

Machine Learning Theory

Least Squares and the Efron-Stein Inequality

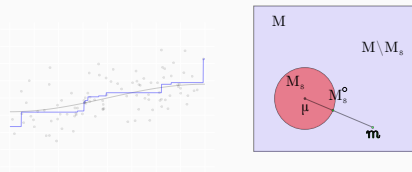
David A. Hirshberg

March 28, 2025

Emory University

Where We Left Things

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



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

$$\|\hat{\mu} - \mu\| < s \quad \text{w.p. } 1 - \delta \quad \text{if} \quad \frac{s^2}{2} \stackrel{(a)}{\geq} \mathbb{E} \max_{m \in \mathcal{M}_s^\circ} \langle \varepsilon, m - \mu \rangle + s\sigma \sqrt{\frac{2M_n}{\delta n}} \quad \text{for } M_n = 1 + 2\log(2n)$$

What We Actually Proved.

$$\|\hat{\mu} - \mu\| < s \quad \text{whenever} \quad \frac{s^2}{2} \stackrel{(b)}{\geq} \max_{m \in \mathcal{M}_s^\circ} \langle \varepsilon, m - \mu \rangle$$

Loose End. w.p. $1 - \delta$, (a) \implies (b). That is, ...

$$\max_{m \in \mathcal{M}_s^\circ} \langle \varepsilon, m - \mu \rangle \leq \mathbb{E} \max_{m \in \mathcal{M}_s^\circ} \langle \varepsilon, m - \mu \rangle + s\sigma \sqrt{\frac{2M_n}{\delta n}} \quad \text{w.p. } 1 - \delta.$$

Our Maximum is Approximately Constant

What we want to show.

$Z = \max_{m \in \mathcal{M}_s^o} \langle \varepsilon, m - \mu \rangle$ satisfies $Z \leq \mathbb{E} Z + s\sigma \sqrt{\frac{2M_n}{\delta n}}$ w.p. $1 - \delta$ for $M_n = 1 + 2 \log(2n)$.

We'll show something a bit stronger.

$$|Z - \mathbb{E} Z| \leq s\sigma \sqrt{\frac{2M_n}{\delta n}} \quad \text{w.p.} \quad 1 - \delta.$$

This is implied by Chebyshev's inequality. A special case of Markov's inequality.

$$\begin{aligned} & P \left\{ |Z - \mathbb{E} Z| \leq \frac{\text{sd}(Z)}{\sqrt{\delta}} \right\} \\ &= P \left\{ |Z - \mathbb{E} Z|^2 \leq \frac{\text{Var}(Z)}{\delta} \right\} \\ &\leq \frac{\mathbb{E} |Z - \mathbb{E} Z|^2}{\frac{\text{Var}(Z)}{\delta}} = \frac{\text{Var}(Z)}{\frac{\text{Var}(Z)}{\delta}} = \delta. \end{aligned}$$

All we need to do is bound the variance. We need to show that ...

$$\frac{\text{sd}(Z)}{\sqrt{\delta}} \leq s\sigma \sqrt{\frac{2M_n}{\delta n}} \quad \text{i.e.} \quad \text{Var}(Z) \leq s^2 \sigma^2 \frac{2M_n}{\delta n}.$$

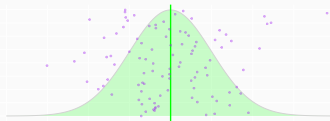
Variance and Independent Copies

$$\text{Var}[Z] = \text{Var}[f(\varepsilon)] \quad \text{for} \quad f(u) = \max_{m \in \mathcal{M}_s^\circ} \sum_{i=1}^n u_i \{m(X_i) - \mu(X_i)\}.$$

- Z is a pretty complicated function of our noise vector ε . To bound its variance, ...
- ...we'll need to think about it a bit differently than you're probably used to.

$$\begin{aligned} \text{Var}[Z] &= \text{E} \left[\{Z - \text{E}[Z]\}^2 \right] \\ &= \frac{1}{2} \text{E} \left[\{Z - \tilde{Z}\}^2 \right] \end{aligned} \quad \text{where } Z \text{ and } \tilde{Z} \text{ are independent and identically distributed.}$$

- It's the mean squared deviation of Z from its expectation.
- And *half* of the mean squared deviation of Z from an *independent copy* of Z .



