# Machine Learning Theory
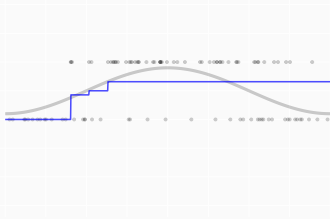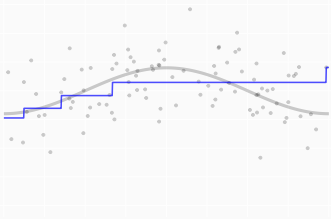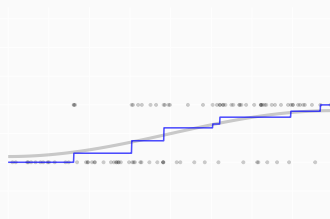
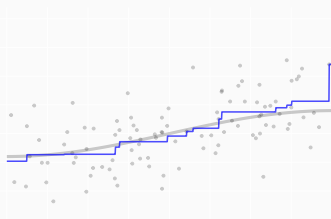Least Squares with Misspecification and Non-Gaussian Noise

David A. Hirshberg

April 1, 2025

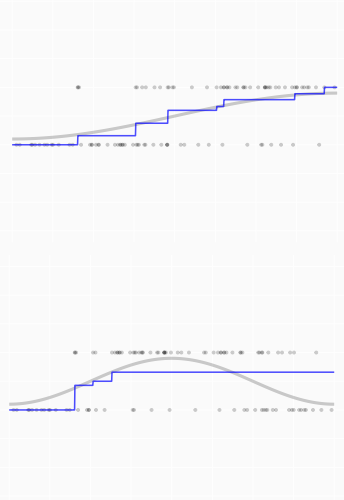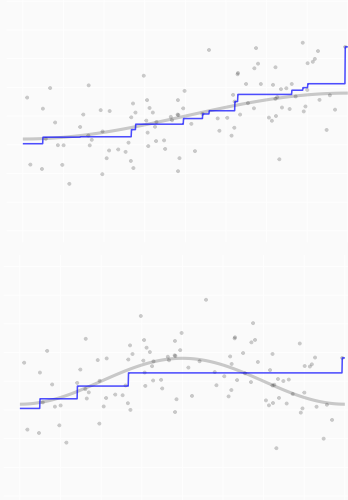Emory University

- The second column is out. We've assumed correct specfication.
- The second row is out. We've assumed normality.

- With misspecification, we estimate the model's best approximation to $\mu$.
- Non-normality doesn't really matter much. We'll look at how it affects our bound.

# Misspecification

- Our error in estimating $\mu$ is bounded by a sum of two terms.

  - The critical radius $s$, i.e., the one satisfying $s^2/2\sigma \geq \mathrm{w}(\mathcal{M}_s^\circ) + s\sqrt{\frac{2M_{n2}^2}{\delta n}}$.
  - The distance from $\mu$ to its best approximation in the model. Or really 3 times that.

  We showed this in the model selection lab using the Cauchy-Schwarz inequality.

- In convex models, we can say more.
  Our error in estimating $\mu^\star$ does not depend on its distance to $\mu$.

## The Argument

For any $\mu^\star \in \mathcal{M}$, we can expand our mean squared error difference as before.

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

But our new 'noise' $\varepsilon_i^\star$ doesn't have mean zero. It's our old noise $\varepsilon_i$, minus something.

$$\varepsilon_i^\star = \{\underbrace{Y_i - \mu(X_i)}_{\varepsilon_i}\} - \{\underbrace{\mu^\star(X_i) - \mu(X_i)}_{\text{something}}\}.$$

So we can think of our mean squared error difference as having three terms:

$$\ell(m) - \ell(\mu^\star) = \|m - \mu^\star\|_{L_2(\mathrm{P_n})}^2 \qquad \qquad \text{squared distance, like before;}$$
$$- \frac{2}{n} \sum_{i=1}^n \varepsilon_i \{m(X_i) - \mu^\star(X_i)\} \qquad \text{a mean zero term, like before;}$$
$$+ \frac{2}{n} \sum_{i=1}^n \{\mu^\star(X_i) - \mu(X_i)\}\{m(X_i) - \mu^\star(X_i)\} \quad \text{and something else.}$$

We can use our argument, ignoring the new term, if that term is always *non-negative.*

Why?

5

$$\ell(m) - \ell(\mu^\star) = \|m - \mu^\star\|_{L_2(P_n)}^2$$
$$- \frac{2}{n} \sum_{i=1}^{n} \varepsilon_i \left\{ m(X_i) - \mu^\star(X_i) \right\}$$
$$+ \frac{2}{n} \sum_{i=1}^{n} \{\mu^\star(X_i) - \mu(X_i)\}\{m(X_i) - \mu^\star(X_i)\}$$



We want to show that if distance from $m$ to $\mu^\star$ is big enough, it wins.

- In particular, it wins in the sense that the loss difference $\ell(m) - \ell(\mu^\star)$ is positive.
- That implies distance from $\hat{\mu}$ to $\mu^\star$ is smaller, as distance doesn't win in that case.

If this new term is non-negative, it helps distance win.

- If the MSE difference is positive when we ignore a non-negative term,
  then it's positive when we don't.

So we want to make sure this new term is non-negative. And we get to choose $\mu^\star$.

- It sounds like we choose what our estimator converges to when we analyze it.
- Obviously we don't really get to do that. It's not really a choice—it's a guess.
- If $\hat{\mu}$ converges to some curve $\mu^\star$, then it can't converge to anything else.

It's the best approximation to $\mu$ in the model.

$$\mu^\star = \operatorname*{argmin}_{m \in \mathcal{M}} \|m - \mu\|^2_{L_2(\mathrm{P_n})}.$$

$$\frac{2}{n} \sum_{i=1}^{n} \{\mu^\star(X_i) - \mu(X_i)\}\{m(X_i) - \mu^\star(X_i)\} = 2\langle \mu^\star - \mu, m - \mu^\star \rangle_{L_2(\mathrm{P_n})}$$

It's proportional to the dot product between two vectors: $\mu \rightarrow \mu^\star$ and $\mu^\star \rightarrow m$.



When the model $\mathcal{M}$ is convex, these vectors are always in the same direction.
That is, this dot product is non-negative for all $m \in \mathcal{M}$.

