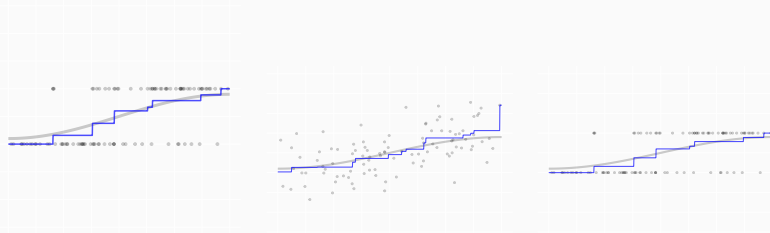


## Non-Gaussian Noise

---

# Review: Probabilistic Classification



Last time, we talked about *probabilistic classification*, i.e. regression with *classification noise*.

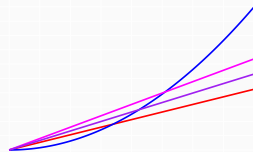
$$\hat{\mu} = \operatorname{argmin}_{m \in \mathcal{M}} \frac{1}{n} \sum_{i=1}^n \{Y_i - m(X_i)\}^2 \quad \text{where} \quad Y_i = \mu(X_i) + \varepsilon_i \quad \text{for} \quad \varepsilon_i = \begin{cases} 1 - \mu(X_i) & \text{w.p. } \mu(X_i) \\ -\mu(X_i) & \text{w.p. } 1 - \mu(X_i) \end{cases}$$

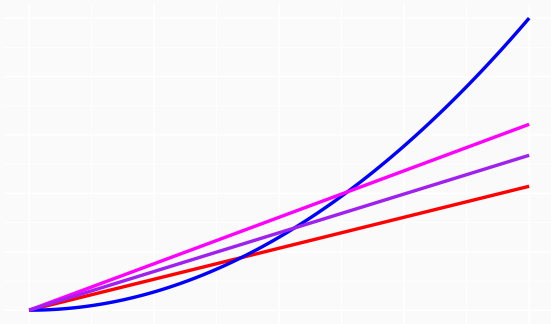
By comparing widths, we showed that this is *easier* than regression with ...

1. random sign noise,  $s_i = \pm 1$  each w.p.  $1/2$ .
2. gaussian noise  $\sigma g_i$  of standard deviation  $\sigma = 1.25$ .

Easier in the sense that our crossing-point argument gives us a better error bound.

$$\frac{s^2}{2} \geq 1.25 w(\mathcal{M}_s) \geq w_s(\mathcal{M}_s) \geq w_\varepsilon(\mathcal{M}_s)$$





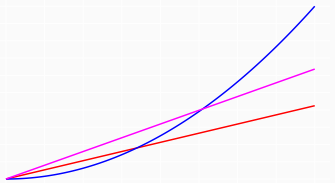
Today, we're going to generalize that result to regression with *any kind of noise*. We'll start with the same abstract bound. It applies no matter how noise is distributed.

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

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

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

# Our Approach



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<sup>1</sup> for our symmetrized noise  $\epsilon - \epsilon'$ . We can bound the impact of this substitution in a model-invariant way.

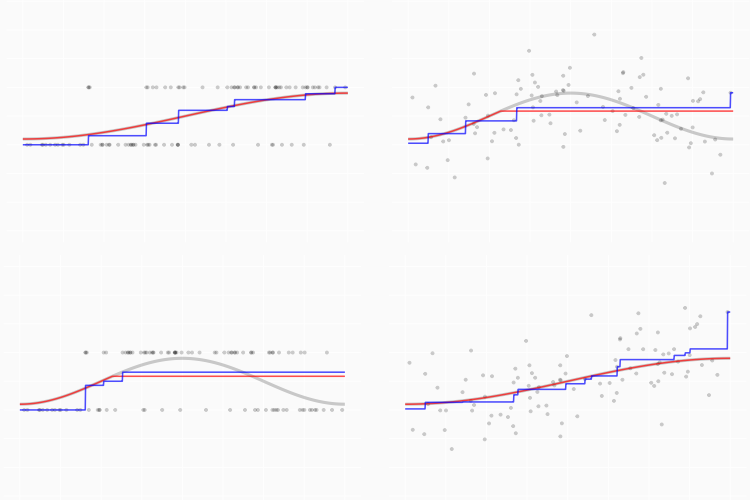
$$w_{\epsilon - \epsilon'}(\cdot) \leq 2M_n w_s(\cdot) \leq \sqrt{2\pi} M_n \times w(\cdot) \quad \text{for} \quad M_n = \mathbb{E} \max_{i \in 1 \dots n} |\epsilon_i|$$