Let's use all this to tackle a simplified version of our problem. We'll lose the max.

Calculate $\text{Var}[f(\varepsilon)]$ for $f(u) = \sum_{i=1}^n u_i$.

We can do this calculation the ‘usual way’ with these independent copies.

$$\begin{aligned}\text{Var}[f(\varepsilon)] &= \frac{1}{2} \text{E} \left[\left\{ \sum_{i=1}^n \varepsilon_i - \sum_{i=1}^n \tilde{\varepsilon}_i \right\}^2 \right] \\ &= \frac{1}{2} \text{E} \left[\left\{ \sum_{i=1}^n (\varepsilon_i - \tilde{\varepsilon}_i) \right\}^2 \right] \\ &= \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n \text{E} [(\varepsilon_i - \tilde{\varepsilon}_i)(\varepsilon_j - \tilde{\varepsilon}_j)] \\ &= \frac{1}{2} \sum_{i=1}^n \text{E} [(\varepsilon_i - \tilde{\varepsilon}_i)^2]\end{aligned}$$

We can use our independent copies to write this more abstractly, keeping everything ‘inside’ our summing function f .

$$\begin{aligned}\varepsilon_i - \tilde{\varepsilon}_i &= (\varepsilon_1 + \dots + \varepsilon_i + \tilde{\varepsilon}_{i-1} + \dots + \tilde{\varepsilon}_n) - (\varepsilon_1 + \dots + \varepsilon_{i-1} + \tilde{\varepsilon}_i + \dots + \tilde{\varepsilon}_n) \\ &= f(\varepsilon^{[i]}) - f(\varepsilon^{[i-1]}) \quad \text{where} \quad \varepsilon^{[i]} = (\varepsilon_1 \quad \varepsilon_2 \quad \dots \quad \varepsilon_i \quad \tilde{\varepsilon}_{i+1} \quad \tilde{\varepsilon}_{i+2} \quad \dots \quad \tilde{\varepsilon}_n)\end{aligned}$$

Calculate $\text{Var}[f(\varepsilon)]$ for $f(u) = \sum_{i=1}^n u_i$.

We can do this calculation the ‘usual way’ with these independent copies.

$$\begin{aligned}
 \text{Var}[f(\varepsilon)] &= \frac{1}{2} \text{E} \left[\left\{ \sum_{i=1}^n \varepsilon_i - \sum_{i=1}^n \tilde{\varepsilon}_i \right\}^2 \right] &= \frac{1}{2} \text{E} \left[\{f(\varepsilon) - f(\tilde{\varepsilon})\}^2 \right] \\
 &= \frac{1}{2} \text{E} \left[\left\{ \sum_{i=1}^n (\varepsilon_i - \tilde{\varepsilon}_i) \right\}^2 \right] \\
 &= \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n \text{E} [(\varepsilon_i - \tilde{\varepsilon}_i)(\varepsilon_j - \tilde{\varepsilon}_j)] \\
 &= \frac{1}{2} \sum_{i=1}^n \text{E} [(\varepsilon_i - \tilde{\varepsilon}_i)^2]
 \end{aligned}$$

We can use our independent copies to write this more abstractly,
keeping everything ‘inside’ our summing function f .

$$\begin{aligned}
 \varepsilon_i - \tilde{\varepsilon}_i &= (\varepsilon_1 + \dots + \varepsilon_i + \tilde{\varepsilon}_{i-1} + \dots + \tilde{\varepsilon}_n) - (\varepsilon_1 + \dots + \varepsilon_{i-1} + \tilde{\varepsilon}_i + \dots + \tilde{\varepsilon}_n) \\
 &= f(\varepsilon^{[i]}) - f(\varepsilon^{[i-1]}) \quad \text{where} \quad \varepsilon^{[i]} = (\varepsilon_1 \quad \varepsilon_2 \quad \dots \quad \varepsilon_i \quad \tilde{\varepsilon}_{i+1} \quad \tilde{\varepsilon}_{i+2} \quad \dots \quad \tilde{\varepsilon}_n)
 \end{aligned}$$

Calculate $\text{Var}[f(\varepsilon)]$ for $f(u) = \sum_{i=1}^n u_i$.

