# Midterm Logistics

## Format

80 minutes, on paper, no calculators or electronics
- Allowed one two-sided index card (handwritten)

## Content

Lectures 1-9, Homeworks 1-4, Groupworks 1-4, feature engineering, and transformations
- You can study with the practice site!

## So… When is it?

Midterm is Tuesday, May 7th in **your** section
- You will receive a randomized seat assignment over the weekend

# Index Card in Depth

## Example Card!

- Posted on Edstem
- On the course website

- We will not print it for you
- You cannot print it either!
- However, you can reference it!
    - And we encourage it!

## Important Notes

- You can write on both sides of the index card
- You must write the index card by hand
    - No typing or writing with a tablet
- You need to find/buy a 4 inch by 6 inch index card
    - OK if 3 in by 5 in
    - Worst case you cut paper to that size

# Problem 1.1 from FA21 Midterm

King Triton just made an Instagram account and has been keeping track of the number of likes his posts have received so far.

His first 7 posts have received a mean of 16 likes; the specific like counts in sorted order are

$$8, 12, 12, 15, 18, 20, 27$$

King Triton wants to predict the number of likes his next post will receive, using a constant prediction rule $h$. For each loss function $L(h, y)$, determine the constant prediction $h^*$ that minimizes empirical risk. If you believe there are multiplie minimizers, specify them all. If you believe you need more information to answer the question or that there is no minimizer, state that clearly. **Give a brief justification for each answer.**

$$L(h, y) = |y - h|$$

empirical risk
w/ absolute loss

: median

$\boxed{15}$

# Problem 1.2 from FA21 Midterm

King Triton just made an Instagram account and has been keeping track of the number of likes his posts have received so far.

His first 7 posts have received a mean of 16 likes; the specific like counts in sorted order are

$$8, 12, 12, 15, 18, 20, 27$$

King Triton wants to predict the number of likes his next post will receive, using a constant prediction rule $h$. For each loss function $L(h, y)$, determine the constant prediction $h^*$ that minimizes empirical risk. If you believe there are multiplie minimizers, specify them all. If you believe you need more information to answer the question or that there is no minimizer, state that clearly. **Give a brief justification for each answer.**

$$L(h, y) = (y - h)^2$$

empirical risk
for square loss:

minimized by
mean.

(16)

King Triton just made an Instagram account and has been keeping track of the number of likes his posts have received so far.

His first 7 posts have received a mean of 16 likes; the specific like counts in sorted order are

$$8, 12, 12, 15, 18, 20, 27$$

King Triton wants to predict the number of likes his next post will receive, using a constant prediction rule $h$. For each loss function $L(h, y)$, determine the constant prediction $h^*$ that minimizes empirical risk. If you believe there are multiplie minimizers, specify them all. If you believe you need more information to answer the question or that there is no minimizer, state that clearly. **Give a brief justification for each answer.**

$$L(h, y) = 4(y - h)^2$$

empirical risk for

scaled squared error

: still minimized by

mean!

16

Nick

# Problem 1.4 from FA21 Midterm

King Triton just made an Instagram account and has been keeping track of the number of likes his posts have received so far.

His first 7 posts have received a mean of 16 likes; the specific like counts in sorted order are

$$8, 12, 12, 15, 18, 20, 27$$

King Triton wants to predict the number of likes his next post will receive, using a constant prediction rule $h$. For each loss function $L(h, y)$, determine the constant prediction $h^*$ that minimizes empirical risk. If you believe there are multiplie minimizers, specify them all. If you believe you need more information to answer the question or that there is no minimizer, state that clearly. **Give a brief justification for each answer.**

$$L(h, y) = \begin{cases} 0 & h = y \\ 100 & h \neq y \end{cases}$$

empirical risk
for scaled 0-1 loss
: minimized by
mode.

Nick

# Problem 1.5 from FA21 Midterm

King Triton just made an Instagram account and has been keeping track of the number of likes his posts have received so far.

His first 7 posts have received a mean of 16 likes; the specific like counts in sorted order are

$$8, 12, 12, 15, 18, 20, 27$$

King Triton wants to predict the number of likes his next post will receive, using a constant prediction rule $h$. For each loss function $L(h, y)$, determine the constant prediction $h^*$ that minimizes empirical risk. If you believe there are multiplie minimizers, specify them all. If you believe you need more information to answer the question or that there is no minimizer, state that clearly. **Give a brief justification for each answer.**

$$L(h, y) = (3y - 4h)^2$$

empirical risk w/ this loss? complicated...

$(3y - 4h)^2$

$\left(3\left(y - \frac{4}{3}h\right)\right)^2$

$3^2\left(y - \frac{4}{3}h\right)^2 \rightarrow 9\left(y - \frac{4}{3}h\right)^2$

does not matter

$\frac{4}{3}h = \text{mean}$

$h^* = \frac{3}{4} \text{mean}$

$\boxed{h = 12}$

want this to be $\bar{x}$

Nick

# Problem 1.6 from FA21 Midterm

King Triton just made an Instagram account and has been keeping track of the number of likes his posts have received so far.

