
DSC 40A - Extra Practice for Final Part 1
Spring 2023

Problem 1. Sloped Mean

Suppose you have a data set y_1, y_2, \dots, y_n with at least three values, $n \geq 3$, and the values are arranged such that $y_1 \leq y_2 \leq \dots \leq y_n$.

We know from class that the mean of the data minimizes mean squared error,

$$R_{sq}(h) = \sum_{i=1}^n (h - y_i)^2.$$

Define a new function that weights larger data points less heavily:

$$S(h) = \left(\sum_{i=1}^{n-2} (h - y_i)^2 \right) + 0.5 \cdot (h - y_{n-1})^2 + 0.1 \cdot (h - y_n)^2.$$

- a) What value of h minimizes $S(h)$? We'll call the value of h that minimizes $S(h)$ the **sloped mean**, since the coefficients of the data values decrease for larger data.

- b) Which do you think is a better hypothesis, the mean or the sloped mean? Is your answer always the same, or does it depend on some property of the data set? Give an example of when you might prefer to use the sloped mean, and when you might prefer the (regular) mean.

Problem 2. Which is bigger? By how much?

Given a data set $y_1 \leq y_2 \leq \dots \leq y_n$, define the following empirical risk functions:

MAE

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

$$R_{\text{sq}}(h) = \frac{1}{n} \sum_{i=1}^n (y_i - h)^2$$

MSE

avg dist of each data point to prediction

Parts (a), (b), and (c) below concern R_{abs} . Parts (d) and (e) concern R_{sq} .

- a) For an arbitrary c with $c < c+1 < y_1$, how does $R_{\text{abs}}(c)$ compare to $R_{\text{abs}}(c+1)$? Can you determine which is bigger, and by how much?

$R_{\text{abs}}(c+1)$: each distance is 1 unit less

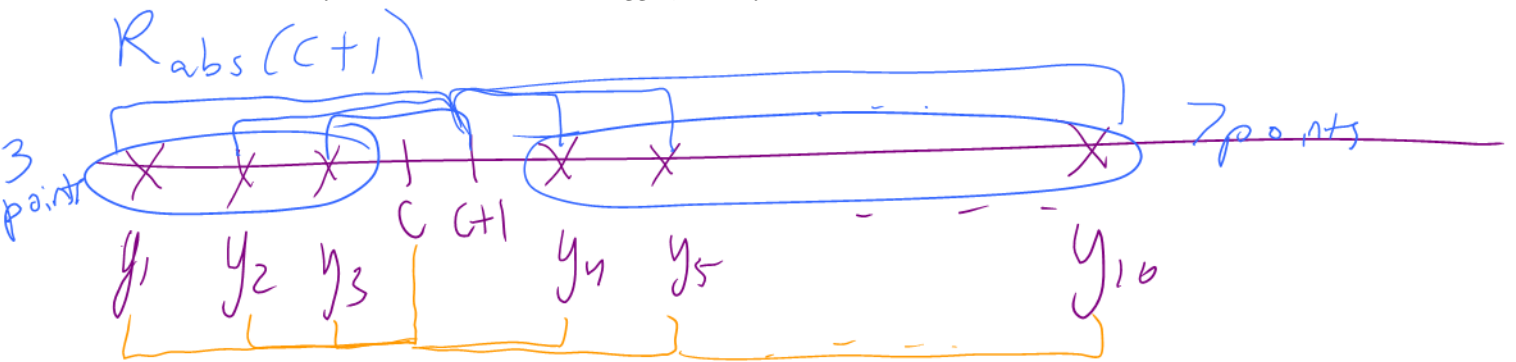


$R_{\text{abs}}(c)$

$$R_{\text{abs}}(c+1) = R_{\text{abs}}(c) - 1$$

- b) For an arbitrary c with $y_n < c < c + 2$, how does $R_{\text{abs}}(c)$ compare to $R_{\text{abs}}(c + 2)$? Can you determine which is bigger, and by how much?

- c) Suppose $n = 10$. For an arbitrary c with $y_3 < c < c+1 < y_4$, how does $R_{\text{abs}}(c)$ compare to $R_{\text{abs}}(c+1)$? Can you determine which is bigger, and by how much?