We can do this calculation the ‘usual way’ with these independent copies.

$$\begin{aligned}
 \text{Var}[f(\varepsilon)] &= \frac{1}{2} \text{E} \left[\left\{ \sum_{i=1}^n \varepsilon_i - \sum_{i=1}^n \tilde{\varepsilon}_i \right\}^2 \right] &&= \frac{1}{2} \text{E} \left[\{f(\varepsilon) - f(\tilde{\varepsilon})\}^2 \right] \\
 &= \frac{1}{2} \text{E} \left[\left\{ \sum_{i=1}^n (\varepsilon_i - \tilde{\varepsilon}_i) \right\}^2 \right] &&= \frac{1}{2} \text{E} \left[\left\{ \sum_{i=1}^n f(\varepsilon^{[i]}) - f(\varepsilon^{[i-1]}) \right\}^2 \right] \\
 &= \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n \text{E} [(\varepsilon_i - \tilde{\varepsilon}_i)(\varepsilon_j - \tilde{\varepsilon}_j)] \\
 &= \frac{1}{2} \sum_{i=1}^n \text{E} [(\varepsilon_i - \tilde{\varepsilon}_i)^2]
 \end{aligned}$$

We can use our independent copies to write this more abstractly,
keeping everything ‘inside’ our summing function f .

$$\begin{aligned}
 \varepsilon_i - \tilde{\varepsilon}_i &= (\varepsilon_1 + \dots + \varepsilon_i + \tilde{\varepsilon}_{i-1} + \dots + \tilde{\varepsilon}_n) - (\varepsilon_1 + \dots + \varepsilon_{i-1} + \tilde{\varepsilon}_i + \dots + \tilde{\varepsilon}_n) \\
 &= f(\varepsilon^{[i]}) - f(\varepsilon^{[i-1]}) \quad \text{where} \quad \varepsilon^{[i]} = (\varepsilon_1 \quad \varepsilon_2 \quad \dots \quad \varepsilon_i \quad \tilde{\varepsilon}_{i+1} \quad \tilde{\varepsilon}_{i+2} \quad \dots \quad \tilde{\varepsilon}_n)
 \end{aligned}$$

Calculate $\text{Var}[f(\varepsilon)]$ for $f(u) = \sum_{i=1}^n u_i$.

We can do this calculation the ‘usual way’ with these independent copies.

$$\begin{aligned}
 \text{Var}[f(\varepsilon)] &= \frac{1}{2} \text{E} \left[\left\{ \sum_{i=1}^n \varepsilon_i - \sum_{i=1}^n \tilde{\varepsilon}_i \right\}^2 \right] &&= \frac{1}{2} \text{E} \left[\{f(\varepsilon) - f(\tilde{\varepsilon})\}^2 \right] \\
 &= \frac{1}{2} \text{E} \left[\left\{ \sum_{i=1}^n (\varepsilon_i - \tilde{\varepsilon}_i) \right\}^2 \right] &&= \frac{1}{2} \text{E} \left[\left\{ \sum_{i=1}^n f(\varepsilon^{[i]}) - f(\varepsilon^{[i-1]}) \right\}^2 \right] \\
 &= \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n \text{E}[(\varepsilon_i - \tilde{\varepsilon}_i)(\varepsilon_j - \tilde{\varepsilon}_j)] &&= \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n \text{E} \left[\{f(\varepsilon^{[i]}) - f(\varepsilon^{[i-1]})\} \{f(\varepsilon^{[j]}) - f(\varepsilon^{[j-1]})\} \right] \\
 &= \frac{1}{2} \sum_{i=1}^n \text{E}[(\varepsilon_i - \tilde{\varepsilon}_i)^2]
 \end{aligned}$$

We can use our independent copies to write this more abstractly,
keeping everything ‘inside’ our summing function f .