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

---

<sup>1</sup>or a random-sign vector

# An Example

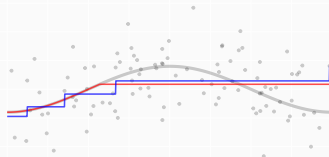


**Figure 1:** real noise  $\rightarrow$  symmetrized noise  $\downarrow$  scaled sign noise  $\leftarrow$  scaled gaussian noise

Here's the same signal with 4 types of noise.

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

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



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

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

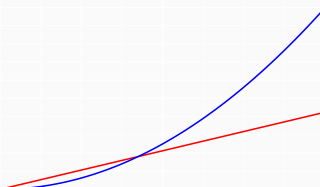
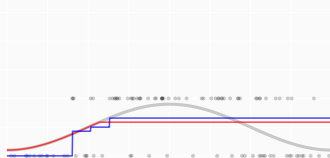
(a) Replacing  $\varepsilon_i$  with  $s_i(\varepsilon_i - \varepsilon'_i)$  is 'free'.

- We stopped here in our classification 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  $2 \times$ .

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

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





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

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

- 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 classification 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.

# Symmetrization, Contraction, and Gaussian Noise

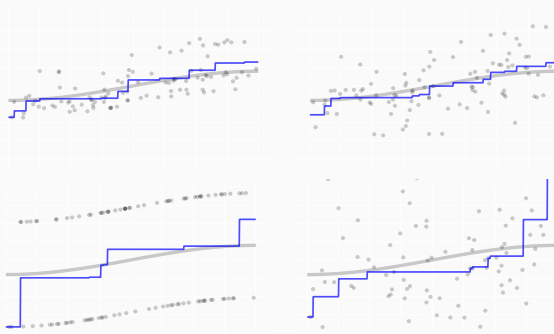


Figure 2: real noise  $\rightarrow$  symmetrized noise  $\downarrow$  scaled sign noise  $\leftarrow$  scaled gaussian noise

After symmetrizing and introducing random signs, i.e. making the substitution

$$\varepsilon_i \rightarrow s_i(\varepsilon_i - \varepsilon'_i),$$

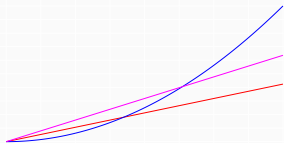
we 'contract out' the symmetrized noise  $\varepsilon - \varepsilon'$  to get a bound in terms of random-sign width.

$$w_\varepsilon(\mathcal{V}) \leq w_{s(\varepsilon - \varepsilon')}(\mathcal{V}) \leq \|\varepsilon - \varepsilon'\|_\infty w_s(\mathcal{V}) \leq \|\varepsilon - \varepsilon'\|_\infty \mathbf{1.25} w(\mathcal{V})$$

We can substitute **1.25 times gaussian width** because that's at least as large as random sign width.

$$\mathbb{E} \max_{v \in \mathcal{V}} \sum_{i=1}^n g_i v_i = \mathbb{E}_s \mathbb{E}_g \max_{v \in \mathcal{V}} \sum_{i=1}^n |g_i| s_i v_i \geq \mathbb{E}_s \max_{v \in \mathcal{V}} \sum_{i=1}^n \mathbb{E}_g |g_i| s_i v_i \underset{\approx 1/1.25}{\geq} \mathbb{E}_s \max_{v \in \mathcal{V}} \sum_{i=1}^n \mathbb{E}_g |g_i| s_i v_i.$$

# Implications for Regression



$$w_\varepsilon(\mathcal{V}) \leq M w_s(\mathcal{V}) \leq 1.25M w(\mathcal{V})$$

$$\text{for } M = E\|\varepsilon - \varepsilon'\|_\infty \leq 2 E\|\varepsilon\|_\infty.$$