# Proof



**Claim.** For any convex set $\mathcal{M}$ in an inner product space, [1]

$$\mu^\star = \operatorname*{argmin}_{m \in \mathcal{M}} \|m - \mu\| \quad \text{satisfies}$$

$$\langle \mu^\star - \mu, \; m - \mu^\star \rangle \geq 0 \quad \text{for all curves} \quad m \in \mathcal{M}.$$

**Proof.** Let $m_\lambda = \lambda(m - \mu^\star) + \mu^\star$.

$$\|m_\lambda - \mu\|^2 = \langle \lambda(m - \mu^\star) + (\mu^\star - \mu), \; \lambda(m - \mu^\star) + (\mu^\star - \mu) \rangle$$
$$= \lambda^2 \|m - \mu^\star\|^2 + \|\mu^\star - \mu\|^2 + 2\lambda \langle m - \mu^\star, \; \mu^\star - \mu \rangle.$$

Because $m_\lambda \in \mathcal{M}$, it follows that this is at least as large as $\|\mu - \mu^\star\|^2$, so

$$0 \leq \lambda^2 \|m - \mu^\star\|^2 + 2\lambda \langle m - \mu^\star, \; \mu^\star - \mu \rangle$$

and therefore, dividing by $\lambda > 0$, that

$$0 \leq \lambda \|m - \mu^\star\|^2 + 2\langle m - \mu^\star, \; \mu^\star - \mu \rangle.$$

<u>Because this holds for arbitrarily small $\lambda > 0$, it must also hold for $\lambda = 0$.</u>

[1] An inner product space is a vector space with a norm $\|u\| = \sqrt{\langle u, u \rangle}$ induced by an inner product $\langle u, v \rangle$.

10

When $\mu^\star \in \mathcal{M}$ isn't the closest point to $\mu$,
these vectors can point in opposite directions.
That is, this dot product can be negative for some $m \in \mathcal{M}$.



The same thing can happen *for the closest point* in a non-convex model.

## Summary

When we use a convex model, the least squares estimator
$\hat{\mu}$ converges to the model's closest point to $\mu$.

- If $\mu$ is in the model, that's $\mu$.
- Otherwise, it's something else.

We can bound our estimator's distance to that closest point $\mu^\star$ just like
we've been bounding distance to $\mu$ when we assumed it was in the model.

$$\|\hat{\mu} - \mu^\star\|_{L_2(\mathrm{P_n})} < s \ \text{ w.p. } \ 1 - \delta \ \text{ if } \ s^2/2\sigma \geq \mathrm{w}(\mathcal{M}_s^\circ) + s\sqrt{2M_n/\delta n}.$$
$$\text{for} \quad \mathcal{M}_s^\circ = \big\{ m \in \mathcal{M} : \|m - \mu^\star\|_{L_2(\mathrm{P_n})} = s \big\} \quad \text{and} \quad M_n = 1 + 2\log(2n).$$

Let's get a feel for what that means by looking at some examples.

**Figure 1:** Increasing Curves.

**Figure 2:** Decreasing Curves.

**Figure 3:** Bounded Variation Curves. $\rho_{\mathrm{TV}} \leq 1$

**Figure 4:** Lipschitz Curves. $\rho_{\mathsf{Lip}} \leq 1$

**Figure 5:** Convex Curves.

**Figure 6:** Concave Curves.

# Non-Gaussian Noise

$$\ell(m) - \ell(\mu^\star) = \|m - \mu^\star\|^2_{L_2(\mathrm{P_n})} \qquad \text{squared distance}$$

$$- \frac{2}{n} \sum_{i=1}^n \varepsilon_i \left\{ m(X_i) - \mu^\star(X_i) \right\} \qquad \text{a mean zero term}$$

$$+ \frac{2}{n} \sum_{i=1}^n \{\mu^\star(X_i) - \mu(X_i)\}\{m(X_i) - \mu^\star(X_i)\} \qquad \text{a non-negative term.}$$



We can bound error using a corresponding *width*, no matter how noise is distributed.

$$\|\hat{\mu} - \mu^\star\|_{L_2(\mathrm{P_n})} < s + 2\sqrt{\frac{2M_{n2}^2}{\delta n}} \quad \text{w.p. } 1 - \delta \text{ for} \quad \frac{s^2}{2} \geq \mathrm{w}_\epsilon(\mathcal{M}_s)$$

$$\text{where} \quad \mathrm{w}_\epsilon(\mathcal{V}) = \mathrm{E} \max_{v \in \mathcal{V}} \langle \epsilon, v \rangle_{L_2(\mathrm{P_n})} \quad \text{and} \quad M_{n2}^2 = \mathrm{E} \max_{i \in 1 \ldots n} \varepsilon_i^2.$$

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

19

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

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

for $\alpha$ depending on $\varepsilon$ but not $\mathcal{M}$ or $s$.

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

1. **Symmetrization**. We'll substitute for $\epsilon_i$ a variant that's symmetric around zero.

$$\epsilon_i \to \epsilon_i - \epsilon_i' \quad \text{where} \quad \epsilon_i' \text{ is an independent copy of } \epsilon_i$$

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

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

$$\mathrm{w}_{\epsilon-\epsilon'}(\cdot) \leq \sqrt{2\pi} M_{n1} \times \mathrm{w}(\cdot) \quad \text{for} \quad M_{n1} = \mathrm{E} \max_{i \in 1\dots n} |\varepsilon_i|$$

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

- If you have a width comparison $\mathrm{w}_\epsilon \leq \alpha \mathrm{w}_\eta$ for some $\alpha \geq 1$.
- This implies a radius comparison $s_\epsilon \leq \alpha s_\eta$ for all convex models $\mathcal{M}$.

$$s_\epsilon = \alpha s_\nu \quad \text{satisfies} \quad \frac{s_\epsilon^2}{2} \geq \mathrm{w}_\varepsilon(\mathcal{M}_{s_\epsilon}) \quad \text{if} \quad \frac{s_\eta^2}{2} \geq \mathrm{w}_\eta(\mathcal{M}_{s_\eta}) \quad \text{for} \quad \text{convex} \quad \mathcal{M}$$

$$\text{and} \quad \mathrm{w}_\varepsilon \leq \alpha \mathrm{w}_\eta \quad \text{for} \quad \alpha \geq 1.$$

- *Interpretation.*
  The noise $\varepsilon$ makes regression at most '$\alpha$ times harder' than the noise $\eta$.
- *This is simplistic and 'lossy'.*
  For most models, our width comparison implies a better radius comparison.