$$\begin{aligned}
 \varepsilon_i - \tilde{\varepsilon}_i &= (\varepsilon_1 + \dots + \varepsilon_i + \tilde{\varepsilon}_{i-1} + \dots + \tilde{\varepsilon}_n) - (\varepsilon_1 + \dots + \varepsilon_{i-1} + \tilde{\varepsilon}_i + \dots + \tilde{\varepsilon}_n) \\
 &= f(\varepsilon^{[i]}) - f(\varepsilon^{[i-1]}) \quad \text{where} \quad \varepsilon^{[i]} = (\varepsilon_1 \quad \varepsilon_2 \quad \dots \quad \varepsilon_i \quad \tilde{\varepsilon}_{i+1} \quad \tilde{\varepsilon}_{i+2} \quad \dots \quad \tilde{\varepsilon}_n)
 \end{aligned}$$

Calculate $\text{Var}[f(\varepsilon)]$ for $f(u) = \sum_{i=1}^n u_i$.

We can do this calculation the ‘usual way’ with these independent copies.

$$\begin{aligned}
 \text{Var}[f(\varepsilon)] &= \frac{1}{2} \text{E} \left[\left\{ \sum_{i=1}^n \varepsilon_i - \sum_{i=1}^n \tilde{\varepsilon}_i \right\}^2 \right] &= \frac{1}{2} \text{E} \left[\{f(\varepsilon) - f(\tilde{\varepsilon})\}^2 \right] \\
 &= \frac{1}{2} \text{E} \left[\left\{ \sum_{i=1}^n (\varepsilon_i - \tilde{\varepsilon}_i) \right\}^2 \right] &= \frac{1}{2} \text{E} \left[\left\{ \sum_{i=1}^n f(\varepsilon^{[i]}) - f(\varepsilon^{[i-1]}) \right\}^2 \right] \\
 &= \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n \text{E} [(\varepsilon_i - \tilde{\varepsilon}_i)(\varepsilon_j - \tilde{\varepsilon}_j)] &= \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n \text{E} \left[\{f(\varepsilon^{[i]}) - f(\varepsilon^{[i-1]})\} \{f(\varepsilon^{[j]}) - f(\varepsilon^{[j-1]})\} \right] \\
 &= \frac{1}{2} \sum_{i=1}^n \text{E} [(\varepsilon_i - \tilde{\varepsilon}_i)^2] &= \frac{1}{2} \sum_{i=1}^n \text{E} \left[\{f(\varepsilon^{[i]}) - f(\varepsilon^{[i-1]})\}^2 \right]
 \end{aligned}$$

We can use our independent copies to write this more abstractly,
keeping everything ‘inside’ our summing function f .

$$\begin{aligned}
 \varepsilon_i - \tilde{\varepsilon}_i &= (\varepsilon_1 + \dots + \varepsilon_i + \tilde{\varepsilon}_{i-1} + \dots + \tilde{\varepsilon}_n) - (\varepsilon_1 + \dots + \varepsilon_{i-1} + \tilde{\varepsilon}_i + \dots + \tilde{\varepsilon}_n) \\
 &= f(\varepsilon^{[i]}) - f(\varepsilon^{[i-1]}) \quad \text{where} \quad \varepsilon^{[i]} = (\varepsilon_1 \quad \varepsilon_2 \quad \dots \quad \varepsilon_i \quad \tilde{\varepsilon}_{i+1} \quad \tilde{\varepsilon}_{i+2} \quad \dots \quad \tilde{\varepsilon}_n)
 \end{aligned}$$

The Variance of Sums: $\text{Var}[f(\varepsilon)]$ for $f(u) = \sum_{i=1}^n u_i$

$$\begin{aligned}\text{Var}[f(\varepsilon)] &= \frac{1}{2} \sum_{i=1}^n \mathbb{E} \left[\left\{ f(\varepsilon^{[i]}) - f(\varepsilon^{[i-1]}) \right\}^2 \right] & \text{for } \varepsilon_j^{[i]} &= \begin{cases} \varepsilon_j & j \leq i \\ \tilde{\varepsilon}_j & j > i \end{cases} \\ &= \frac{1}{2} \sum_{i=1}^n \mathbb{E} \left[\left\{ f(\varepsilon) - f(\varepsilon^{(i)}) \right\}^2 \right] & \text{for } \varepsilon_j^{(i)} &= \begin{cases} \tilde{\varepsilon}_i & j = i \\ \varepsilon_j & j \neq i \end{cases}\end{aligned}$$

We can derive the (simpler) second formula from the one we've just worked out. Here's the argument.

