Non-Gaussian Noise

Review: Probabilistic Classification



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

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

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 σg_i of standard deviation $\sigma = 1.25$.

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

$$\frac{s^2}{2} \ge 1.25 \operatorname{w}(\mathcal{M}_s) \stackrel{\Longrightarrow}{\ge} \operatorname{w}_s(\mathcal{M}_s) \stackrel{\Longrightarrow}{\ge} \operatorname{w}_{\varepsilon}(\mathcal{M}_s)$$



Starting Point



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.

$$\begin{split} \|\hat{\mu} - \mu^{\star}\|_{L_{2}(\mathbf{P}_{n})} &< s + 2\sqrt{\frac{2\Sigma_{n}}{\delta n}} \quad \text{w.p. } 1 - \delta \text{ for } \frac{s^{2}}{2} \geq \mathbf{w}_{\boldsymbol{\epsilon}}(\mathcal{M}_{s}) \\ \text{where } \mathbf{w}_{\boldsymbol{\epsilon}}(\mathcal{V}) &= \mathbf{E}\max_{\boldsymbol{v}\in\mathcal{V}}\langle\boldsymbol{\epsilon},\boldsymbol{v}\rangle_{L_{2}(\mathbf{P}_{n})} \text{ and } \Sigma_{n} = \mathbf{E}\max_{\boldsymbol{i}\in1\dots,n}\varepsilon_{\boldsymbol{i}}^{2}. \end{split}$$

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

Our Approach



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

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

for α depending on ε but not \mathcal{M} or s.

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

1. Symmetrization. We'll substitute for ϵ_i a variant that's symmetric around zero.

 $\epsilon_i
ightarrow \epsilon_i - \epsilon_i'$ where ϵ_i' is an independent copy of ϵ_i

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

2. Contraction. We'll substitute a gaussian vector¹ for our symmetrized noise $\epsilon - \epsilon'$. We can bound the impact of this substitution in a model-invariant way.

$$\mathbf{w}_{\epsilon-\epsilon'}(\cdot) \leq 2M_n \, \mathbf{w}_s(\cdot) \leq \sqrt{2\pi} M_n \times \mathbf{w}(\cdot) \quad \text{for} \quad M_n = \mathbf{E} \max_{i \in 1...n} |\varepsilon_i|$$

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

¹or a random-sign vector

An Example



4

Symmetrization

$$\begin{split} \mathbf{w}_{\varepsilon}(\mathcal{V}) &\leq \mathbf{w}_{s(\varepsilon-\varepsilon')}(\mathcal{V}) \leq 2 \, \mathbf{w}_{s\varepsilon}(\mathcal{V}) \\ \mathbf{E} \max_{v \in \mathcal{V}} \sum_{i=1}^{n} \varepsilon_{i} v_{i} &= \mathbf{E} \max_{v \in \mathcal{V}} \sum_{i=1}^{n} (\varepsilon_{i} - \mathbf{E} \, \varepsilon'_{i}) v_{i} \\ &\stackrel{(a)}{\leq} \mathbf{E} \, \mathbf{E}' \max_{v \in \mathcal{V}} \sum_{i=1}^{n} (\varepsilon_{i} - \varepsilon'_{i}) v_{i} \\ &= \mathbf{E}_{s} \, \mathbf{E} \, \mathbf{E}' \max_{v \in \mathcal{V}} \sum_{i=1}^{n} s_{i} (\varepsilon_{i} - \varepsilon'_{i}) v_{i} \\ &\stackrel{(b)}{\leq} \mathbf{E}_{s} \, \mathbf{E} \max_{v \in \mathcal{V}} \sum_{i=1}^{n} s_{i} \varepsilon_{i} + \mathbf{E}_{s} \, \mathbf{E}' \max_{v \in \mathcal{V}} \sum_{i=1}^{n} s_{i} \varepsilon'_{i} v_{i} = 2 \, \mathbf{E}_{s} \, \mathbf{E} \max_{v \in \mathcal{V}} \sum_{i=1}^{n} \varepsilon_{i} s_{i} v_{i}. \end{split}$$