## Proof: Width Comparisons imply Radius Comparisons

**Claim.** If $w_\varepsilon \leq \alpha\, w_\eta$ for $\alpha \geq 1$, then for any convex model $\mathcal{M}$, the critical radius using noise $\varepsilon$ is at most $\alpha$ times the critical radius using noise $\eta$, i.e.

$$\frac{(\alpha s)^2}{2} \geq w_\varepsilon(\mathcal{M}_{\alpha s}) \quad \text{if} \quad \frac{s^2}{2} \geq w_\eta(\mathcal{M}_s) \quad \text{and} \quad w_\varepsilon \leq \alpha\, w_\eta \quad \text{for} \quad \alpha \geq 1.$$

**Proof.** If $s^2/2 \geq w_\eta(\mathcal{M}_s)$, then

## Proof: Width Comparisons imply Radius Comparisons

**Claim.** If $\mathrm{w}_\varepsilon \leq \alpha \, \mathrm{w}_\eta$ for $\alpha \geq 1$, then for any convex model $\mathcal{M}$, the critical radius using noise $\varepsilon$ is at most $\alpha$ times the critical radius using noise $\eta$, i.e.

$$\frac{(\alpha s)^2}{2} \geq \mathrm{w}_\varepsilon(\mathcal{M}_{\alpha s}) \quad \text{if} \quad \frac{s^2}{2} \geq \mathrm{w}_\eta(\mathcal{M}_s) \quad \text{and} \quad \mathrm{w}_\varepsilon \leq \alpha \, \mathrm{w}_\eta \quad \text{for} \quad \alpha \geq 1.$$

**Proof.** If $s^2/2 \geq \mathrm{w}_\eta(\mathcal{M}_s)$, then

$$
\begin{aligned}
\alpha s/2 &\geq \alpha \, \mathrm{w}_\eta(\mathcal{M}_s)/s && \text{multiplying both sides by } \alpha/s \\
&\geq \alpha \, \mathrm{w}_\eta(\mathcal{M}_{\alpha s})/(\alpha s) && \text{using sublinearity of } f(s) = \mathrm{w}_\eta(\mathcal{M}_s) \\
&\geq \mathrm{w}_\varepsilon(\mathcal{M}_{\alpha s})/(\alpha s) && \text{using our premise } \alpha \, \mathrm{w}_\eta \geq \mathrm{w}_\varepsilon.
\end{aligned}
$$

Multiplying both sides by $\alpha s$, we get our claim.

**Where we are.** We have a bound that depends on the model $\mathcal{M}$ and the distribution of the noise $\varepsilon$ in a complex and entangled way.

$$\|\hat{\mu} - \mu^\star\|_{L_2(\mathrm{P_n})} < s_\epsilon + 2\sqrt{\frac{2M_{n2}^2}{\delta n}} \quad \text{w.p. } 1 - \delta \text{ for} \quad \frac{s_\epsilon^2}{2} \geq \mathrm{w}_\epsilon(\mathcal{M}_{s_\epsilon})$$

$$\text{where} \quad \mathrm{w}_\epsilon(\mathcal{V}) = \mathrm{E} \max_{v \in \mathcal{V}} \langle \epsilon, v \rangle_{L_2(\mathrm{P_n})} \quad \text{and} \quad M_{n2}^2 = \mathrm{E} \max_{i \in 1 \ldots n} \varepsilon_i^2.$$

**Where we're going.** We'll derive a bound that depends on the model $\mathcal{M}$ and the distribution of the noise $\varepsilon$ in simpler and disentangled way.

$$\leq \sqrt{2\pi} M_{n2} \left\{ \mathrm{w}(\mathcal{M}) + \sqrt{1/(\delta n)} \right\}$$

$$\|\hat{\mu} - \mu^\star\|_{L_2(\mathrm{P_n})} < \sqrt{2\pi} M_{n1} \, \mathrm{w}(\mathcal{M}) + 2\sqrt{\frac{2M_{n2}^2}{\delta n}} \leq \quad \text{w.p. } 1 - \delta \text{ for} \quad \frac{s^2}{2} \geq \mathrm{w}(\mathcal{M}_s)$$

$$\text{where} \quad \mathrm{w}(\mathcal{V}) = \mathrm{E} \max_{v \in \mathcal{V}} \langle g, v \rangle_{L_2(\mathrm{P_n})} \quad \text{and} \quad M_{n1} = \mathrm{E} \max_{i \in 1 \ldots n} |\varepsilon_i|.$$

23

## Non-Gaussian Noise

Example: Probabilistic Classification

Figure 7: classification noise → symmetrized classification noise → random-sign noise

Suppose we have independent *binary observations*.

$$Y_i = \begin{cases} 1 & \text{with conditional probability } \mu(X_i) \\ 0 & \text{otherwise} \end{cases}$$

$$= \mu(X_i) + \varepsilon_i \quad \text{for} \quad \varepsilon_i = \begin{cases} 1 - \mu(X_i) & \text{with conditional probability } \mu(X_i) \\ -\mu(X_i) & \text{with conditional probability } 1 - \mu(X_i) \end{cases}.$$

Note that this *classification noise* $\varepsilon_i$ has conditional mean zero.

$$\mathrm{E}[\varepsilon_i \mid X_i] = \mu(X_i)\{1 - \mu(X_i)\} + \{1 - \mu(X_i)\}\{-\mu(X_i)\} = 0.$$

**Figure 7:** classification noise → symmetrized classification noise → random-sign noise

What we need to bound is *classification-noise width*

$$\mathrm{w}_\epsilon(\mathcal{V}) = \frac{1}{n} \, \mathrm{E} \max_{v \in \mathcal{V}} \sum_{i=1}^{n} \varepsilon_i v_i.$$

We'll show it's no bigger than a version with *symmetrized noise*.

$$\varepsilon_i - \varepsilon_i' = \begin{cases} +1 & \text{when } \varepsilon_i = 1 - \mu(X_i), \ \varepsilon_i' = \mu(X_i) \\ -1 & \text{when } \varepsilon_i = \mu(X_i), \ \varepsilon_i' = 1 - \mu(X_i) \\ 0 & \text{when } \varepsilon_i = \varepsilon_i' \end{cases}$$

**Figure 7:** classification noise → symmetrized classification noise → random-sign noise

And we'll show that *this* is no bigger than a version with *random sign noise*

$$\mathrm{w}_\epsilon(\mathcal{V}) \leq \mathrm{w}_{\epsilon-\epsilon'}(\mathcal{V}) \leq \mathrm{w}_s(\mathcal{V}) \quad \text{where} \quad s_i = \pm 1 \text{ w.p. } 1/2.$$

The trick will be multiplying the symmetrized noise by a random sign.
It's already symmetric, so that doesn't change its distribution.

$$\varepsilon_i - \varepsilon_i' \overset{dist}{=} s_i(\varepsilon_i - \varepsilon_i')$$

Then we'll *contract out* the symmetrized noise, leaving the random sign. You'll see.

We bound our maximum in terms of one involving symmetric noise.
We'll work with an *independent copy* $\varepsilon'$ of our noise vector $\varepsilon$.

$$\mathrm{E}_\varepsilon \max_{v \in \mathcal{V}} \sum_{i=1}^{n} \varepsilon_i v_i \overset{(a)}{=} \mathrm{E}_\varepsilon \max_{v \in \mathcal{V}} \sum_{i=1}^{n} (\varepsilon_i - \mathrm{E}_{\varepsilon'} \, \varepsilon_i') v_i$$

$$\overset{(b)}{=} \mathrm{E}_\varepsilon \max_{v \in \mathcal{V}} \mathrm{E}_{\varepsilon'} \sum_{i=1}^{n} (\varepsilon_i - \varepsilon_i') v_i$$

$$\overset{(c)}{\leq} \mathrm{E}_\varepsilon \mathrm{E}_{\varepsilon'} \max_{v \in \mathcal{V}} \sum_{i=1}^{n} (\varepsilon_i - \varepsilon_i') v_i.$$

Why do these steps work?

