

Machine Learning Theory

Lecture 7: Least Squares and Population Mean Squared Error

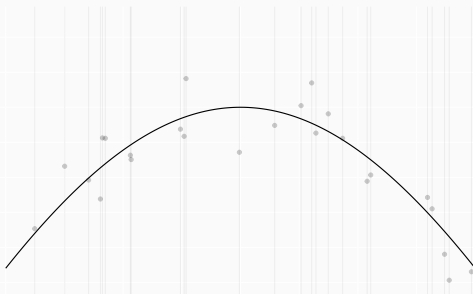
David A. Hirshberg

May 24, 2024

Emory University

Least squares with gaussian noise

We observe $Y_i = \mu(X_i) + \epsilon_i$ for $\epsilon_i \stackrel{iid}{\sim} N(0, \sigma^2)$.



We're estimating the curve $\mu(x)$.

We've been bounding sample mean squared error.

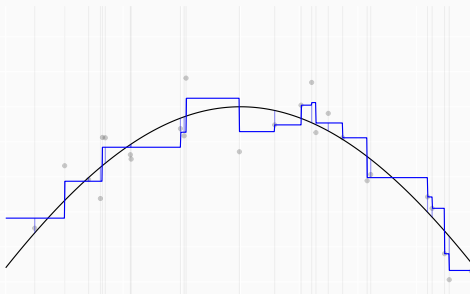
$$\|\hat{\mu} - \mu\|_{L_2(P_n)}^2 < s^2 \quad \text{with high probability.}$$

It's the average of the squared errors we make at the points $X_1 \dots X_n$.

What does this tell us?

This statement tells us $\hat{\mu}$ and μ are close, on average, at the observed X_i .

$$\frac{1}{n} \sum_{i=1}^n \{\mu(X_i) - \hat{\mu}(X_i)\}^2 < s^2.$$



It doesn't tell us they're close in the gaps between the X_i .

Today, we'll address that. We're going to bound population mean squared error.

$$\|\hat{\mu} - \mu\|_{L_2(\mathbb{P})}^2 < s^2 \quad \text{with high probability.}$$

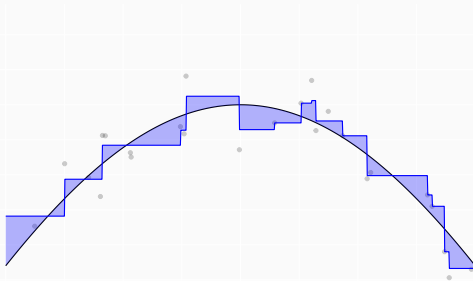
What Population Mean Squared Error Is

It's the mean squared error we make a random point X_{i+1} distributed like $X_1 \dots X_n$.

$$\|\hat{\mu} - \mu\|_{L_2(P)}^2 = \mathbb{E}_{X_{i+1}} [\{\hat{\mu}(X_{i+1}) - \mu(X_{i+1})\}^2]$$

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

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

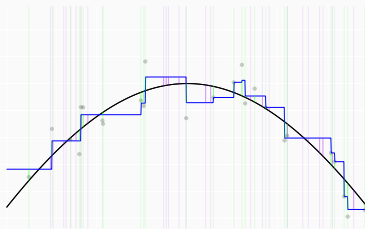


Why we care about Population Mean Squared Error: Generalization

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

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

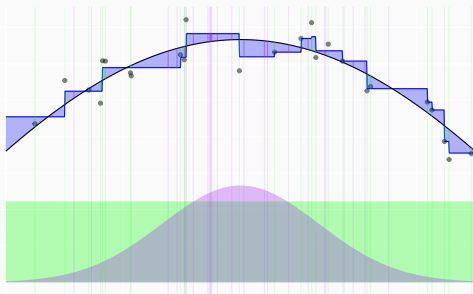
This can be a bit different from accuracy on our original sample $X_1 \dots X_n$.