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

$$\begin{split} \mathbf{E} \max_{v \in \mathcal{V}} \sum_{i=1}^{n} \varepsilon_{i} v_{i} &= \mathbf{E} \max_{v \in \mathcal{V}} \sum_{i=1}^{n} (\varepsilon_{i} - \mathbf{E} \varepsilon_{i}') v_{i} \\ &\stackrel{(a)}{\leq} \mathbf{E} \mathbf{E}' \max_{v \in \mathcal{V}} \sum_{i=1}^{n} (\varepsilon_{i} - \varepsilon_{i}') v_{i} \\ &= \mathbf{E}_{s} \mathbf{E} \mathbf{E}' \max_{v \in \mathcal{V}} \sum_{i=1}^{n} s_{i} (\varepsilon_{i} - \varepsilon_{i}') v_{i} \\ &\stackrel{(b)}{\leq} \mathbf{E}_{s} \mathbf{E} \max_{v \in \mathcal{V}} \sum_{i=1}^{n} s_{i} \varepsilon_{i} + \mathbf{E}_{s} \mathbf{E}' \max_{v \in \mathcal{V}} \sum_{i=1}^{n} s_{i} \varepsilon_{i}' v_{i} = 2 \mathbf{E}_{s} \mathbf{E} \max_{v \in \mathcal{V}} \sum_{i=1}^{n} \varepsilon_{i} s_{i} v_{i}. \end{split}$$

(a) Replacing ε_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 $arepsilon_i$ again.

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

Contraction

$$\begin{split} \mathbf{w}_{\eta}(\mathcal{V}) &= \mathbf{w}_{s\eta}(\mathcal{V}) \leq \mathbf{E} \|\eta\|_{\infty} \, \mathbf{w}_{\eta}(\mathcal{V}) \quad \text{if} \quad \eta \stackrel{dist}{=} -\eta. \\ \mathbf{E}_{s} \, \mathbf{E}_{\eta} \max_{v \in \mathcal{V}} \sum_{i=1}^{n} \eta_{i} s_{i} v_{i} \leq \mathbf{E}_{\eta} \max_{\substack{u \in \mathbb{R}^{n} \\ |u_{i}| \leq \|\eta\|_{\infty}}} \mathbf{E}_{s} \max_{v \in \mathcal{V}} \sum_{i=1}^{n} u_{i} s_{i} v_{i} \\ &= \mathbf{E}_{\eta} \|\eta\|_{\infty} \max_{u \in [-1,1]^{n}} \mathbf{E}_{s} \max_{v \in \mathcal{V}} \sum_{i=1}^{n} u_{i} s_{i} v_{i} \\ &= \mathbf{E}_{\eta} \|\eta\|_{\infty} \times \max_{u \in \{-1,1\}^{n}} \mathbf{E}_{s} \max_{v \in \mathcal{V}} \sum_{i=1}^{n} u_{i} s_{i} v_{i} \\ &= \mathbf{E}_{\eta} \|\eta\|_{\infty} \times \mathbf{E}_{s} \max_{v \in \mathcal{V}} \sum_{i=1}^{n} s_{i} v_{i} \end{split}$$



Contraction

$$\begin{split} \mathbf{w}_{\eta}(\mathcal{V}) &= \mathbf{w}_{s\eta}(\mathcal{V}) \leq \mathbf{E} \|\eta\|_{\infty} \, \mathbf{w}_{\eta}(\mathcal{V}) \quad \text{if} \quad \eta \stackrel{dist}{=} -\eta. \\ \mathbf{E}_{s} \, \mathbf{E}_{\eta} \, \max_{v \in \mathcal{V}} \sum_{i=1}^{n} \eta_{i} s_{i} v_{i} \leq \mathbf{E}_{\eta} \, \max_{\substack{u \in \mathbb{R}^{n} \\ |u_{i}| \leq ||\eta||_{\infty}}} \mathbf{E}_{s} \, \max_{v \in \mathcal{V}} \sum_{i=1}^{n} u_{i} s_{i} v_{i} \\ &= \mathbf{E}_{\eta} \|\eta\|_{\infty} \, \max_{u \in [-1,1]^{n}} \mathbf{E}_{s} \, \max_{v \in \mathcal{V}} \sum_{i=1}^{n} u_{i} s_{i} v_{i} \\ &= \mathbf{E}_{\eta} \|\eta\|_{\infty} \times \, \max_{u \in \{-1,1\}^{n}} \mathbf{E}_{s} \, \max_{v \in \mathcal{V}} \sum_{i=1}^{n} u_{i} s_{i} v_{i} \\ &= \mathbf{E}_{\eta} \|\eta\|_{\infty} \times \mathbf{E}_{s} \, \max_{v \in \mathcal{V}} \sum_{i=1}^{n} s_{i} v_{i} \end{split}$$