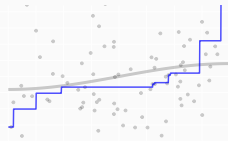
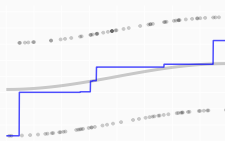
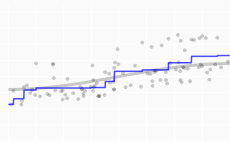
In terms of our crossing-point bounds, regression with arbitrary independent noise,

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

is no harder than with scaled random sign noise or with gaussian noise

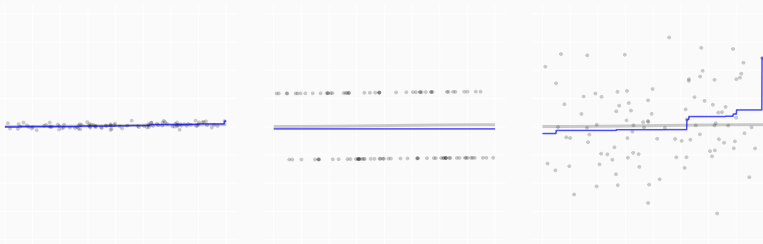
$$\text{i.e. } Y_i = \mu(X_i) + Ms_i \quad \text{for } s_i = \begin{cases} +1 & \text{w.p. } 1/2 \\ -1 & \text{w.p. } 1/2 \end{cases}$$

$$\text{or } Y_i = \mu(X_i) + 1.25Mg_i \quad \text{for } g_i \sim N(0, 1)$$



The **scale factor** is  $2 \times$  the expected magnitude of our noise vector's largest element.

# Lost Precision



**Figure 3:** standard gaussian noise  $\rightarrow$  scaled random sign noise  $\rightarrow$  scaled gaussian noise

- This isn't the absolute best bound we can get.
- For example, if we start with standard gaussian noise, we lose ...
- ...a factor of roughly  $7\sqrt{\log(2n)}$  going to random sign width and back.

$$w_\varepsilon(\mathcal{V}) \leq 2 w_{s\varepsilon}(\mathcal{V}) \leq 2 \times 2 \sqrt{2 \log(2n)} w_s(\mathcal{V}) \leq 4 \sqrt{2 \log(n)} \times \sqrt{\frac{\pi}{2}} w_\varepsilon(\mathcal{V}) \approx 7 \sqrt{\log(2n)} w_\varepsilon(\mathcal{V}).$$

$\geq \mathbb{E} \|\varepsilon\|_\infty$   $= \mathbb{E} |g_1|$

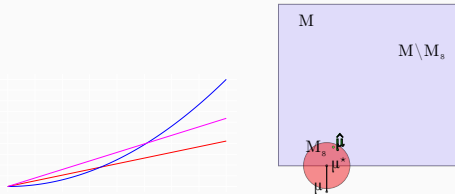
- 'Symmetrization' cost us a factor of 2.
- Contraction costs us a factor of  $\mathbb{E} \max_{i \leq n} |\varepsilon_i| \leq 2\sqrt{2 \log(2n)}$ . (See HW Appendix B)
- Converting random signs back to gaussians costs us a factor of  $\sqrt{\frac{\pi}{2}} \approx 1.25$ .

We're in the right ballpark. For sample sizes  $n$  between 50 and 50 million, that factor is between 15 and 30. But if we want a more precise error bound, we need to be a little more careful.

## Sampling

---

# What We've Done



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

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

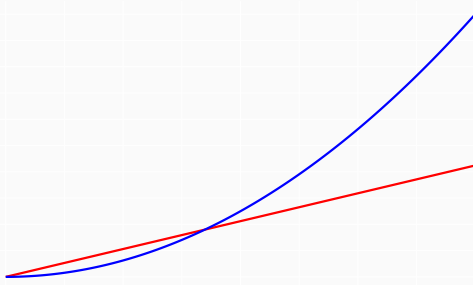
- It depends on the model's size through the *critical radius* of random-sign width.  
 $s$  satisfying  $s^2/2 \geq w_s(\mathcal{M}_s)$  for  $\mathcal{M}_s = \{m \in \mathcal{M} : \|m - \mu^*\|_{L_2(\mathbf{P}_n)} \leq s\}$ 
  - This is a one-number summary of the random-sign width of neighborhoods ...
  - ...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  $X_1 \dots X_n$  is are drawn independently from some distribution  $\mathbf{P}$ . Think sampling with replacement from a population.
- We'll bound the *population root mean squared error*  $\|\hat{\mu} - \mu^*\|_{L_2(\mathbf{P})}$ .