- TV regression spends its 'variation budget' jumping to fit on the original sample.
- Between observations, it doesn't know whether it should jump or not.
- So we 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

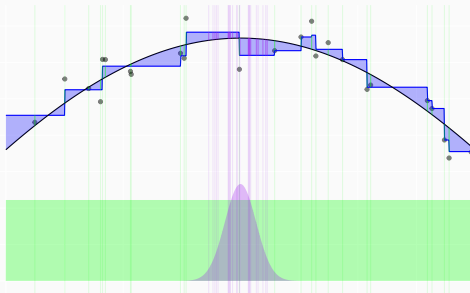
If we're interested in average accuracy for new points from a different distribution Q , we can bound this by comparing this distribution's density to that of our observations.



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

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' .



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

This is what we want, more or less, to estimate a treatment effect using a discontinuity. We want the value of the curve just left and right of the discontinuity.

Bounding Population Mean Squared Error

It's like bounding Sample Mean Squared Error

To bound Sample MSE

- We showed that with high probability the least squares estimator $\hat{\mu}$ must be in a neighborhood of the curve μ .
- Or, more generally, the model's best approximation to it, $\hat{\mu}^*$.
- Both of these things were defined in terms of distance on the sample.

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

To bound Population MSE

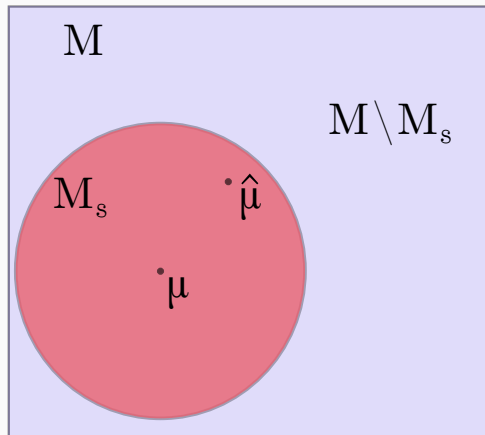
- We do exactly the same thing with different definitions.
- We define our neighborhood and best approximation using population distance.

$$\mathcal{M}_s = \{m \in \mathcal{M} : \|m - \mu^*\|_{L_2(P)} \leq s\}$$
$$\mu^* = \operatorname{argmin}_{m \in \mathcal{M}} \|m - \mu\|_{L_2(P)}.$$

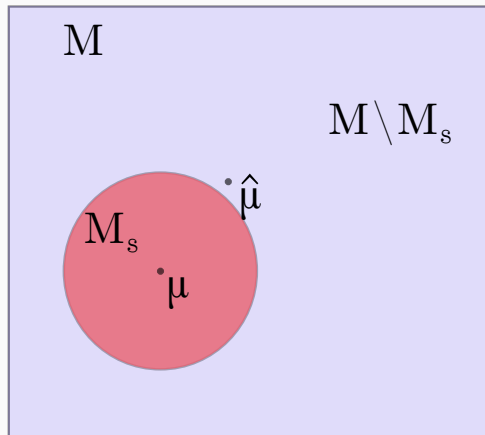
Warning: We're using the same notation for different but analogous things.
Context should tell us what's what.

We'll focus on the case that the curve μ is in the model \mathcal{M} .
We can generalize to the case that it isn't just like we did for sample MSE.

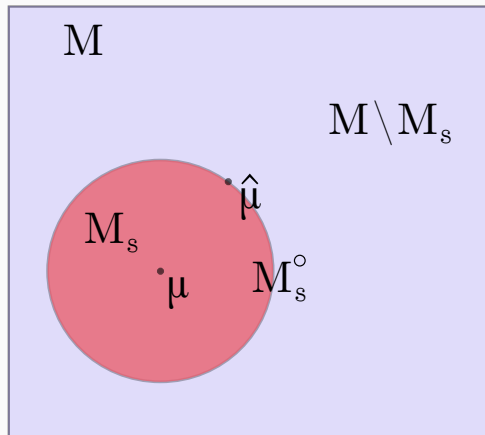
We'll show our estimator is in a neighborhood of μ .