(a) $\mathrm{E}_{\varepsilon'} \, \varepsilon_i' = 0.$

We bound our maximum in terms of one involving symmetric noise.
We'll work with an *independent copy* $\varepsilon'$ of our noise vector $\varepsilon$.

$$\mathrm{E}_{\varepsilon} \max_{v \in \mathcal{V}} \sum_{i=1}^{n} \varepsilon_i v_i \overset{(a)}{=} \mathrm{E}_{\varepsilon} \max_{v \in \mathcal{V}} \sum_{i=1}^{n} (\varepsilon_i - \mathrm{E}_{\varepsilon'} \, \varepsilon_i') v_i$$

$$\overset{(b)}{=} \mathrm{E}_{\varepsilon} \max_{v \in \mathcal{V}} \mathrm{E}_{\varepsilon'} \sum_{i=1}^{n} (\varepsilon_i - \varepsilon_i') v_i$$

$$\overset{(c)}{\leq} \mathrm{E}_{\varepsilon} \, \mathrm{E}_{\varepsilon'} \max_{v \in \mathcal{V}} \sum_{i=1}^{n} (\varepsilon_i - \varepsilon_i') v_i.$$

Why do these steps work?

(a) $\mathrm{E}_{\varepsilon'} \, \varepsilon_i' = 0$.

(b) Expectation is linear.

25

We bound our maximum in terms of one involving symmetric noise.
We'll work with an *independent copy* $\varepsilon'$ of our noise vector $\varepsilon$.

$$
\begin{aligned}
\mathrm{E}_\varepsilon \max_{v \in \mathcal{V}} \sum_{i=1}^n \varepsilon_i v_i
&\overset{(a)}{=} \mathrm{E}_\varepsilon \max_{v \in \mathcal{V}} \sum_{i=1}^n (\varepsilon_i - \mathrm{E}_{\varepsilon'}\, \varepsilon_i') v_i \\
&\overset{(b)}{=} \mathrm{E}_\varepsilon \max_{v \in \mathcal{V}} \mathrm{E}_{\varepsilon'} \sum_{i=1}^n (\varepsilon_i - \varepsilon_i') v_i \\
&\overset{(c)}{\leq} \mathrm{E}_\varepsilon \, \mathrm{E}_{\varepsilon'} \max_{v \in \mathcal{V}} \sum_{i=1}^n (\varepsilon_i - \varepsilon_i') v_i.
\end{aligned}
$$

Why do these steps work?

(a) $\mathrm{E}_{\varepsilon'}\, \varepsilon_i' = 0$.

(b) Expectation is linear.

(c) Maximizing the average gives us something smaller than averaging the maxima.

25

We bound our maximum in terms of one involving symmetric noise.
We'll work with an *independent copy* $\varepsilon'$ of our noise vector $\varepsilon$.

$$
\begin{aligned}
\mathrm{E}_\varepsilon \max_{v \in \mathcal{V}} \sum_{i=1}^n \varepsilon_i v_i
&\overset{(a)}{=} \mathrm{E}_\varepsilon \max_{v \in \mathcal{V}} \sum_{i=1}^n (\varepsilon_i - \mathrm{E}_{\varepsilon'}\, \varepsilon_i') v_i \\
&\overset{(b)}{=} \mathrm{E}_\varepsilon \max_{v \in \mathcal{V}}\, \mathrm{E}_{\varepsilon'} \sum_{i=1}^n (\varepsilon_i - \varepsilon_i') v_i \\
&\overset{(c)}{\le} \mathrm{E}_\varepsilon\, \mathrm{E}_{\varepsilon'} \max_{v \in \mathcal{V}} \sum_{i=1}^n (\varepsilon_i - \varepsilon_i') v_i.
\end{aligned}
$$

Why do these steps work?

(a) $\mathrm{E}_{\varepsilon'}\, \varepsilon_i' = 0$.

(b) Expectation is linear.

(c) Maximizing the average gives us something smaller than averaging the maxima.
  · In (c), we choose the maximizing $v \in \mathcal{V}$ *for each $\varepsilon'$.*

We bound our maximum in terms of one involving symmetric noise.
We'll work with an *independent copy* $\varepsilon'$ of our noise vector $\varepsilon$.

$$\mathrm{E}_\varepsilon \max_{v \in \mathcal{V}} \sum_{i=1}^{n} \varepsilon_i v_i \overset{(a)}{=} \mathrm{E}_\varepsilon \max_{v \in \mathcal{V}} \sum_{i=1}^{n} (\varepsilon_i - \mathrm{E}_{\varepsilon'} \varepsilon_i') v_i$$

$$\overset{(b)}{=} \mathrm{E}_\varepsilon \max_{v \in \mathcal{V}} \mathrm{E}_{\varepsilon'} \sum_{i=1}^{n} (\varepsilon_i - \varepsilon_i') v_i$$

$$\overset{(c)}{\leq} \mathrm{E}_\varepsilon \, \mathrm{E}_{\varepsilon'} \max_{v \in \mathcal{V}} \sum_{i=1}^{n} (\varepsilon_i - \varepsilon_i') v_i.$$

Why do these steps work?

(a) $\mathrm{E}_{\varepsilon'} \varepsilon_i' = 0$.
(b) Expectation is linear.
(c) Maximizing the average gives us something smaller than averaging the maxima.
  · In (c), we choose the maximizing $v \in \mathcal{V}$ *for each $\varepsilon'$*.
  · If we wanted to choose the same one each time, like we do in (b), we could.

We introduce independent random signs $s_i = \pm 1$ w.p. $1/2$, changing nothing.

$$\mathrm{E}_\varepsilon \, \mathrm{E}_{\varepsilon'} \max_{v \in \mathcal{V}} \sum_{i=1}^n (\varepsilon_i - \varepsilon'_i) v_i = \mathrm{E}_s \, \mathrm{E}_\varepsilon \, \mathrm{E}_{\varepsilon'} \max_{v \in \mathcal{V}} \sum_{i=1}^n s_i (\varepsilon_i - \varepsilon'_i) v_i.$$

Why does this change nothing?

We introduce independent random signs $s_i = \pm 1$ w.p. $1/2$, changing nothing.

$$\mathrm{E}_\varepsilon \, \mathrm{E}_{\varepsilon'} \max_{v \in \mathcal{V}} \sum_{i=1}^n (\varepsilon_i - \varepsilon_i') v_i = \mathrm{E}_s \, \mathrm{E}_\varepsilon \, \mathrm{E}_{\varepsilon'} \max_{v \in \mathcal{V}} \sum_{i=1}^n s_i (\varepsilon_i - \varepsilon_i') v_i.$$

Why does this change nothing?

