

# Machine Learning Theory

## Bounded Variation Regression

---

David A. Hirshberg

January 24, 2025

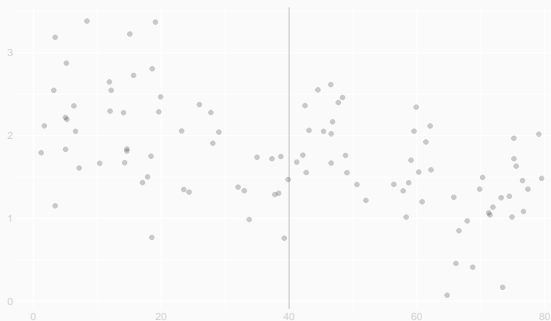
Emory University

## Review

---

## A Fictionalized RDD Example

- Question. Are smaller classes better for 5th graders?
- Data. A state caps class sizes at 40.
  - When there are  $x \leq 40$  5th graders enrolled in a school, they run one class of size  $x$ .
  - When there are  $x > 40$ , they run two classes of average size  $x/2$ .

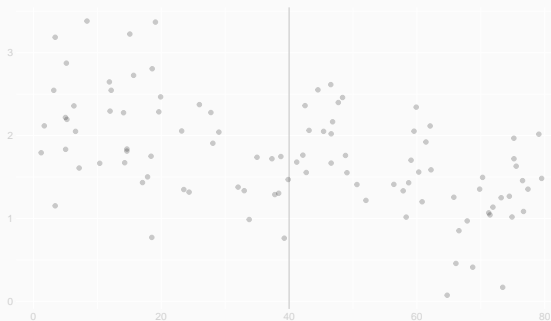


This is a fake-data version of a study of Angrist and Lavy [1999]. The state is Israel.

- It has been simplified to make our discussion easier.
- Real schools sometimes had more than 80 5th-graders enrolled.
- And they didn't follow this cap perfectly.

## A Fictionalized RDD Example

- Question. Are smaller classes better for 5th graders?
- Data. A state caps class sizes at 40.
  - When there are  $x \leq 40$  5th graders enrolled in a school, they run one class of size  $x$ .
  - When there are  $x > 40$ , they run two classes of average size  $x/2$ .



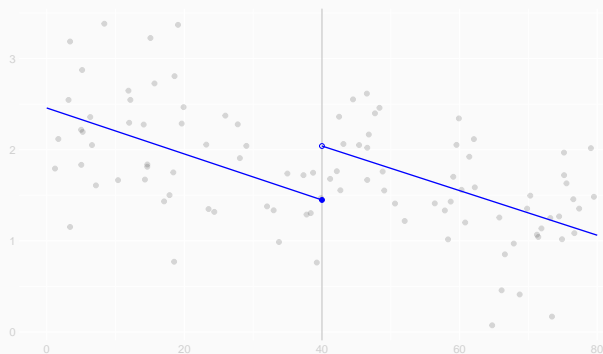
We can use this to estimate the average effect of having 20 vs. 40 students/class.

$$\text{effect} = \mu(40+) - \mu(40-) \quad \text{where}$$

$$\mu(x) = E[\text{avg. test score}_i \mid \text{enrolled 5th graders}_i = x].$$

All we have to do is estimate  $\mu(x)$  just to the left and to the right of 40.

# The Simple Approach



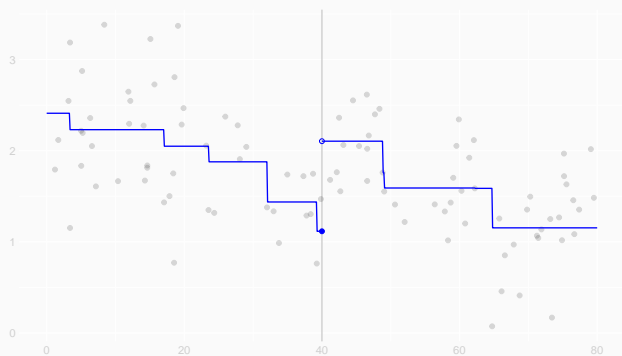
$$\widehat{\text{effect}} = \hat{\mu}_{\text{right}}(40) - \hat{\mu}_{\text{left}}(40) = 0.59$$

where

$$\hat{\mu}_{\text{left}} = \operatorname{argmin}_{m \in \mathcal{M}} \sum_{i: X_i \leq 40} \{Y_i - m(X_i)\}^2 \quad \text{and} \quad \hat{\mu}_{\text{right}} = \operatorname{argmin}_{m \in \mathcal{M}} \sum_{i: X_i > 40} \{Y_i - m(X_i)\}^2.$$

$$\text{for } \mathcal{M} = \{m(x) = b_0 + b_1 x : b \in \mathbb{R}^2\}$$

# The Simple Approach



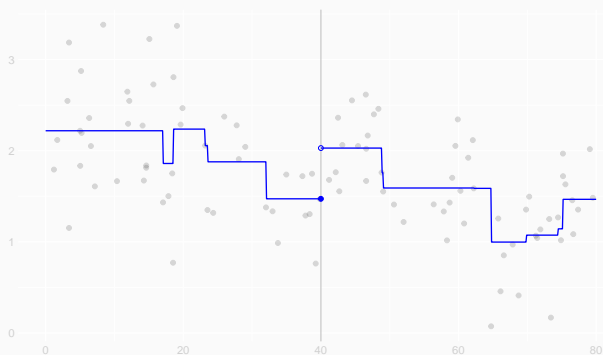
$$\widehat{\text{effect}} = \hat{\mu}_{\text{right}}(40) - \hat{\mu}_{\text{left}}(40) = 0.99$$

where

$$\hat{\mu}_{\text{left}} = \operatorname{argmin}_{m \in \mathcal{M}} \sum_{i: X_i \leq 40} \{Y_i - m(X_i)\}^2 \quad \text{and} \quad \hat{\mu}_{\text{right}} = \operatorname{argmin}_{m \in \mathcal{M}} \sum_{i: X_i > 40} \{Y_i - m(X_i)\}^2.$$

for  $\mathcal{M} = \{\text{decreasing } m\}$

# The Simple Approach



$$\widehat{\text{effect}} = \hat{\mu}_{\text{right}}(40) - \hat{\mu}_{\text{left}}(40) = 0.56$$

where

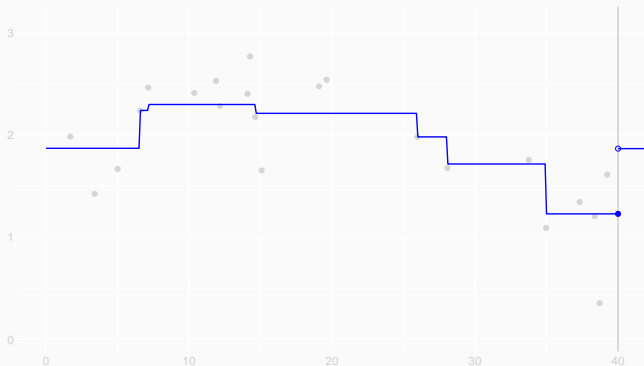
$$\hat{\mu}_{\text{left}} = \operatorname{argmin}_{m \in \mathcal{M}} \sum_{i: X_i \leq 40} \{Y_i - m(X_i)\}^2 \quad \text{and} \quad \hat{\mu}_{\text{right}} = \operatorname{argmin}_{m \in \mathcal{M}} \sum_{i: X_i > 40} \{Y_i - m(X_i)\}^2.$$

$$\text{for } \mathcal{M} = \{m : \int |m'(x)| \leq B\} \quad \text{where } B = 1.5$$

## The Bounded Variation Model

---





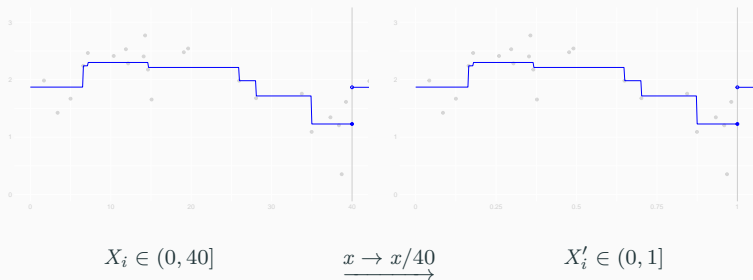
$$\hat{\mu} = \operatorname{argmin}_{m \in \mathcal{M}} \frac{1}{n} \sum_{i=1}^n \{Y_i - m(X_i)\}^2 \quad \text{where} \quad \mathcal{M} = \{m : \int |m'(x)| \leq B\}.$$

This is too vague to actually implement.

1. We haven't specified the domain we're integrating over.
2. Nor have we said what this means for non-differentiable functions.

Let's fix that.

The integration domain is easy. It's a matter of convention. We shift and scale our  $X_i$  into the unit interval  $[0, 1]$ . That's our integration domain.