$R_{\text{abs}}(c)$

3 points on left: $c+1$ is 1 unit further from each

7 points on right: $c+1$ is 1 unit closer to each

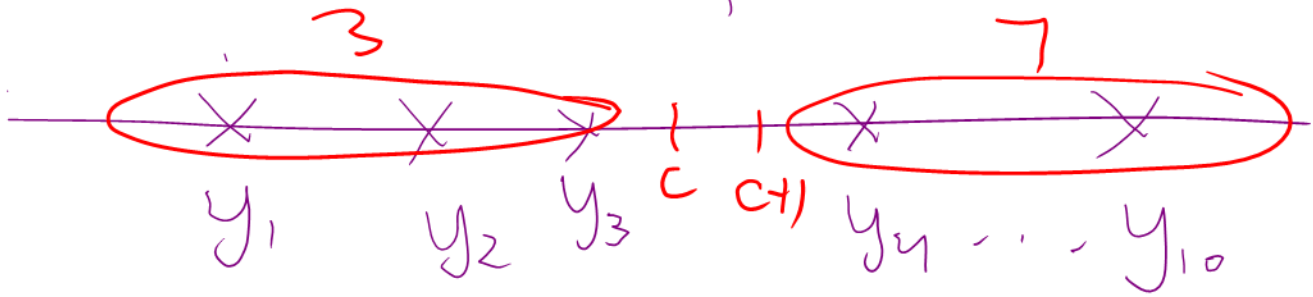
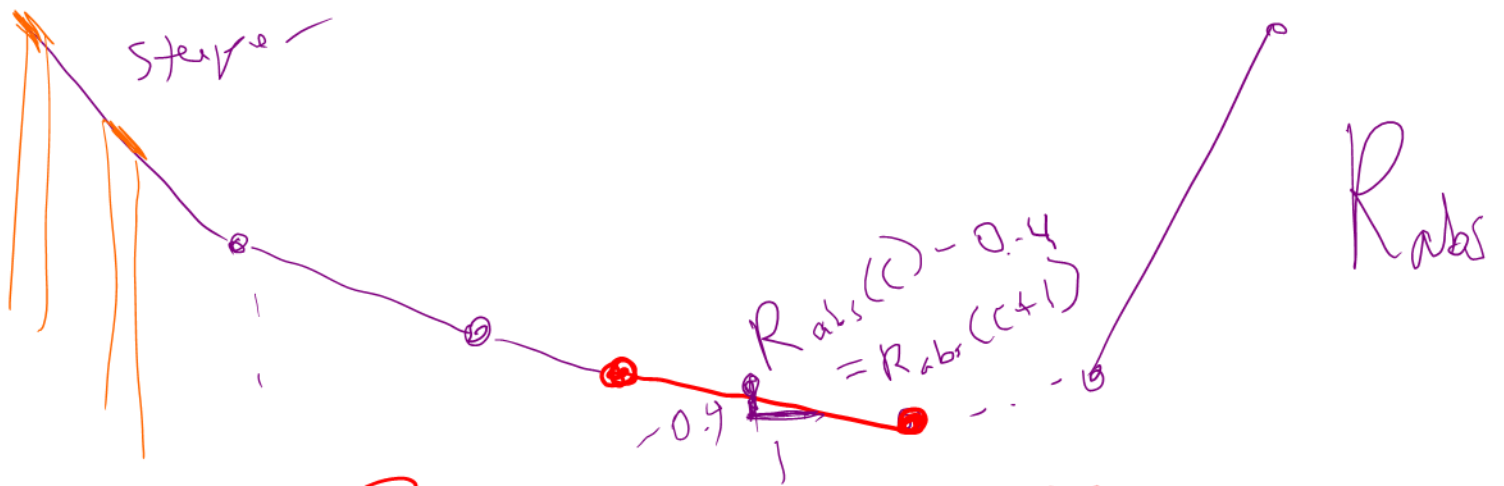
$$R_{\text{abs}}(c+1) = \frac{1}{n} \sum_{i=1}^n |y_i - (c+1)|$$

$$\approx \frac{1}{n} \left(\begin{array}{l} \text{total dist} \\ \text{of } c+1 \text{ to} \\ \text{each } y_i \end{array} \right)$$

$$= \frac{1}{10} \left(\begin{array}{l} \text{total dist of } c \text{ to} \\ \text{each } y_i \end{array} \right) \underbrace{+ 3 - 7}_{-4}$$

$$= \frac{1}{10} \left(\begin{array}{l} \text{total dist} \\ \text{of } c \text{ to} \\ \text{each } y_i \end{array} \right) - \frac{4}{10}$$

