqrt{\Sigma_n} \quad \text{and} \quad 2 \leq \sqrt{2\pi}$$

## Non-Gaussian Noise

---

Example: Probabilistic Classification

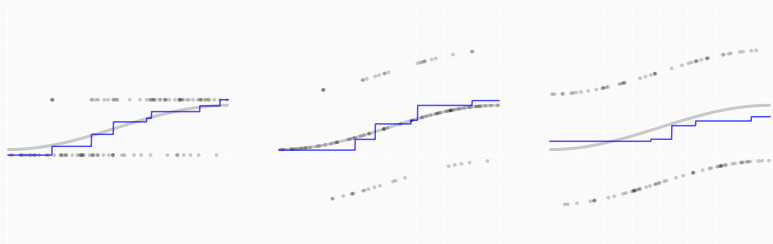


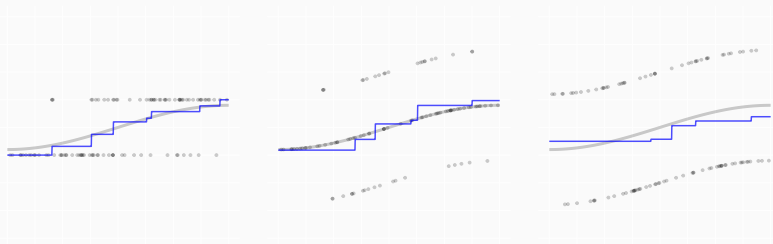
Figure 11: classification noise  $\rightarrow$  symmetrized classification noise  $\rightarrow$  random-sign noise

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.

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



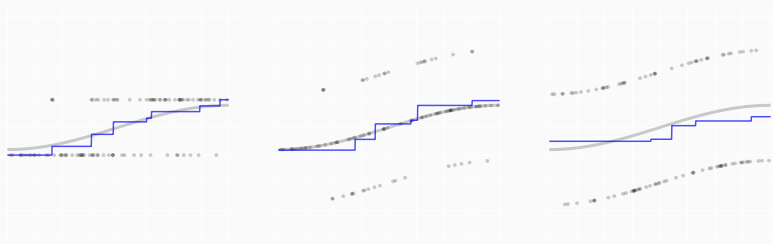
**Figure 11:** classification noise  $\rightarrow$  symmetrized classification noise  $\rightarrow$  random-sign noise

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*.

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



**Figure 11:** 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.

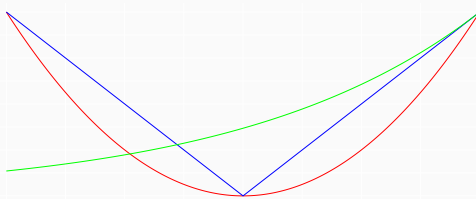
## Step 3

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.



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

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



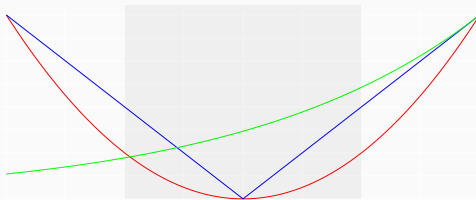
## Step 3

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.

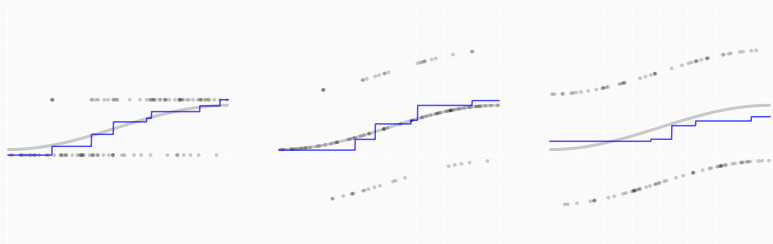
$$\begin{aligned}\max_{u \in \{-1, 1\}^n} f(u) &= \max_{u \in \{-1, 1\}^n} \mathbb{E}_s \max_{v \in \mathcal{V}} \sum_{i=1}^n s_i u_i v_i \\ &= \mathbb{E}_s \max_{v \in \mathcal{V}} \sum_{i=1}^n s_i v_i.\end{aligned}$$

Why?

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

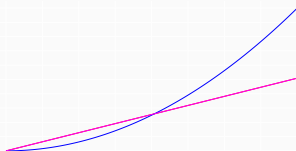
- For  $u_i \in \{-1, 1\}$ , the distributions of  $u_i$  and  $s_i u_i$  are the same.
- So the distribution of the sum, and its maximum, are the same at every corner  $u$ .
- Including the vector of all ones  $u = (1, 1, \dots, 1)$ .