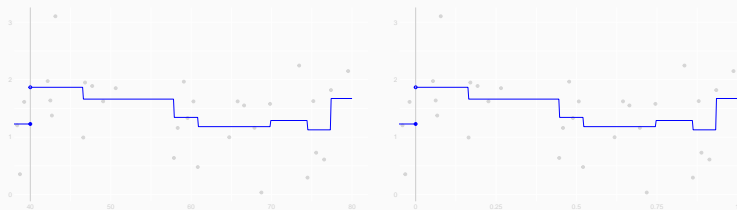
What's nice about this is that it's an average.

$$\int_0^1 |m'(x)| dx = E|m'(\tilde{X})| \text{ for } \tilde{X} \text{ uniformly distributed in } [0, 1]$$

Averages are easy to think about and compare to other things.

$$E|m'(\tilde{X})| \leq \sqrt{E m'(\tilde{X})^2} \leq \max_x |m'(x)|$$

The integration domain is easy. It's a matter of convention. We shift and scale our  $X_i$  into the unit interval  $[0, 1]$ . That's our integration domain.



$$X_i \in (40, 80]$$

$$x \rightarrow \frac{(x - 40)}{40}$$

$$\tilde{X}_i \in (0, 1]$$

What's nice about this is that it's an average.

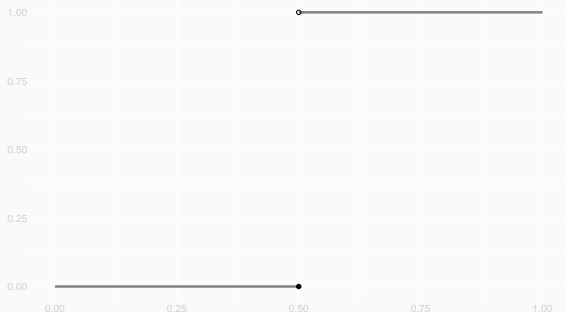
$$\int_0^1 |m'(x)| dx = \mathbb{E}|m'(\tilde{X})| \text{ for } \tilde{X} \text{ uniformly distributed in } [0, 1]$$

Averages are easy to think about and compare to other things.

$$\mathbb{E}|m'(\tilde{X})| \leq \sqrt{\mathbb{E} m'(\tilde{X})^2} \leq \max_x |m'(x)|$$

The nondifferentiable thing is subtler.  
Think about this problem—an infinite-data noiseless regression.

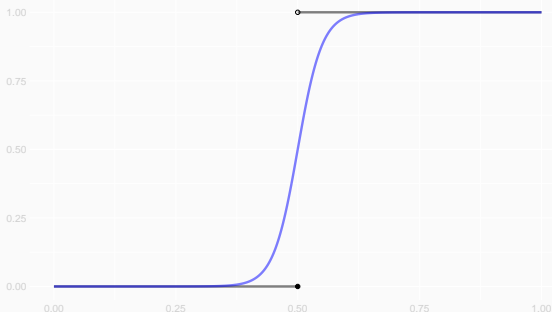
$$\hat{\mu} = \operatorname{argmin}_{m \in \mathcal{M}} \int_0^1 \{\mu(x) - m(x)\}^2 \quad \text{where} \quad \mathcal{M} = \left\{ \text{differentiable } m : \int_0^1 |m'(x)| \leq B \right\}.$$



Here's  $\mu$ .  
What is the solution  $\hat{\mu}$ ?

The nondifferentiable thing is subtler.  
Think about this problem—an infinite-data noiseless regression.

$$\hat{\mu} = \operatorname{argmin}_{m \in \mathcal{M}} \int_0^1 \{\mu(x) - m(x)\}^2 \quad \text{where} \quad \mathcal{M} = \left\{ \text{differentiable } m : \int_0^1 |m'(x)| \leq B \right\}.$$

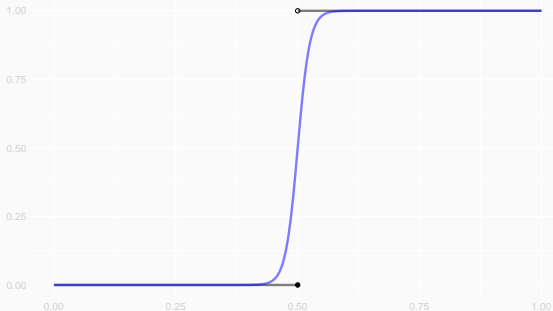


It's not **this function**  $m$ .  
We can find a better fit.

The nondifferentiable thing is subtler.

Think about this problem—an infinite-data noiseless regression.

$$\hat{\mu} = \operatorname{argmin}_{m \in \mathcal{M}} \int_0^1 \{\mu(x) - m(x)\}^2 \quad \text{where} \quad \mathcal{M} = \left\{ \text{differentiable } m : \int_0^1 |m'(x)| \leq B \right\}.$$



This is better, but we can keep going.

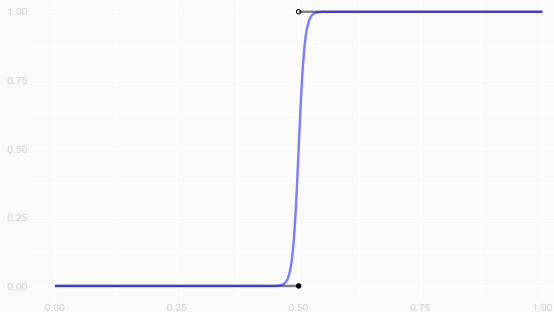
We can get arbitrarily close.

$$\mu = \lim_k m_k \quad \text{for} \quad m_k \in \mathcal{M}.$$

The nondifferentiable thing is subtler.

Think about this problem—an infinite-data noiseless regression.

$$\hat{\mu} = \operatorname{argmin}_{m \in \mathcal{M}} \int_0^1 \{\mu(x) - m(x)\}^2 \quad \text{where} \quad \mathcal{M} = \left\{ \text{differentiable } m : \int_0^1 |m'(x)| \leq B \right\}.$$



This is better, but we can keep going.

We can get arbitrarily close.

$$\mu = \lim_k m_k \quad \text{for} \quad m_k \in \mathcal{M}.$$

$$m_k(x) = \frac{1}{1 + e^{-k(x-.5)}}$$

The nondifferentiable thing is subtler.

Think about this problem—an infinite-data noiseless regression.

$$\hat{\mu} = \operatorname{argmin}_{m \in \mathcal{M}} \int_0^1 \{\mu(x) - m(x)\}^2 \quad \text{where} \quad \mathcal{M} = \left\{ \text{differentiable } m : \int_0^1 |m'(x)| \leq B \right\}.$$



This is better, but we can keep going.

We can get arbitrarily close.

$$\mu = \lim_k m_k \quad \text{for} \quad m_k \in \mathcal{M}.$$

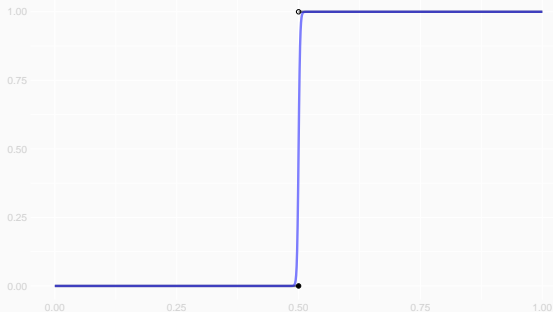
$$m_k(x) = \frac{1}{1 + e^{-k(x-.5)}}$$



The nondifferentiable thing is subtler.

Think about this problem—an infinite-data noiseless regression.

$$\hat{\mu} = \operatorname{argmin}_{m \in \mathcal{M}} \int_0^1 \{\mu(x) - m(x)\}^2 \quad \text{where} \quad \mathcal{M} = \left\{ \text{differentiable } m : \int_0^1 |m'(x)| \leq B \right\}.$$



This is better, but we can keep going.

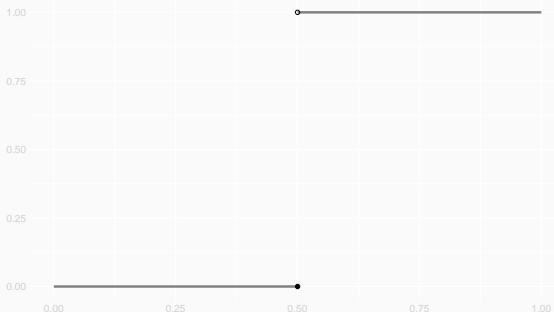
We can get arbitrarily close.

$$\mu = \lim_k m_k \quad \text{for} \quad m_k \in \mathcal{M}.$$

$$m_k(x) = \frac{1}{1 + e^{-k(x-.5)}}$$

The nondifferentiable thing is subtler.  
Think about this problem—an infinite-data noiseless regression.

$$\hat{\mu} = \operatorname{argmin}_{m \in \mathcal{M}} \int_0^1 \{\mu(x) - m(x)\}^2 \quad \text{where} \quad \mathcal{M} = \left\{ \text{differentiable } m : \int_0^1 |m'(x)| \leq B \right\}.$$



The **only reasonable answer** is  $\hat{\mu} = \mu$   
But  $\mu$  isn't in our model  $\mathcal{M}$ .

The argmin  $\hat{\mu}$  **doesn't exist**.

$$\hat{\mu} = \operatorname{argmin}_{m \in \mathcal{M}} \int_0^1 \{\mu(x) - m(x)\}^2 \quad \text{where} \quad \mathcal{M} = \left\{ \text{differentiable } m : \int_0^1 |m'(x)| \leq B \right\}.$$

Mean squared error does not have a minimum **in this model**.

But being pedantic doesn't get us far. We know  $\mu$  **should** be the solution.

We need a way to define our model so that it can be.