- The pair of vectors $\varepsilon^{[i]}, \varepsilon^{[i-1]}$ have the same joint distribution as $\varepsilon, \varepsilon^{(i)}$.
- It follows that any functions of those pairs,

$$\text{e.g. } f(\varepsilon^{[i]}) - f(\varepsilon^{[i-1]}) \quad \text{and} \quad f(\varepsilon) - f(\varepsilon^{(i)}),$$

have the same distribution. And therefore the same expectation.

How do we know our pairs have the same distribution?

- The first vectors, $\varepsilon^{[i]}$ and ε , have the same distribution.
- To get the second vector from the first, we do the same thing. We replace the i th component with an independent copy.

The Efron-Stein inequality: $\text{Var}[f(\varepsilon)]$ for arbitrary f

$$\text{Var}[f(\varepsilon)] \leq \frac{1}{2} \sum_{i=1}^n \mathbb{E} \left[\left\{ f(\varepsilon) - f(\varepsilon^{(i)}) \right\}^2 \right] \quad \text{for} \quad \varepsilon_j^{(i)} = \begin{cases} \tilde{\varepsilon}_i & j = i \\ \varepsilon_j & j \neq i \end{cases}$$

- Something very cool happens when we write things this way.
 - What we've derived isn't just a new formula for the variance of a sum.
 - It's a variance bound for *any function* of a vector of independent random variables.
- We call this the *Efron-Stein inequality*.
- There's an equivalent 'positive part' version that's sometimes easier to use.

$$\text{Var}[f(\varepsilon)] \leq \sum_{i=1}^n \mathbb{E} \left[\left\{ f(\varepsilon) - f(\varepsilon^{(i)}) \right\}_+^2 \right] \quad \text{for} \quad \{z\}_+ = \max\{z, 0\}.$$

- This is nice because $f(x) = \{x\}_+^2$ is increasing (whereas $f(x) = x^2$ is not).
- And that means we can substitute an upper bound for what's inside it.

$$\text{Var}[f(\varepsilon)] \leq \sum_{i=1}^n \mathbb{E} \{ F_i \}_+^2 \leq \sum_{i=1}^n \mathbb{E} F_i^2 \quad \text{for} \quad F_i \geq f(\varepsilon) - f(\varepsilon^{(i)}).$$

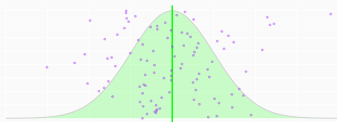
The 'Positive Part' Efron-Stein inequality

$$\begin{aligned}\text{Var}[f(\varepsilon)] &\leq \frac{1}{2} \sum_{i=1}^n \mathbb{E} \left[\left\{ f(\varepsilon) - f(\varepsilon^{(i)}) \right\}^2 \right] \\ &= \sum_{i=1}^n \mathbb{E} \left[\left\{ f(\varepsilon) - f(\varepsilon^{(i)}) \right\}_+^2 \right] \quad \text{for } \{z\}_+ = \max\{z, 0\}.\end{aligned}$$

- What's changed from the first formula to the second?
 - The differences on the right have been replaced with their *positive parts*.
 - We've lost the $\frac{1}{2}$ to compensate.
- Why is this equivalent? Symmetry.
- For any random variable S with a symmetric distribution¹, $\mathbb{E} S^2 = 2 \mathbb{E}\{S\}_+^2$.

Proof.

$$\begin{aligned}S^2 &= \{S\}_+^2 + \{-S\}_+^2 \\ &= \mathbb{E}\{S\}_+^2 + \mathbb{E}\{-S\}_+^2 \\ &= 2 \mathbb{E}\{S\}_+^2.\end{aligned}$$



□

¹A random variable S has a symmetric distribution if S and $-S$ have the same distribution.