His first 7 posts have received a mean of 16 likes; the specific like counts in sorted order are
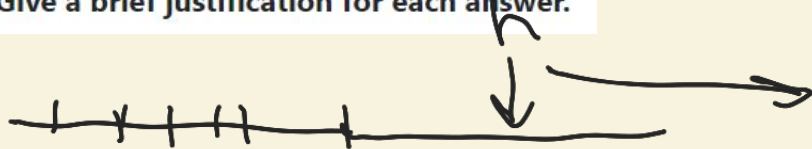
$$8, 12, 12, 15, 18, 20, 27$$

King Triton wants to predict the number of likes his next post will receive, using a constant prediction rule $h$. For each loss function $L(h, y)$, determine the constant prediction $h^*$ that minimizes empirical risk. If you believe there are multiplie minimizers, specify them all. If you believe you need more information to answer the question or that there is no minimizer, state that clearly. **Give a brief justification for each answer.**

$$L(h, y) = (y - h)^3$$

No minimizer

as h increases,
the Loss approaches
$-\infty$

Nick

**Problem 1.**

The table below shows the softness and color of several different avocados, which we want to use to predict their ripeness. Each variable is measured on a scale of 1 to 5. For softness, 5 is softest, for color, 5 is darkest, and for ripeness, 5 is ripest.

| Avocado | Softness | Color | Ripeness |
|---------|----------|-------|----------|
| 1 | 3 | 4 | 2.5 |
| 2 | 1 | 2 | 2 |
| 3 | 4 | 5 | 5 |

Suppose we have decided on the following hypothesis function: given an avocado's softness and color, we average these numbers to produce a predicted ripeness.

**f)** Write down the *design matrix, X*.

$$n = 3$$
$$2 + 1 = 3$$
$$\boxed{3 \times 3}$$

$$X = \begin{bmatrix} 1 & 3 & 4 \\ 1 & 1 & 2 \\ 1 & 4 & 5 \end{bmatrix}$$

$$3 \times 1$$

$$y = \begin{bmatrix} 2.5 \\ 2 \\ 5 \end{bmatrix}$$

Harshi

**Problem 1.**

The table below shows the softness and color of several different avocados, which we want to use to predict their ripeness. Each variable is measured on a scale of 1 to 5. For softness, 5 is softest, for color, 5 is darkest, and for ripeness, 5 is ripest.

| Avocado | Softness | Color | Ripeness |
|---|---|---|---|
| 1 | 3 | 4 | 2.5 |
| 2 | 1 | 2 | 2 |
| 3 | 4 | 5 | 5 |

$$H(x) = w_0 + w_1 x_1 + \ldots + w_d x_d$$

$$f(x) = 0 + \frac{1}{2} x_1 + \frac{1}{2} x_2$$
softness    color

Suppose we have decided on the following hypothesis function: given an avocado's softness and color, we average these numbers to produce a predicted ripeness.

**g)** Write down the *parameter vector*, $\vec{w}$ that corresponds to this particular choice of hypothesis function. The parameter vector should have three components, one for the bias, and one for each of the

$w \ (d+1) \times 1$

$x \ n \times (d+1)$

$$w = \begin{bmatrix} w_0 \\ w_1 \\ w_2 \end{bmatrix} = \begin{bmatrix} 0 \\ 0.5 \\ 0.5 \end{bmatrix}$$
$3 \times 1$

a hypothesis func, not $w^*$

prediction
$= \frac{\text{softness} + \text{color}}{2}$
$= \frac{1}{2} \text{softness} + \frac{1}{2} \text{color}$

**Harshi**

$X$ 3×3

$w$ 3×1

(3×3) (3)(1)

$X$      $w$

3×1

**Problem 1.**

The table below shows the softness and color of several different avocados, which we want to use to predict their ripeness. Each variable is measured on a scale of 1 to 5. For softness, 5 is softest, for color, 5 is darkest, and for ripeness, 5 is ripest.

$X = n \times (d+1)$

$w = (d+1) \times 1$

prod : $n \times 1$

| Avocado | Softness | Color | Ripeness |
|---|---|---|---|
| 1 | 3 | 4 | 2.5 |
| 2 | 1 | 2 | 2 |
| 3 | 4 | 5 | 5 |

$H(x) =$

prediction:
$\begin{bmatrix} 3.5 \\ 1.5 \\ 4.5 \end{bmatrix}$

Suppose we have decided on the following hypothesis function: given an avocado's softness and color, we average these numbers to produce a predicted ripeness.

**h)** Check that the entries of $X\vec{w}$ are the predicted ripenesses you found above.

3×3

$X = \begin{bmatrix} 1 & 3 & 4 \\ 1 & 1 & 2 \\ 1 & 4 & 5 \end{bmatrix}$

$w = \begin{bmatrix} 0 \\ 0.5 \\ 0.5 \end{bmatrix}$

Av1   $1 \times 0 + 3 \times 0.5 + 4 \times 0.5 = \dfrac{3+4}{2}$   $\begin{bmatrix} 3.5 \end{bmatrix}$

Av2   $1 \times 0 + 1 \times 0.5 + 2 \times 0.5 = \dfrac{1+2}{2} = 1.5$

Av3   $1 \times 0 + 4 \times 0.5 + 5 \times 0.5 = \dfrac{4+5}{2} = 4.5$

3×1

**Harshi**

Consider the dataset shown below.

| $x^{(1)}$ | $x^{(2)}$ | $x^{(3)}$ | $y$ |
|-----------|-----------|-----------|-----|
| 0 | 6 | 8 | -5 |
| 3 | 4 | 5 | 7 |
| 5 | -1 | -3 | 4 |
| 0 | 2 | 1 | 2 |

We want to use multiple regression to fit a prediction rule of the form

$$H(x^{(1)}, x^{(2)}, x^{(3)}) = w_0 + w_1 x^{(1)} x^{(3)} + w_2 (x^{(2)} - x^{(3)})^2$$

Write down the design matrix $X$ and observation vector $\vec{y}$ for this scenario. No justification needed.