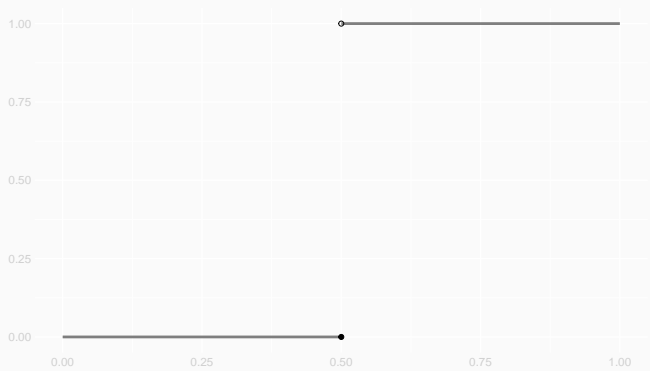
In a sense, we don't need to change much.

$$\hat{\mu} = \operatorname{argmin}_{m \in \mathcal{M}} \int_0^1 \{\mu(x) - m(x)\}^2 \quad \text{where} \quad \mathcal{M} = \left\{ m : \int_0^1 |m'(x)| \leq B \right\}.$$

All we've got to do is define the integral  $\int_0^1 |m'(x)|$  for non-differentiable functions.

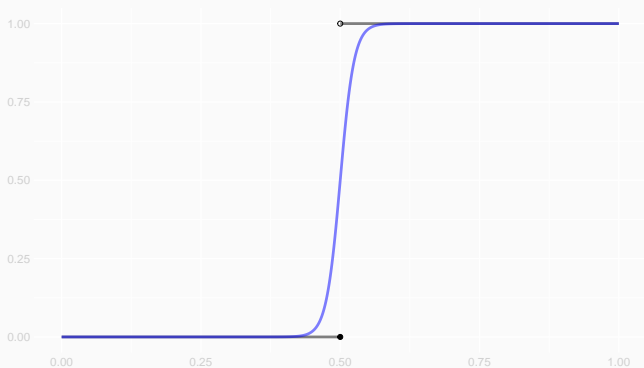
We know what we want from our definition.

$$\int_0^1 |m'(x)| dx = ?$$



We know what we want from our definition.

$$\int_0^1 |m'(x)| dx \approx \int_0^1 |m'_k(x)| dx$$



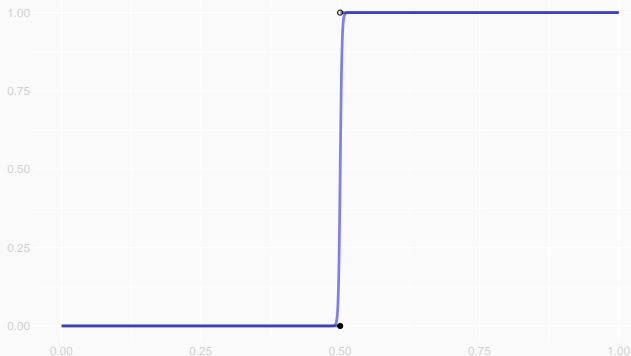
where

$$m_k(x) \approx m(x).$$

The rest is calculus.

We know what we want from our definition.

$$\int_0^1 |m'(x)| dx = \lim_k \int_0^1 |m'_k(x)| dx$$



where

$$m_k(x) \rightarrow m(x).$$

The rest is calculus.

If we have a dense partition  $0 = x_1 < \dots < x_k = 1$  in  $[0, 1]$ , then:

$$\begin{aligned}\int_0^1 |m'(x)| &\approx \sum_{j=1}^k |m'(x_j)|(x_{j+1} - x_j) && \text{a sum approximates our integral} \\ &\approx \sum_{j=1}^k \left| \frac{m(x_{j+1}) - m(x_j)}{x_{j+1} - x_j} \right| (x_{j+1} - x_j) && \text{a slope approximates our derivative} \\ &= \sum_{j=1}^k |m(x_{j+1}) - m(x_j)| && \text{and we've got no calculus stuff left.}\end{aligned}$$

If we take denser and denser partitions ( $\max_j x_{j+1} - x_j \rightarrow 0$ ), then this is exact.

So we've got a definition.

$$\int_0^1 |m'(x)| dx = \lim \sum_{j=1}^k |m(x_{j+1}) - m(x_j)| \quad \text{as} \quad \max_j x_{j+1} - x_j \rightarrow 0.$$

$$\int_0^1 |m'(x)| dx = \lim \sum_{j=1}^k |m(x_{j+1}) - m(x_j)| \quad \text{as} \quad \max_j x_{j+1} - x_j \rightarrow 0.$$

This definition is a bit complicated. It involves a weird kind of limit. You'd have to think about whether that limit exists. So we use another.

$$\int_0^1 |m'(x)| dx = \sup_{\substack{\text{finite partitions} \\ 0=x_1 < \dots < x_k=1}} \sum_{j=1}^k |m(x_{j+1}) - m(x_j)|.$$

We've replaced the *limit over denser partitions* with a *supremum over all partitions*.

Is that ok? Are these the same?



$$\sup_{\substack{\text{finite partitions} \\ 0=x_0 < \dots < x_k=1}} \sum_{j=1}^k |m(x_{j+1}) - m(x_j)| \stackrel{?}{=} \lim_{j=1}^k |m(x_{j+1}) - m(x_j)| \text{ as } \max_j x_{j+1} - x_j \rightarrow 0.$$

I'm not being totally rigorous here. But yes, they're the same.

The sum can only get bigger if we *refine* our partition by adding intermediate points.

$$\begin{aligned} |m(x_{j+1}) - m(x_j)| &= |m(x_{j+1}) - m(\tilde{x}) + m(\tilde{x}) - m(x_j)| \\ &\leq |m(x_{j+1}) - m(\tilde{x})| + |m(\tilde{x}) - m(x_j)| \quad \text{for } \tilde{x} \in (x_j, x_{j+1}). \end{aligned}$$

That's the triangle inequality in action.

And it means that our supremum is *always* a limit for increasingly dense partitions.

It doesn't however, mean we can't stop at some point.

There may be a finite partition that gives us the maximal sum.

Thinking about these partitions can help us understand our 'integral'.

# The Bounded Variation Model

$$\mathcal{M} = \{m : \rho_{TV}(m) \leq B\}$$

where

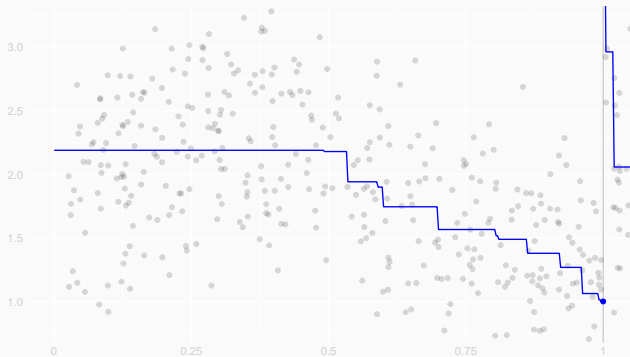
$$\begin{aligned} \rho_{TV}(m) &= \int_0^1 |m'(x)| dx && \text{for differentiable } m \\ &= \sup_{\substack{\text{finite partitions} \\ 0=x_0 < \dots < x_k=1}} \sum_{j=1}^k |m(x_{j+1}) - m(x_j)| && \text{in complete generality.} \end{aligned}$$

It's a **ball** in the **total variation seminorm**  $\rho_{TV}$ .

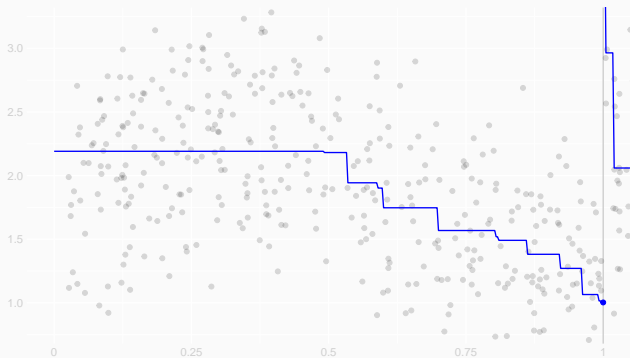
How do we know that it's a seminorm?

## Understanding Total Variation

---



Here's a monotone regression estimator  $\hat{\mu}$ .  
What is its total variation  $\rho_{TV}(\hat{\mu})$ ?

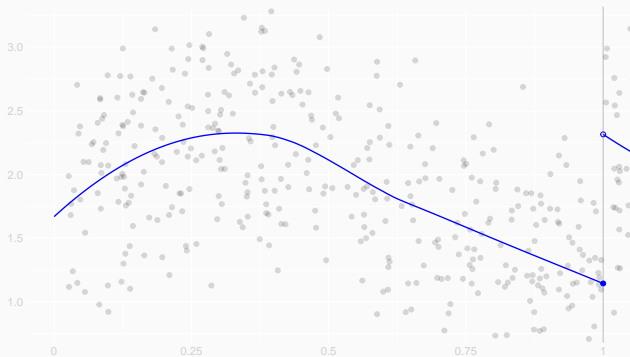


Here's a monotone regression estimator  $\hat{\mu}$ .  
 What is its total variation  $\rho_{TV}(\hat{\mu})$ ?