# 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.

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



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

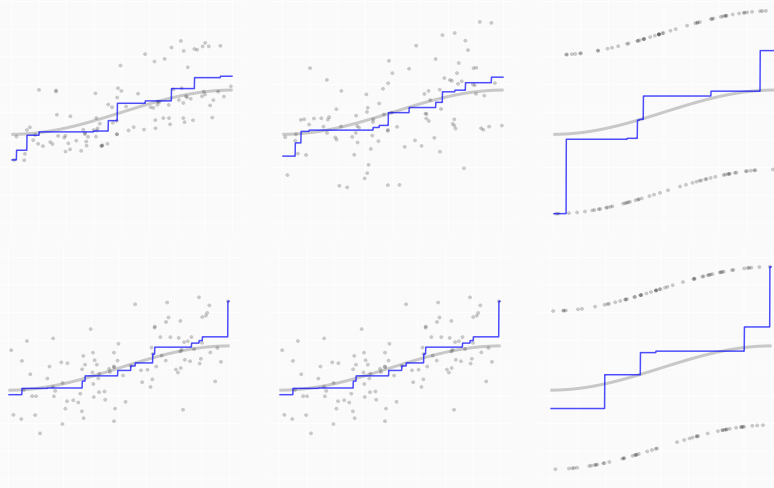


Figure 12: real noise  $\rightarrow$  symmetrized noise  $\rightarrow$  scaled sign noise



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

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

(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.

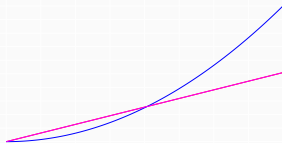
(b) Replacing  $\varepsilon_i$  with  $s_i \varepsilon_i$  increases width by at most  $2 \times$ .

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

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

- 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



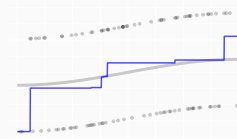
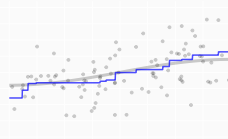
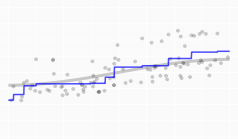
$$w_{\varepsilon}(\mathcal{V}) \leq \mathbb{E} \|\varepsilon_i - \varepsilon'_i\|_{\infty} w_s(\mathcal{V}) \leq 2 \mathbb{E} \|\varepsilon_i\|_{\infty} w_s(\mathcal{V})$$

Regression with arbitrary independent noise, i.e.

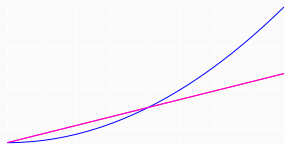
$$Y_i = \mu(X_i) + \varepsilon_i \quad \text{where} \quad \varepsilon_1 \dots \varepsilon_n \text{ are independent,}$$

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

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



# The Symmetric Case



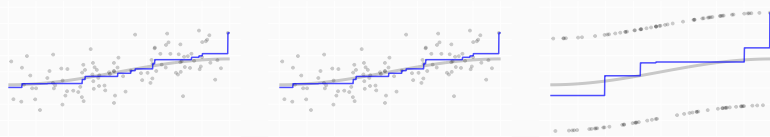
$$w_\varepsilon(\mathcal{V}) \leq \mathbb{E}\|\varepsilon_i\|_\infty w_s(\mathcal{V})$$

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.

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



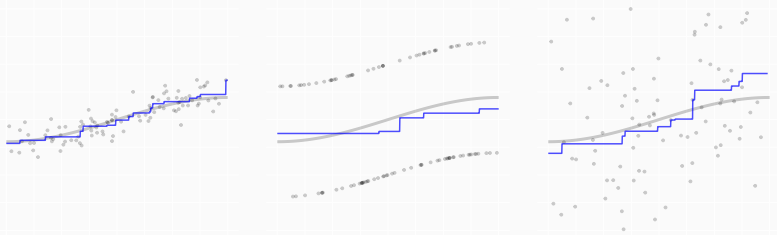
**Figure 13:** real noise  $\rightarrow$  symmetrized noise  $\rightarrow$  scaled sign noise

