

# Machine Learning Theory

## Lecture 5: Least Squares with Misspecification and Non-Gaussian Noise

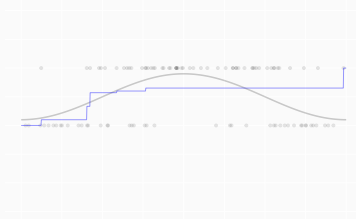
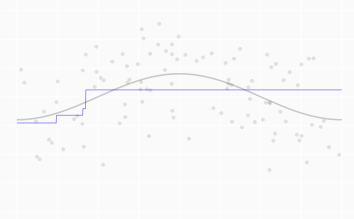
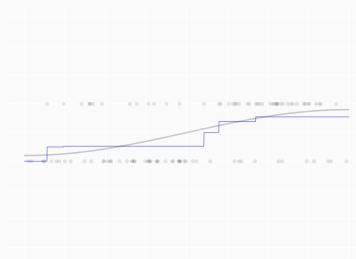
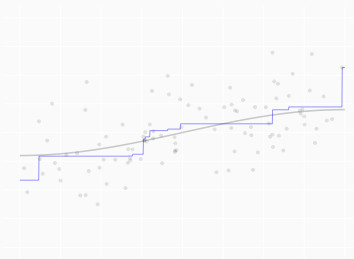
---

David A. Hirshberg

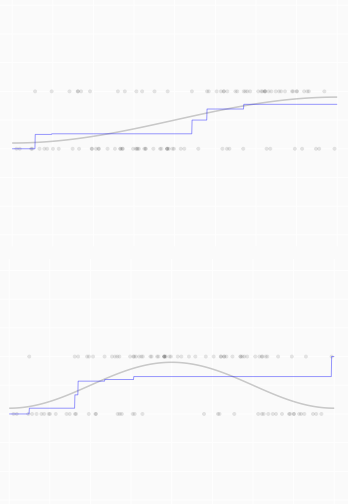
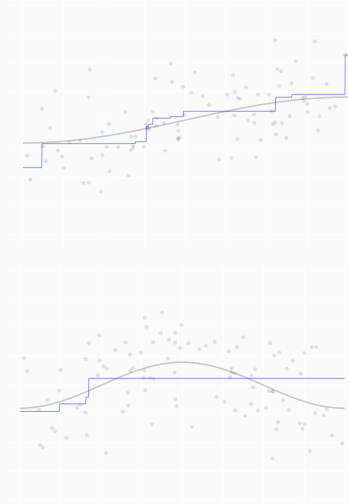
May 24, 2024

Emory University

# When Does Our Theory Apply?

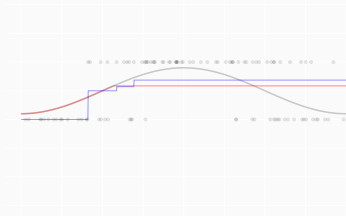
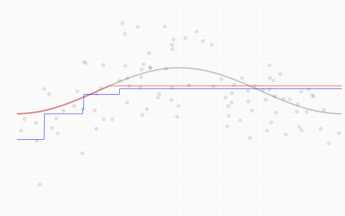
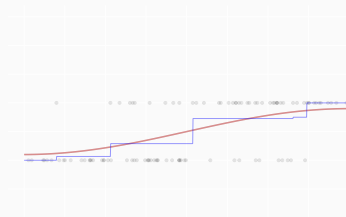
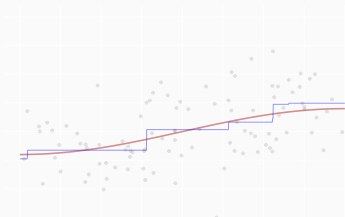


# When Does Our Theory Apply?



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

# Today, We Fix That



# Today, We Fix That

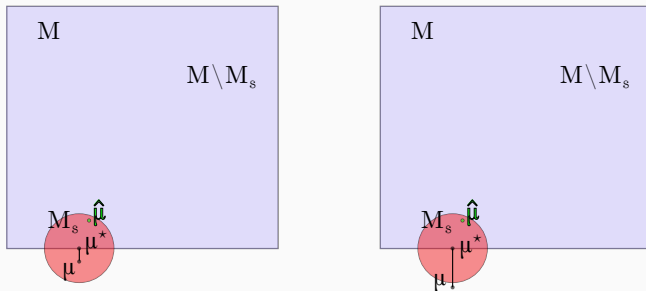


- With misspecification, we estimate the model's **best approximation** to  $\mu$ .
- Non-normality doesn't really matter much. We'll look at how it affects our bound.