The Variance of our Maximum

$$\text{Var}[f(\varepsilon)] \leq \sum_{i=1}^n \mathbb{E} \left[\left\{ f(\varepsilon) - f(\varepsilon^{(i)}) \right\}_+^2 \right] \quad \text{for } f(x) = \max_{m \in \mathcal{M}_s^\circ} \langle x, m - \mu \rangle$$

What do the terms on the right look like?

$$\begin{aligned} f(\varepsilon) - f(\varepsilon^{(i)}) &= \max_{m \in \mathcal{M}_s^\circ} \langle \varepsilon, m - \mu \rangle - \max_{m \in \mathcal{M}_s^\circ} \langle \varepsilon^{(i)}, m - \mu \rangle \\ &\leq \langle \varepsilon, \hat{m} - \mu \rangle - \langle \varepsilon^{(i)}, \hat{m} - \mu \rangle \quad \text{for } \hat{m} = \operatorname{argmax}_{m \in \mathcal{M}_s^\circ} \langle \varepsilon, m - \mu \rangle \\ &= \langle \varepsilon - \varepsilon^{(i)}, \hat{m} - \mu \rangle = \frac{1}{n} \{ \hat{m}(X_i) - \mu(X_i) \} (\varepsilon_i - \tilde{\varepsilon}_i). \end{aligned}$$

Plugging in these bounds, we get ...

$$\begin{aligned} \text{Var}[f(\varepsilon)] &\leq \frac{1}{n} \times \mathbb{E} \frac{1}{n} \sum_{i=1}^n \{ \hat{m}(X_i) - \mu(X_i) \}_{U_i}^2 (\varepsilon_i - \tilde{\varepsilon}_i)_{V_i}^2 &&= \frac{1}{n} \times \mathbb{E} \langle U, V \rangle_{L_2(P_n)} \\ &= \frac{1}{n} \times \frac{1}{n} \sum_{i=1}^n \{ \hat{m}(X_i) - \mu(X_i) \}^2 \mathbb{E} \max_{i \in 1 \dots n} (\varepsilon_i - \tilde{\varepsilon}_i)^2 &&= \frac{1}{n} \times \mathbb{E} \|U\|_{L_1(P_n)} \|V\|_{L_\infty(P_n)} \\ &= \frac{1}{n} \times s^2 \times \mathbb{E} \max_{i \in 1 \dots n} (\varepsilon_i - \tilde{\varepsilon}_i)^2 \\ &\leq \frac{1}{n} \times s^2 \times 2\sigma^2 M_n \quad \text{for } M_n = 1 + 2 \log(2n).^2 \end{aligned}$$

² M_n bounds the maximum of the squares of n independent standard normals. Scaling by $2\sigma^2$ gives a bound for normals with variance $\text{Var}[\varepsilon_i - \tilde{\varepsilon}_i] = 2\sigma^2$.

A Proof of the Efron-Stein inequality

$$\begin{aligned}\text{Var}[f(\varepsilon)] &= \mathbb{E} f(\varepsilon)^2 - \{\mathbb{E} f(\varepsilon)\}^2 \\&= \mathbb{E} f(\varepsilon)^2 - \mathbb{E} f(\varepsilon) \mathbb{E} f(\tilde{\varepsilon}) \\&= \mathbb{E} f(\varepsilon) \{f(\varepsilon) - \mathbb{E} f(\varepsilon)\} \\&= \mathbb{E} f(\varepsilon) \left\{ \sum_{i=1}^n f(\varepsilon^{[i]}) - f(\varepsilon^{[i-1]}) \right\} \\&= \sum_{i=1}^n \mathbb{E} f(\varepsilon) \left\{ f(\varepsilon^{[i]}) - f(\varepsilon^{[i-1]}) \right\} \quad \text{where} \quad \varepsilon_j^{[i]} = \begin{cases} \varepsilon_j & j \leq i \\ \tilde{\varepsilon}_j & j > i \end{cases}\end{aligned}$$

The Swapping Trick

$$\text{Var} [f(\varepsilon)] = \sum_{i=1}^n \mathbb{E} f(\varepsilon) \left\{ f(\varepsilon^{[i]}) - f(\varepsilon^{[i-1]}) \right\} \quad \text{where} \quad \varepsilon_j^{[i]} = \begin{cases} \varepsilon_j & j \leq i \\ \tilde{\varepsilon}_j & j > i \end{cases}$$