$X$   $4 \times 3$

$$X = \begin{bmatrix} 1 & 0 & 4 \\ 1 & 15 & 1 \\ 1 & -15 & 4 \\ 1 & 0 & 1 \end{bmatrix} \qquad y = \begin{bmatrix} -5 \\ 7 \\ 4 \\ 2 \end{bmatrix}$$

$d = 2$

$d + 1 = 3$

$w_0 + w_1 \boxed{b_1} + w_2 \boxed{b_2}$

$x^1 x^3 \qquad (x^2 - x^3)^2$

For the $X$ and $\vec{y}$ that you have written down, let $\vec{w}$ be the optimal parameter vector, which comes from solving the normal equations $\underline{X^T X \vec{w} = X^T \vec{y}}$. Let $\underline{\vec{e} = \vec{y} - X\vec{w}}$ be the error vector, and let $e_i$ be the $i$th component of this error vector. Show that

$$4e_1 + e_2 + 4e_3 + e_4 = 0.$$

$$\boxed{X^T e = 0}$$

$$X = \begin{bmatrix} 1 & 0 & 4 \\ 1 & 15 & 1 \\ 1 & -15 & 4 \\ 1 & 0 & 1 \end{bmatrix}$$

$$X^T = \begin{bmatrix} 1 & 1 & 1 & 1 \\ 0 & 15 & -15 & 0 \\ 4 & 1 & 4 & 1 \end{bmatrix}$$

$$X^T e = 0$$

$$X^T e = \begin{bmatrix} 1 & 1 & 1 & 1 \\ 0 & 15 & -15 & 0 \\ 4 & 1 & 4 & 1 \end{bmatrix} \begin{bmatrix} e_1 \\ e_2 \\ e_3 \\ e_4 \end{bmatrix} = 0$$

$$\rightarrow 4e_1 + e_2 + 4e_3 + e_4 = 0$$

$$= \begin{bmatrix} e_1 + e_2 + e_3 + e_4 \\ 15e_2 - 15e_3 \\ 4e_1 + e_2 + 4e_3 + e_4 \end{bmatrix} = \begin{bmatrix} 0 \\ 0 \\ 0 \end{bmatrix}$$

Harshi

$$X^T X w = X^T y$$

$$\vec{e} = \vec{y} - X\vec{w}$$

$$X^T e = 0 \qquad X \text{ orth } e$$

$$0 = X^T y - X^T X w$$

$$0 = X^T (y - Xw)$$

$$0 = X^T e$$

In the National Football League (NFL), the highest paid position is the quarterback. The job of the quarterback on an (American) football team is, among other things, to throw (or "pass") the football to other players, who then score "touchdowns." Each time a quarterback throws the ball to another player and that other player scores, we say the quarterback made a "touchdown pass."

Suppose that we have access to a dataset containing information about a random sample of 50 quarterbacks. For each quarterback, we have the number of touchdown passes they threw, along with their salary in 2023. In the 2023 dataset, the number of touchdown passes for all quarterbacks has a mean of 17 and a standard deviation of 3.

We minimize mean squared error to fit a linear hypothesis function, $H(x) = w_0 + w_1 x$, to this dataset. We will use the hypothesis function to help other players predict their 2023 salary in millions of dollars ($y$) based on their number of touchdown passes ($x$).

**a)** 🥑🥑🥑🥑 CJ Stroud was one of the quarterbacks in our 2023 dataset. Suppose that in 2023, he had 26 touchdown passes and his salary was only 10 million, the smallest salary in our sample.

In 2024, Stroud signed a new contract based on his performance. In 2024, he again threw 26 touchdowns, but his salary shot up to 50 million!

Suppose we create two linear hypothesis functions, one using the dataset from 2023 when Stroud had a salary of 10 million and another using the dataset from 2024 when Stroud had a salary of 50 million. Assume that all other players threw the same amount of touchdowns and had the same salary in both datasets. **That is, only this one data point is different between these two datasets.**

Suppose the optimal slope and intercept fit on the first dataset (2023) are $w_1^*$ and $w_0^*$, respectively, and the optimal slope and intercept fit on the second dataset (2024) are $w_1'$ and $w_0'$, respectively.

What is the difference between the new slope and the old slope? **That is, what is $w_1' - w_1^*$?** The answer you get should be a number with no variables.

**Note**: Since we want to salary in millions of dollars, use 10 instead of 10,000,000 for Stroud's salary in 2023.

$$w_1^* = \frac{\sum_{i=1}^{n}(x_i - \bar{x})y_i}{\sum_{i=1}^{n}(x_i - \bar{x})^2}$$

Yosen

$$w_1^* = \frac{\sum\limits_{i=1}^{n}(x_i - \overline{x})y_i}{\sum\limits_{i=1}^{n}(x_i - \overline{x})^2}$$