- We can 'contract out' any symmetrically distributed noise vector η by ...
 - 1. multiplying in independent random signs s_i . Symmetry $\implies s_i \eta_i \stackrel{dist}{=} \eta_i$.
 - 2. maximizing over a cube containing η .
- 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

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

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

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

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

$$\operatorname{E}\max_{v\in\mathcal{V}}\sum_{i=1}^n g_iv_i = \operatorname{E}_s\operatorname{E}_g \max_{v\in\mathcal{V}}\sum_{i=1}^n |g_i| \; s_iv_i \geq \operatorname{E}_s \max_{v\in\mathcal{V}}\sum_{i=1}^n \operatorname{E}_g |g_i| \; s_iv_i.$$

7

Implications for Regression



$$\begin{split} \mathbf{v}_{\varepsilon}(\mathcal{V}) &\leq M \, \mathbf{w}_{s}(\mathcal{V}) \leq 1.25 M \, \mathbf{w}(\mathcal{V}) \\ \text{for} \quad \mathbf{M} = \mathbf{E} \| \varepsilon - \varepsilon' \|_{\infty} \leq 2 \, \mathbf{E} \| \varepsilon \|_{\infty} \end{split}$$

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

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

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

i.e.
$$Y_i = \mu(X_i) + Ms_i$$
 for $s_i = \begin{cases} +1 & \text{w.p. } 1/2 \\ -1 & \text{w.p. } 1/2 \end{cases}$
or $= \mu(X_i) + 1.25Ma_i$ for $a_i \sim N(0, 1)$



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



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 ...
- \cdot ...a factor of roughly $7\sqrt{\log(2n)}$ going to random sign width and back.

$$\|\mathbf{w}_{arepsilon}(\mathcal{V}) \leq 2 \, \mathbf{w}_{sarepsilon}(\mathcal{V}) \leq 2 imes 2 \sqrt{2 \log(2n)} \, \mathbf{w}_{s}(\mathcal{V}) \leq 4 \sqrt{2 \log(n)} imes \sqrt{rac{\pi}{2}} \, \mathbf{w}_{arepsilon}(\mathcal{V}) pprox 7 \sqrt{\log(2n)} \, \mathbf{w}_{arepsilon}(\mathcal{V}).$$

- (a) 'Symmetrization' cost us a factor of 2.
- (b) Contraction costs us a factor of $\operatorname{Emax}_{i \leq n} |\varepsilon_i| \leq 2\sqrt{2\log(2n)}$. (See HW Appendix B)
- (c) 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



We have a bound that's valid for any signal μ and any vector of independent noise ε .

$$\|\hat{\mu} - \mu^{\star}\|_{L_{2}(\mathbf{P}_{n})} < 2\sqrt{\Sigma_{n}} \left(s + \sqrt{\frac{2}{\delta n}}\right) \quad \text{w.p. } 1 - \delta \text{ for } \quad \frac{s^{2}}{2} \ge \mathbf{w}_{s}(\mathcal{M}_{s})$$

- It depends on the model's size through the *critical radius* of random-sign width.
 - s satisfying $s^2/2 \ge w_s(\mathcal{M}_s)$ for $\mathcal{M}_s = \{m \in \mathcal{M} : \|m \mu^\star\|_{L_2(\mathcal{P}_n)} \le s\}$
 - This is a one-number summary of the random-sign width of neighborhoods ...
 - \cdot ...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.

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

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



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 P. Think sampling with replacement from a population.
- We'll bound the population root mean squared error $\|\hat{\mu} \mu^{\star}\|_{L_2(\mathbf{P})}$.