$$= R_{\text{abs}}(c) - \frac{4}{10}$$



formula for slope of R_{abs} at h

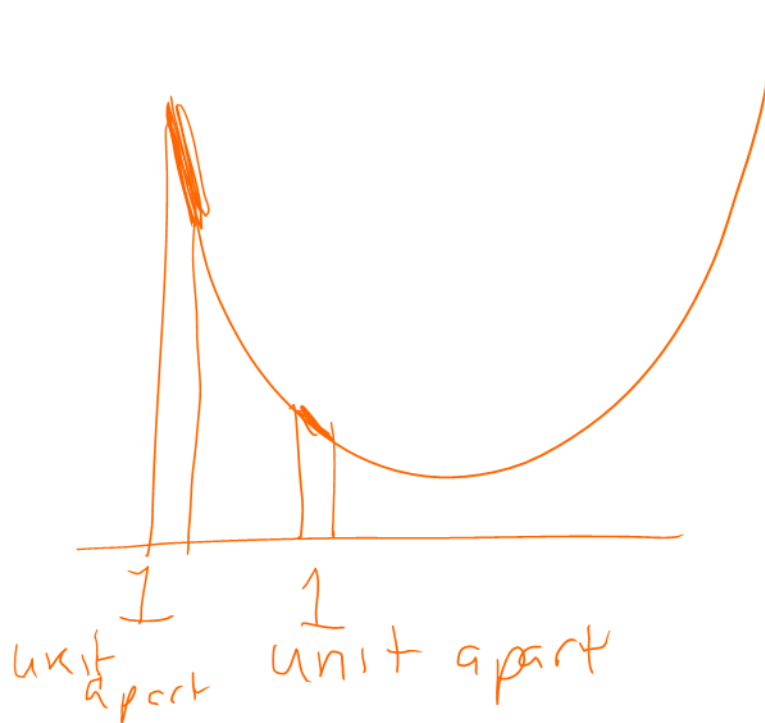
$$\frac{1}{n} \left(\begin{array}{l} \# y_i < h \\ \text{(left)} \end{array} - \begin{array}{l} \# y_i > h \\ \text{(right)} \end{array} \right)$$

$$\frac{1}{10} (3 - 7) = -0.4$$

d) For an arbitrary c with $c < y_1$, how does $R_{sq}(c)$ compare to $R_{sq}(c-1)$? Can you determine which is bigger, and by how much?

$$\underline{R_{sq}(c-1) > R_{sq}(c)}$$

R_{sq}



- e) For an arbitrary c with $c > y_n$, how does $R_{\text{sq}}(c)$ compare to $R_{\text{sq}}(c+1)$? Can you determine which is bigger, and by how much?

Problem 3. Matrix, Vector, Scalar, or Nonsense?

Suppose M is an $m \times n$ matrix, v is a vector in \mathbb{R}^n , and s is a scalar. Determine whether each of the following quantities is a matrix, vector, scalar, or nonsense (undefined).

a) Mv

b) vM

c) v^2

d) $M^T M$

e) MM^T

f) $v^T Mv$

g) $(sMv) \cdot (sMv)$

h) $(sv^T M^T)^T$

i) $v^T M^T Mv$

j) $vv^T + M^T M$

Problem 4. Orthogonality

- a) Is it possible for a vector to be orthogonal to itself?

- b) Show that if \vec{u} is orthogonal to both \vec{v} and \vec{w} , then \vec{u} is also orthogonal to any linear combination of \vec{v} and \vec{w} , $\alpha\vec{v} + \beta\vec{w}$.

- c) Show that if $A^T \vec{b} = \mathbf{0}$, then \vec{b} is orthogonal to the **column space** of A , which is the space of all linear combinations of the columns of A .

$$SD(x) = \sqrt{\frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^2}$$

Problem 5. Regression

Suppose you have a dataset

$$\{(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n)\}$$

where the standard deviation of the x -values, $SD(x)$, is twice the standard deviation of the y -values, $SD(y)$.
Let

$$y = a + bx$$

} Usual way

be the regression line with x as the predictor variable and y as the response variable. Let

$$x = c + dy$$

} Another way

be the regression line with y as the predictor variable and x as the response variable. **Express b in terms of d .**

$$b = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sum_{i=1}^n (x_i - \bar{x})^2}$$

$$d = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sum_{i=1}^n (y_i - \bar{y})^2}$$