To do it, we'll rule out the possibility that it's outside it.



And it's enough to rule out the possibility it's on the boundary.



Why? A Review

There are two forces working against each other when we're doing least squares.

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

- The first term, squared distance, tells us we don't want something far from μ .
- The second term, which has mean zero and scale \propto distance, obscures that.
- Once we're far enough away, the first term wins out. It grows faster with distance.

The distance involved is Sample MSE, but we can switch to Population MSE.

The difference between the two is just another mean-zero term.

$$\begin{aligned} \|m - \mu\|_{L_2(P)}^2 - \|m - \mu\|_{L_2(P_n)}^2 &= \frac{1}{n} \sum_{i=1}^n \mathbb{E} \{m(X_i) - \mu(X_i)\}^2 - \{m(X_i) - \mu(X_i)\}^2 \\ &= \frac{1}{n} \sum_{i=1}^n (\mathbb{E} - 1) \{m(X_i) - \mu(X_i)\}^2 \quad \text{in compact notation.} \end{aligned}$$

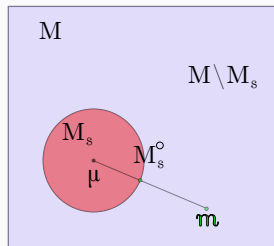
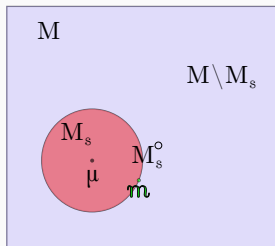
$$\ell(m) - \ell(\mu) = \|m - \mu\|_{L_2(P)}^2 - \frac{1}{n} \sum_{i=1}^n \left[\underbrace{2\varepsilon_i \{m(X_i) - \mu(X_i)\}}_{\text{old mean zero term}} + \underbrace{(\mathbb{E} - 1) \{m(X_i) - \mu(X_i)\}^2}_{\text{new mean zero term}} \right]$$

The Argument

$$\ell(m) - \ell(\mu) = \|m - \mu\|_{L_2(P)}^2 - \frac{1}{n} \sum_{i=1}^n \left[\underbrace{2\varepsilon_i \{m(X_i) - \mu(X_i)\}}_{\text{old mean zero term}} - \underbrace{(\mathbb{E} - 1) \{m(X_i) - \mu(X_i)\}^2}_{\text{new mean zero term}} \right]$$

For a certain *distance* s , squared distance exceeds this mean zero term for all curves in the model at that distance, i.e., curves $m \in \mathcal{M}_s^o$.

- That rules out the possibility that our minimizer is *any* curve at that distance.
- And any curve at a greater distance. It's representative is ruled out, and so is it.



So all we've got to do is characterize a radius s for which this happens.

The Characterization

$$\ell(m) - \ell(\mu) = s^2 - \frac{1}{n} \sum_{i=1}^n [2\varepsilon_i \{m(X_i) - \mu(X_i)\} + (E-1) \{m(X_i) - \mu(X_i)\}^2]$$

exceeds zero for all curves m in the neighborhood boundary \mathcal{M}_s° if ...

$$\frac{s^2}{2} > \max_{m \in \mathcal{M}_s^\circ} \frac{2}{n} \sum_{i=1}^n \varepsilon_i \{m(X_i) - \mu(X_i)\}$$

and

$$\frac{s^2}{2} > \max_{m \in \mathcal{M}_s^\circ} \frac{1}{n} \sum_{i=1}^n (E-1) \{m(X_i) - \mu(X_i)\}^2.$$

Here we're using half of s^2 to bound each mean zero term.

We know what we need for the old one.

It's enough that s^2 exceeds a multiple of its mean.

$$\frac{s^2}{2} \geq 2c_\delta w_\varepsilon(\mathcal{M}_s^\circ) = c_\delta E \max_{m \in \mathcal{M}_s^\circ} \frac{2}{n} \sum_{i=1}^n \varepsilon_i \{m(X_i) - \mu(X_i)\}.$$

To find a radius s satisfying this, we bound the width of our neighborhood boundary.