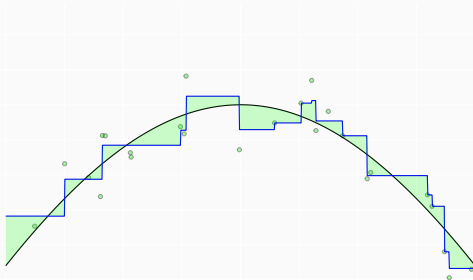
# What Population Mean Squared Error Is

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  $X_i$ .

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



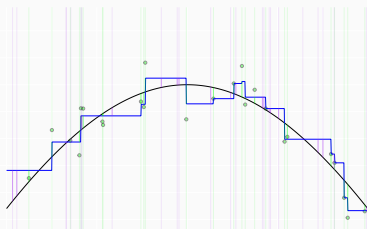


# Why we care about Population Mean Squared Error: Generalization

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.

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

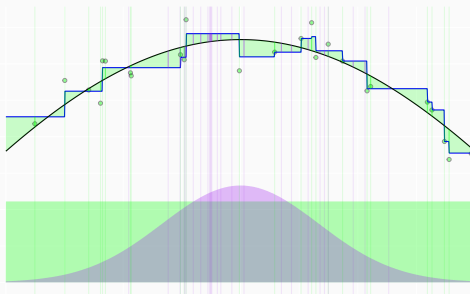
This can be a bit different from accuracy on our original sample  $X_1 \dots X_n$ .



- 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

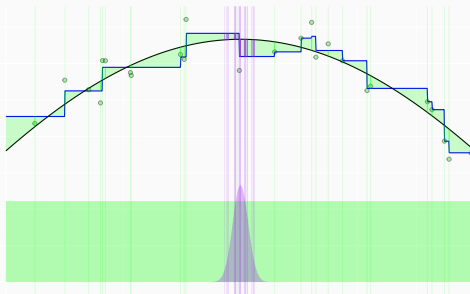
If we're interested in average accuracy for new points from a different distribution  $Q$ , we can bound this by comparing this distribution's density to that of our observations.



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

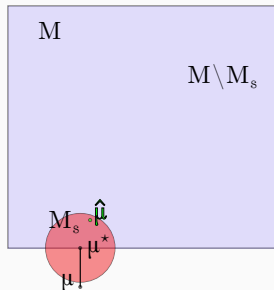
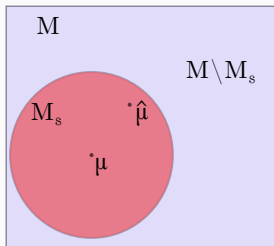
# Why we care about Population Mean Squared Error: Generalization

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.  

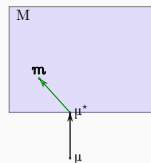
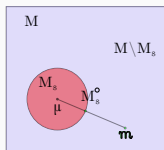
$$\hat{\mu} \in \mathcal{M}_s \text{ for } \mathcal{M}_s = \{m \in \mathcal{M} : \|m - \mu^*\|_{L_2(\mathbb{P})} \leq s\} \text{ for } \mu^* = \operatorname{argmin}_{m \in \mathcal{M}} \|m - \mu\|_{L_2(\mathbb{P})}$$
- We'll do this using essentially the same argument we used to bound sample MSE.
  1. 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  $V = \{m(X_1) - \mu^*(X_1) \dots m(X_n) - \mu^*(X_n) : m \in \mathcal{M}_s^\circ\}$ ,