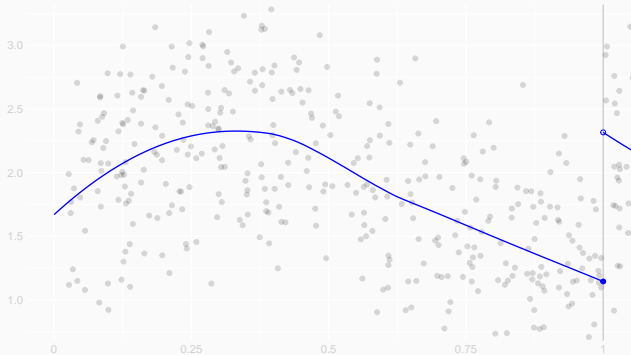
It's  $\hat{\mu}(0) - \hat{\mu}(1) \approx 1.25$ .

$$\sum_{j=1}^k |\hat{\mu}(x_{j+1}) - \hat{\mu}(x_j)| = \sum_{j=1}^k \hat{\mu}(x_j) - \hat{\mu}(x_{j+1}) = \hat{\mu}(x_0) - \hat{\mu}(x_k).$$

For a monotone function, it doesn't matter what partition we sum over.  
 The sum is always the same.



Here's another estimator  $\hat{\mu}$ . It's a *local polynomial regression estimator*.  
What is its total variation  $\rho_{TV}(\hat{\mu})$ ?



Here's another estimator  $\hat{\mu}$ . It's a *local polynomial regression estimator*.  
 What is its total variation  $\rho_{TV}(\hat{\mu})$ ?

It's  $|\hat{\mu}(0) - \hat{\mu}(x_*)| + |\hat{\mu}(x_*) - \hat{\mu}(1)| \approx 1.75$  for  $x_* = \operatorname{argmax}_x \hat{\mu}(x)$ .

$$\sum_{j=1}^k |\hat{\mu}(x_{j+1}) - \hat{\mu}(x_j)| = \sum_{j=1}^{j_*-1} \hat{\mu}(x_{j+1}) - \hat{\mu}(x_j) + \sum_{j=j_*}^k \hat{\mu}(x_j) - \hat{\mu}(x_{j+1})$$

for  $x_{j_*} = x_*$  because  $m$  is increasing on  $[0, x_*]$  and decreasing on  $[x_*, 1]$ .

As long as our partition contains all local extrema, the sum is always the same.

$$\begin{aligned}\rho_{TV}(m) &= |m(1) - m(0)| && \text{if } m \text{ is monotone (increasing or decreasing) on } [0, 1] \\ &= \sum_j |m(x_{j+1}) - m(x_j)| && \text{if } m \text{ is monotone on each interval } [x_j, x_{j+1}].\end{aligned}$$



## Try these!

What is  $\rho_{TV}(m)$  for ...

1.  $m(x) = x$

2.  $m(x) = x^2$

3.  $m(x) = e^x$

4.  $m(x) = \sin(\pi x)$

5.  $m(x) = \begin{cases} 0 & \text{if } x < 1 \\ 1 & \text{if } x = 1 \end{cases}$

6.  $m(x) = \sin(1/x)$

We'll look at some more interesting cases in this week's homework.

## Why the Bounded Variation Model is Useful

---

## How does a constraint on total variation keep us from overfitting?

Hint.

- A monotonicity constraint helps because noise jumps up and down.
- To fit noise would violate monotonicity.

## How does a constraint on total variation keep us from overfitting?

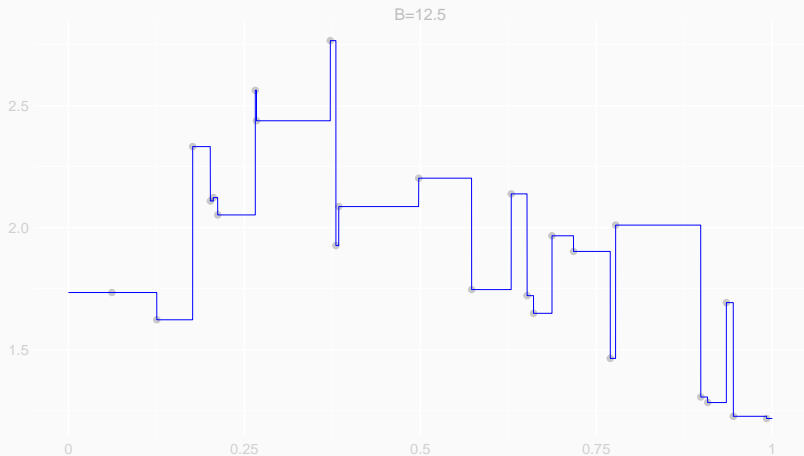
Hint.

- A monotonicity constraint helps because noise jumps up and down.
- To fit noise would violate monotonicity.

Answer.

- A TV constraint helps because noise jumps up and down *a lot*.
- To fit it, we'd need a curve  $m$  with huge total variation  $\rho_{TV}(m)$ .

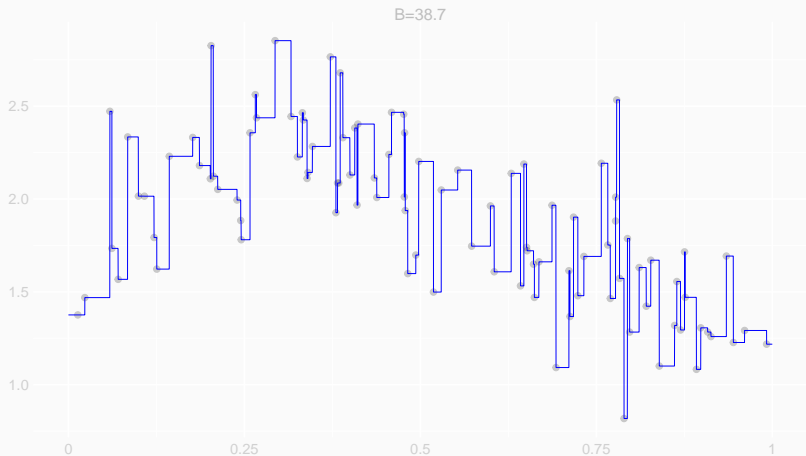
To fit perfectly, we need to increase our variation budget with sample size. Fast.



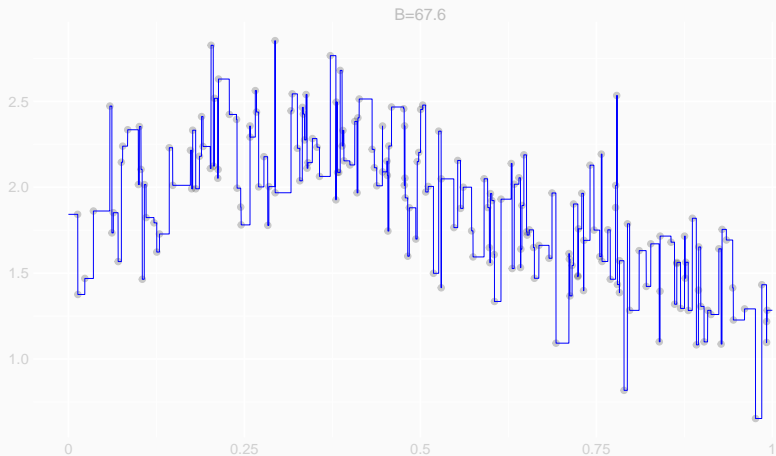
To fit perfectly, we need to increase our variation budget with sample size. Fast.



To fit perfectly, we need to increase our variation budget with sample size. Fast.

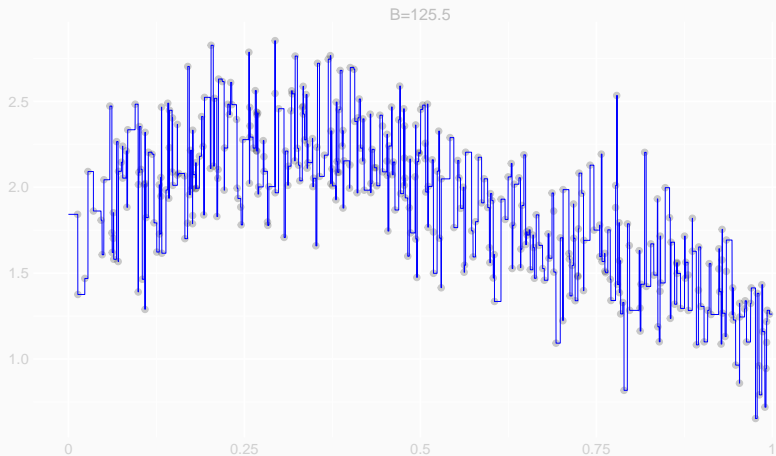


To fit perfectly, we need to increase our variation budget with sample size. Fast.