To deal with the new part, we do the same. We need to show this.

$$\frac{s^2}{2} \gtrsim c_\delta \mathbb{E} \max_{m \in \mathcal{M}_s^\circ} \frac{1}{n} \sum_{i=1}^n (\mathbb{E} - 1) \{m(X_i) - \mu(X_i)\}^2$$

Can we find a radius s satisfying this essentially the same way?

To deal with the new part, we do the same. We need to show this.

$$\frac{s^2}{2} \gtrsim c_\delta \mathbb{E} \max_{m \in \mathcal{M}_s^\circ} \frac{1}{n} \sum_{i=1}^n (\mathbb{E} - 1) \{m(X_i) - \mu(X_i)\}^2$$

Can we find a radius s satisfying this essentially the same way?

Yes. And we can understand why using *symmetrization*.

- Let's think about independent copy $X'_1 \dots X'_n$ of our observed points.
- They have the same distribution, so it's equivalent to center on the copy's mean.
- But they're independent, so we have a little more latitude in terms of arithmetic.

Step 1

The mean (E_X) over $X_1 \dots X_n$ is equivalent to a mean ($E_{X'}$) over our copy.

$$\begin{aligned} & E_X \max_{m \in \mathcal{M}_s^\circ} \sum_{i=1}^n [(E_X - 1) \{m(X_i) - \mu(X_i)\}^2] \\ &= E_X \max_{m \in \mathcal{M}_s^\circ} \sum_{i=1}^n [E_X \{m(X_i) - \mu(X_i)\}^2 - \{m(X_i) - \mu(X_i)\}^2] \\ &= E_X \max_{m \in \mathcal{M}_s^\circ} \sum_{i=1}^n [E_{X'} \{m(X'_i) - \mu(X'_i)\}^2 - \{m(X_i) - \mu(X_i)\}^2] \\ &= E_X \max_{m \in \mathcal{M}_s^\circ} E_{X'} \sum_{i=1}^n [\{m(X'_i) - \mu(X'_i)\}^2 - \{m(X_i) - \mu(X_i)\}^2]. \end{aligned}$$

Step 2

The mean of a maximum is larger than the maximum of a mean.

$$\begin{aligned} & \mathbb{E}_X \max_{m \in \mathcal{M}_s^\circ} \mathbb{E}_{X'} \sum_{i=1}^n [\{m(X'_i) - \mu(X'_i)\}^2 - \{m(X_i) - \mu(X_i)\}^2] \\ & \leq \mathbb{E}_X \mathbb{E}_{X'} \max_{m \in \mathcal{M}_s^\circ} \sum_{i=1}^n [\{m(X'_i) - \mu(X'_i)\}^2 - \{m(X_i) - \mu(X_i)\}^2]. \end{aligned}$$

Step 3

Terms as symmetric, so multiplying by *signs* $s_i = \pm 1$ doesn't change anything.

$$\begin{aligned} & \mathbb{E}_X \mathbb{E}_{X'} \max_{m \in \mathcal{M}_s^\circ} \sum_{i=1}^n [\{m(X'_i) - \mu(X'_i)\}^2 - \{m(X_i) - \mu(X_i)\}^2] \\ &= \mathbb{E}_X \mathbb{E}_{X'} \max_{m \in \mathcal{M}_s^\circ} \sum_{i=1}^n s_i [\{m(X'_i) - \mu(X'_i)\}^2 - \{m(X_i) - \mu(X_i)\}^2]. \end{aligned}$$

where our signs $s_i = \pm 1$ each with probability $1/2$.