## Misspecification

---

## What happens when $\mu$ isn't in the model?



- Our error in estimating  $\mu$  is bounded by a sum of two terms.
  - The critical radius  $s$ , i.e., the one satisfying  $s^2 \geq 2\sigma c_\delta w(\mathcal{M}_s)$ .
  - 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  $\mu^*$  does not depend on its distance to  $\mu$ .

# The Argument

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

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

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

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

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

$$\begin{aligned} \ell(m) - \ell(\mu^*) &= \|m - \mu^*\|_{L_2(\mathbf{P}_n)}^2 && \text{squared distance, like before;} \\ &- \frac{2}{n} \sum_{i=1}^n \varepsilon_i \{m(X_i) - \mu^*(X_i)\} && \text{a mean zero term, like before;} \\ &+ \frac{2}{n} \sum_{i=1}^n \{\mu^*(X_i) - \mu(X_i)\} \{m(X_i) - \mu^*(X_i)\} && \text{and something else.} \end{aligned}$$

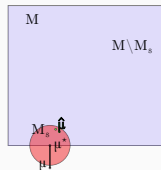
We can use our argument, ignoring the new term, if that term is always *non-negative*.

Why?



Why.

$$\begin{aligned}\ell(m) - \ell(\mu^*) &= \|m - \mu^*\|_{L_2(P_n)}^2 \\ &\quad - \frac{2}{n} \sum_{i=1}^n \varepsilon_i \{m(X_i) - \mu^*(X_i)\} \\ &\quad + \frac{2}{n} \sum_{i=1}^n \{\mu^*(X_i) - \mu(X_i)\} \{m(X_i) - \mu^*(X_i)\}\end{aligned}$$



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

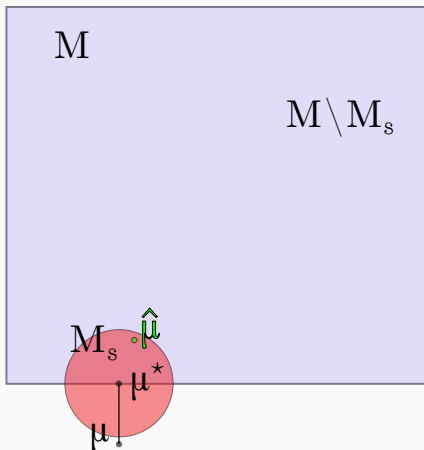
- In particular, it wins in the sense that the loss difference  $\ell(m) - \ell(\mu^*)$  is positive.
- That implies distance from  $\hat{\mu}$  to  $\mu^*$  is smaller, as distance doesn't win in that case.

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.

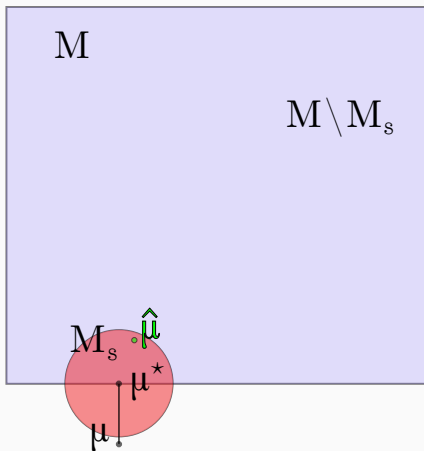
So we want to make sure this new term is non-negative. And we get to choose  $\mu^*$ .

## This sounds weird



- 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  $\mu^*$ , then it can't converge to anything else.

## The right choice



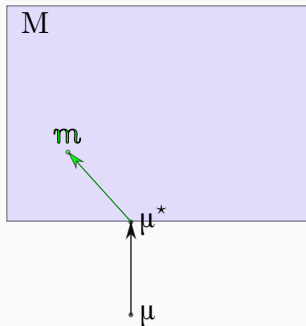
It's the best approximation to  $\mu$  in the model.

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

With this choice, the new term is always non-negative

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

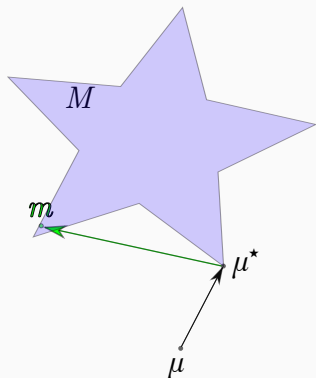
It's proportional to the dot product between two vectors:  $\mu \rightarrow \mu^*$  and  $\mu^* \rightarrow m$ .



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 for Homework!

## That's not true for other choices

When  $\mu^* \in \mathcal{M}$  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.

