

DSC 40A

Theoretical Foundations of Data Science I

In This Video


Which prediction minimizes the mean error?

Recommended Reading

Course Notes: Chapter 1, Section 1

The Best Prediction

- ▶ We want the best prediction, h^* .
- ▶ Goal: find h that minimizes the mean error:

$$R(h) = \frac{1}{n} \sum_{i=1}^n |y_i - h|$$


- ▶ This is an optimization problem.

Question

Can we use calculus to minimize R ?

Minimizing with Calculus

- Calculus: take derivative, set equal to zero, solve.

$$R(h) = \left(\frac{1}{n}\right) \sum_{i=1}^n |h - y_i|$$

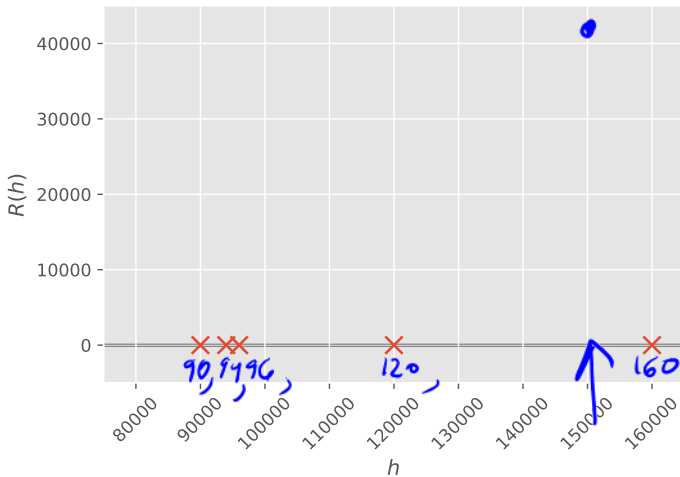
$$R'(h) = \frac{1}{n} \left(\sum_{i=1}^n \frac{d}{dh} (|h - y_i|) \right)$$

problem:
not differentiable
because
has cusp



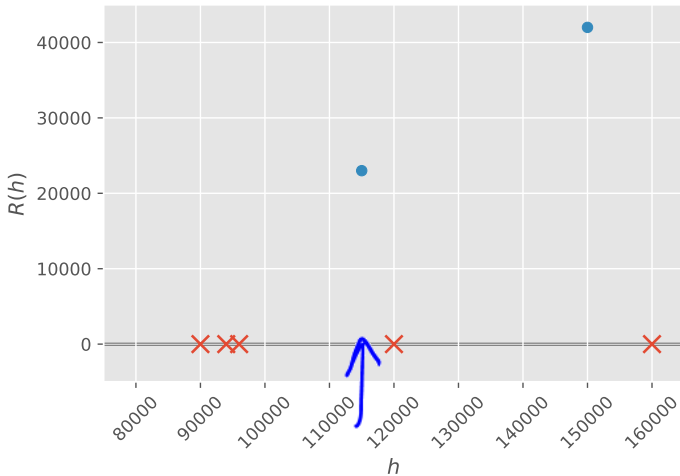
Can't be minimized with calc.

Plotting the Mean Error



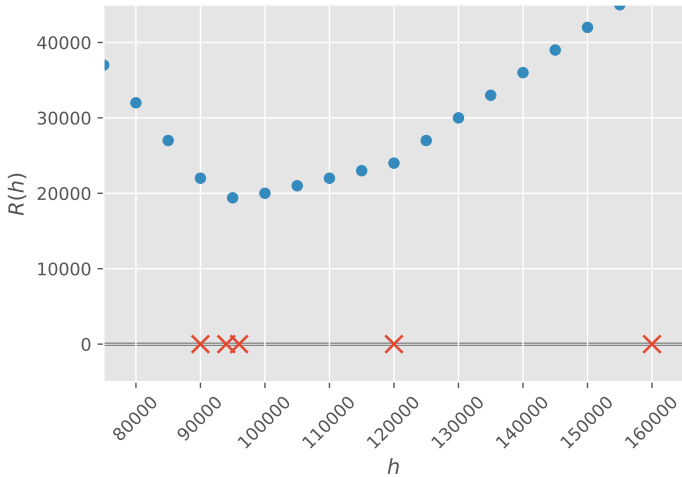
Recall: $R(\underline{150,000}) = 42,000$ ← mean error