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

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

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

$$\|\hat{\mu} - \mu^{\star}\|_{L_{2}(\mathbf{P})}^{2} = \int {\{\hat{\mu}(x) - \mu^{\star}(x)\}^{2} p(x) dx} \quad \text{if} \quad X_{i} \quad \text{has the density} \quad p(x).$$



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.

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

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.



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

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'.



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

Same Argument, Different Neighborhood



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

$$\hat{\mu} \in \mathcal{M}_s \quad \text{for} \quad \mathcal{M}_s = \{ m \in \mathcal{M} \, : \, \|m - \mu^\star\|_{L_2(\mathbf{P})} \leq s \} \quad \text{for} \quad \mu^\star = \operatorname*{argmin}_{m \in \mathcal{M}} \|m - \mu\|_{L_2(\mathbf{P})}$$

- 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 μ^{\star} '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.

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

Reduction to a Maximal Inequality

$$\ell(m) - \ell(\mu^{*}) = \frac{1}{n} \sum_{i=1}^{n} Z_{i}(m) := \{m(X_{i}) - \mu^{*}(X_{i})\}^{2} - 2 \{Y_{i} - \mu^{*}(X_{i})\}\{m(X_{i}) - \mu^{*}(X_{i})\}$$
$$= E Z_{i}(m) + \frac{1}{n} \sum_{i=1}^{n} Z_{i}(m) - E Z_{i}(m).$$
Convexity Helps
as Usual.

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

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

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

$$- \operatorname{E}\left[\left\{Y_{i} - \mu^{*}(X_{i})\right\}\left\{m(X_{i}) - \mu^{*}(X_{i})\right\}\right] = - \operatorname{E}\left[\left\{\operatorname{E}\left[Y_{i} \mid X_{i}\right] - \mu^{*}(X_{i})\right\}\left\{m(X_{i}) - \mu^{*}(X_{i})\right\}\right]$$
$$= \left\langle\mu^{*} - \mu, m - \mu^{*}\right\rangle_{L_{2}(P)} \ge 0 \quad \text{for all} \quad m \in \mathcal{M}$$
It follows that $m \in \mathcal{M}_{s}$ if $m \in \mathcal{M}_{s}$ and $s^{2} > \max_{m \in \mathcal{M}_{s}^{o}} \frac{1}{n} \sum_{i=1}^{n} Z_{i}(m) - \operatorname{E} Z_{i}(m)$

Bounding the New Maximum

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

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

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

- (a) Write the centers $E Z_i(v)$ 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.

$$\begin{split} n \times \mathbf{E} \, \bar{Z} \stackrel{(a)}{=} \mathbf{E}_Z \max_{m \in \mathcal{M}_S^0} \mathbf{E}_{Z'} \sum_{i=1}^n \left\{ Z_i(m) - Z'_i(m) \right\} \\ \stackrel{(b)}{\leq} \mathbf{E}_Z \, \mathbf{E}_{Z'} \, \mathbb{E}_s \max_{m \in \mathcal{M}_S^0} \sum_{i=1}^n s_i \left\{ Z_i(m) - Z'_i(m) \right\} \\ \stackrel{(c)}{\leq} \mathbf{E}_Z \, \mathbf{E}_{Z'} \, \mathbf{E}_s \max_{m,m' \in \mathcal{M}_S^0} \sum_{i=1}^n s_i Z_i(m) + (-s_i) Z_i(m') \\ = 2 \, \mathbf{E}_Z \, \mathbf{E}_s \max_{m \in \mathcal{M}_S^0} \sum_{i=1}^n s_i Z_i(m) \end{split}$$

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

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

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 $\mathcal{M}_s - \mu^*$.