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  $\mu^*$  just like we've been bounding distance to  $\mu$  when we assumed it was in the model.

$$\|\hat{\mu} - \mu^*\|_{L_2(\mathbf{P}_n)} < s \text{ w.p. } 1 - \delta \text{ if } s^2 \geq 2\sigma c_\delta w(\mathcal{M}_s)$$

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

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

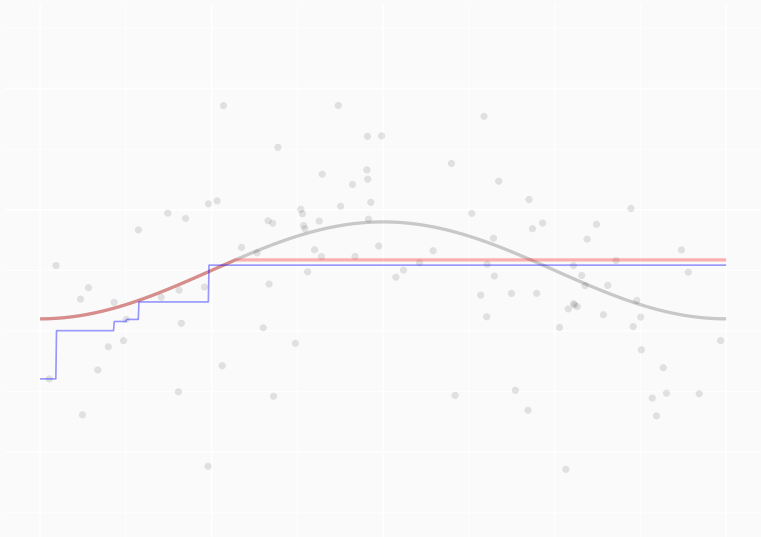
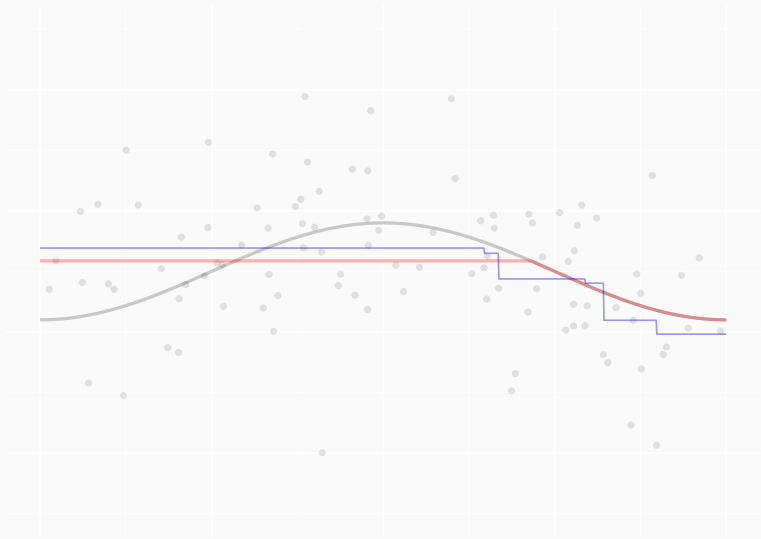


Figure 1: Increasing Curves.