· Because the inner mean ($\mathrm{E}_\varepsilon \, \mathrm{E}_{\varepsilon'}$) doesn't depend on the signs $s_i$.

· That's because $\varepsilon_i$ and $\varepsilon_i'$ have the same distribution.

· And this implies $(\varepsilon_i - \varepsilon_i')$ and $(\varepsilon_i' - \varepsilon) = -(\varepsilon_i - \varepsilon_i')$ do, too.
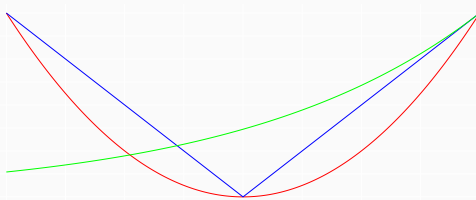
We swap the order of our averages and think about the inner average as a *function* of our vector of symmetric noise.

$$\mathrm{E}_s\,\mathrm{E}_\varepsilon\,\mathrm{E}_{\varepsilon'}\max_{v\in\mathcal{V}}\sum_{i=1}^n s_i(\varepsilon_i-\varepsilon_i')v_i = \mathrm{E}_\varepsilon\,\mathrm{E}_{\varepsilon'}\,\mathrm{E}_s\max_{v\in\mathcal{V}}\sum_{i=1}^n s_i(\varepsilon_i-\varepsilon_i')v_i$$

$$= \mathrm{E}_\varepsilon\,\mathrm{E}_{\varepsilon'}\,f(\varepsilon-\varepsilon') \quad \text{for} \quad f(u)=\mathrm{E}_s\max_{v\in\mathcal{V}}\sum_{i=1}^n s_i u_i v_i.$$

This function $f$ is convex.

What does that mean? These, for example, are all convex.



$$f\{(1-\lambda)a+\lambda b\} \leq (1-\lambda)f(a)+\lambda f(b) \quad \text{for} \quad \lambda\in[0,1]. \quad \text{That's Convexity}$$

We swap the order of our averages and think about the
inner average as a *function* of our vector of symmetric noise.

$$\mathrm{E}_s \, \mathrm{E}_\varepsilon \, \mathrm{E}_{\varepsilon'} \max_{v \in \mathcal{V}} \sum_{i=1}^n s_i(\varepsilon_i - \varepsilon_i')v_i = \mathrm{E}_\varepsilon \, \mathrm{E}_{\varepsilon'} \, \mathrm{E}_s \max_{v \in \mathcal{V}} \sum_{i=1}^n s_i(\varepsilon_i - \varepsilon_i')v_i$$

$$= \mathrm{E}_\varepsilon \, \mathrm{E}_{\varepsilon'} \, f(\varepsilon - \varepsilon') \quad \text{for} \quad f(u) = \mathrm{E}_s \max_{v \in \mathcal{V}} \sum_{i=1}^n s_i u_i v_i.$$

This function $f$ is convex.

How do we know? Maximizing each term is better than maximizing their sum.

$$\begin{aligned}
f\{(1-\lambda)a + \lambda b\} &= \mathrm{E}_s \max_{v \in \mathcal{V}} \left\{ (1-\lambda)\sum_{i=1}^n s_i a_i v_i + \lambda \sum_{i=1}^n s_i b_i v_i \right\} \\
&\leq \mathrm{E}_s \left\{ \max_{v \in \mathcal{V}} (1-\lambda)\sum_{i=1}^n s_i a_i v_i + \max_{v \in \mathcal{V}} \lambda \sum_{i=1}^n s_i b_i v_i \right\} \\
&= (1-\lambda) \, \mathrm{E}_s \max_{v \in \mathcal{V}} \sum_{i=1}^n s_i a_i v_i + \lambda \, \mathrm{E}_s \max_{v \in \mathcal{V}} \lambda \sum_{i=1}^n s_i b_i v_i \\
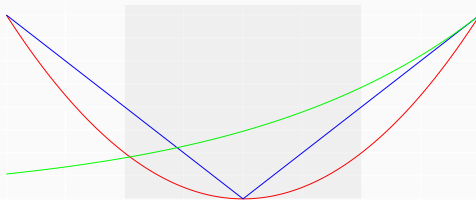&= (1-\lambda)f(a) + \lambda f(b).
\end{aligned}$$

We swap the order of our averages and think about the inner average as a *function* of our vector of symmetric noise.

$$\mathrm{E}_s \, \mathrm{E}_\varepsilon \, \mathrm{E}_{\varepsilon'} \max_{v \in \mathcal{V}} \sum_{i=1}^{n} s_i(\varepsilon_i - \varepsilon_i')v_i = \mathrm{E}_\varepsilon \, \mathrm{E}_{\varepsilon'} \, \mathrm{E}_s \max_{v \in \mathcal{V}} \sum_{i=1}^{n} s_i(\varepsilon_i - \varepsilon_i')v_i$$

$$= \mathrm{E}_\varepsilon \, \mathrm{E}_{\varepsilon'} \, f(\varepsilon - \varepsilon') \quad \text{for} \quad f(u) = \mathrm{E}_s \max_{v \in \mathcal{V}} \sum_{i=1}^{n} s_i u_i v_i.$$

This function $f$ is convex.

Why does this matter? The max of a convex function over a cube occurs at a corner.



What cube?

The vector of symmetric noise, $\varepsilon - \varepsilon'$, is in the *unit cube* $[-1, 1]^n$.

$$\varepsilon_i - \varepsilon_i' = \begin{cases} 0 & \text{when } \varepsilon_i = \varepsilon_i' \\ +1 & \text{when } \varepsilon_i = 1 - \mu(X_i), \ \varepsilon_i' = \mu(X_i) \\ -1 & \text{when } \varepsilon_i = \mu(X_i), \ \varepsilon_i' = 1 - \mu(X_i). \end{cases}$$

The average over this random vector is bounded by the maximum over the cube it's in.

$$\mathrm{E}_\varepsilon \, \mathrm{E}_{\varepsilon'} \, \mathrm{E}_s \max_{v \in \mathcal{V}} \sum_{i=1}^n s_i (\varepsilon_i - \varepsilon_i') v_i \leq \max_{u \in [-1,1]^n} \mathrm{E}_s \max_{v \in \mathcal{V}} \sum_{i=1}^n s_i u_i v_i$$

$$= \max_{u \in [-1,1]^n} f(u) \quad \text{max over the cube}$$

$$= \max_{u \in \{-1,1\}^n} f(u) \quad \text{max over its corners}$$

We characterize this maximum over corners. Remember what $f$ is.

$$\max_{u \in \{-1,1\}^n} f(u) = \max_{u \in \{-1,1\}^n} \mathrm{E}_s \max_{v \in \mathcal{V}} \sum_{i=1}^{n} s_i u_i v_i$$

$$= \mathrm{E}_s \max_{v \in \mathcal{V}} \sum_{i=1}^{n} s_i v_i.$$

Why?
Hint. What's the distribution of $s_i$? And $s_i u_i$ for $u_i \in \{-1,1\}$?

