Machine Learning Theory

Least Squares and the Efron-Stein Inequality

David A. Hirshberg May 4, 2025

Emory University

Where We Left Things

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



Claim. When $Y_i = \mu(X_i) + \varepsilon_i$ for $\mu \in \mathcal{M}$ and $\varepsilon_i \stackrel{iid}{\sim} N(0, \sigma^2)$, $\|\hat{\mu} - \mu\| < s \quad \text{w.p. } 1 - \delta \text{ if } \quad \frac{s^2}{2} \stackrel{(a)}{\geq} \operatorname{E} \max_{m \in \mathcal{M}_s} \langle \varepsilon, \ m - \mu \rangle + s\sigma \sqrt{\frac{2\Sigma_n}{\delta n}} \text{ for } \Sigma_n = 1 + 2\log(2n).$

What We Actually Proved.

$$\|\hat{\mu} - \mu\| < s$$
 whenever $\frac{s^2}{2} \ge \max_{m \in \mathcal{M}_{\mathbb{Q}}^{\circ}} \langle \varepsilon, m - \mu \rangle$

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

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

What we want to show.

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

We'll show something a bit stronger.

$$|Z - \operatorname{E} Z| < s\sigma \sqrt{\frac{2\Sigma_n}{\delta n}} \quad \text{w.p.} \quad 1 - \delta.$$

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

$$P\left\{|Z - \operatorname{E} Z| \ge \frac{\operatorname{sd}(Z)}{\sqrt{\delta}}\right\} = P\left\{|Z - \operatorname{E} Z|^2 \ge \frac{\operatorname{Var}(Z)}{\delta}\right\} \le \frac{\operatorname{E}|Z - \operatorname{E} Z|^2}{\frac{\operatorname{Var}(Z)}{\delta}} = \frac{\operatorname{Var}(Z)}{\frac{\operatorname{Var}(Z)}{\delta}} = \delta$$

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

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

Variance and Independent Copies

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

- 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.

$$Var[Z] = E\left[\{Z - E[Z]\}^2\right]$$
$$= \frac{1}{2} E\left[\{Z - \tilde{Z}\}^2\right]$$

where Z and \tilde{Z} are independent and identically distributed.

- It's the mean squared deviation of Z from its expectation.
- \cdot 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
$$\operatorname{Var}[f(\varepsilon)]$$
 for $f(u) = \sum_{i=1}^{n} u_i$

$$\operatorname{Var}\left[f(\varepsilon)\right] = \frac{1}{2} \operatorname{E}\left[\left\{\sum_{i=1}^{n} \varepsilon_{i} - \sum_{i=1}^{n} \tilde{\varepsilon}_{i}\right\}^{2}\right]$$
$$= \frac{1}{2} \operatorname{E}\left[\left\{\sum_{i=1}^{n} (\varepsilon_{i} - \tilde{\varepsilon}_{i})\right\}^{2}\right]$$
$$= \frac{1}{2} \sum_{i=1}^{n} \sum_{j=1}^{n} \operatorname{E}\left[(\varepsilon_{i} - \tilde{\varepsilon}_{i})(\varepsilon_{j} - \tilde{\varepsilon}_{j})\right]$$
$$= \frac{1}{2} \sum_{i=1}^{n} \operatorname{E}\left[(\varepsilon_{i} - \tilde{\varepsilon}_{i})^{2}\right]$$

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

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

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

$$\operatorname{Var}\left[f(\varepsilon)\right] = \frac{1}{2} \operatorname{E}\left[\left\{\sum_{i=1}^{n} \varepsilon_{i} - \sum_{i=1}^{n} \tilde{\varepsilon}_{i}\right\}^{2}\right] = \frac{1}{2} \operatorname{E}\left[\left\{f(\varepsilon) - f(\tilde{\varepsilon})\right\}^{2}\right]$$
$$= \frac{1}{2} \operatorname{E}\left[\left\{\sum_{i=1}^{n} (\varepsilon_{i} - \tilde{\varepsilon}_{i})\right\}^{2}\right]$$
$$= \frac{1}{2} \sum_{i=1}^{n} \sum_{j=1}^{n} \operatorname{E}\left[(\varepsilon_{i} - \tilde{\varepsilon}_{i})(\varepsilon_{j} - \tilde{\varepsilon}_{j})\right]$$
$$= \frac{1}{2} \sum_{i=1}^{n} \operatorname{E}\left[(\varepsilon_{i} - \tilde{\varepsilon}_{i})^{2}\right]$$