Plotting the Mean Error

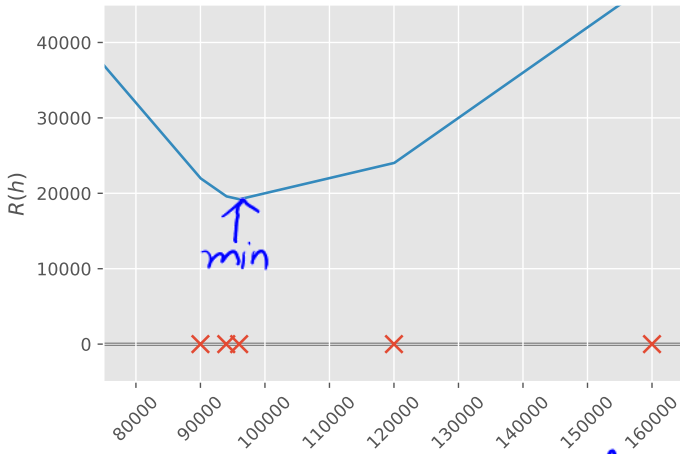


Recall: $R(115,000) = 23,000$

Plotting the Mean Error



Plotting the Mean Error



1) continuous

2) made of line segments (piecewise linear)

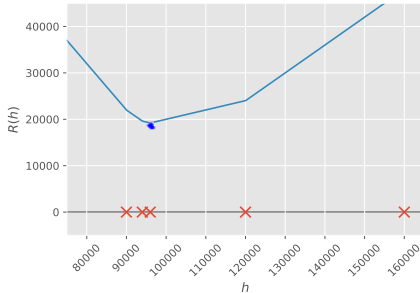
Question

A local minimum occurs when the slope of a function goes from _____ . Select all that apply.

- A) positive to negative
- B) negative to positive
- C) positive to zero
- D) negative to zero



Goal



- ▶ Find where slope of R goes from negative to non-negative.
- ▶ Want a formula for the slope of R at h .

pos or
zero

or

Sums of Linear Functions

► Let

$$f_1(x) = 3x + 7 \quad f_2(x) = 5x - 4 \quad f_3(x) = -2x - 8$$

► What is the slope of $f(x) = f_1(x) + f_2(x) + f_3(x)$?

$$3x + 5x - 2x + c$$

$$6x + c$$

↑
slope is 6

Sums of Absolute Values

► Let

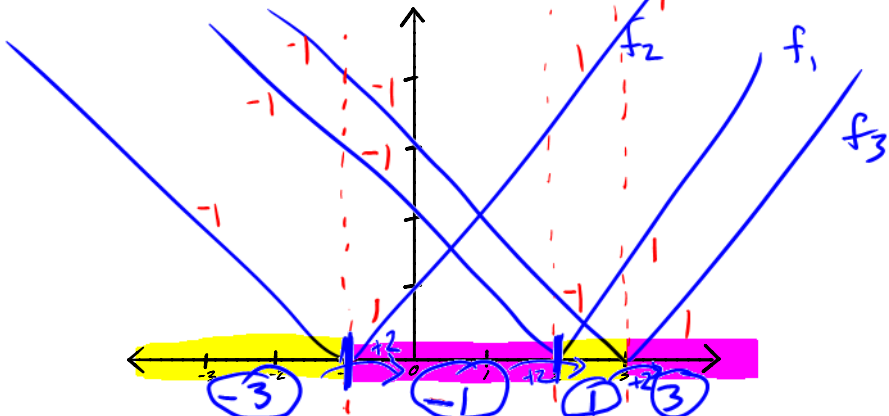
$$-1 \sqrt{\quad}$$

$$f_1(x) = |x - \underbrace{2}_{\substack{\circledast \\ 2}}$$

$$f_2(x) = |x + \underbrace{1}_{\substack{\circledast \\ -1}}$$

$$f_3(x) = |x - \underbrace{3}_{\substack{\circledast \\ 3}}$$

► What is the slope of $f(x) = f_1(x) + f_2(x) + f_3(x)$?