Stroud: player $j$

$$W_1^* = \frac{\sum\limits_{i \neq j}^{n}(x_i - \overline{x})y_i + (x_j - \overline{x})y_j}{\sum\limits_{i=1}^{n}(x_i - \overline{x})^2}$$

$$W_1^{*\prime} = \frac{\sum\limits_{i \neq j}^{n}(x_i - \overline{x})y_i + (x_j - \overline{x})y_j'}{\sum\limits_{i=1}^{n}(x_i - \overline{x})^2}$$

$$W_1^{*\prime} - W_1^* = \frac{\left(\sum\limits_{i \neq j}^{n}(x_i - \overline{x})y_i + (x_j - \overline{x})y_j'\right) - \left(\sum\limits_{i \neq j}^{n}(x_i - \overline{x})y_i + (x_j - \overline{x})y_j\right)}{\sum\limits_{i=1}^{n}(x_i - \overline{x})^2}$$

$$= \frac{(x_j - \overline{x})(y_j' - y_j)}{\sum\limits_{i=1}^{n}(x_i - \overline{x})^2}$$

$$= \frac{(x_j - \bar{x})(y_j' - y_j)}{n\sigma^2}$$

$$\sigma^2 = \frac{\sum\limits_{i=1}^{n}(x_i - \bar{x})^2}{n}$$

$$x_j = 26$$

$$= \frac{(26 - 17)(50 \sim 10)}{50(3)^2}$$

$$= \frac{9 \cdot 40}{50 \cdot 9}$$

$$= \frac{4}{5}$$

**b)** 🥑🥑🥑 Let $H^*(x)$ be the linear hypothesis function fit on the 2023 dataset (i.e. $H^*(x) = w_0^* + w_1^* x$) and $H'(x)$ be the linear hypothesis function fit on the 2024 dataset (i.e. $H'(x) = w_0' + w_1' x$).
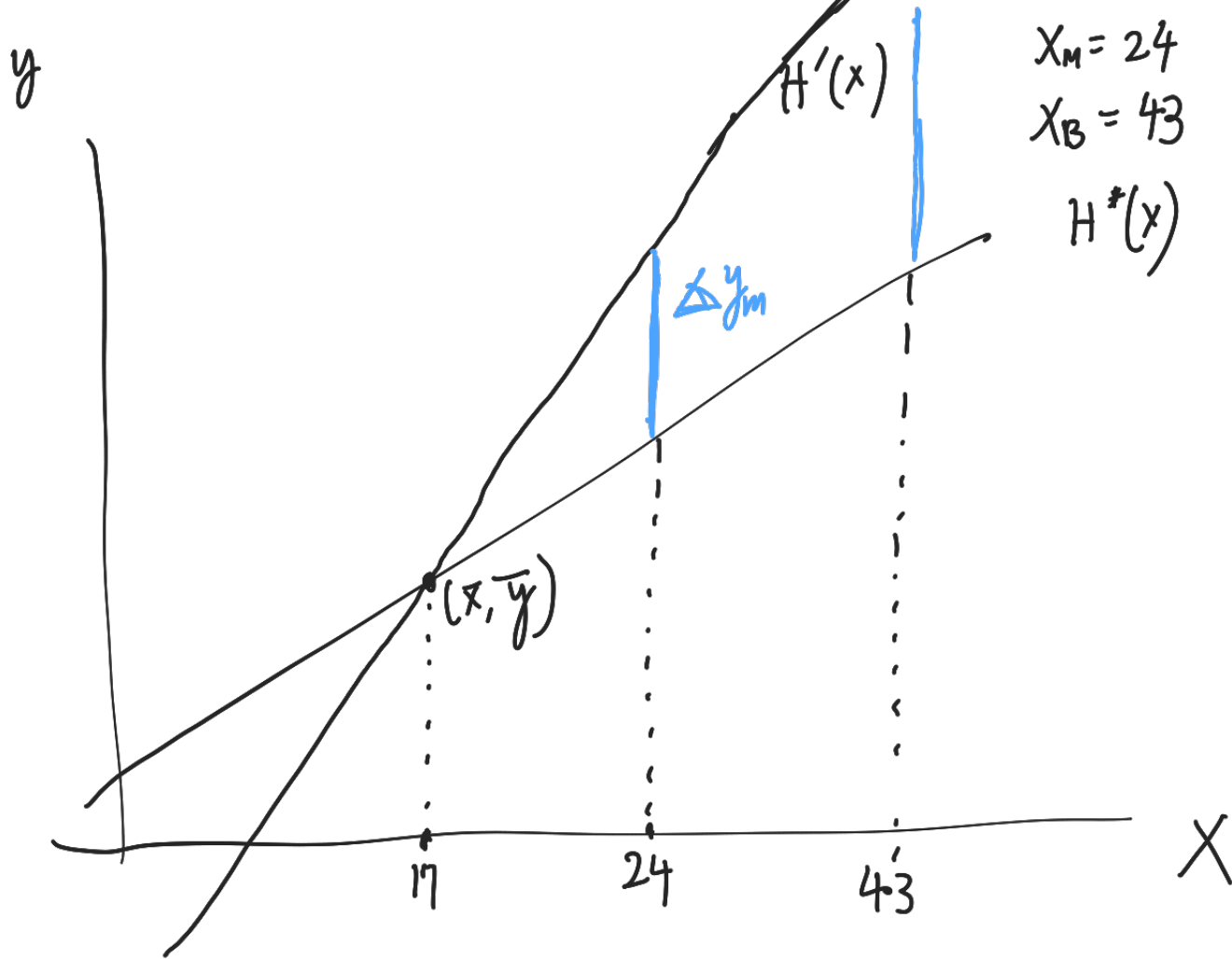
Consider two other players, Davis Mills and Tom Brady, neither of whom were part of our original sample in 2023. Suppose that in 2021, Mills had 24 touchdowns and Brady had 43 touchdowns.

Both Mills and Brady want to try and use one of our linear hypothesis functions to predict their salary for next year.

Suppose they both first use $H^*(x)$ to determine their predicted yields as per the first rule (when Stroud had a salary of 10 million). Then, they both then use $H'(x)$ to determine predicted yields as per the second rule (when Stroud had a salary of 50 million).

Whose prediction changed more by switching from $H^*(x)$ to $H'(x)$ – Mills' or Brady's?