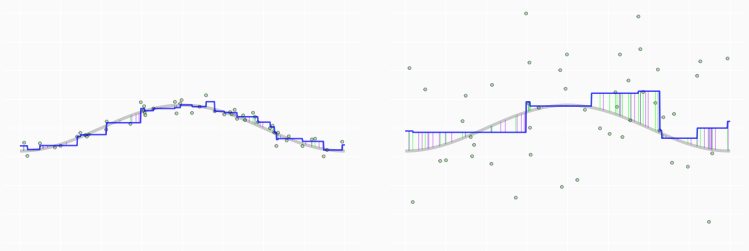
$$\text{that's } \mathbb{E}_s \max_{m \in \mathcal{M}_s^\circ} \sum_{i=1}^n s_i \psi_i\{m(X_i) - \mu^*(X_i)\} \leq L \max_{m \in \mathcal{M}_s^\circ} \sum_{i=1}^n s_i \{m(X_i) - \mu^*(X_i)\}$$

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

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

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

# Interpretation



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

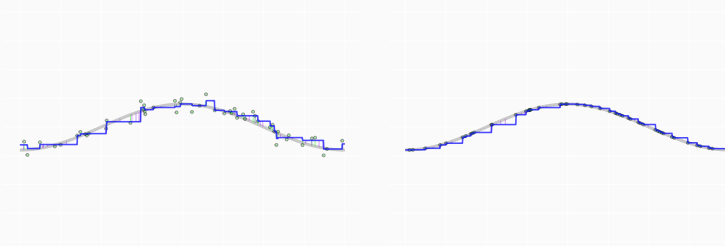
if  $\frac{s^2}{2} \geq \mathbb{E} \text{w}_s(\mathcal{M}_s)$  and  $\|m - \mu\|_\infty \leq B$

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

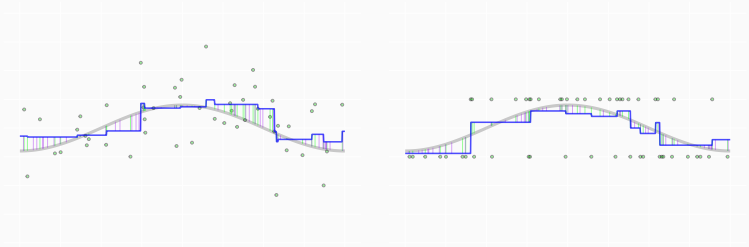
## Appendices

---

## Appendices

---

### Boundedness

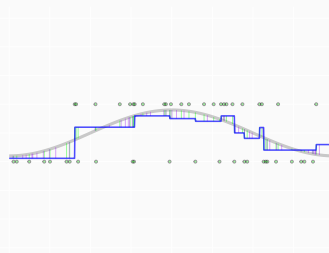
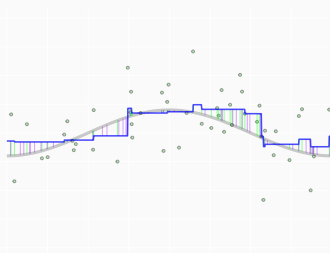


Our Population MSE bound introduces a new consideration: boundedness of  $\|m - \mu\|_\infty$  in neighborhoods of  $\mu^*$ .

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

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

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



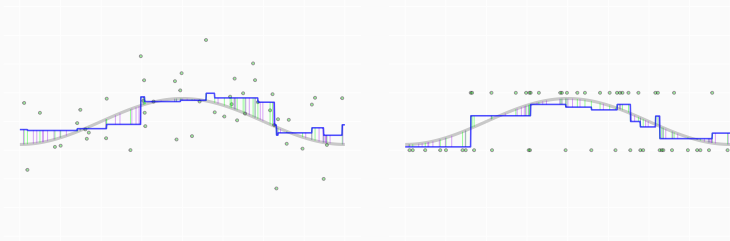
**Option 1.** Baking it into the Model.

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

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

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

$$\implies$$



## Option 2. Arguing Based on Bounded Data.

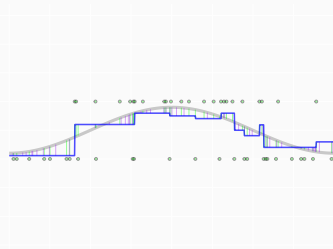
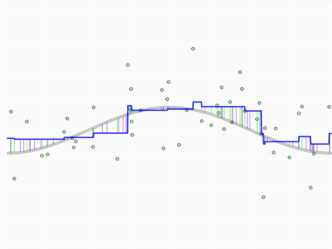
In many models, you can show that  $\hat{\mu}$  will be within the range of the data.

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

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.