**Figure 2:** Decreasing Curves.



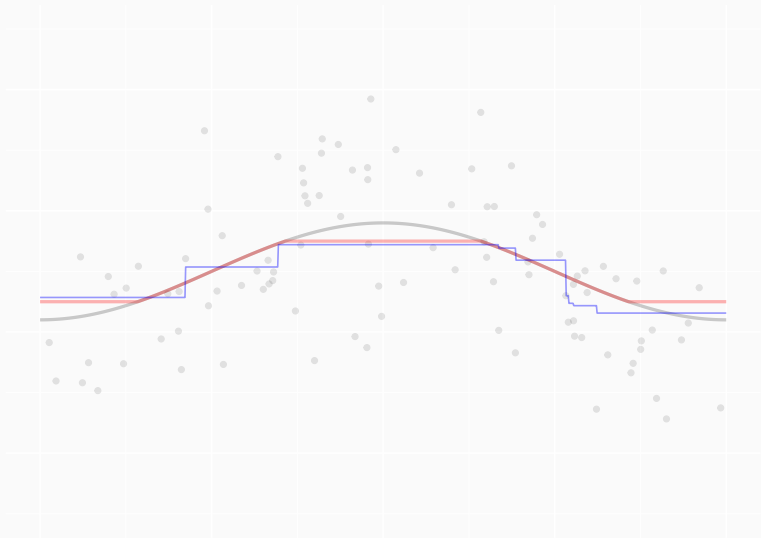


Figure 3: Bounded Variation Curves.  $\rho_{TV} \leq 1$

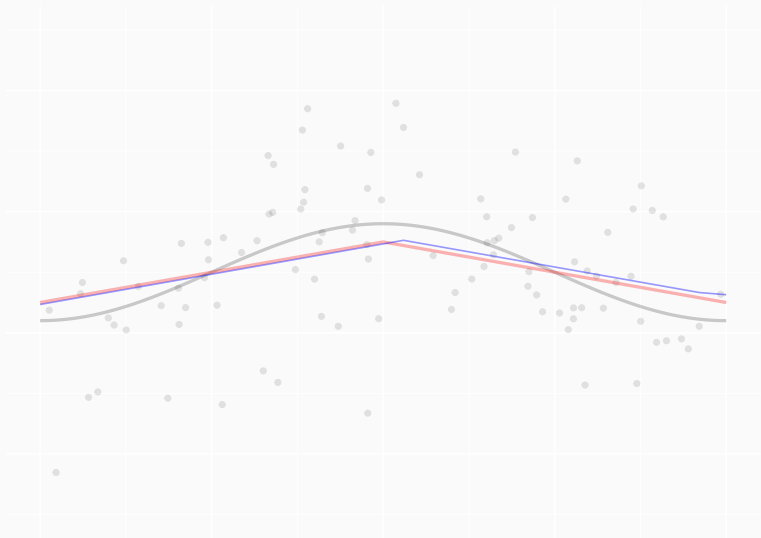


Figure 4: Lipschitz Curves.  $\rho_{\text{Lip}} \leq 1$

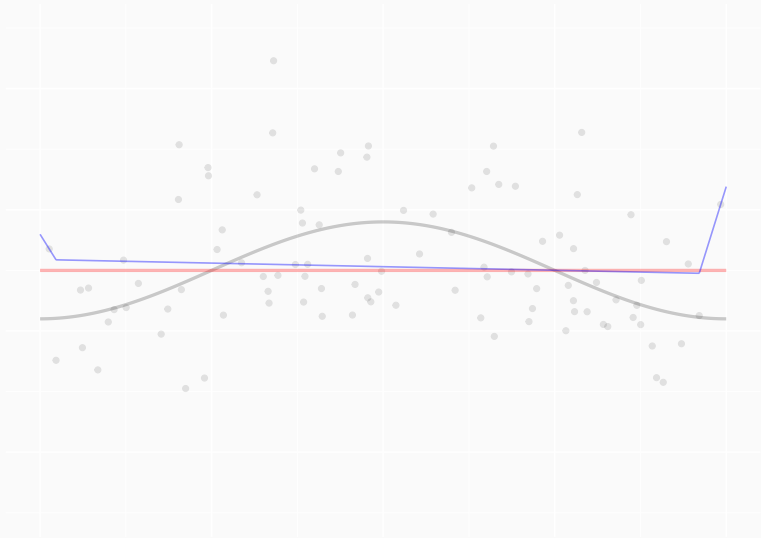
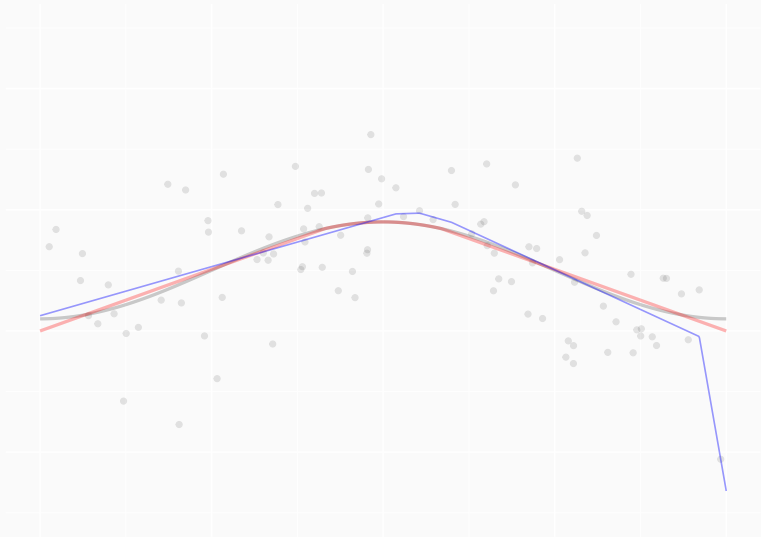


Figure 5: Convex Curves.