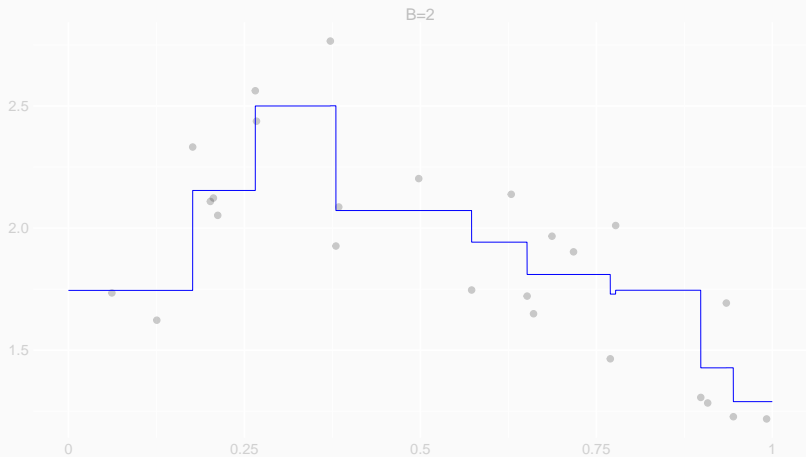




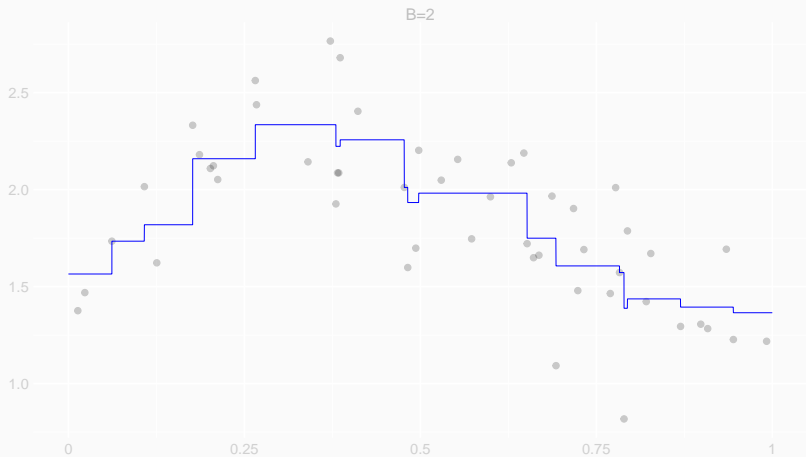
To fit perfectly, we need to increase our variation budget with sample size. Fast.



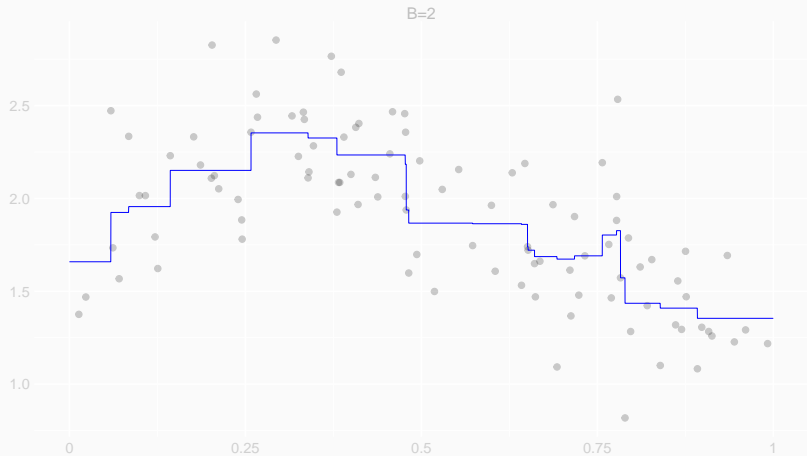
To fit the overall trend, we don't.



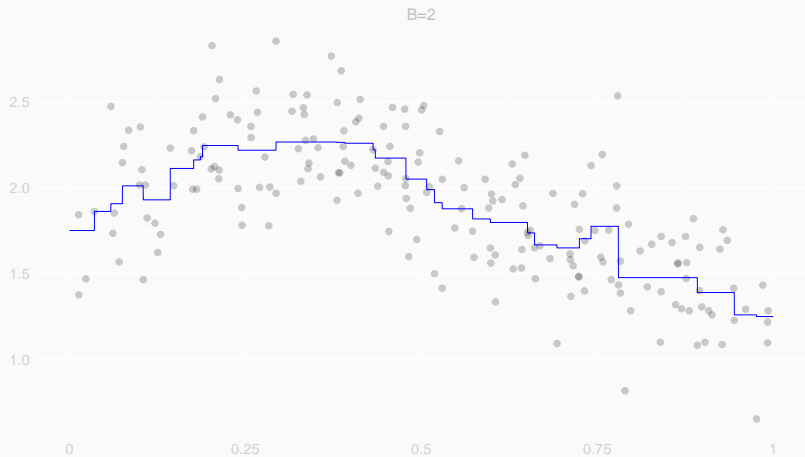
To fit the overall trend, we don't.



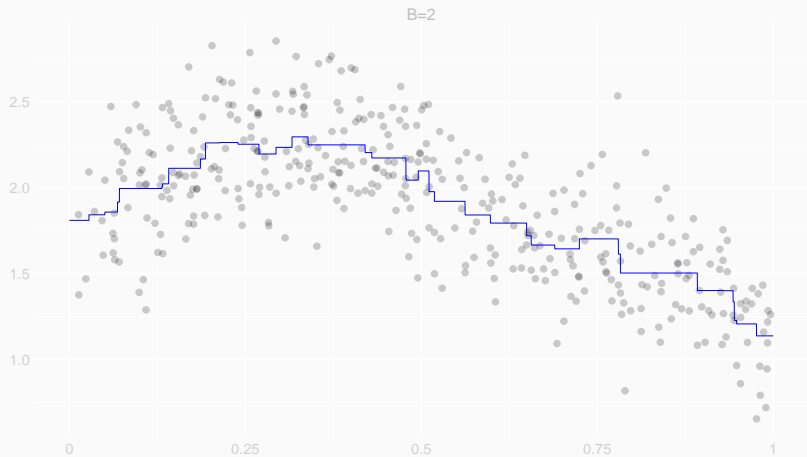
To fit the overall trend, we don't.



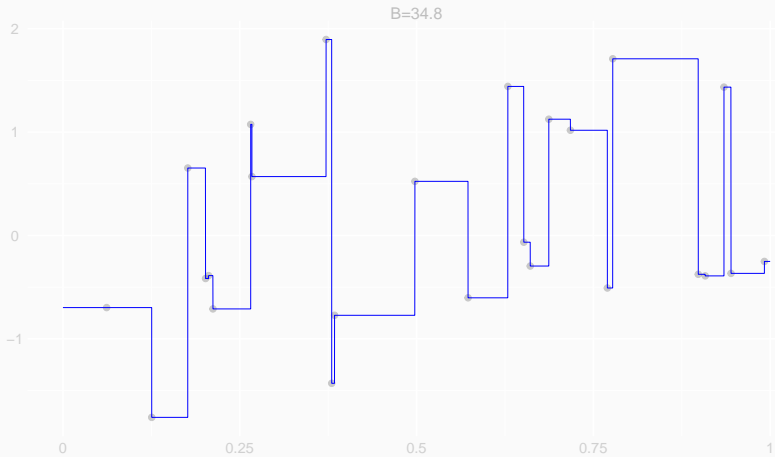
To fit the overall trend, we don't.



To fit the overall trend, we don't.

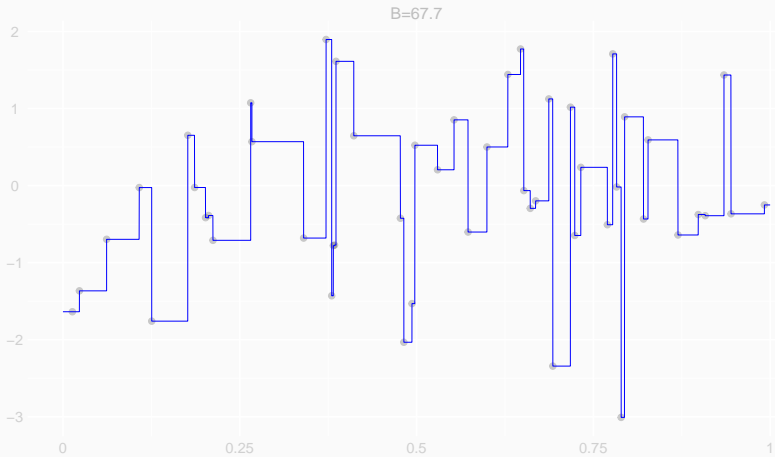


It takes a similar increase to fit pure gaussian noise.



$$Y_i \sim N(0, 1)$$

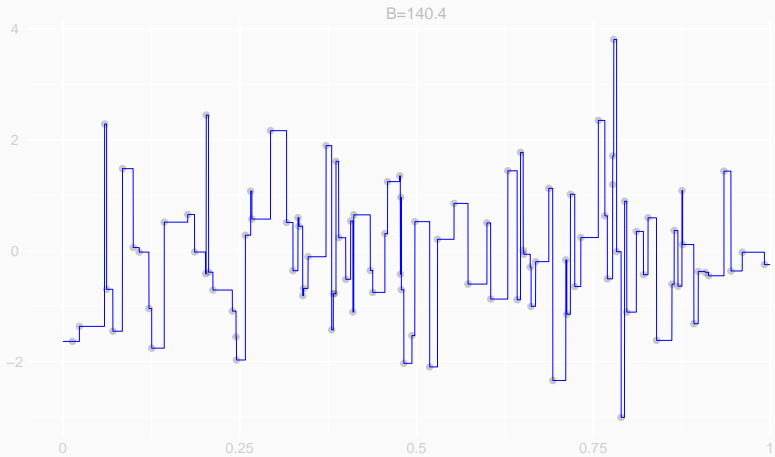
It takes a similar increase to fit pure gaussian noise.



$$Y_i \sim N(0, 1)$$

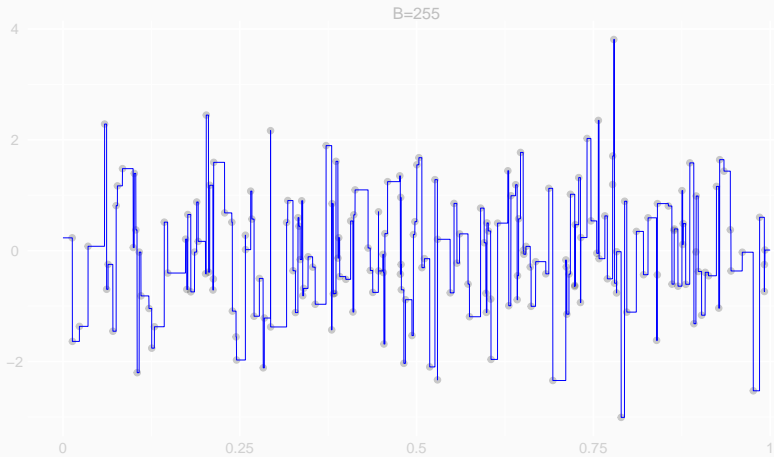


It takes a similar increase to fit pure gaussian noise.