<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

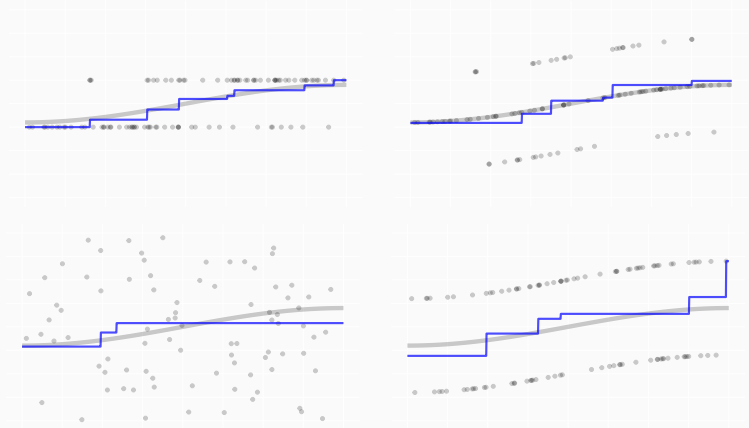


Figure 14: real noise  $\rightarrow$  symmetrized noise  $\downarrow$  scaled sign noise  $\leftarrow$  scaled gaussian noise

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

$$\leq \sqrt{2\pi} \approx 2.5 \times \mathbb{E} \|\varepsilon\|_{\infty}$$

# Implications for Regression

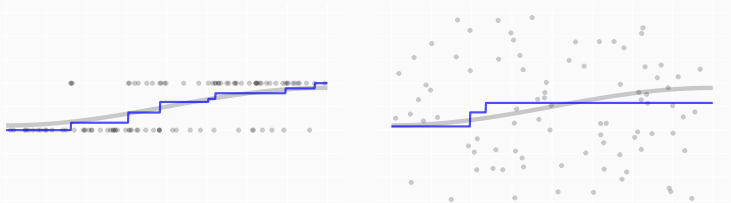


Figure 15: real noise  $\rightarrow$  scaled gaussian noise

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.

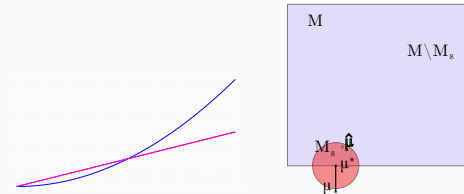
This means we don't have to bound a million different kinds of widths for each model.  
We can bound random-sign width or gaussian width. Whichever is easier.



## Sampling

---

# What We've Done



We have a bound that's valid for any signal  $\mu$  and any vector of independent noise  $\varepsilon$ .

$$\|\hat{\mu} - \mu^*\|_{L_2(P_n)} < 2\sqrt{\Sigma_n} \left( s + \sqrt{\frac{2}{\delta n}} \right) \quad \text{w.p. } 1 - \delta \quad \text{for } \frac{s^2}{2} \geq w_s(\mathcal{M}_s)$$

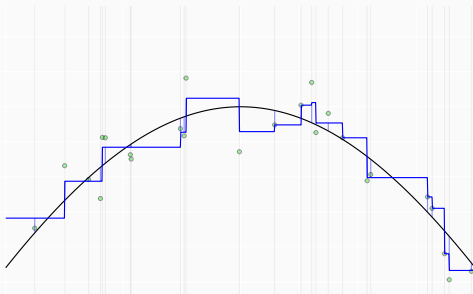
- It depends on the model's size through the *critical radius* of random-sign width.  
 $s$  satisfying  $s^2/2 \geq w_s(\mathcal{M}_s)$  for  $\mathcal{M}_s = \{m \in \mathcal{M} : \|m - \mu^*\|_{L_2(P_n)} \leq s\}$ 
  - This is a one-number summary of the random-sign width of neighborhoods ...
  - ...of the model's best approximation to the signal. It's the summary that matters.
- It depends on the noise's size through the expected maximum square.

$$\Sigma_n = \mathbb{E} \max_{i \in 1 \dots n} |\varepsilon_i|^2$$

# What does this tell us?

Bounds like this say how close  $\hat{\mu}$  and  $\mu^*$  are, on average, on our sample  $X_1 \dots X_n$ .

$$\frac{1}{n} \sum_{i=1}^n \{\hat{\mu}(X_i) - \mu^*(X_i)\} < \dots$$



It doesn't tell us how close they are in the gaps between those points.

- Let's think about what happens when  $X_1 \dots X_n$  are drawn independently from some distribution  $\mathbf{P}$ . Think sampling with replacement from a population.
- We'll bound the *population root mean squared error*  $\|\hat{\mu} - \mu^*\|_{L_2(\mathbf{P})}$ .