$$n \times \mathbf{E} Z \leq 2 \mathbf{E} \mathbf{E}_{s} \max_{m \in \mathcal{M}_{s}^{\circ}} \sum_{i=1}^{n} s_{i} Z_{i}(m) = \{m(X_{i}) - \mu^{*}(X_{i})\}^{2} - 2 \{Y_{i} - \mu^{*}(X_{i})\}\{m(X_{i}) - \mu^{*}(X_{i})\}$$
$$\leq 4 \mathbf{E} \left\{ \max_{m \in \mathcal{M}_{s}^{\circ}} ||m - \mu||_{L_{\infty}(\mathbf{Pn})} + ||\varepsilon||_{L_{\infty}(\mathbf{Pn})} \right\} \mathbf{E}_{s} \max_{m \in \mathcal{M}_{s}^{\circ}} \sum_{i=1}^{n} s_{i} \{m(X_{i}) - \mu^{*}(X_{i})\}$$

We've compared that to the width of the neighborhood itself using ... Lemma (Lipschitz Comparison)

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

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

that's
$$E_s \max_{m \in \mathcal{M}_s^0} \sum_{i=1}^n s_i \psi_i \{ m(X_i) - \mu^*(X_i) \} \le L \max_{m \in \mathcal{M}_s^0} \sum_{i=1}^n s_i \{ m(X_i) - \mu^*(X_i) \}$$

where $L = \max_i \max_{m \in \mathcal{M}_s^0} |\psi_i' \{ m(X_i) - \mu^*(X_i) \} |$
 $= \max_i \max_{m \in \mathcal{M}_s^0} |2\{ m(X_i) - \mu^*(X_i) \} - 2\{ Y_i - \mu^*(X_i) \} |$
 $\le 2 \max_{m \in \mathcal{M}_s^0} ||m - \mu||_{L_{\infty}(P_n)} + 2||\varepsilon||_{L_{\infty}(P_n)}.$

Interpretation



$$\begin{split} \|\hat{\mu} - \mu^{\star}\|_{L_{2}(\mathbf{P})} &\leq s \times 2\left\{\sqrt{\Sigma_{n}} + B\right\} + \sqrt{\frac{\operatorname{Var}(\bar{Z})}{\delta}} \quad \text{w.p. } 1 - \delta \\ \text{if } \frac{s^{2}}{2} &\geq \operatorname{Ew}_{s}(\mathcal{M}_{s}) \quad \text{and} \quad \|m - \mu\|_{\infty} \leq B \end{split}$$

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

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

Left: With little noise, our estimator μ̂ fits substantially better at the sample points X_i.
 Right: With more, it doesn't. The observations are far enough from μ 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 ($\Sigma_n = 0$),

$$\begin{aligned} \|\hat{\mu} - \mu\|_{L_2(\mathbf{P})} &\leq s \times 2\left\{\sqrt{\Sigma_n} + B\right\} + \sqrt{\frac{\operatorname{Var}(\bar{Z})}{\delta}} \quad \text{w.p. } 1 - \delta \\ \text{if } \frac{s^2}{2} &\geq \operatorname{Ew}_s(\mathcal{M}_s) \quad \text{and} \quad \|m - \mu\|_{\infty} \leq B \end{aligned}$$

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

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 μ between the sample points $X_1 \ldots X_n$, ...
- ...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 μ^{\star} .

$$\begin{split} \|\hat{\mu} - \mu\|_{L_2(\mathbf{P})} &\leq s \times 2\left\{\sqrt{\Sigma_n} + B\right\} + \sqrt{\frac{\operatorname{Var}(\bar{Z})}{\delta}} \quad \text{w.p. } 1 - \delta \\ & \text{if } \frac{s^2}{2} \geq \operatorname{Ew}_s(\mathcal{M}_s) \quad \text{and} \quad \|m - \mu\|_{\infty} \leq B \end{split}$$

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 $\hat{\mu}$ is will be within the range of the data.

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

This is true, in particular, for Monotone and Bounded Variation Regression. We can add this constraint to our model when doing our analysis.

$$\begin{split} \|\hat{\mu} - \mu^{\star}\|_{L_{2}(\mathbf{P})} &< s \quad \text{if} \quad \ell(m) - \ell(\mu^{\star}) > 0 \quad \text{for all} \quad m \in \mathcal{M} \dots \\ \dots \text{ with } \|m\|_{\infty} \leq B \text{ and } \|m - \mu^{\star}\|_{L_{2}(\mathbf{P})} \geq s \end{split}$$



The are other options.