$y$

$H'(x)$

$x_M = 24$
$x_B = 43$

$H^\#(x)$

$\Delta y_m$

$(x, \bar{y})$

17   24   43

$X$

c) 🥑🥑 In this problem, we'll consider how our answer to part (b) might have been different if Stroud had fewer touchdowns in both 2023 and 2024.

- Suppose Stroud instead had 17 touchdowns in both 2023 and 2024. If his salary increased from 2023 to 2024, and everyone else's data stayed the same, which slope would be larger: $H^*(x)$ or $H'(x)$?

- Suppose Stroud instead had 10 touchdowns in both 2023 and 2024. If his salary increased from 2023 to 2024, and everyone else's data stayed the same, which slope would be larger: $H^*(x)$ or $H'(x)$?

You don't have to actually calculate the new slopes, but given the information in the problem and the work you've already done, you should be able to answer the questions and give brief justification.

Stroud: player $j$

$$W_1^* = \frac{\sum_{i \neq j}^{n} (X_i - \bar{X}) y_i + (X_j - \bar{X}) y_j}{\sum_{i=1}^{n} (X_i - \bar{X})^2}$$

$$W_1^{*\,\prime} = \frac{\sum_{i \neq j}^{n} (X_i - \bar{X}) y_i + (X_j - \bar{X}) y_j'}{\sum_{i=1}^{n} (X_i - \bar{X})^2}$$

$$W_1^{*\prime} - W_1^* = \frac{\left(\sum_{i \neq j}^{n} (X_i - \bar{X}) y_i + (X_j - \bar{X}) y_j'\right) - \left(\sum_{i \neq j}^{n} (X_i - \bar{X}) y_i + (X_j - \bar{X}) y_j\right)}{\sum_{i=1}^{n} (X_i - \bar{X})^2}$$

$$= \frac{(X_j - \bar{X})(y_j' - y_j)}{\sum_{i=1}^{n} (X_i - \bar{X})^2}$$

1.

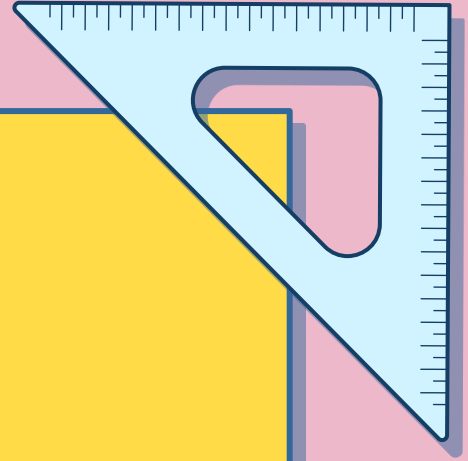$x_j' \; \cancel{\bigcirc} = 17, \quad \cancel{\bigcirc}\!\!_{x_j} = 10 \qquad x_j - \bar{x}$

$y_j = 9$

$$= \frac{(\underset{10}{\cancel{x}} - \underset{17}{\cancel{x}})( \; y_j' - y_j)}{\underset{i=1}{\overset{n}{\sum}} (X_i - \bar{x})^2}$$

$\overset{50-10}{}$

$\longrightarrow n\sigma^2$

2. $\quad X_j' = 10$

# 5 min Break Time!

Consider a dataset that consists of $y_1, \cdots, y_n$. In class, we used calculus to minimize mean squared error, $R_{sq}(h) = \frac{1}{n} \sum_{i=1}^{n} (h - y_i)^2$. In this problem, we want you to apply the same approach to a slightly different loss function defined below:

$$L_{\text{midterm}}(y, h) = (\alpha y - h)^2 + \lambda h$$

Write down the empiricial risk $R_{\text{midterm}}(h)$ by using the above loss function.

$$\star\ R(h) = \frac{1}{n} \sum_{i=1}^{n} L(y, h) \qquad R(h) = \frac{1}{n} \sum_{i=1}^{n} \left[ (\alpha y - h)^2 + \lambda h \right]$$

$$R(h) = \frac{1}{n} \sum_{i=1}^{n} (\alpha y - h)^2 + (\lambda h n) \rightarrow \frac{1}{n} \sum_{i=1}^{n} (\alpha y - h)^2 + \lambda h$$

Consider a dataset that consists of $y_1, \cdots, y_n$. In class, we used calculus to minimize mean squared error, $R_{sq}(h) = \frac{1}{n} \sum_{i=1}^{n} (h - y_i)^2$. In this problem, we want you to apply the same approach to a slightly different loss function defined below:

$$L_{\text{midterm}}(y, h) = \underbrace{(\alpha y - h)^2} + \lambda h$$

The mean of dataset is $\bar{y}$, i.e. $\bar{y} = \frac{1}{n} \sum_{i=1}^{n} y_i$. Find $h^*$ that minimizes $R_{\text{midterm}}(h)$ using calculus. Your result should be in terms of $\bar{y}$, $\alpha$ and $\lambda$.

$$\frac{\partial}{\partial h} \left[ f(g(x)) \right] = f'(g(x)) \, g'(x) \qquad R(h) = \frac{1}{n} \sum_{i=1}^{n} (\alpha y - h)^2 + \lambda h$$

$$f(x) = x^2 \qquad f'(x) = 2x$$

$$= \frac{2}{n} \sum_{i=1}^{n} (\alpha y_i - h) + \lambda R$$

$$(\alpha y - h)^2$$

$$g(x) = \alpha y - h \qquad g'(x) = -1$$

$$-2\alpha \left( \frac{1}{n} \sum_{i=1}^{n} y_i \right) + \frac{2}{n}(hn) + \lambda$$