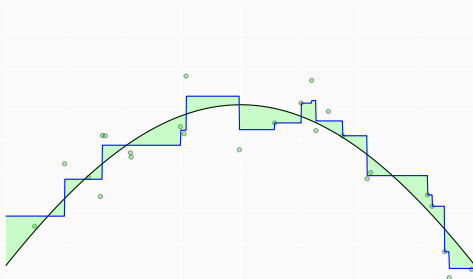
# What Population Mean Squared Error Is

It's the mean squared error we make at random point  $X'$  distributed like  $X_1 \dots X_n$ .

$$\|\hat{\mu} - \mu^*\|_{L_2(\mathbf{P})}^2 = \mathbb{E}_{X'} [\{\hat{\mu}(X') - \mu^*(X')\}^2]$$

That's the integral of the squared distance between the two curves,  
multiplied by the density of  $X_i$ .

$$\|\hat{\mu} - \mu^*\|_{L_2(\mathbf{P})}^2 = \int \{\hat{\mu}(x) - \mu^*(x)\}^2 p(x) dx \quad \text{if } X_i \text{ has the density } p(x).$$

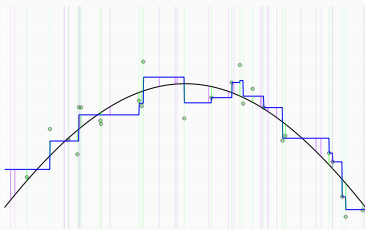


# Why we care about Population Mean Squared Error: Generalization

If we're interested in average accuracy for a bunch of new points  $X'_1 \dots X'_{n'}$ , distributed like  $X_1 \dots X_n$ , that's more or less exactly what it is.

$$\|\hat{\mu} - \mu\|_{L_2(P)}^2 = E_{X'} [\{\hat{\mu}(X') - \mu(X')\}^2] \stackrel{LLN}{\approx} \frac{1}{n'} \sum_{i=1}^{n'} \{\hat{\mu}(X'_i) - \mu(X'_i)\}^2.$$

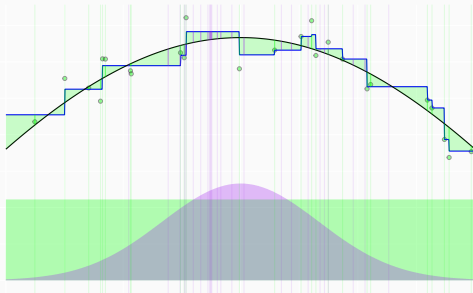
This can be a bit different from accuracy on our original sample  $X_1 \dots X_n$ .



- BV regression spends its 'variation budget' jumping to fit on the original sample.
- Between those points, it doesn't know whether it should jump or not.
  - So we can get larger error at our new points.
  - It's usually not much larger, but sometimes it is. We'll see why.

## Why we care about Population Mean Squared Error: Generalization

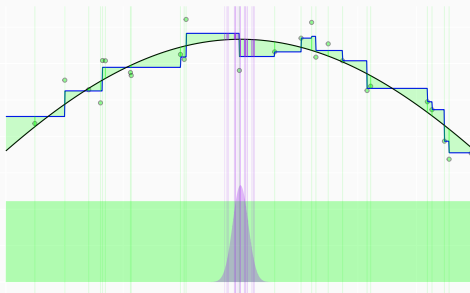
If we're interested in average accuracy for new points from a different distribution  $Q$ , we can bound this by comparing this distribution's density to that of our observations.



$$\begin{aligned}\frac{1}{n'} \sum_{i=1}^{n'} \{\hat{\mu}(X'_i) - \mu(X'_i)\}^2 &\approx \|\hat{\mu} - \mu\|_{L_2(Q)}^2 = \int \{\hat{\mu}(x) - \mu(x)\}^2 \frac{q(x)}{p(x)} p(x) dx \\ &\leq \max_x \frac{q(x)}{p(x)} \|\hat{\mu} - \mu\|_{L_2(P)}^2.\end{aligned}$$

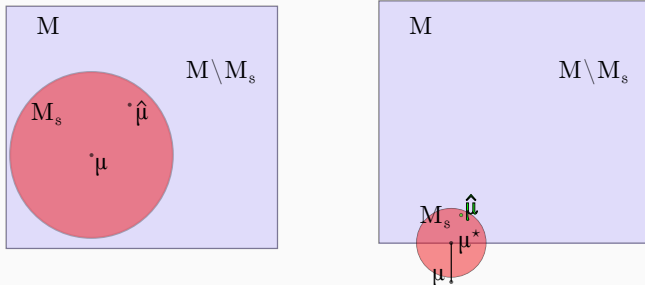
# Why we care about Population Mean Squared Error: Generalization

If we're interested in accuracy at a specific point  $x'$ , we can think of this new distribution  $Q$  as a little bump around  $x'$ .



$$\{\hat{\mu}(x') - \mu(x')\}^2 \approx \|\hat{\mu} - \mu\|_{L_2(Q_\epsilon)}^2 \quad \text{for} \quad Q = N(x', \epsilon^2).$$

## Same Argument, Different Neighborhood



- We want to show that  $\hat{\mu}$  is in a *population-distance neighborhood* of  $\mu$ .
- Or, if we've chosen the model wrong, at least its best population-distance approximation.  