The Slope of the Mean Error

$R(h)$ is a sum of absolute value functions (times $\frac{1}{n}$):

$$R(h) = \frac{1}{n} (|h - y_1| + |h - y_2| + \dots + |h - y_n|)$$

$$R(h) = \frac{1}{n} \sum_{i=1}^n |h - y_i|$$

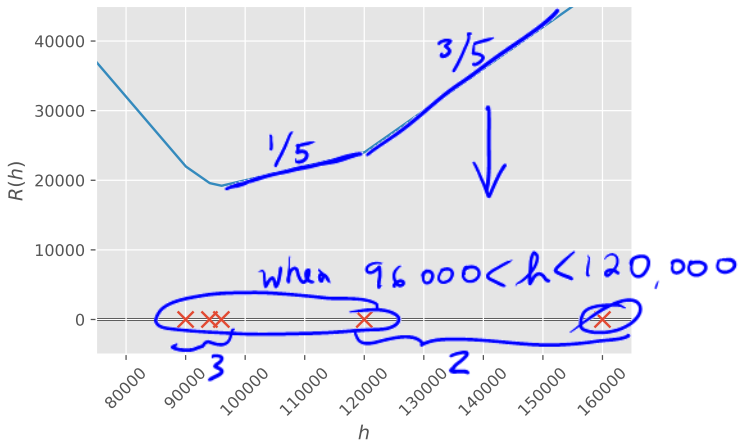

$$= \frac{1}{n} \left(\sum_{\substack{y_i < h \\ \text{pos}}} |h - y_i| + \sum_{\substack{y_i > h \\ \text{neg}}} |h - y_i| + \sum_{\substack{y_i = h \\ 0}} |h - y_i| \right)$$

$$= \frac{1}{n} \left(\underbrace{\sum_{y_i < h} (h - y_i)}_{\substack{\text{slope of} \\ \text{each term} \\ \text{is } 1}} + \sum_{y_i > h} \underbrace{-\frac{y_i - h}{(h - y_i)}}_{\substack{\text{slope of} \\ \text{each term} \\ \text{is } -1}} \right)$$

The Slope of the Mean Error

The slope of R at h is:

$$\frac{1}{n} [(\# \text{ of } y_i\text{'s} < h) - (\# \text{ of } y_i\text{'s} > h)]$$



Where the Slope's Sign Changes

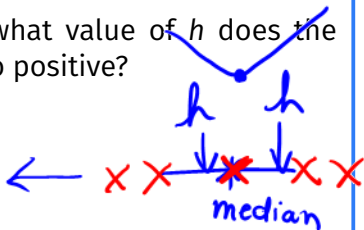
The slope of R at h is:

$$\frac{1}{n} \cdot [(\# \text{ of } y_i\text{'s} < h) - (\# \text{ of } y_i\text{'s} > h)]$$

Question

Suppose that n is odd. At what value of h does the slope of R go from negative to positive?

- A) $h = \text{mean of } y_1, \dots, y_n$
- B) $h = \text{median of } y_1, \dots, y_n$
- C) $h = \text{mode of } y_1, \dots, y_n$



Summary: The Median Minimizes the Mean Error

- ▶ Our problem was: find h^* which minimizes the mean error,

$$R(h) = \frac{1}{n} \sum_{i=1}^n |y_i - h|.$$

- ▶ The answer is: $\text{Median}(y_1, \dots, y_n)$.
- ▶ The **best prediction**¹ is the **median**.
- ▶ **Next time:** We consider a different measure of error that is differentiable.

¹in terms of mean error