Zoe

derivative
$$R(h) = -2x\bar{y} + 2h + \lambda$$

$$0 = -2xy + 2h + \lambda$$

$$\frac{2xy - \lambda = 2h}{2} \rightarrow h^* = x\bar{y} - \frac{\lambda}{2}$$

---

6.2

$$\vec{w}^* = \begin{bmatrix} 2000 \\ 10\,000 \\ -1000 \end{bmatrix} \begin{matrix} w_0 \\ w_1 \\ w_2 \end{matrix}$$ Carats: 0.65
length: 4 cm

predicted price $= w_0 + w_1(\text{carat}) + w_2(\text{length})$

$$2000 + 10\,000\underset{carat}{(0.65)} + (-1000)\underset{length}{(4)}$$

$$2000 + 6500 \quad -4000 = 4500$$

Billy's aunt owns a jewellery store, and gives him data on $5000$ of the diamonds in her store. For each diamond, we have:

- **carat**: the weight of the diamond, in carats
- **length**: the length of the diamond, in centimeters
- **width**: the width of the diamond, in centimeters
- **price**: the value of the diamond, in dollars

The first 5 rows of the 5000-row dataset are shown below:

| | carat | length | width | price |
|---|---|---|---|---|
| 1 | 0.40 | 4.81 | 4.76 | 1323 |
| 2 | 1.04 | 6.58 | 6.53 | 5102 |
| 3 | 0.40 | 4.74 | 4.76 | 696 |
| 4 | 0.40 | 4.67 | 4.65 | 798 |
| 5 | 0.50 | 4.90 | 4.95 | 987 |

Billy has enlisted our help in predicting the price of a diamond given various other features.

Suppose we want to fit a linear prediction rule that uses two features, carat and length, to predict price. Specifically, our prediction rule will be of the form

$$\text{predicted price} = \underset{}{w_0} + w_1 \cdot \underset{x_1}{\underline{\text{carat}}} + w_2 \cdot \underset{x_2}{\underline{\text{length}}}$$

each row $= \begin{bmatrix} 1 & \text{carat} & \text{length} \end{bmatrix}$

$\underset{5000}{\underbrace{}}$

We will use least squares to find $\vec{w}^* = \begin{bmatrix} w_0^* \\ w_1^* \\ w_2^* \end{bmatrix}$.

Write out the first 5 rows of the design matrix, $X$. Your matrix should not have any variables in it.

$$X = \begin{bmatrix} 1 & 0.4 & 4.81 \\ 1 & 1.04 & 6.58 \\ 1 & 0.4 & 4.74 \\ 1 & 0.4 & 4.67 \\ 1 & 0.5 & 4.90 \end{bmatrix} r \times c$$

$r = 5 \quad c = 3$

$(d+1) = 3$

$d = 2$

Zoe

Billy's aunt owns a jewellery store, and gives him data on $5000$ of the diamonds in her store. For each diamond, we have:

- **carat**: the weight of the diamond, in carats
- **length**: the length of the diamond, in centimeters
- **width**: the width of the diamond, in centimeters
- **price**: the value of the diamond, in dollars

The first 5 rows of the 5000-row dataset are shown below:

| carat | length | width | price |
|-------|--------|-------|-------|
| 0.40 | 4.81 | 4.76 | 1323 |
| 1.04 | 6.58 | 6.53 | 5102 |
| 0.40 | 4.74 | 4.76 | 696 |
| 0.40 | 4.67 | 4.65 | 798 |
| 0.50 | 4.90 | 4.95 | 987 |

Billy has enlisted our help in predicting the price of a diamond given various other features.

Suppose $\vec{e} = \begin{bmatrix} e_1 \\ e_2 \\ \dots \\ e_n \end{bmatrix}$ is the error/residual vector, defined as

$$\vec{e} = \vec{y} - X\vec{w}^*$$

where $\vec{y}$ is the observation vector containing the prices for each diamond.

For each of the following quantities, state whether they are guaranteed to be equal to 0 the scalar, $\vec{0}$ the vector of all 0s, or neither. No justification is necessary.

- $\sum_{i=1}^{n} e_i$
- $||\vec{y} - X\vec{w}^*||^2$
- $X^T X \vec{w}^*$
- $2X^T X \vec{w}^* - 2X^T \vec{y}$

Zoe

$$\sum_{i=1}^{n} e_i$$

↑ sum of entries

guaranteed to equal 0

**Your answer in HW4**

2b

$$\vec{e} = \vec{y} - X\vec{w}$$

$\|\vec{e}\|^2$ ⌐ length of error vector

$$\|\vec{y} - X\vec{w}^*\|^2$$

$$\vec{y} = X\vec{w}^* = 0$$

$2 - 2 = 0$

└ zero because same val

actual    pred

Not guaranteed if $\vec{y} \neq X\vec{w}^*$

$X^T X \vec{w}^*$

$\vec{w}^* = (X^T X)^{-1} X^T y$       Not guaranteed

$\underbrace{X^T X} \vec{w} = X^T \vec{y}$

$X^T X \vec{w} = 0$ ?