$$\hat{\mu} \in \mathcal{M}_s \text{ for } \mathcal{M}_s = \{m \in \mathcal{M} : \|m - \mu^*\|_{L_2(\mathbb{P})} \leq s\} \text{ for } \mu^* = \operatorname{argmin}_{m \in \mathcal{M}} \|m - \mu\|_{L_2(\mathbb{P})}$$
- We'll do this using essentially the same argument we used to bound sample MSE.
  1. We know that the  $\hat{\mu}$ 's squared error loss is at least as good as  $\mu^*$ 's.
  2. We find a radius  $s$  for which every curve with this property is in the neighborhood  $\mathcal{M}_s$ .
- It amounts to showing the loss difference  $\ell(m) - \ell(\mu^*)$  is positive outside this neighborhood.

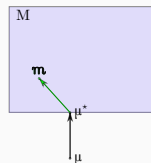
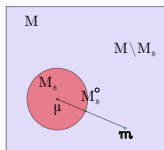
$$m \in \mathcal{M}_s \text{ if } m \in \mathcal{M}_s \text{ and } \ell(m) - \ell(\mu^*) > 0 \text{ for all } m \in \mathcal{M} \setminus \mathcal{M}_s$$



# Reduction to a Maximal Inequality

$$\begin{aligned}\ell(m) - \ell(\mu^*) &= \frac{1}{n} \sum_{i=1}^n Z_i(m) := \{m(X_i) - \mu^*(X_i)\}^2 - 2 \{Y_i - \mu^*(X_i)\} \{m(X_i) - \mu^*(X_i)\} \\ &= E Z_i(m) + \frac{1}{n} \sum_{i=1}^n Z_i(m) - E Z_i(m).\end{aligned}$$

Convexity Helps  
as Usual.



1. The loss difference is positive outside the neighborhood if it's positive on its boundary.

$$m \in \mathcal{M}_s \text{ if } m \in \mathcal{M}_s \text{ and } \ell(m) - \ell(\mu^*) > 0 \text{ for all } m \in \mathcal{M} \setminus \mathcal{M}_s^\circ$$

2. The projection theorem tells us an unwanted term in  $E Z_i(m)$  is non-negative.

$$\begin{aligned}-E \left[ \{Y_i - \mu^*(X_i)\} \{m(X_i) - \mu^*(X_i)\} \right] &= -E \left[ \left\{ E[Y_i | X_i] - \mu^*(X_i) \right\} \{m(X_i) - \mu^*(X_i)\} \right] \\ &= \langle \mu^* - \mu, m - \mu^* \rangle_{L_2(P)} \geq 0 \quad \text{for all } m \in \mathcal{M}\end{aligned}$$

$$\text{It follows that } m \in \mathcal{M}_s \text{ if } m \in \mathcal{M}_s \text{ and } s^2 > \max_{m \in \mathcal{M}_s^\circ} \frac{1}{n} \sum_{i=1}^n Z_i(m) - E Z_i(m)$$

# Bounding the New Maximum

We show this maximum is **approximately constant**, i.e. close to its expectation.

$$\bar{Z} := \max_{m \in \mathcal{M}_s^\circ} \frac{1}{n} \sum_{i=1}^n Z_i(m) - \mathbb{E} Z_i(m) \quad \text{satisfies} \quad \bar{Z} \leq \mathbb{E} \bar{Z} + \sqrt{\frac{\text{Var}(\bar{Z})}{\delta n}} \quad \text{w.p. } 1 - \delta$$

We use **symmetrization** to bound its expectation in terms of random-sign width.

- (a) Write the centers  $\mathbb{E} Z_i(v)$  in terms of an independent copy of our sample.
- (b) Compare the result to a maximum of an average of symmetric random variables.
- (c) Introduce random signs and compare to two copies of a simpler maximum.

$$\begin{aligned} n \times \mathbb{E} \bar{Z} &\stackrel{(a)}{=} \mathbb{E}_Z \max_{m \in \mathcal{M}_s^\circ} \mathbb{E}_{Z'} \sum_{i=1}^n \{Z_i(m) - Z'_i(m)\} \\ &\stackrel{(b)}{\leq} \mathbb{E}_Z \mathbb{E}_{Z'} \mathbb{E}_s \max_{m \in \mathcal{M}_s^\circ} \sum_{i=1}^n s_i \{Z_i(m) - Z'_i(m)\} \\ &\stackrel{(c)}{\leq} \mathbb{E}_Z \mathbb{E}_{Z'} \mathbb{E}_s \max_{m, m' \in \mathcal{M}_s^\circ} \sum_{i=1}^n s_i Z_i(m) + (-s_i) Z_i(m') \\ &= 2 \mathbb{E}_Z \mathbb{E}_s \max_{m \in \mathcal{M}_s^\circ} \sum_{i=1}^n s_i Z_i(m) \end{aligned}$$