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

$$\begin{split} \varepsilon_{i} - \tilde{\varepsilon}_{i} &= (\varepsilon_{1} + \ldots + \varepsilon_{i} + \tilde{\varepsilon}_{i-1} + \ldots + \tilde{\varepsilon}_{n}) - (\varepsilon_{1} + \ldots + \varepsilon_{i-1} + \tilde{\varepsilon}_{i} + \ldots + \tilde{\varepsilon}_{n}) \\ &= f\left(\varepsilon^{[i]}\right) - f\left(\varepsilon^{[i-1]}\right) \quad \text{where} \quad \varepsilon^{[i]} = \left(\varepsilon_{1} \quad \varepsilon_{2} \quad \ldots \quad \varepsilon_{i} \quad \tilde{\varepsilon}_{i+1} \quad \tilde{\varepsilon}_{i+2} \quad \ldots \quad \tilde{\varepsilon}_{n}\right) \end{split}$$

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

$$\operatorname{Var}\left[f(\varepsilon)\right] = \frac{1}{2} \operatorname{E}\left[\left\{\sum_{i=1}^{n} \varepsilon_{i} - \sum_{i=1}^{n} \tilde{\varepsilon}_{i}\right\}^{2}\right] = \frac{1}{2} \operatorname{E}\left[\left\{f(\varepsilon) - f(\tilde{\varepsilon})\right\}^{2}\right]$$
$$= \frac{1}{2} \operatorname{E}\left[\left\{\sum_{i=1}^{n} (\varepsilon_{i} - \tilde{\varepsilon}_{i})\right\}^{2}\right] = \frac{1}{2} \operatorname{E}\left[\left\{\sum_{i=1}^{n} f\left(\varepsilon^{[i]}\right) - f\left(\varepsilon^{[i-1]}\right)\right\}^{2}\right]$$
$$= \frac{1}{2} \sum_{i=1}^{n} \sum_{j=1}^{n} \operatorname{E}\left[(\varepsilon_{i} - \tilde{\varepsilon}_{i})(\varepsilon_{j} - \tilde{\varepsilon}_{j})\right]$$
$$= \frac{1}{2} \sum_{i=1}^{n} \operatorname{E}\left[(\varepsilon_{i} - \tilde{\varepsilon}_{i})^{2}\right]$$

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

$$\begin{split} \varepsilon_{i} - \tilde{\varepsilon}_{i} &= (\varepsilon_{1} + \ldots + \varepsilon_{i} + \tilde{\varepsilon}_{i-1} + \ldots + \tilde{\varepsilon}_{n}) - (\varepsilon_{1} + \ldots + \varepsilon_{i-1} + \tilde{\varepsilon}_{i} + \ldots + \tilde{\varepsilon}_{n}) \\ &= f\left(\varepsilon^{[i]}\right) - f\left(\varepsilon^{[i-1]}\right) \quad \text{where} \quad \varepsilon^{[i]} = \left(\varepsilon_{1} \quad \varepsilon_{2} \quad \ldots \quad \varepsilon_{i} \quad \tilde{\varepsilon}_{i+1} \quad \tilde{\varepsilon}_{i+2} \quad \ldots \quad \tilde{\varepsilon}_{n}\right) \end{split}$$

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

$$\operatorname{Var}\left[f(\varepsilon)\right] = \frac{1}{2} \operatorname{E}\left[\left\{\sum_{i=1}^{n} \varepsilon_{i} - \sum_{i=1}^{n} \tilde{\varepsilon}_{i}\right\}^{2}\right] = \frac{1}{2} \operatorname{E}\left[\left\{f(\varepsilon) - f(\tilde{\varepsilon})\right\}^{2}\right]$$
$$= \frac{1}{2} \operatorname{E}\left[\left\{\sum_{i=1}^{n} (\varepsilon_{i} - \tilde{\varepsilon}_{i})\right\}^{2}\right] = \frac{1}{2} \operatorname{E}\left[\left\{\sum_{i=1}^{n} f\left(\varepsilon^{[i]}\right) - f\left(\varepsilon^{[i-1]}\right)\right\}^{2}\right]$$
$$= \frac{1}{2} \sum_{i=1}^{n} \sum_{j=1}^{n} \operatorname{E}\left[(\varepsilon_{i} - \tilde{\varepsilon}_{i})(\varepsilon_{j} - \tilde{\varepsilon}_{j})\right] = \frac{1}{2} \sum_{i=1}^{n} \sum_{j=1}^{n} \operatorname{E}\left[\left\{f\left(\varepsilon^{[i]}\right) - f\left(\varepsilon^{[i-1]}\right)\right\}\left\{f\left(\varepsilon^{[j]}\right)\right\}$$
$$= \frac{1}{2} \sum_{i=1}^{n} \operatorname{E}\left[(\varepsilon_{i} - \tilde{\varepsilon}_{i})^{2}\right]$$