**Figure 6:** Concave Curves.

## Non-Gaussian Noise

---

$$\begin{aligned}\ell(m) - \ell(\mu^*) &= \|m - \mu^*\|_{L_2(P_n)}^2 && \text{squared distance} \\ &- \frac{2}{n} \sum_{i=1}^n \epsilon_i \{m(X_i) - \mu^*(X_i)\} && \text{a mean zero term} \\ &+ \frac{2}{n} \sum_{i=1}^n \{\mu^*(X_i) - \mu(X_i)\} \{m(X_i) - \mu^*(X_i)\} && \text{a non-negative term.}\end{aligned}$$

- No matter what noise vector  $\epsilon$  we have, we have an error bound determined by the *width* associated with it.

$$\|\hat{\mu} - \mu^*\|_{L_2(P_n)} < s \quad \text{w.p. } 1 - \delta \quad \text{for } s^2 \geq 2c_\delta w_\epsilon(\mathcal{M}_s^\circ)$$

where  $w_\epsilon(\mathcal{V}) = \mathbb{E} \max_{v \in \mathcal{V}} \langle \epsilon, v \rangle_{L_2(P_n)}.$

- To take advantage of our gaussian width calculations, we'll bound this width in terms of gaussian width.
- At the heart of it are two ideas called *symmetrization* and *contraction*.
- We'll substitute for  $\epsilon_i$  a variant that's symmetric around zero.

$$\epsilon_i \rightarrow \epsilon_i - \epsilon'_i \quad \text{where } \epsilon'_i \text{ is an independent copy of } \epsilon_i$$

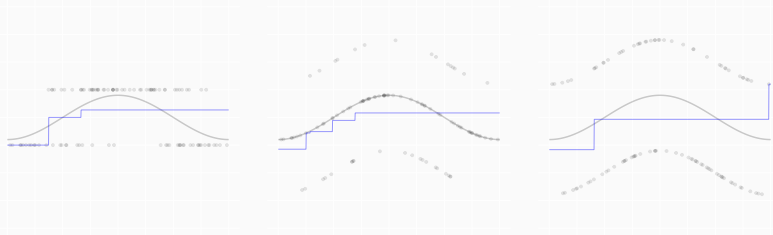
- We'll show this substitution doesn't change much. That's symmetrization.
- Substituting a gaussian often doesn't change much either. That's contraction.

Non-Gaussian Noise

---

Probabilistic Classification

# The Setting



Suppose we have independent *binary observations*.

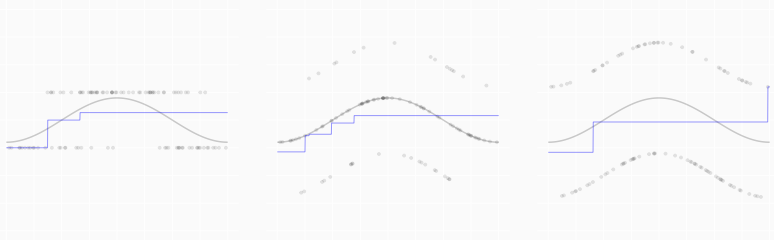
$$Y_i = \begin{cases} 1 & \text{with conditional probability } \mu(X_i) \\ 0 & \text{otherwise} \end{cases}$$
$$= \mu(X_i) + \varepsilon_i \quad \text{for} \quad \varepsilon_i = \begin{cases} 1 - \mu(X_i) & \text{with conditional probability } \mu(X_i) \\ -\mu(X_i) & \text{with conditional probability } 1 - \mu(X_i) \end{cases}.$$

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

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



# The Setting



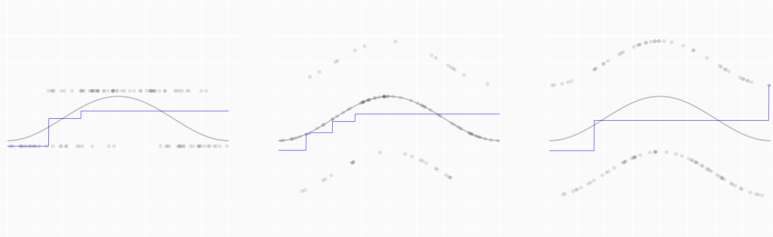
What we need to bound is *classification-noise width*

$$w_{\epsilon}(\mathcal{V}) = \frac{1}{n} \mathbb{E} \max_{v \in \mathcal{V}} \sum_{i=1}^n \epsilon_i v_i.$$

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

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

# The Setting



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

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

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

$$\epsilon_i - \epsilon'_i \stackrel{\text{dist}}{=} s_i(\epsilon_i - \epsilon'_i)$$

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

## Step 1

We bound our maximum in terms of one involving symmetric noise.

- We'll work with a conditionally independent copy  $\varepsilon'_1 \dots \varepsilon'_n$  of our noise.
- It has the same distribution as  $\varepsilon_1 \dots \varepsilon_n$  conditional on  $X_1 \dots X_n$ .

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

Why do these steps work?

(a)  $\mathbb{E}_{\varepsilon'} \varepsilon'_i = 0$ .

## Step 1

We bound our maximum in terms of one involving symmetric noise.

- We'll work with a conditionally independent copy  $\varepsilon'_1 \dots \varepsilon'_n$  of our noise.
- It has the same distribution as  $\varepsilon_1 \dots \varepsilon_n$  conditional on  $X_1 \dots X_n$ .