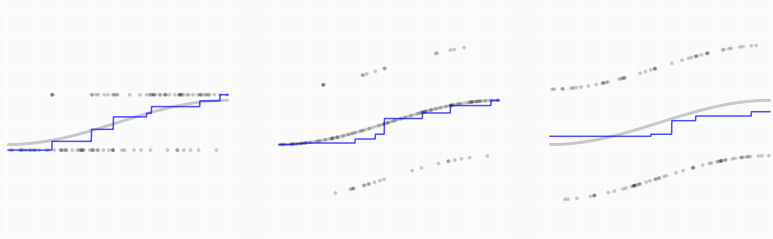We characterize this maximum over corners. Remember what $f$ is.

$$\max_{u \in \{-1,1\}^n} f(u) = \max_{u \in \{-1,1\}^n} \mathrm{E}_s \max_{v \in \mathcal{V}} \sum_{i=1}^{n} s_i u_i v_i$$

$$= \mathrm{E}_s \max_{v \in \mathcal{V}} \sum_{i=1}^{n} s_i v_i.$$

Why?

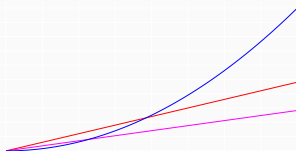Hint. What's the distribution of $s_i$? And $s_i u_i$ for $u_i \in \{-1,1\}$?

- For $u_i \in \{-1,1\}$, the distributions of $u_i$ and $s_i u_i$ are the same.
- So the distribution of the sum, and its maximum, are the same at every corner $u$.
- Including the vector of all ones $u = (1,1,\dots,1)$.

classification noise width $\leq$ symmetrized classification noise width $\leq$ random sign width
This means probabilistic classification is *easier* than regression
with random sign noise. Or, at least, that we get a better bound.

$$\frac{s^2}{2} \geq \mathrm{w}_s(\mathcal{M}_s) \quad \text{and} \quad \mathrm{w}_s(\mathcal{M}_s) \geq \mathrm{w}_\varepsilon(\mathcal{M}_s) \quad \implies \quad \frac{s^2}{2} \geq \mathrm{w}_\varepsilon(\mathcal{M}_s)$$

People call random sign width, or something like it, *Rademacher Complexity*.

$$\text{Rademacher Complexity}(\mathcal{V}) = \underset{v \in \mathcal{V}}{\text{E} \max} \langle s, v \rangle_{L_2(\text{P}_n)} \quad \text{for} \quad \text{i.i.d.} \quad s_i = \begin{cases} +1 & \text{w.p.} \ 1/2 \\ -1 & \text{w.p.} \ 1/2 \end{cases}$$

$$\underset{\text{w}_s(\mathcal{V})}{}$$

$$\text{or maybe} = \underset{v \in \mathcal{V}}{\text{E} \max} \left| \langle s, v \rangle_{L_2(\text{P}_n)} \right|$$

- This second definition is the same if $\mathcal{V}$ is symmetric, i.e. $v \in \mathcal{V} \implies -v \in \mathcal{V}$.
- Otherwise, it can be a little bigger.
    - At most $2\times$ bigger. Prove it!
    - Use the bound $\max a, b \leq a + b$ and the symmetry of $s$'s distribution.

31

# Non-Gaussian Noise

The General Case

**Figure 8:** real noise → symmetrized noise → scaled sign noise

## Symmetrization

$$\mathrm{w}_\varepsilon(\mathcal{V}) \leq \mathrm{w}_{s(\varepsilon-\varepsilon')}(\mathcal{V}) \leq 2\,\mathrm{w}_{s\varepsilon}(\mathcal{V})$$

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

(a) Replacing $\varepsilon_i$ with $s_i(\varepsilon_i - \varepsilon_i')$ is 'free'.
  - We stopped here in our 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.
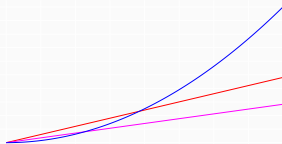
(b) Replacing $\varepsilon_i$ with $s_i\varepsilon_i$ increases width by at most $2\times$.

$$\mathrm{w}_\eta(\mathcal{V}) = \mathrm{w}_{s\eta}(\mathcal{V}) \le \mathrm{E}\|\eta\|_\infty \, \mathrm{w}_\eta(\mathcal{V}) \quad \text{if} \quad \eta \overset{dist}{=} -\eta.$$

$$\mathrm{E}_s \, \mathrm{E}_\eta \max_{v \in \mathcal{V}} \sum_{i=1}^n \eta_i s_i v_i \le \mathrm{E}_\eta \max_{\substack{u \in \mathbb{R}^n \\ |u_i| \le \|\eta\|_\infty}} \mathrm{E}_s \max_{v \in \mathcal{V}} \sum_{i=1}^n u_i s_i v_i$$

$$= \mathrm{E}_\eta \|\eta\|_\infty \max_{u \in [-1,1]^n} \mathrm{E}_s \max_{v \in \mathcal{V}} \sum_{i=1}^n u_i s_i v_i$$

$$= \mathrm{E}_\eta \|\eta\|_\infty \times \max_{u \in [-1,1]^n} \mathrm{E}_s \max_{v \in \mathcal{V}} \sum_{i=1}^n u_i s_i v_i$$

$$= \mathrm{E}_\eta \|\eta\|_\infty \times \mathrm{E}_s \max_{v \in \mathcal{V}} \sum_{i=1}^n s_i v_i$$

- We can 'contract out' any symmetrically distributed noise vector $\eta$ by ...
  1. multiplying in independent random signs $s_i$. Symmetry $\implies s_i \eta_i \overset{dist}{=} \eta_i$.
  2. maximizing over a cube containing $\eta$.
- We just have to use a big enough cube.
  - In our example, $\eta = \varepsilon - \varepsilon'$ was in the unit cube $[-1, 1]^n$ deterministically.
  - Generally, we maximize over a random cube $[-\|\eta\|_\infty, \|\eta\|_\infty]^n$.
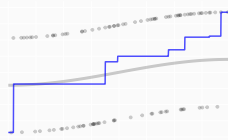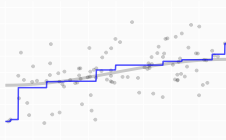  - And we can pull out the cube's radius $\|\eta\|_\infty$ as a multiplicative factor.

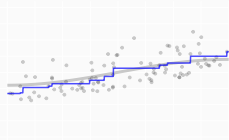$$\mathrm{w}_\varepsilon(\mathcal{V}) \le \mathrm{E}\|\varepsilon_i - \varepsilon_i'\|_\infty \, \mathrm{w}_s(\mathcal{V}) \le 2\,\mathrm{E}\|\varepsilon_i\|_\infty \, \mathrm{w}_s(\mathcal{V})$$

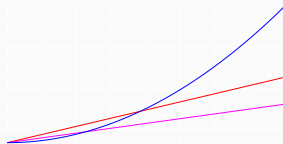Regression with arbitrary independent noise, i.e.