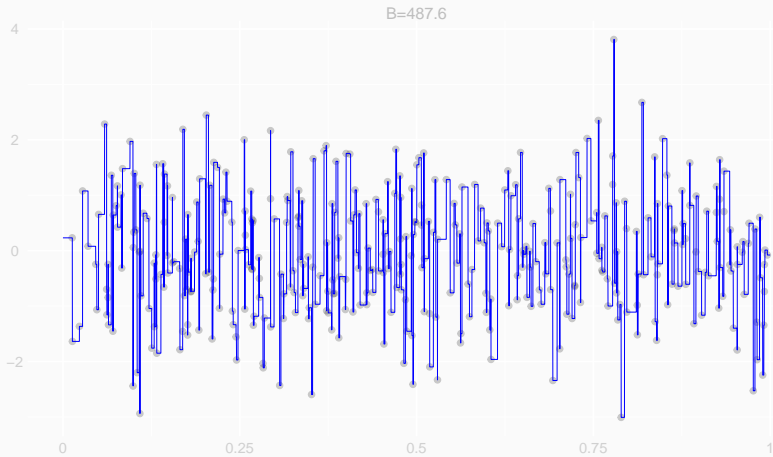
$$Y_i \sim N(0, 1)$$

It takes a similar increase to fit pure gaussian noise.



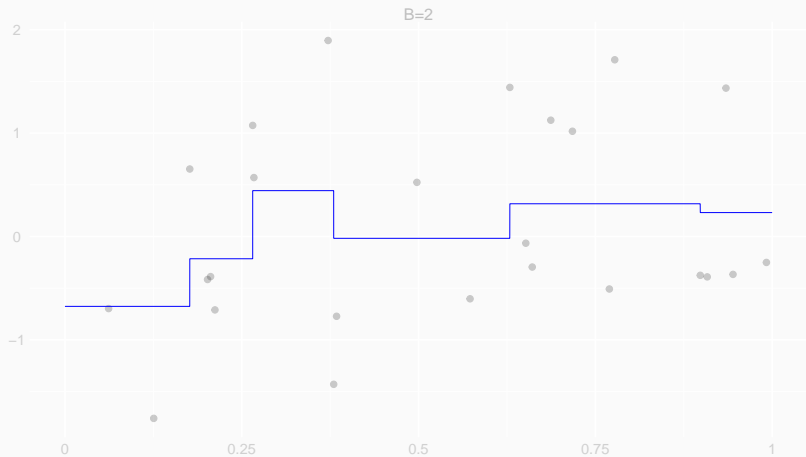
$$Y_i \sim N(0, 1)$$

It takes a similar increase to fit pure gaussian noise.



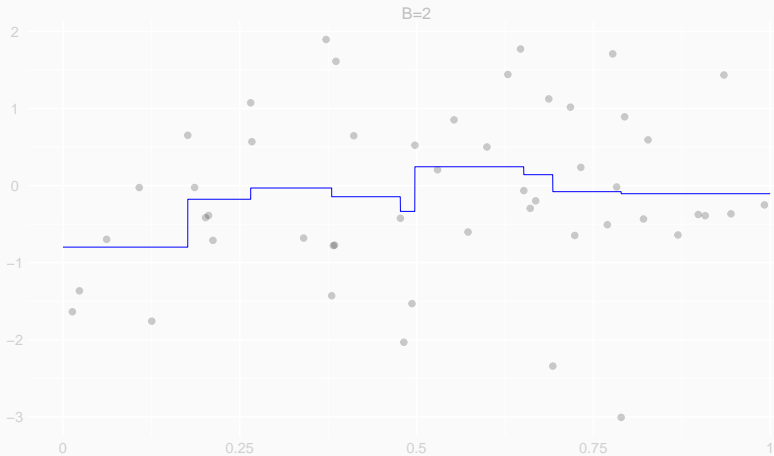
$$Y_i \sim N(0, 1)$$

And we can find that there's no trend if we keep our budget small.



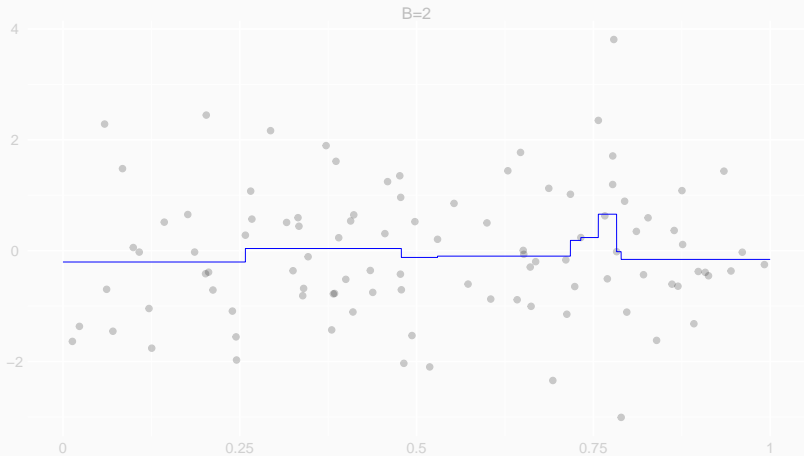
$$Y_i \sim N(0, 1)$$

And we can find that there's no trend if we keep our budget small.



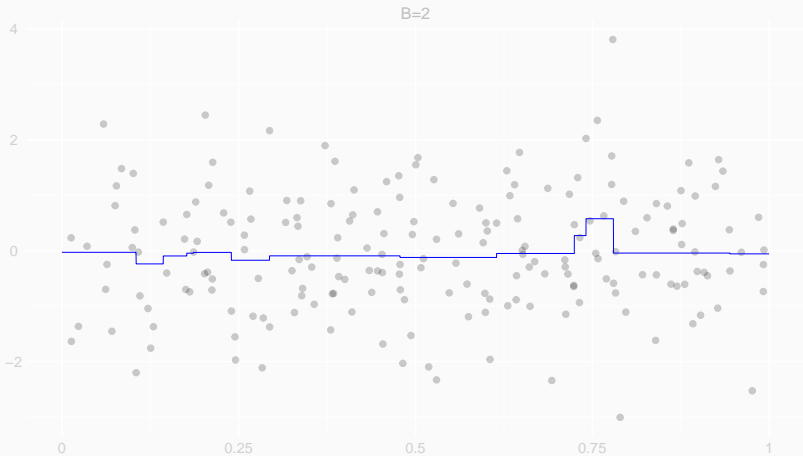
$$Y_i \sim N(0, 1)$$

And we can find that there's no trend if we keep our budget small.



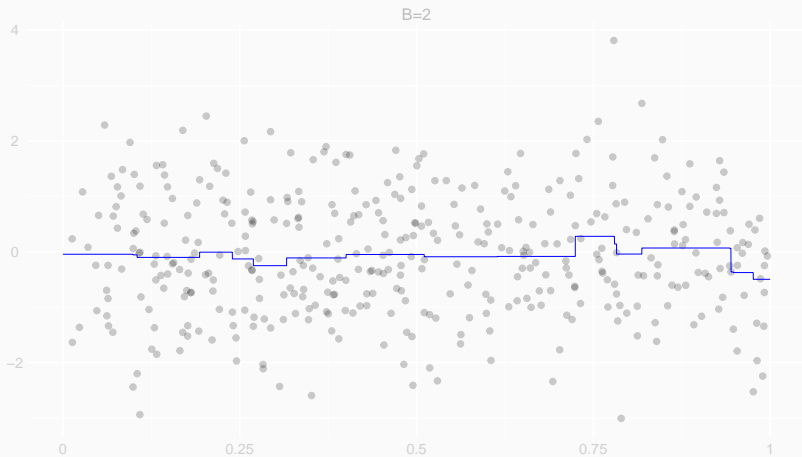
$$Y_i \sim N(0, 1)$$

And we can find that there's no trend if we keep our budget small.



$$Y_i \sim N(0, 1)$$

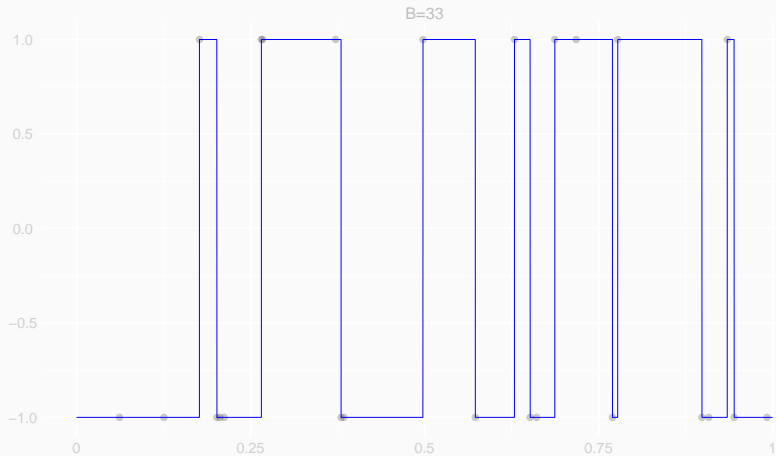
And we can find that there's no trend if we keep our budget small.