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

Why do these steps work?

- (a)  $\mathbb{E}_{\varepsilon'} \varepsilon'_i = 0$ .
- (b) Expectation is linear.

## Step 1

We bound our maximum in terms of one involving symmetric noise.

- We'll work with a conditionally independent copy  $\varepsilon'_1 \dots \varepsilon'_n$  of our noise.
- It has the same distribution as  $\varepsilon_1 \dots \varepsilon_n$  conditional on  $X_1 \dots X_n$ .

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

Why do these steps work?

- (a)  $\mathbb{E}_{\varepsilon'} \varepsilon'_i = 0$ .
- (b) Expectation is linear.
- (c) Maximizing the average gives us something smaller than averaging the maxima.

## Step 1

We bound our maximum in terms of one involving symmetric noise.

- We'll work with a conditionally independent copy  $\varepsilon'_1 \dots \varepsilon'_n$  of our noise.
- It has the same distribution as  $\varepsilon_1 \dots \varepsilon_n$  conditional on  $X_1 \dots X_n$ .

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

Why do these steps work?

- (a)  $\mathbb{E}_{\varepsilon'} \varepsilon'_i = 0$ .
- (b) Expectation is linear.
- (c) Maximizing the average gives us something smaller than averaging the maxima.
  - In the second, we choose the maximizing  $v \in \mathcal{V}$  for each  $\varepsilon'$ .

## Step 1

We bound our maximum in terms of one involving symmetric noise.

- We'll work with a conditionally independent copy  $\varepsilon'_1 \dots \varepsilon'_n$  of our noise.
- It has the same distribution as  $\varepsilon_1 \dots \varepsilon_n$  conditional on  $X_1 \dots X_n$ .

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

Why do these steps work?

- (a)  $\mathbb{E}_{\varepsilon'} \varepsilon'_i = 0$ .
- (b) Expectation is linear.
- (c) Maximizing the average gives us something smaller than averaging the maxima.
  - In the second, 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 the first, we could.

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

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

Why does this change nothing?



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

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

Why does this change nothing?

- Because the inner mean  $(\mathbb{E}_{\varepsilon} \mathbb{E}_{\varepsilon'})$  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  $-(\varepsilon_i - \varepsilon'_i)$  do, too.

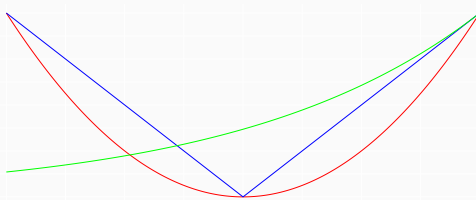
## Step 3

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

$$\begin{aligned} \mathbb{E}_s \mathbb{E}_\varepsilon \mathbb{E}_{\varepsilon'} \max_{v \in \mathcal{V}} \sum_{i=1}^n s_i (\varepsilon_i - \varepsilon'_i) v_i &= \mathbb{E}_\varepsilon \mathbb{E}_{\varepsilon'} \mathbb{E}_s \max_{v \in \mathcal{V}} \sum_{i=1}^n s_i (\varepsilon_i - \varepsilon'_i) v_i \\ &= \mathbb{E}_\varepsilon \mathbb{E}_{\varepsilon'} f(\varepsilon - \varepsilon') \quad \text{for} \quad f(u) = \mathbb{E}_s \max_{v \in \mathcal{V}} \sum_{i=1}^n s_i u_i v_i. \end{aligned}$$

This function  $f$  is convex.

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



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

## Step 3

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

$$\begin{aligned} \mathbb{E}_s \mathbb{E}_\varepsilon \mathbb{E}_{\varepsilon'} \max_{v \in \mathcal{V}} \sum_{i=1}^n s_i (\varepsilon_i - \varepsilon'_i) v_i &= \mathbb{E}_\varepsilon \mathbb{E}_{\varepsilon'} \mathbb{E}_s \max_{v \in \mathcal{V}} \sum_{i=1}^n s_i (\varepsilon_i - \varepsilon'_i) v_i \\ &= \mathbb{E}_\varepsilon \mathbb{E}_{\varepsilon'} f(\varepsilon - \varepsilon') \quad \text{for} \quad f(u) = \mathbb{E}_s \max_{v \in \mathcal{V}} \sum_{i=1}^n s_i u_i v_i. \end{aligned}$$

This function  $f$  is convex.

How do we know? Maximizing two things separately is better than maximizing their sum.