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

$$\begin{split} \varepsilon_{i} - \tilde{\varepsilon}_{i} &= (\varepsilon_{1} + \ldots + \varepsilon_{i} + \tilde{\varepsilon}_{i-1} + \ldots + \tilde{\varepsilon}_{n}) - (\varepsilon_{1} + \ldots + \varepsilon_{i-1} + \tilde{\varepsilon}_{i} + \ldots + \tilde{\varepsilon}_{n}) \\ &= f\left(\varepsilon^{[i]}\right) - f\left(\varepsilon^{[i-1]}\right) \quad \text{where} \quad \varepsilon^{[i]} = \left(\varepsilon_{1} \quad \varepsilon_{2} \quad \ldots \quad \varepsilon_{i} \quad \tilde{\varepsilon}_{i+1} \quad \tilde{\varepsilon}_{i+2} \quad \ldots \quad \tilde{\varepsilon}_{n}\right) \end{split}$$

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

$$\operatorname{Var}\left[f(\varepsilon)\right] = \frac{1}{2} \operatorname{E}\left[\left\{\sum_{i=1}^{n} \varepsilon_{i} - \sum_{i=1}^{n} \tilde{\varepsilon}_{i}\right\}^{2}\right] = \frac{1}{2} \operatorname{E}\left[\left\{f(\varepsilon) - f(\tilde{\varepsilon})\right\}^{2}\right]$$
$$= \frac{1}{2} \operatorname{E}\left[\left\{\sum_{i=1}^{n} (\varepsilon_{i} - \tilde{\varepsilon}_{i})\right\}^{2}\right] = \frac{1}{2} \operatorname{E}\left[\left\{\sum_{i=1}^{n} f\left(\varepsilon^{[i]}\right) - f\left(\varepsilon^{[i-1]}\right)\right\}^{2}\right]$$
$$= \frac{1}{2} \sum_{i=1}^{n} \sum_{j=1}^{n} \operatorname{E}\left[(\varepsilon_{i} - \tilde{\varepsilon}_{i})(\varepsilon_{j} - \tilde{\varepsilon}_{j})\right] = \frac{1}{2} \sum_{i=1}^{n} \sum_{j=1}^{n} \operatorname{E}\left[\left\{f\left(\varepsilon^{[i]}\right) - f\left(\varepsilon^{[i-1]}\right)\right\}\left\{f\left(\varepsilon^{[j]}\right)\right\}^{2}\right]$$
$$= \frac{1}{2} \sum_{i=1}^{n} \operatorname{E}\left[(\varepsilon_{i} - \tilde{\varepsilon}_{i})^{2}\right] = \frac{1}{2} \sum_{i=1}^{n} \operatorname{E}\left[\left\{f\left(\varepsilon^{[i]}\right) - f\left(\varepsilon^{[i-1]}\right)\right\}^{2}\right]$$

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

$$\begin{split} \varepsilon_{i} - \tilde{\varepsilon}_{i} &= (\varepsilon_{1} + \ldots + \varepsilon_{i} + \tilde{\varepsilon}_{i-1} + \ldots + \tilde{\varepsilon}_{n}) - (\varepsilon_{1} + \ldots + \varepsilon_{i-1} + \tilde{\varepsilon}_{i} + \ldots + \tilde{\varepsilon}_{n}) \\ &= f\left(\varepsilon^{[i]}\right) - f\left(\varepsilon^{[i-1]}\right) \quad \text{where} \quad \varepsilon^{[i]} = \left(\varepsilon_{1} \quad \varepsilon_{2} \quad \ldots \quad \varepsilon_{i} \quad \tilde{\varepsilon}_{i+1} \quad \tilde{\varepsilon}_{i+2} \quad \ldots \quad \tilde{\varepsilon}_{n}\right) \end{split}$$

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

$$\begin{aligned} \operatorname{Var}[f(\varepsilon)] &= \frac{1}{2} \sum_{i=1}^{n} \operatorname{E}\left[\left\{f\left(\varepsilon^{[i]}\right) - f\left(\varepsilon^{[i-1]}\right)\right\}^{2}\right] & \text{for } \varepsilon_{j}^{[i]} = \begin{cases} \varepsilon_{j} & j \leq i\\ \widetilde{\varepsilon}_{j} & j > i \end{cases} \\ &= \frac{1}{2} \sum_{i=1}^{n} \operatorname{E}\left[\left\{f(\varepsilon) - f\left(\varepsilon^{(i)}\right)\right\}^{2}\right] & \text{for } \varepsilon_{j}^{(i)} = \begin{cases} \widetilde{\varepsilon}_{i} & j = i\\ \varepsilon_{j} & j \neq i \end{cases} \end{aligned}$$

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,

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

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.