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

is no harder than with scaled-up random sign noise, i.e.

$$Y_i = \mu(X_i) + Ms_i \quad \text{for} \quad M = \mathrm{E}\|\varepsilon_i - \varepsilon_i'\|_\infty \quad \text{and} \quad s_i = \begin{cases} +1 & \text{w.p. } 1/2 \\ -1 & \text{w.p. } 1/2 \end{cases}.$$

$$\mathrm{w}_\varepsilon(\mathcal{V}) \le \mathrm{E}\|\varepsilon_i\|_\infty \, \mathrm{w}_s(\mathcal{V})$$

Regression with arbitrary independent *symmetric* noise, i.e.

$$Y_i = \mu(X_i) + \varepsilon_i \quad \text{where} \quad \varepsilon_1 \ldots \varepsilon_n \text{ are independent with } \varepsilon_i \stackrel{dist}{=} -\varepsilon_i,$$

is no harder than with scaled-up random sign noise, i.e.

$$Y_i = \mu(X_i) + M s_i \quad \text{for}^2 \quad M = \mathrm{E}\|\varepsilon_i\|_\infty \quad \text{and} \quad s_i = \begin{cases} +1 & \text{w.p. } 1/2 \\ -1 & \text{w.p. } 1/2 \end{cases}.$$
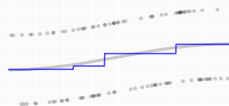


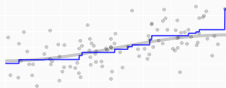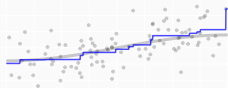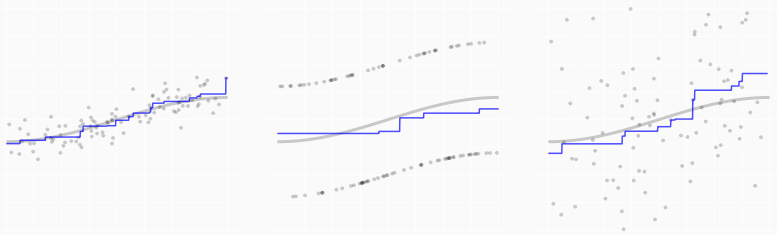**Figure 9:** real noise $\rightarrow$ symmetrized noise $\rightarrow$ scaled sign noise

---

[2] $M = \mathrm{E}\|\varepsilon_i\|_\infty \le 2\sigma\sqrt{2\log(2n)}$ for $\varepsilon_i \sim N(0, \sigma^2)$. See Appendix B of the Gaussian Width Homework.

## Non-Gaussian Noise

Comparison to the Gaussian Case

- So far, we've bounded arbitrary-noise width in terms of random-sign width.
- But often, it's easier to understand gaussian width. That's good enough.[3]

$$\frac{1}{2\sqrt{\log(2n)}} \, \mathrm{w}_g(\mathcal{V}) \leq \mathrm{w}_s(\mathcal{V}) \leq \sqrt{\frac{\pi}{2}} \, \mathrm{w}_g(\mathcal{V})$$

$\underset{\approx.2 \text{ for } n=100}{\phantom{x}}$ $\qquad\qquad\qquad\qquad$ $\underset{\approx 1.25}{\phantom{x}}$

- We just saw it can't be that much bigger than random-sign width.
- And we can show it's at least 4/5 as big.

$$\mathrm{E}\max_{v\in\mathcal{V}}\sum_{i=1}^{n} g_i v_i = \mathrm{E}_s\,\mathrm{E}_g\max_{v\in\mathcal{V}}\sum_{i=1}^{n}|g_i|s_i v_i \geq \mathrm{E}_s\max_{v\in\mathcal{V}}\sum_{i=1}^{n}\underset{=\sqrt{\frac{2}{\pi}}}{\mathrm{E}_g|g_i|}s_i v_i.$$

<hr>

[3] We can show $.125\,\mathrm{w}_g(\mathcal{V}) \leq \mathrm{w}_s(\mathcal{V}) \leq 1.25\,\mathrm{w}_g(\mathcal{V})$ for $n \leq 10$ trillion by bounding $\mathrm{E}\|g\|_\infty$ more carefully.
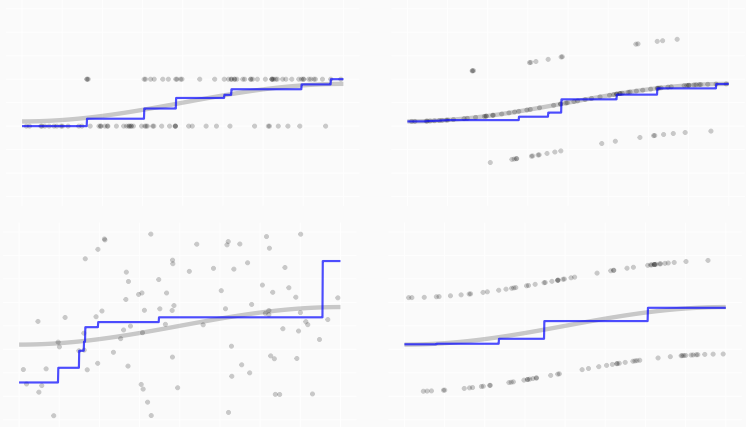
**Figure 10:** real noise $\rightarrow$ symmetrized noise $\downarrow$ scaled sign noise $\leftarrow$ scaled gaussian noise

$$\mathrm{w}_\varepsilon(\mathcal{V}) \leq \underset{\leq 2\,\mathrm{E}\|\varepsilon\|_\infty}{\mathrm{w}_{\varepsilon-\varepsilon'}(\mathcal{V})} \leq \mathrm{E}\|\varepsilon - \varepsilon'\|_\infty \quad \mathrm{w}_s(\mathcal{V}) \leq \underset{\leq\sqrt{2\pi}\approx 2.5\times\mathrm{E}\|\varepsilon\|_\infty}{\sqrt{\frac{\pi}{2}}\,\mathrm{E}\|\varepsilon - \varepsilon'\|_\infty} \quad \mathrm{w}_g(\mathcal{V})$$

**Figure 11:** real noise $\rightarrow$ scaled gaussian noise

For any noise vector $\varepsilon$ with independent components $\varepsilon_i$,

$$\mathrm{w}_\varepsilon(\mathcal{V}) \leq 2\,\mathrm{E}\|\varepsilon\|_\infty \cdot \mathrm{w}_s(\mathcal{V}) \leq \sqrt{2\pi}\,\mathrm{E}\|\varepsilon\|_\infty \cdot \mathrm{w}_g(\mathcal{V}).$$

- We can bound the width $\mathrm{w}_\varepsilon$ in terms of
  1. random-sign width
  2. the maximum absolute value of $\varepsilon$'s components.
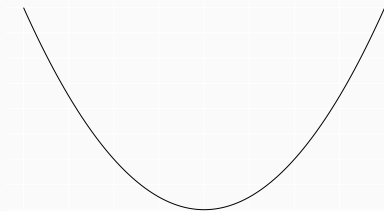- And we can bound random-sign width in terms of gaussian width.