Look at

$$\frac{b}{d} = \frac{\sum_{i=1}^n (y_i - \bar{y})^2}{\sum_{i=1}^n (x_i - \bar{x})^2}$$

$$\Rightarrow \frac{b}{d} = \frac{n \cdot (SD(y))^2}{n \cdot (SD(x))^2}$$

$$\Rightarrow \frac{b}{d} = \frac{(SD(y))^2}{(2 \cdot SD(y))^2}$$

$$\Rightarrow \frac{b}{d} = \frac{\cancel{SD(y)}^2}{4 \cdot \cancel{SD(y)}^2}$$

$$\Rightarrow \frac{b}{d} = \frac{1}{4}$$

$$\Rightarrow \boxed{b = d/4}$$

$$b = r \cdot \frac{SD(y)}{SD(x)}, \quad d = r \cdot \frac{SD(x)}{SD(y)}$$

given $SD(x) = 2 SD(y)$

$$b = r \cdot \frac{\cancel{SD(y)}}{2 \cancel{SD(y)}}$$

$$d = r \cdot \frac{2 \cancel{SD(y)}}{\cancel{SD(y)}}$$

$$b = \frac{r}{2}$$

$$d = 2r$$

$$\frac{b}{d} = \frac{\frac{r}{2}}{2r} = \frac{\cancel{r}}{4\cancel{r}} = \frac{1}{4}$$

$$\boxed{b = d/4}$$

Problem 6. Farmfluencer

Billy the avocado farmer heard about the success of 72 year-old Gerald Stratford's viral gardening videos on Twitter and Instagram. After witnessing Gerald turn into the so-called [King of Big Veg](#) overnight, Billy is feeling inspired to up his social media game (he's also feeling a little bit jealous).

Billy is new to Instagram and is trying to understand how people gain followers. In particular, he wants to be able to predict the number of followers, y , based on these features:

- number of people they follow, $x^{(1)}$
 - number of years since first post, $x^{(2)}$
 - average number of posts per day, $x^{(3)}$
- a) Suppose Billy has access to a large data set of Instagram accounts, and he uses multiple regression on this data to fit a linear prediction rule of the form

$$H(\vec{x}) = w_0 + w_1x^{(1)} + w_2x^{(2)} + w_3x^{(3)}.$$

What does w_2 represent in terms of Instagram followers?

- b) What if instead of the number of years since the first post, $x^{(2)}$, Billy instead uses the number of days since the first post, $x^{(4)}$. Now he uses multiple regression to fit a prediction rule of the form

$$H'(\vec{x}) = w'_0 + w'_1x^{(1)} + w'_3x^{(3)} + w'_4x^{(4)}.$$

How do the parameters of this prediction rule (w'_0, w'_1, w'_3, w'_4) compare to the parameters of original prediction rule (w_0, w_1, w_2, w_3) ?

Problem 7. Changing the Prediction Rule

Suppose we have a dataset consisting of variables $x^{(1)}, x^{(2)}$, and y . We use multiple regression to fit a prediction rule of the form

$$H(x^{(1)}, x^{(2)}) = w_0 + w_1(x^{(1)} + x^{(2)}) + w_2 x^{(1)} x^{(2)} + w_3(x^{(1)} + 1)(x^{(2)} + 1) \quad (1)$$

and then again use multiple regression to fit a different prediction rule of the form

$$H'(x^{(1)}, x^{(2)}) = w'_0 + w'_1 x^{(1)} + w'_2 x^{(2)} + w'_3 x^{(1)} x^{(2)} \quad (2)$$

Which form, (1) or (2), will yield a prediction rule with lower mean squared error? Justify your answer.

can have different coeff for $x^{(1)}, x^{(2)}$

H:

$$w_0 + w_1 x^{(1)} + w_1 x^{(2)} + w_2 x^{(1)} x^{(2)}$$

$$+ w_3 (x^{(1)} x^{(2)} + x^{(1)} + x^{(2)} + 1)$$

$$= (w_0 + w_3) + (w_1 + w_3) x^{(1)} + (w_1 + w_3) x^{(2)}$$

$$+ (w_2 + w_3) x^{(1)} x^{(2)}$$

$$= C_0 + C_1 x^{(1)} + C_1 x^{(2)} + C_2 x^{(1)} x^{(2)}$$

forced to same coeff for $x^{(1)}$ and $x^{(2)}$

$$MSE(H'_{(2)}) \leq MSE(H)$$