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

- 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.

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

• This is nice because $f(x) = \{x\} + 2$ is increasing (whereas $f(x) = x^2$ is not).

 $\cdot\,$ And that means we can substitute an upper bound for what's inside it.

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

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

- 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, $ES^2 = 2E\{S\}^2_+$.

Proof.

$$S^{2} = \{S\}_{+}^{2} + \{-S\}_{+}^{2}$$

= $E\{S\}_{+}^{2} + E\{-S\}_{+}^{2}$
= $2E\{S\}_{+}^{2}$.

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

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

What do the terms on the right look like?

$$f(\varepsilon) - f(\varepsilon^{(i)}) = \max_{m \in \mathcal{M}_{s}^{\circ}} \langle \varepsilon, \ m - \mu \rangle - \max_{m \in \mathcal{M}_{s}^{\circ}} \langle \varepsilon^{(i)}, \ m - \mu \rangle$$
$$\leq \langle \varepsilon, \ \hat{m} - \mu \rangle - \left\langle \varepsilon^{(i)}, \ \hat{m} - \mu \right\rangle \quad \text{for} \quad \hat{m} = \operatorname*{argmax}_{m \in \mathcal{M}_{s}^{\circ}} \langle \varepsilon, \ m - \mu \rangle$$
$$= \left\langle \varepsilon - \varepsilon^{(i)}, \ \hat{m} - \mu \right\rangle = \frac{1}{n} \{ \hat{m}(X_{i}) - \mu(X_{i}) \} (\varepsilon_{i} - \tilde{\varepsilon}_{i}).$$

Plugging in these bounds, we get ...

$$\begin{aligned} \operatorname{Var}[f(\varepsilon)] &\leq \frac{1}{n} \times \operatorname{E} \frac{1}{n} \sum_{i=1}^{n} \left\{ \hat{m}(X_{i}) - \mu(X_{i}) \right\}^{2} (\varepsilon_{i} - \tilde{\varepsilon}_{i})^{2} &= \frac{1}{n} \times \operatorname{E} \langle U, V \rangle_{L_{2}(\operatorname{Pn})} \\ &= \frac{1}{n} \times \frac{1}{n} \sum_{i=1}^{n} \left\{ \hat{m}(X_{i}) - \mu(X_{i}) \right\}^{2} \operatorname{E} \max_{i \in 1...n} (\varepsilon_{i} - \tilde{\varepsilon}_{i})^{2} &= \frac{1}{n} \times \operatorname{E} \| U \|_{L_{1}(\operatorname{Pn})} \| V \|_{L_{\infty}(\operatorname{Pn})} \\ &= \frac{1}{n} \times s^{2} \times \operatorname{E} \max_{i \in 1...n} (\varepsilon_{i} - \tilde{\varepsilon}_{i})^{2} \\ &\leq \frac{1}{n} \times s^{2} \times 2\sigma^{2} \Sigma_{n} \quad \text{for} \quad \Sigma_{n} = 1 + 2\log(2n).^{2} \end{aligned}$$

 ${}^{2}\Sigma_{n}$ bounds the maximum of the squares of n independent standard normals. Scaling by $2\sigma^{2}$ gives a bound for normals with variance $\operatorname{Var}[\varepsilon_{i} - \tilde{\varepsilon}_{i}] = 2\sigma^{2}$. A Proof of the Efron-Stein inequality

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

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

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

$$\begin{split} f(\varepsilon) \Big\{ f\Big(\varepsilon^{[i]}\Big) - f\Big(\varepsilon^{[i-1]}\Big) \Big\} &\to f(\varepsilon^{(i)}) \Big\{ f\Big(\varepsilon^{[i-1]}\Big) - f\Big(\tilde{\varepsilon}^{[i]}\Big) \Big\} \quad \text{for} \quad \tilde{\varepsilon}_j^{(i)} = \begin{cases} \tilde{\varepsilon}_i & j = i\\ \varepsilon_j & j \neq i \end{cases} \\ &= -f(\varepsilon^{(i)}) \Big\{ f\Big(\varepsilon^{[i]}\Big) - f\Big(\varepsilon^{[i-1]}\Big) \Big\}. \end{split}$$

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

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

Finishing Up

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

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 $\max_{m\in\mathcal{M}_s^\circ}\langle\varepsilon,\ m-\mu\rangle$ is based on Example 3.6 in Chapter 3.
 - The bound M_n on $\operatorname{Emax}_{i\in 1...n} \varepsilon_i^2$ for $\varepsilon_i \stackrel{iid}{\sim} N(0,1)$ is from Lemma 11.3 in Chapter 11.