This means we don't have to bound a million different kinds of widths for each model.
We can bound random-sign width or gaussian width. Whichever is easier.

# Background: Convex Functions Are Maximized At Extreme Points

## Definition

A function $f$ is convex if *secants* lie above the curve.

$$f\{(1 - \lambda)a + \lambda b\} \leq (1 - \lambda)f(a) + \lambda f(b) \quad \text{for} \quad \lambda \in [0, 1]$$



We can give this a *probabilistic interpretation* for a random variable $Z_\lambda$.

$$f(\mathrm{E}\, Z_\lambda) \leq \mathrm{E}\, f(Z_\lambda) \quad \text{where} \quad Z_\lambda =$$

## Definition

A function $f$ is convex if *secants* lie above the curve.

$$f\{(1 - \lambda)a + \lambda b\} \leq (1 - \lambda)f(a) + \lambda f(b) \quad \text{for} \quad \lambda \in [0, 1]$$
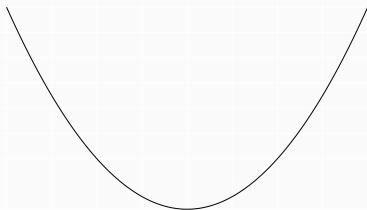


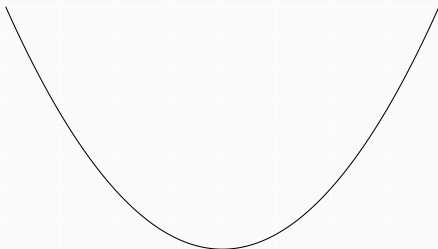We can give this a *probabilistic interpretation* for a random variable $Z_\lambda$.

$$f(\mathrm{E}\, Z_\lambda) \leq \mathrm{E}\, f(Z_\lambda) \quad \text{where} \quad Z_\lambda = \begin{cases} a & \text{w.p. } 1 - \lambda \\ b & \text{w.p. } \lambda \end{cases}$$

## Jensen's Inequality

In fact, this is true all random variables $Z$.
If $f$ is convex, its mean value exceeds its value at the mean.

$$f(\operatorname{E} Z) \leq \operatorname{E} f(Z)$$

That's called Jensen's Inequality.



You can prove it for discrete random variables via induction.

## Jensen's Inequality Proof

#### Base case.
It's true for random variables taking on 2 values.

$$f(\lambda_1 z_1 + \lambda_2 z_2) \leq \lambda_1 f(z_1) + \lambda_2 f(z_2) \quad \text{if} \quad \lambda_1, \lambda_2 \geq 0 \quad \text{satisfy} \quad \lambda_1 + \lambda_2 = 1$$
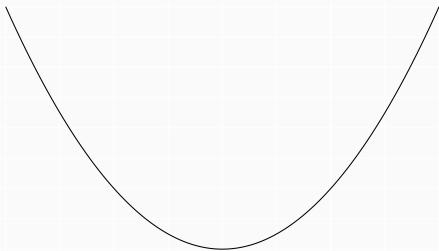
#### Inductive Step.
We'll show that if it's true for random variables taking on
$n-1$ values, then it's also true for ones taking on $n$ values.

$$\begin{aligned}
f\left\{\sum_{i=1}^{n} \lambda_i z_i\right\} &= f\left\{(1-\lambda_n)\left(\sum_{i=1}^{n-1} \frac{\lambda_i}{1-\lambda_n} z_i\right) + \lambda_n z_n\right\} \\
&\leq (1-\lambda_n) f\left(\sum_{i=1}^{n-1} \frac{\lambda_i}{1-\lambda_n} z_i\right) + \lambda_n f(z_n) \\
&\leq (1-\lambda_n) \sum_{i=1}^{n-1} \frac{\lambda_i}{1-\lambda_n} f(z_i) + \lambda_n f(z_n) \\
&= \sum_{i=1}^{n-1} \lambda_i f(z_i) + \lambda_n f(z_n)
\end{aligned}$$

## Maxima of Convex Functions

Convex functions have no local maxima.



That means the maximum of a convex function over an interval occurs at an endpoint.

**Proof.**

$$\max_{x\in[a,b]} f(x) = \max_{\lambda\in[0,1]} f\{\underbrace{(1-\lambda)a + \lambda b}_{f(\mathrm{E}\,Z_\lambda)}\} \leq \max_{\lambda\in[0,1]} \underbrace{(1-\lambda)f(a) + \lambda f(b)}_{\mathrm{E}\,f(Z_\lambda)} = \max\{f(a), f(b)\}$$

This is essentially true in higher dimensions as well.
We just need the right generalizations of *interval* and its *endpoints*.

The natural generalizations a *convex polytope* and its *extreme points*.

### Definitions.

A **convex polytope** is the set of all weighted averages of some set of vectors $u_1 \ldots u_K$.

$$\mathcal{U} = \left\{ \sum_i \lambda_i u_i \ : \ \lambda \in \Lambda \right\} \quad \text{where} \quad \Lambda = \left\{ \lambda \ : \ \lambda_i \geq 0 \ \text{for all } i \ \text{and} \ \sum_i \lambda_i = 1 \right\}$$

Its **extreme points** are the subset of these vectors that are not redundant.
That is, they're the ones we cannot write as weighted averages of the others.

### Examples.

- A triangle is the set of weighted averages of its three vertices, its extreme points.
- A square is the set of weighted averages of its four vertices, its extreme points.
- A cube in $\mathbb{R}^n$ is the set of weighted averages of its $2^n$ vertices, its extreme points.

The maximum of a convex function over a convex polytope
occurs at an extreme point.

### Proof.
It's more-or-less the same as the one-dimensional case.
We apply Jensen's inequality to a *random extreme point* $Z_\lambda$.

$$\max_{u \in \mathcal{U}} f(u) = \max_{\lambda \in \Lambda} \underbrace{f\left(\sum_i \lambda_i u_i\right)}_{f(\mathrm{E}\, Z_\lambda)} \leq \max_{\lambda \in \Lambda} \underbrace{\sum_i \lambda_i f(u_i)}_{\mathrm{E}\, f(Z_\lambda)} \leq \max_i f(u_i)$$

where

$$Z_\lambda = \begin{cases} u_1 & \text{w.p. } \lambda_1 \\ \vdots & \vdots \\ u_K & \text{w.p. } \lambda_K \end{cases}$$