We can use the **Efron-Stein inequality** to bound the variance. Come back and try it later!

$$\begin{aligned} \text{Var}(\bar{Z}) &\stackrel{\text{why?}}{\leq} \frac{1}{n^2} \sum_{i=1}^n \mathbb{E} \{Z_i(\hat{m}) - Z'_i(\hat{m})\}_+^2 \quad \text{for} \quad \hat{m} = \operatorname{argmax}_{m \in \mathcal{M}_s^\circ} \sum_{i=1}^n Z_i(m) - \mathbb{E} Z_i(m) \\ &\leq \dots \end{aligned}$$

# Contracting Out Lipschitz Functions

What we get is  $2 \times$  the expected random-sign width of some set of vectors, but it's not just the set of the vectors in our neighborhood  $\mathcal{M}_s - \mu^*$ .

$$\begin{aligned} n \times \mathbb{E} Z &\leq 2 \mathbb{E} \mathbb{E}_s \max_{m \in \mathcal{M}_s^\circ} \sum_{i=1}^n s_i Z_i(m) = \{m(X_i) - \mu^*(X_i)\}^2 - 2 \{Y_i - \mu^*(X_i)\} \{m(X_i) - \mu^*(X_i)\} \\ &\leq 4 \mathbb{E} \left\{ \max_{m \in \mathcal{M}_s^\circ} \|m - \mu\|_{L_\infty(\mathcal{P}_n)} + \|\varepsilon\|_{L_\infty(\mathcal{P}_n)} \right\} \mathbb{E}_s \max_{m \in \mathcal{M}_s^\circ} \sum_{i=1}^n s_i \{m(X_i) - \mu^*(X_i)\} \end{aligned}$$

We've compared that to the width of the neighborhood itself using ...

**Lemma (Lipschitz Comparison)**

$$\mathbb{E}_s \max_{v \in \mathcal{V}} \sum_{i=1}^n s_i \psi_i(v_i) \leq L \mathbb{E}_s \max_{v \in \mathcal{V}} \sum_{i=1}^n s_i v_i \text{ if } |\psi_i(u_i) - \psi_i(v_i)| \leq L |u_i - v_i| \text{ for all } u, v \in \mathcal{V}.$$

For  $\psi_i(v) = v_i^2 - 2\{Y_i - \mu^*(X_i)\}v_i$  and  $V = \{m(X_1) - \mu^*(X_1) \dots m(X_n) - \mu^*(X_n) : m \in \mathcal{M}_s^\circ\}$ ,

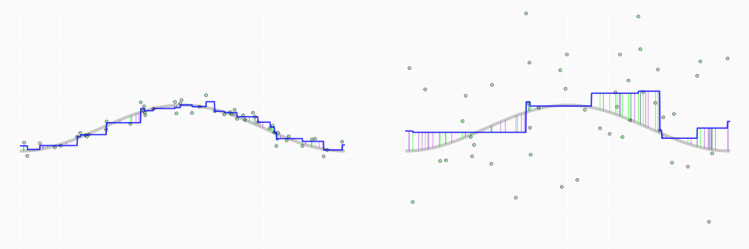
$$\text{that's } \mathbb{E}_s \max_{m \in \mathcal{M}_s^\circ} \sum_{i=1}^n s_i \psi_i\{m(X_i) - \mu^*(X_i)\} \leq L \max_{m \in \mathcal{M}_s^\circ} \sum_{i=1}^n s_i \{m(X_i) - \mu^*(X_i)\}$$

$$\text{where } L = \max_i \max_{m \in \mathcal{M}_s^\circ} |\psi'_i\{m(X_i) - \mu^*(X_i)\}|$$

$$= \max_i \max_{m \in \mathcal{M}_s^\circ} |2\{m(X_i) - \mu^*(X_i)\} - 2\{Y_i - \mu^*(X_i)\}|$$

$$\leq 2 \max_{m \in \mathcal{M}_s^\circ} \|m - \mu\|_{L_\infty(\mathcal{P}_n)} + 2\|\varepsilon\|_{L_\infty(\mathcal{P}_n)}.$$