- Think of the i th term as a function of ε : $g_i(\varepsilon) = f(\varepsilon) \{ f(\varepsilon^{[i]}) - f(\varepsilon^{[i-1]}) \}$.
- Swapping $\varepsilon_i \rightarrow \tilde{\varepsilon}_i$ doesn't change the distribution of ε .
- So it doesn't change the distribution — or expectation — of $g_i(\varepsilon)$.

$$\begin{aligned} f(\varepsilon) \left\{ f(\varepsilon^{[i]}) - f(\varepsilon^{[i-1]}) \right\} &\xrightarrow{A_i = g_i(\varepsilon)} \xrightarrow{B_i = g_i(\varepsilon^{(i)})} f(\varepsilon^{(i)}) \left\{ f(\varepsilon^{[i-1]}) - f(\tilde{\varepsilon}^{[i]}) \right\} \quad \text{for} \quad \tilde{\varepsilon}_j^{(i)} = \begin{cases} \tilde{\varepsilon}_i & j = i \\ \varepsilon_j & j \neq i \end{cases} \\ &= -f(\varepsilon^{(i)}) \left\{ f(\varepsilon^{[i]}) - f(\varepsilon^{[i-1]}) \right\}. \end{aligned}$$

Because $A_i = B_i = (A_i + B_i)/2$, it follows that $\text{Var} [f(\varepsilon)] = \frac{1}{2} \sum_{i=1}^n \mathbb{E}[A_i + B_i]$ where

$$A_i + B_i = \left\{ f(\varepsilon) - f(\varepsilon^{(i)}) \right\} \left\{ f(\varepsilon^{[i]}) - f(\varepsilon^{[i-1]}) \right\}$$

$$\begin{aligned}\text{Var}[f(\varepsilon)] &= \frac{1}{2} \sum_{i=1}^n \mathbb{E} \left[\left\{ f(\varepsilon) - f(\varepsilon^{(i)}) \right\} \left\{ f(\varepsilon^{[i]}) - f(\varepsilon^{[i-1]}) \right\} \right] \\ &\stackrel{(a)}{\leq} \frac{1}{2} \sqrt{ \sum_{i=1}^n \mathbb{E} \left[\left\{ f(\varepsilon) - f(\varepsilon^{(i)}) \right\}^2 \right] \sum_{i=1}^n \mathbb{E} \left[\left\{ f(\varepsilon^{[i]}) - f(\varepsilon^{[i-1]}) \right\}^2 \right] } \\ &\stackrel{(b)}{=} \frac{1}{2} \sqrt{ \left\{ \sum_{i=1}^n \mathbb{E} \left[\left\{ f(\varepsilon) - f(\varepsilon^{(i)}) \right\}^2 \right] \right\}^2 } \\ &= \frac{1}{2} \sum_{i=1}^n \mathbb{E} \left[\left\{ f(\varepsilon) - f(\varepsilon^{(i)}) \right\}^2 \right].\end{aligned}$$

The rest boils down to

- (a) Using the $\langle \cdot, \cdot \rangle_{L_2(\mathbf{P})}$ Cauchy-Schwarz bound on each term in the sum.
- (b) Our observation, from a few slides back, that $\{f(\varepsilon) - f(\varepsilon^{(i)})\}^2$ and $\{f(\varepsilon^{[i]}) - f(\varepsilon^{[i-1]})\}^2$ have the same distribution.

- Sourav Chatterjee's class Stein's method and applications.
 - The proof of the Efron-Stein inequality is based on lecture 10.
- Boucheron, Lugosi, and Massart's *Concentration inequalities: A nonasymptotic theory of independence*.
 - The bound on the variance of the maximum $\max_{m \in \mathcal{M}_s^o} \langle \varepsilon, m - \mu \rangle$ is based on Example 3.6 in Chapter 3.
 - The bound M_n on $\mathbb{E} \max_{i \in 1 \dots n} \varepsilon_i^2$ for $\varepsilon_i \stackrel{iid}{\sim} N(0, 1)$ is from Lemma 11.3 in Chapter 11.