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

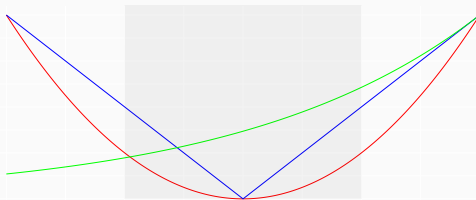
## Step 3

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

$$\begin{aligned} \mathbb{E}_s \mathbb{E}_\varepsilon \mathbb{E}_{\varepsilon'} \max_{v \in \mathcal{V}} \sum_{i=1}^n s_i (\varepsilon_i - \varepsilon'_i) v_i &= \mathbb{E}_\varepsilon \mathbb{E}_{\varepsilon'} \mathbb{E}_s \max_{v \in \mathcal{V}} \sum_{i=1}^n s_i (\varepsilon_i - \varepsilon'_i) v_i \\ &= \mathbb{E}_\varepsilon \mathbb{E}_{\varepsilon'} f(\varepsilon - \varepsilon') \quad \text{for} \quad f(u) = \mathbb{E}_s \max_{v \in \mathcal{V}} \sum_{i=1}^n s_i u_i v_i. \end{aligned}$$

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.

$$\begin{aligned} \mathbb{E}_{\varepsilon} \mathbb{E}_{\varepsilon'} \mathbb{E}_s \max_{v \in \mathcal{V}} \sum_{i=1}^n s_i (\varepsilon_i - \varepsilon'_i) v_i &\leq \max_{u \in [-1, 1]^n} \mathbb{E}_s \max_{v \in \mathcal{V}} \sum_{i=1}^n s_i u_i v_i \\ &= \max_{u \in [-1, 1]^n} f(u) \quad \text{max over the cube} \\ &= \max_{u \in \{-1, 1\}^n} f(u) \quad \text{max over its corners} \end{aligned}$$

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

$$\begin{aligned}\max_{u \in \{-1,1\}^n} f(u) &= \max_{u \in \{-1,1\}^n} \mathbb{E}_s \max_{v \in \mathcal{V}} \sum_{i=1}^n s_i u_i v_i \\ &= \mathbb{E}_s \max_{v \in \mathcal{V}} \sum_{i=1}^n s_i v_i.\end{aligned}$$

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.

$$\begin{aligned}\max_{u \in \{-1,1\}^n} f(u) &= \max_{u \in \{-1,1\}^n} \mathbb{E}_s \max_{v \in \mathcal{V}} \sum_{i=1}^n s_i u_i v_i \\ &= \mathbb{E}_s \max_{v \in \mathcal{V}} \sum_{i=1}^n s_i v_i.\end{aligned}$$

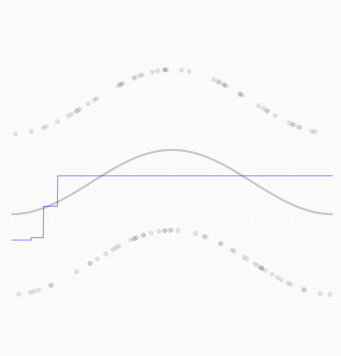
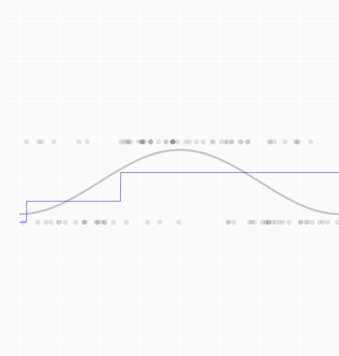
Why?

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

They're the same.

So the distribution of the sum—and its maximum—is the same at every corner  $u$ .  
Including the vector of all ones.

# Summary



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.

$$s^2 \geq 2c_\delta w_s(\mathcal{M}_s^\circ) \quad \text{and} \quad w_s(\mathcal{M}_s^\circ) \geq w_\varepsilon(\mathcal{M}_s^\circ) \quad \implies \quad s^2 \geq 2c_\delta w_\varepsilon(\mathcal{M}_s^\circ)$$



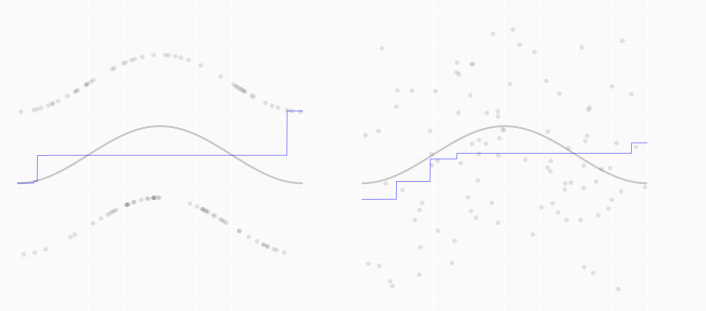
People call this random sign width, or something like it, *Rademacher Complexity*.