$$Y_i \sim N(0, 1)$$

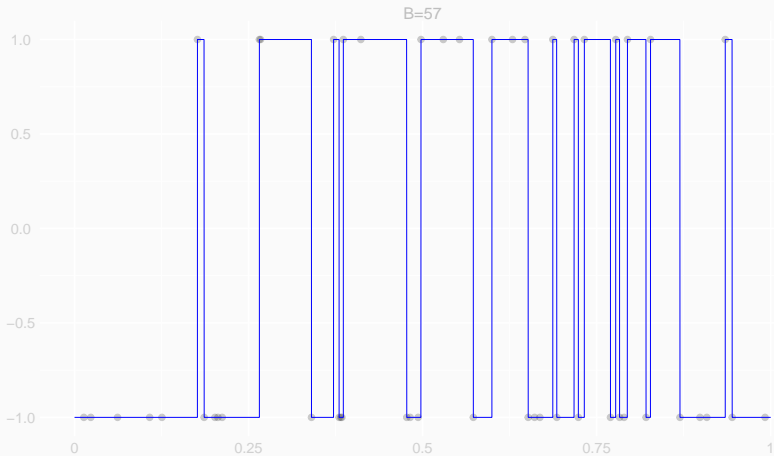


Same deal for fitting random signs  $\pm 1$ . Coin flips.



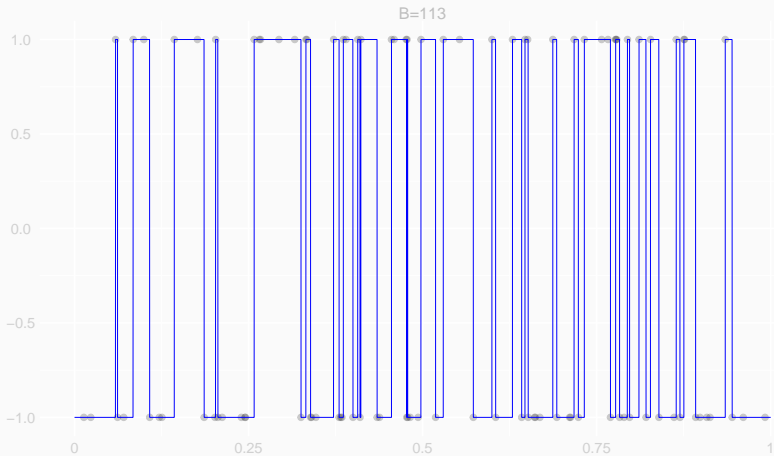
$Y_i = \pm 1$  each with probability  $1/2$

Same deal for fitting random signs  $\pm 1$ . Coin flips.



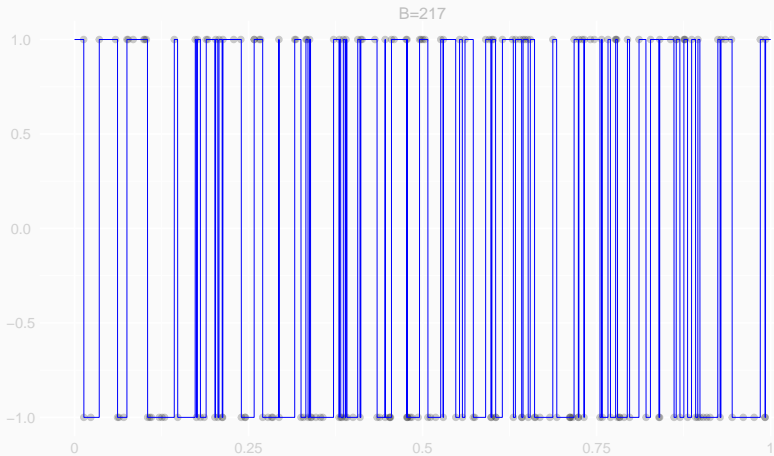
$Y_i = \pm 1$  each with probability  $1/2$

Same deal for fitting random signs  $\pm 1$ . Coin flips.



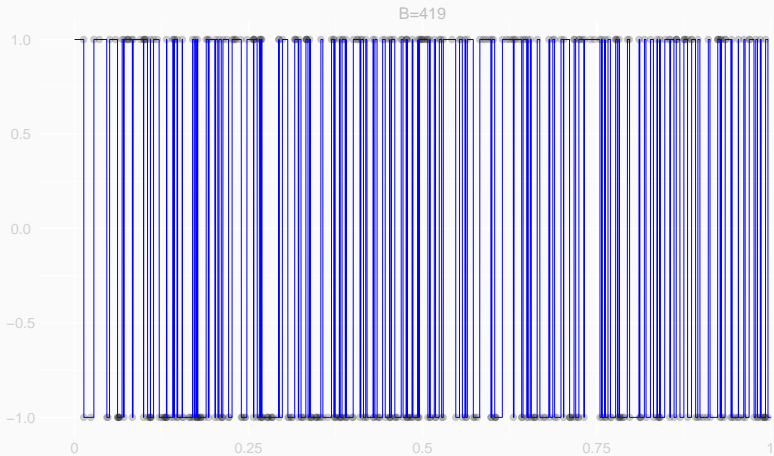
$Y_i = \pm 1$  each with probability  $1/2$

Same deal for fitting random signs  $\pm 1$ . Coin flips.



$Y_i = \pm 1$  each with probability  $1/2$

Same deal for fitting random signs  $\pm 1$ . Coin flips.



$Y_i = \pm 1$  each with probability  $1/2$

## The price of noise.

If a curve fits noise perfectly, i.e. if  $m(X_i) = \varepsilon_i$ , then what do we know about  $\rho_{TV}(m)$ ?

## The price of noise.

If a curve fits noise perfectly, i.e. if  $m(X_i) = \varepsilon_i$ , then what do we know about  $\rho_{TV}(m)$ ?

$$\rho_{TV}(m) \geq \sum_{i=1}^{n-1} |m(X_{i+1}) - m(X_i)| = \sum_{i=1}^{n-1} |\varepsilon_{i+1} - \varepsilon_i|.$$

Why  $\geq$  and not  $=$ ?

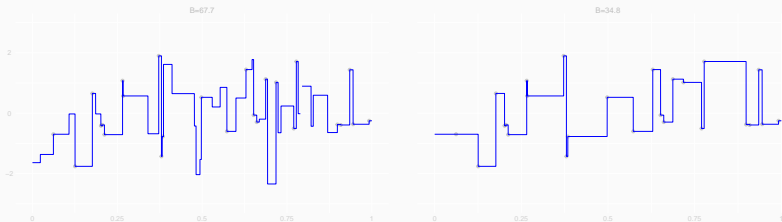
# The price of noise.

If a curve fits noise perfectly, i.e. if  $m(X_i) = \varepsilon_i$ , then what do we know about  $\rho_{TV}(m)$ ?

$$\rho_{TV}(m) \geq \sum_{i=1}^{n-1} |m(X_{i+1}) - m(X_i)| = \sum_{i=1}^{n-1} |\varepsilon_{i+1} - \varepsilon_i|.$$

Why  $\geq$  and not  $=$ ?

The observations are just one partition, so all this tells us is a lower bound.  
We didn't say what this curve does between the observations.



It might jump for no reason. It might not.  
If it doesn't, this is an equality.



What if our noise  $\varepsilon_i$  is a coin flip.

$\varepsilon_i = \pm 1$  each with probability  $1/2$ .

How much total variation would we need for a perfect fit?

First, translate this into a question about the noise.

What if our noise  $\varepsilon_i$  is a coin flip.

$\varepsilon_i = \pm 1$  each with probability  $1/2$ .

How much total variation would we need for a perfect fit?

First, translate this into a question about the noise.

What is  $\sum_{i=1}^{n-1} |\varepsilon_{i+1} - \varepsilon_i|$ ?

What if our noise  $\varepsilon_i$  is a coin flip.

$\varepsilon_i = \pm 1$  each with probability  $1/2$ .

How much total variation would we need for a perfect fit?

First, translate this into a question about the noise.

What is  $\sum_{i=1}^{n-1} |\varepsilon_{i+1} - \varepsilon_i|$  ?

In the worst case, it's  $(n - 1) \times 2$ .

- We'd get that if we flipped sign every time.
- Heads, Tails, Heads, Tails, ...
- $|1 - -1| + |-1 - 1| + |1 - -1| + \dots$  ( $n - 1$  times).
- Not very random-looking. Not very likely.

## The price of noise. Coin flip edition.

What if our noise  $\varepsilon_i$  is a coin flip.

$\varepsilon_i = \pm 1$  each with probability  $1/2$ .

How much total variation would we need for a perfect fit?

First, translate this into a question about the noise.

What is  $\sum_{i=1}^{n-1} |\varepsilon_{i+1} - \varepsilon_i|$ ?