Step 4

Each copy, multiplied by this random sign, is centered. Maximize each separately.

$$\begin{aligned} & \mathbb{E}_X \mathbb{E}_{X'} \max_{m \in \mathcal{M}_s^\circ} \sum_{i=1}^n s_i [\{m(X'_i) - \mu(X'_i)\}^2 - \{m(X_i) - \mu(X_i)\}^2] \\ & \leq \mathbb{E}_{X'} \max_{m \in \mathcal{M}_s^\circ} \sum_{i=1}^n s_i \{m(X'_i) - \mu(X'_i)\}^2 + \mathbb{E}_X \max_{m \in \mathcal{M}_s^\circ} \sum_{i=1}^n s_i \{m(X_i) - \mu(X_i)\}^2. \end{aligned}$$

This leaves us with two copies of the same thing.

Summary

We pay a factor of two to substitute a random sign s_i for centering ($E - 1$).

$$\begin{aligned} & E \max_{m \in \mathcal{M}_s^\circ} \frac{1}{n} \sum_{i=1}^n (E - 1) \{m(X_i) - \mu(X_i)\}^2 \\ & \leq 2 E \max_{m \in \mathcal{M}_s^\circ} \frac{1}{n} \sum_{i=1}^n s_i \{m(X_i) - \mu(X_i)\}^2. \end{aligned}$$

What this Leaves Us With. In the Gaussian-noise Case.

What we're left with is like the gaussian width we used to bound the old term.

$$\mathbb{E} \max_{m \in \mathcal{M}_s^\circ} \frac{2}{n} \sum_{i=1}^n s_i \{m(X_i) - \mu(X_i)\}^2 \quad \text{new term bound}$$

$$\mathbb{E} \max_{m \in \mathcal{M}_s^\circ} \frac{2\sigma}{n} \sum_{i=1}^n g_i \{m(X_i) - \mu(X_i)\} \quad \text{old term bound.}$$

But there are a few differences.

1. We've got a random sign s_i instead of standard normals g_i .
2. We've got squared differences $(m - \mu)^2$ instead of $m - \mu$.
3. It doesn't scale with the noise level σ .

For the most part, the third difference is the only one that matters.

Why?

- A mean like this with random signs is never much bigger than one with gaussians.
- The squaring only matters if there's some X_i where $m - \mu$ is big.

Both of these are implied by a *contraction inequality* for random-sign width.
A new one. We'll investigate this, using simulation as well as theory, in this week's lab.

Preview: Consequences in the Gaussian-noise Case.

In effect, when we're interested in population MSE

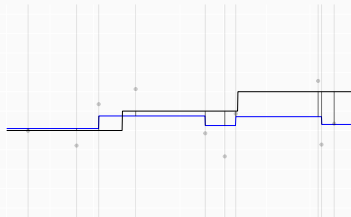
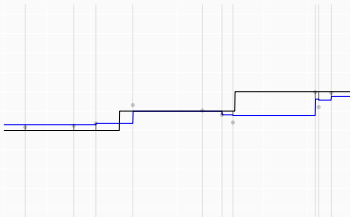
- random sampling from the population acts like gaussian noise
- that noise's standard deviation σ_{samp} is no larger than $2 \times$ the extreme values of the functions in our centered neighborhood.

$$\begin{aligned} \|\hat{\mu} - \mu\|_{L_2(\mathbf{P})} \leq s \quad & \text{with probability } 1 - \delta \quad \text{if} \quad s^2 \geq 2c_\delta (\sigma_{\text{samp}} + \sigma_{\text{noise}}) w(\mathcal{M}_s^\circ) \\ \text{where} \quad \sigma_{\text{samp}} &:= 2 \max_{m \in \mathcal{M}_s^\circ} \|m - \mu\|_\infty \\ \text{and} \quad \sigma_{\text{noise}} &:= \max_i \sqrt{\mathbb{E}\{Y_i - \mu(X_i)\}^2} \end{aligned}$$

When we're interested in sample MSE, we don't have that. Therefore

- sample MSE \approx population MSE unless there's relatively little noise.
- sample MSE \lesssim population MSE generally.

Low Noise vs High Noise



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 μ that we can't estimate it all that precisely even where we have some data.