# Interpretation



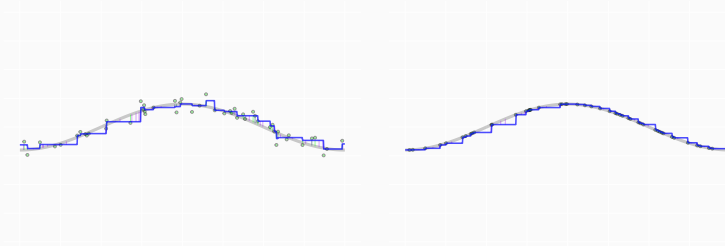
$$\|\hat{\mu} - \mu^*\|_{L_2(\mathcal{P})} \leq s \times 2\left\{\sqrt{\Sigma_n} + B\right\} + \sqrt{\frac{\text{Var}(\bar{Z})}{\delta}} \quad \text{w.p. } 1 - \delta$$

if  $\frac{s^2}{2} \geq \mathbb{E} \text{w}_s(\mathcal{M}_s)$  and  $\|m - \mu\|_\infty \leq B$

This is the bound we'd get on sample MSE with additional scaled random-sign noise,

$$\text{i.e. if we'd observed } Y_i = \mu(X_i) + \varepsilon_i + B s_i$$

Left: With little noise, our estimator  $\hat{\mu}$  fits substantially better at the sample points  $X_i$ .  
Right: With more, it doesn't. The observations are far enough from  $\mu$  that we can't estimate it all that precisely even where we have some data.



**Signal Recovery** is regression without any noise at all. In that case ( $\Sigma_n = 0$ ),

$$\|\hat{\mu} - \mu\|_{L_2(\mathcal{P})} \leq s \times 2 \left\{ \sqrt{\Sigma_n} + B \right\} + \sqrt{\frac{\text{Var}(\bar{Z})}{\delta}} \quad \text{w.p. } 1 - \delta$$

$$\text{if } \frac{s^2}{2} \geq \mathbb{E} \text{w}_s(\mathcal{M}_s) \quad \text{and} \quad \|m - \mu\|_{\infty} \leq B$$

This is the bound we'd get on sample MSE with *only* scaled random-sign noise.

$$\text{i.e. if we'd observed } Y_i = \mu(X_i) + \varepsilon_i + Bs_i$$

- This is an extreme case of the low-noise regime. And it's still hard.
- When you want to estimate  $\mu$  between the sample points  $X_1 \dots X_n, \dots$
- ...what you want to see obscured by bounded 'sampling noise'  $\in [-B, B]$ .

Chapter 6 of Talagrand's Upper and Lower Bounds for Stochastic Processes.

- Random Signs vs. Gaussians: Proposition 6.22
- Contraction: Lemma 6.4.5
- Lipschitz Contraction: Theorem 6.5.1

## Appendices

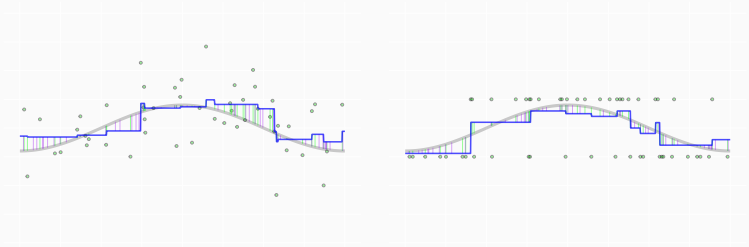
---

## Appendices

---

### Boundedness



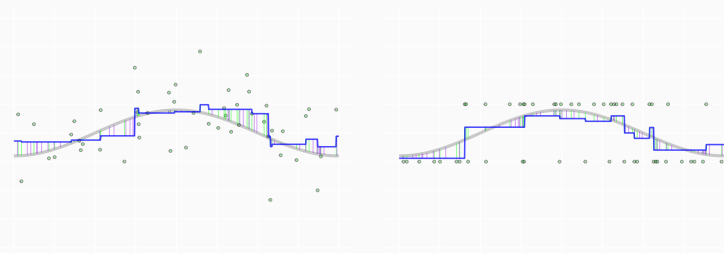


Our Population MSE bound introduces a new consideration: boundedness of  $\|m - \mu\|_\infty$  in neighborhoods of  $\mu^*$ .

$$\|\hat{\mu} - \mu\|_{L_2(\mathbb{P})} \leq s \times 2 \left\{ \sqrt{\Sigma_n} + B \right\} + \sqrt{\frac{\text{Var}(\bar{Z})}{\delta}} \quad \text{w.p. } 1 - \delta$$

$$\text{if } \frac{s^2}{2} \geq \mathbb{E} \mathbf{w}_s(\mathcal{M}_s) \quad \text{and} \quad \|m - \mu\|_\infty \leq B$$

Getting a bound  $B$  can take a bit of work. There are options.