$\vec{y}$ orthogonal
to every column of X

$2 X^T X \vec{w}^* - 2 X^T \vec{y}$ ⟵

$\underline{hint}$:     $X^T X \vec{w} = X^T \vec{y}$

$2 ( X^T X \vec{w} - X^T \vec{y} = \vec{0} )$

$2 X^T X \vec{w} - 2 X^T \vec{y} = \vec{0}$

$\begin{bmatrix} \phantom{x} \\ \phantom{x} \end{bmatrix}$ $-$ $\begin{bmatrix} \phantom{x} \\ \phantom{x} \end{bmatrix}$ $=$ $\begin{bmatrix} \phantom{x} \\ \phantom{x} \end{bmatrix}$
           3×1              3×1           3×1

yes! guaranteed
to be $\vec{0}$

$\begin{bmatrix} \phantom{x} \\ \phantom{x} \end{bmatrix}$ $\begin{bmatrix} \phantom{x} \\ \phantom{x} \end{bmatrix}$ $\begin{bmatrix} \phantom{x} \\ \phantom{x} \end{bmatrix}$ $\begin{bmatrix} \phantom{x} \\ \phantom{x} \end{bmatrix}$

d+1 xn      ~~n~~ ×d×1      ~~d~~×1 xn   =   d+1 xn

Suppose we introduce two more features:

- width alone, and
- area, which is defined as length times width

Suppose we also decide to remove the intercept term of our prediction rule. With all of these changes, our prediction rule is now

**no $w_0$?**

predicted price $= w_1 \cdot$ carat $+ w_2 \cdot$ length $+ w_3 \cdot$ width $+ w_4 \cdot$ (length $\cdot$ width)

**$0+$**

- Write out just the first 2 rows of the design matrix $X$ for this new prediction rule. You do **not** need to simplify the numbers in your matrix, it is fine if they involve the multiplication symbol.
- Is the optimal coefficient for carat, $w_1^*$, for this new prediction rule guaranteed to be equal to 10000, the optimal coefficient for carat in our original prediction rule? No justification is necessary.

| carat | length | width | price |
|---|---|---|---|
| 0.40 | 4.81 | 4.76 | 1323 |
| 1.04 | 6.58 | 6.53 | 5102 |
| 0.40 | 4.74 | 4.76 | 696 |
| 0.40 | 4.67 | 4.65 | 798 |
| 0.50 | 4.90 | 4.95 | 987 |

**each row of X:**

$[\text{carat} \quad \text{length} \quad \text{width} \quad \ell \cdot w]$

$$X = \begin{bmatrix} 0.4 & 4.81 & 4.76 & 4.81(4.76) \\ 1.04 & 6.58 & 6.53 & 6.58(6.53) \end{bmatrix}$$

$$w^* = \begin{bmatrix} w_1 \\ w_2 \\ w_3 \\ w_4 \end{bmatrix}$$

$w_1 = 1000$    $6.2$

$$\text{old } w = \begin{bmatrix} 2000 \\ 10000 \\ -1000 \end{bmatrix} \begin{matrix} \leftarrow w_0 \\ \leftarrow w_1 \\ \leftarrow w_2 \end{matrix}$$

$$\text{old } \begin{bmatrix} w_1 \\ w_2 \\ w_3 \\ w_4 \end{bmatrix} \quad \text{new} \begin{bmatrix} w_1 \\ w_2 \\ w_3 \\ w_4 \\ w_5 \end{bmatrix} \neq 0$$

**Zoe**

## Problem 4. Slippery Slope

In Lecture 2, we found that $h^* = \text{Median}(y_1, y_2, ..., y_n)$ is the constant prediction that minimizes mean absolute error:

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

Suppose that we have a dataset of numbers $y_1, y_2, ..., y_n$ such that $n$ **is odd** and the values are arranged in increasing order. That is, $y_1 \leq y_2, \leq ... \leq y_n$.

**Note: Parts (a) and (b) are independent of each other.**

a) 🥑🥑🥑🥑 Suppose that $R_{\text{abs}}(\alpha) = V$, where $V$ is the minimum value of $R_{\text{abs}}(h)$ and $\alpha$ is one of the numbers in our dataset.

Let $\alpha + \beta$ be the smallest value greater than $\alpha$ in our dataset, where $\beta > 0$. Another way of thinking about this is that $\beta = (\text{smallest value greater than } \alpha) - \alpha$.

Suppose we modify our dataset by replacing the value $\alpha$ with the value $\alpha + \beta + 1$. In our new dataset of $n$ values, what is the new minimum value of $R_{\text{abs}}(h)$ and at what value of $h$ is it minimized? Your answers to both parts should only involve the variables $V$, $\alpha$, $\beta$, $n$, and/or one or more constants.

Absolute error

Ratio $(\alpha) = V \rightarrow$ min value

$\alpha \rightarrow$ median      old

$\alpha + \beta$    $\beta > 0$

new

$$\alpha \rightarrow \alpha + \beta + 1$$
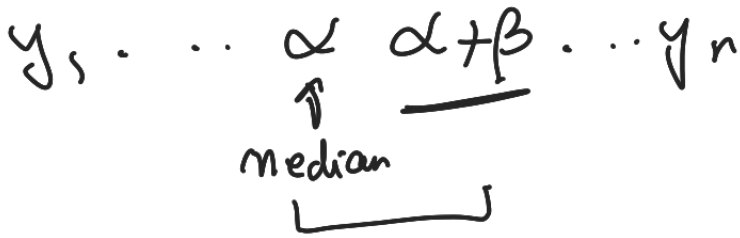
median ?

Ratio $(h)$ new value ?