$$\text{Rademacher Complexity}(\mathcal{V}) = \mathbb{E} \max_{v \in \mathcal{V}} \langle s, v \rangle_{L_2(\mathbf{P}_n)}$$

$$\text{or maybe } = \mathbb{E} \max_{v \in \mathcal{V}} |\langle s, v \rangle_{L_2(\mathbf{P}_n)}| \quad \text{which is slightly different.}$$

## Comparison To Gaussian Width

---



Regression with random sign noise isn't much harder than with gaussian noise.  
Random sign width can't be much bigger than gaussian width.

$$w_s(\mathcal{V}) \leq \sqrt{\frac{\pi}{2}} w_g(\mathcal{V}) \approx 1.25 w_g(\mathcal{V}).$$

$$\begin{aligned} \mathbb{E} \max_{v \in \mathcal{V}} \sum_{i=1}^n g_i v_i &= \mathbb{E}_s \mathbb{E}_g \max_{v \in \mathcal{V}} \sum_{i=1}^n |g_i| s_i v_i \\ &\geq \mathbb{E}_s \max_{v \in \mathcal{V}} \mathbb{E}_g \sum_{i=1}^n |g_i| s_i v_i = \mathbb{E}_s \max_{v \in \mathcal{V}} \sum_{i=1}^n \mathbb{E} |g_i| s_i v_i = \mathbb{E}_s \max_{v \in \mathcal{V}} \sum_{i=1}^n \sqrt{\frac{2}{\pi}} s_i v_i. \end{aligned}$$

We've got the tools to show it isn't much smaller, either.

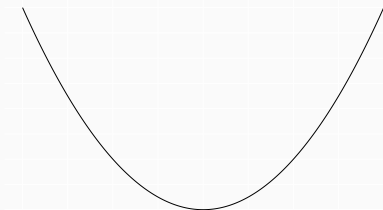
- We'll prove that. It's a factor  $\propto \sqrt{\log(n)}$  smaller at most.
- We'll bound widths for arbitrary noise in terms of random sign width.
  - And therefore in terms of gaussian width.
  - It'll be almost exactly the same symmetrization and contraction argument.
- We'll talk about *sampling*, too. We'll bound population mean squared error.
  - We'll see that sampling isn't much more problematic than random sign noise.
  - That's a symmetrization and contraction thing, too.
  - But we'll need a subtler contraction argument.

Background: Convex Functions Are  
Maximized At Extreme Points

---

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

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



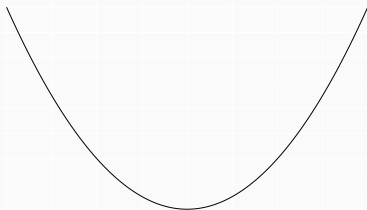
We can give this a *probabilistic interpretation* for a random variable  $Z_\lambda$ .

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

## Definition

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

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



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

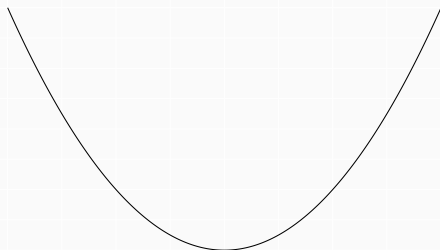
$$f(\mathbb{E} Z_\lambda) \leq \mathbb{E} f(Z_\lambda) \quad \text{where} \quad Z_\lambda = \begin{cases} a & \text{w.p. } 1 - \lambda \\ b & \text{w.p. } \lambda \end{cases}$$

# 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.

$$f(E Z) \leq E f(Z)$$

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.

$$f(\lambda_1 z_1 + \lambda_2 z_2) \leq \lambda_1 f(z_1) + \lambda_2 f(z_2) \quad \text{if} \quad \lambda_1, \lambda_2 \geq 0 \quad \text{satisfy} \quad \lambda_1 + \lambda_2 = 1$$

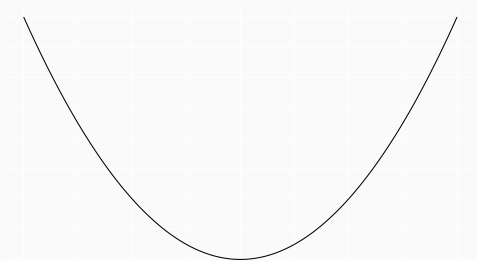
## Inductive Step.

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.**

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

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  $u_1 \dots u_K$ .

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

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.

# Maxima of Convex Functions over Polytopes

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

**Proof.**

It's more-or-less the same as the one-dimensional case. We use Jensen's inequality.

$$\max_{u \in \mathcal{U}} f(u) = \max_{\lambda \in \Lambda} f\left(\sum_i \lambda_i u_i\right) \leq \max_{\lambda \in \Lambda} \sum_i \lambda_i f(u_i) \leq \max_i f(u_i)$$