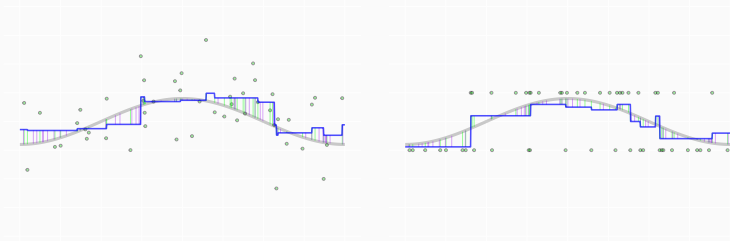
**Option 1.** Baking it into the Model.

$$\mathcal{M} = \{m : \|m\|_{\infty} \leq B \text{ and } \rho_{TV}(m) \leq B\}$$

$$\implies \|m - \mu\|_{\infty} \leq \|m\|_{\infty} + \|\mu\|_{\infty} \leq B + \|\mu\|_{\infty}$$

$$\mathcal{M} = \{m : m(0) = 0 \text{ and } \rho_{TV}(m) \leq B\}$$

$$\implies$$



## Option 2. Arguing Based on Bounded Data.

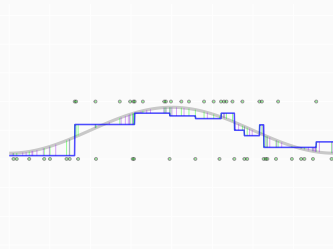
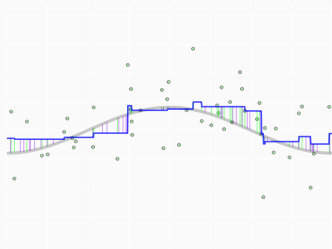
In many models, you can show that  $\hat{\mu}$  will be within the range of the data.

$$\text{i.e. } \min_{i \leq n} Y_i \leq \hat{\mu}(x) \leq \max_{i \leq n} Y_i$$

This is true, in particular, for Monotone and Bounded Variation Regression.

We can add this constraint to our model when doing our analysis.

$$\begin{aligned} & \|\hat{\mu} - \mu^*\|_{L_2(\mathbf{P})} < s \quad \text{if} \quad \ell(m) - \ell(\mu^*) > 0 \quad \text{for all} \quad m \in \mathcal{M} \dots \\ & \dots \text{ with } \|m\|_{\infty} \leq B \quad \text{and} \quad \|m - \mu^*\|_{L_2(\mathbf{P})} \geq s \end{aligned}$$



The are other options.

## Appendices

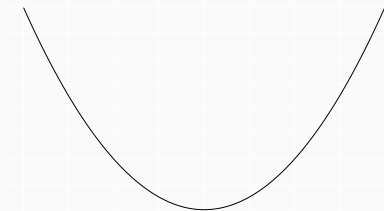
---

Convex Functions Are Maximized At  
Extreme Points

## Definition

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

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



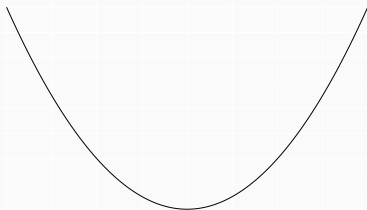
We can give this a *probabilistic interpretation* for a random variable  $Z_\lambda$ .

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

## 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  $Z_\lambda$ .

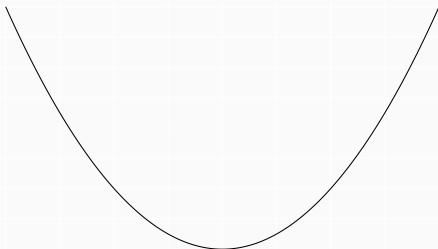
$$f(\mathbb{E} Z_\lambda) \leq \mathbb{E} f(Z_\lambda) \quad \text{where} \quad Z_\lambda = \begin{cases} a & \text{w.p. } 1 - \lambda \\ b & \text{w.p. } \lambda \end{cases}$$

# 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.

$$f(E Z) \leq E f(Z)$$

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.

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

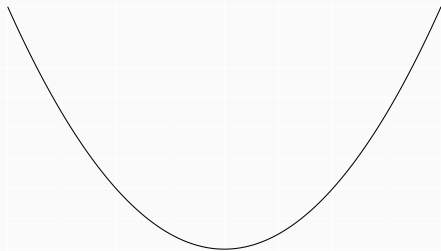
## Inductive Step.

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 apply Jensen's inequality to a *random extreme point*  $Z_\lambda$ .

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

$f(\mathbb{E} Z_\lambda) \qquad \mathbb{E} f(Z_\lambda)$

where

$$Z_\lambda = \begin{cases} u_1 & \text{w.p. } \lambda_1 \\ \vdots & \vdots \\ u_K & \text{w.p. } \lambda_K \end{cases}$$