In the average case, it's half that.

$$|\varepsilon_{i+1} - \varepsilon_i| = \begin{cases} 0 & \text{with probability } 1/2 \\ 2 & \text{with probability } 1/2 \end{cases}$$

and therefore

$$\mathbb{E}|\varepsilon_{i+1} - \varepsilon_i| = \frac{0 + 2}{2} = 1 \quad \text{and} \quad \mathbb{E} \sum_{i=1}^{n-1} |\varepsilon_{i+1} - \varepsilon_i| = (n-1) \times 1.$$

What if our noise  $\varepsilon_i$  is standard normal?

i.e. what is  $\sum_{i=1}^{n-1} |\varepsilon_{i+1} - \varepsilon_i|$  for  $\varepsilon_i \sim N(0, 1)$  ?

Is it bigger or smaller?

What if our noise  $\varepsilon_i$  is standard normal?

i.e. what is  $\sum_{i=1}^{n-1} |\varepsilon_{i+1} - \varepsilon_i|$  for  $\varepsilon_i \sim N(0, 1)$  ?

Is it bigger or smaller?

In the average case, it's  $(n - 1) \times \sqrt{2} \mathbf{E}|\varepsilon_i| \approx (n - 1) \times 1.1$ .

$$|\varepsilon_{i+1} - \varepsilon_i| \sim |N(0, 2)| = \sqrt{2}|N(0, 1)|$$

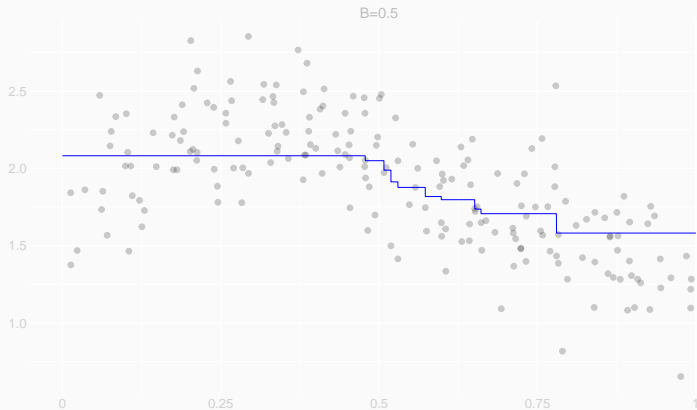
and therefore

$$\mathbf{E}|\varepsilon_{i+1} - \varepsilon_i| = \sqrt{2} \mathbf{E}|\varepsilon_i| \quad \text{and} \quad \mathbf{E} \sum_{i=1}^{n-1} |\varepsilon_{i+1} - \varepsilon_i| = (n - 1) \times \sqrt{2} \mathbf{E}|\varepsilon_i|.$$

$= \sqrt{\frac{2}{\pi}}$

We don't tend to use variation budgets anywhere near that big.  
If  $\mu$  were wildly discontinuous, we might be a bit stuck even if we knew it exactly.

But we do have to choose a budget somehow.  
And if we go too small, that's not good either.



In lab, we'll talk about how to choose it automatically.

## BV Regression in Action.

Maybe you've heard of it. Or something like it.

---



*Total Variation Denoising is BV Regression in 2D.*

Original



Noisy image



Denoised image



People use it all the time.

This image, by MAL, is licensed under CC BY-SA 3.0.

*Total Variation Denoising* is BV Regression in 2D.



They used it to capture the first-ever image of a black hole.  
We'll talk about and implement it in a few weeks.

# The *Highly Adaptive Lasso* is a generalization of BV regression that works well even with high dimensional data.

Google Scholar highly adaptive lasso

About 32,100 results (0.83 sec)

Articles

Any time  
Since 2022  
Since 2021  
Since 2016  
Custom range...

Sort by relevance  
Sort by date

Any type  
Review articles

include parents  
 include citations

Create alert

**The highly adaptive lasso estimator** [pdf] [jeep.oreg](#)  
D. Berstner, M. Van Der Laan - 2016 IEEE International ... - [ieeexplore.ieee.org](#)  
Find it @ Emory  
Estimation of a regression function is a common goal of statistical learning. We propose a novel nonparametric regression estimator that, in contrast to many existing methods, does not rely on local smoothness assumptions nor is it constructed using local smoothing ...  
☆ Save ☆ Cite Cited by 80 Related articles All 5 versions

**A generally efficient targeted minimum loss based estimator based on the highly adaptive lasso** [pdf] [nih.gov](#)  
M. Van Der Laan - The International journal of biostatistics, 2017 - [degruyter.com](#)  
Suppose we observe  $n$  independent and identically distributed observations of a finite dimensional bounded random variable. This article is concerned with the construction of an efficient targeted minimum loss-based estimator (TMLE) of a pathwise differentiable target ...  
☆ Save ☆ Cite Cited by 53 Related articles All 10 versions

**hal9001: Scalable highly adaptive lasso regression in R** [pdf] [teoyj.org](#)  
N. Ishida, J. G. Glynn, M. Van Der Laan - Journal of Open Source ... - 2020 - [jeep.oreg](#)  
Find it @ Emory  
The `hal9001` R package provides a computationally efficient implementation of the **highly adaptive lasso** (HAL), a flexible nonparametric regression and machine learning algorithm endowed with several theoretically convenient properties. `hal9001` puts an implementation ...  
☆ Save ☆ Cite Cited by 12 Related articles All 4 versions 36

**Efficient estimation of pathwise differentiable target parameters with the under-smoothed highly adaptive lasso** [pdf] [arxiv.org](#)  
M. Van Der Laan, D. Berstner, V. Cal - arXiv preprint arXiv:1908.09507, 2019 - [arxiv.org](#)  
We consider estimation of a functional parameter of a realistically modelled data distribution based on observing independent and identically distributed observations. We define an  $S$ -th order **Highly Adaptive Lasso** Minimum Loss Estimator (Spline-HAL-MLE) of a ...  
☆ Save ☆ Cite Cited by 18 Related articles All 2 versions 36

**Robust inference on the average treatment effect using the outcome highly adaptive lasso** [pdf] [wiley.com](#)  
D. J. Berstner, M. Van Der Laan - Biometrics, 2020 - [Wiley Online Library](#)  
Find it @ Emory  
Many estimators of the average effect of a treatment on an outcome require estimation of the propensity score, the outcome regression, or both. It is often beneficial to utilize flexible techniques, such as semiparametric regression or machine learning, to estimate these ...  
☆ Save ☆ Cite Cited by 9 Related articles All 8 versions

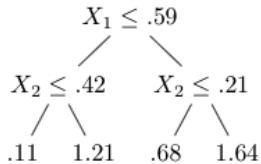
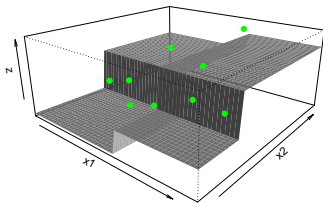
**Highly adaptive lasso (hal)**  
M. Van Der Laan, D. Berstner - Targeted Learning in Data Science, 2018 - Springer  
In this chapter, we define a general nonparametric estimator of ad-variate functional valued parameter  $\psi_0$ . This parameter is defined as a minimizer of an expectation of a loss function  $L(\psi; \mathcal{D})$  that is guaranteed to converge to the true  $\psi_0$  at a rate faster than  $n^{-1/4}$ , for all ...  
☆ Save ☆ Cite Cited by 4 Related articles

**Uniform consistency of the highly adaptive lasso estimator of infinite dimensional parameters** [pdf] [arxiv.org](#)  
M. Van Der Laan, A. E. Bibaud - arXiv preprint arXiv:1709.06226, 2017 - [arxiv.org](#)  
Consider the case that we observe  $S \times n$  independent and identically distributed copies of a random variable with a probability distribution known to be an element of a specified statistical model. We are interested in estimating an infinite dimensional target parameter ...  
☆ Save ☆ Cite Cited by 8 Related articles All 2 versions 36

**Nonparametric bootstrap inference for the targeted highly adaptive LASSO** [pdf] [arxiv.org](#)

It's popular with one of the biggest causal inference groups out there. We'll touch on this, too.

Random forests are a computationally-efficient approximation to BV regression.



- Forests, too, are constant except for jumps.
- TV regression uses the jumps that minimize mean squared error subject to a constraint on the curve's total variation  $\rho_{TV}(\hat{\mu})$ .
- This is fast enough in 1D or 2D, but can get slow high dimensions.
- Random forests use a greedy algorithm that tries to minimize MSE with constraints that are similar in intention and effect.

## Variations

---

# There's more than one way to be a smooth curve.

1. Often, we think the curve we're looking for shouldn't jump around all the time.
2. Sometimes, we think it shouldn't jump at all.  
And that it shouldn't even get close to making big jumps.
3. Maybe we'll even think it shouldn't get steep at all.

These are all different ways to be smooth. Each is stronger than the last.

$$\max_{x \in [0,1]} |m'(x)| \leq B \implies \sqrt{\int_0^1 |m'(x)|^2} \leq B \implies \int_0^1 |m'(x)| \leq B.$$

Lipschitz smooth (3)                      Sobolev smooth (2)                      Bounded Variation (1)

In many ways these kinds of models are similar, but they all have their quirks. Which you use makes a difference, especially for behavior near the data's edges.

- We'll start talking about Lipschitz smoothness in this week's homework.
- And we'll add Sobolev smoothness to the mix later.

## References

---

Joshua D Angrist and Victor Lavy. Using maimonides' rule to estimate the effect of class size on scholastic achievement. *The Quarterly journal of economics*, 114(2): 533–575, 1999.