__old__

$y_1 \ldots \quad \underset{\underset{\text{median}}{\uparrow}}{\alpha} \quad \underline{\alpha + \beta} \ldots \quad y_n$

__new__

$y_1 \ldots \underline{\alpha + \beta} \quad \alpha + \beta + 1 \ldots y_n$

new median    is  $\underline{\alpha + \beta}$

new minimizer $\longrightarrow$

$$\underline{\text{old}}$$

$$y_s \cdots \cdots \underset{\uparrow}{\alpha} \quad \underline{\alpha + \beta} \cdots y_n$$

$$\text{median}$$

imagine

new

$$\boxed{\alpha} \boxed{\quad} + \beta$$

$$y_1 \cdots \underline{\alpha + \beta} \underset{\uparrow}{\underline{\alpha + \beta + 1}} \cdots y_n$$

median

old left    old dev sum
olds        total of n values.

$\frac{n-1}{2}$ value to the left

$\beta_1 + \beta$
$y_2 + \beta$
$\cdots y$ just before new median
$\left(\frac{n-1}{2}\right)$ of them

$0 \rightarrow$ middle or median position
values to the right
$\rightarrow$ moved $\beta$ closer to each
vlads

olads $\rightarrow$ old.
olads $+ \beta \left(\frac{n-1}{2}\right) \rightarrow$ new

$\frac{\beta(n-1)}{2}$

change:  left values

$-$ olads.

vlads $- \beta\left(\frac{n-1}{2} - 1\right)$    $-\beta\left(\frac{n-1}{2}\right) - 1$

$\xrightarrow{\text{those many.}}$  $\left(\frac{n-1}{2}\right) - 1$

new sum of
abs dev $\frac{1}{h} \left( V_n + \frac{\beta(n-1)}{2} + 0 + (1 - \beta) + \left( -\beta\left(\frac{(n-1)}{2} - 1\right) \right) \right)$
ol

left vals    med   val imm
                   r of med

val to right
except imm

$$\text{old} \quad \alpha + \beta - \alpha \rightarrow \underline{\beta}$$

$$\text{new} \quad (\alpha + \beta + 1) - (\alpha + \beta) \rightarrow \underline{1}$$

$$1 - \beta$$

$$f_{abs}(old) = V$$
$$\alpha$$

new
$$V_n \int V + \frac{1}{n}$$

## Problem 4. Slippery Slope

In Lecture 2, we found that $h^* = \text{Median}(y_1, y_2, ..., y_n)$ is the constant prediction that minimizes mean absolute error:

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

Suppose that we have a dataset of numbers $y_1, y_2, ..., y_n$ such that $n$ **is odd** and the values are arranged in increasing order. That is, $y_1 \leq y_2, \leq ... \leq y_n$.

**Note: Parts (a) and (b) are independent of each other.**

**b)** 🥑🥑🥑 Let $y_a$ and $y_b$ be two values in our dataset such that $y_a < y_b$ and that the slope of $R_{\text{abs}}(h)$ is the same between $h = y_a$ and $h = y_b$. Specifically, let $d$ be the slope of $R_{\text{abs}}(h)$ between $y_a$ and $y_b$.

Suppose we introduce a new value $q$ to our dataset such that $q > y_b$. In our new dataset of $n + 1$ values, the slope of $R_{\text{abs}}(h)$ is still the same between $h = y_a$ and $h = y_b$, but it's no longer equal to $d$. What is the slope of $R_{\text{abs}}(h)$ between $h = y_a$ and $h = y_b$ in our new dataset? Your answer should depend on $d$, $n$, $q$, and/or one or more constants.

Harshi

$f_{abs} \rightarrow$ absolute loss.

$$y_a < y_b$$

slope $b/w = d$

new datapoint $q$

such that $q > y_b$.

$$\frac{d}{dh} f_{abs}(h) = \frac{1}{n}(\text{\# pts to left of } h - \text{\# pts to right of } h) \qquad \text{※} \underline{\text{Lecture 2}}$$

↑

@ pt $c$     $c$ $b/w$ $y_a$ and $y_b$     slope is $d$

slope
old $d$ $\rightarrow$ $d = \frac{1}{n}(\text{\# pts left of } c - \text{\# pts to the right of } c)$

$\rightarrow nd = (\text{\# pts left of } c - \text{\# pts of right of } c)$

new
slope $\rightarrow d' = \frac{1}{n+1}(\text{\#pts left of } c - (\text{\#pts right of } c + 1))$

$$d' = \frac{1}{n+1} \left( \text{\#pts left of } c - \text{\# pts right of } c - 1 \right)$$

$$d' = \frac{1}{n+1} \left( nd - 1 \right) \implies \boxed{\frac{nd - 1}{n + 1}}$$

Albert collected $400$ data points from a radiation detector. Each data point contains $3$ features: feature $A$, feature $B$ and feature $C$. The true particle energy $E$ is also reported. Albert wants to design a linear regression algorithm to predict the energy $E$ of each particle, given a combination of one or more of feature $A$, $B$, and $C$. As the first step, Albert calculated the correlation coefficients among $A$, $B$, $C$ and $E$. He wrote it down in the following table, where each cell of the table represents the correlaton of two terms:

|   | $A$ | $B$ | $C$ | $E$ |
|---|------|------|------|------|
| $A$ | 1 | -0.99 | 0.13 | 0.8 |
| $B$ | -0.99 | 1 | 0.25 | -0.95 |
| $C$ | 0.13 | 0.25 | 1 | 0.72 |
| $E$ | 0.8 | -0.95 | 0.72 | 1 |

Albert wants to start with a simple model: fitting only a single feature to obtain the true energy (i.e. $y = w_0 + w_1 x$). Which feature should he choose as $x$ to get the lowest mean square error?

- ○ $A$
- ● $B$
- ○ $C$

Albert collected $400$ data points from a radiation detector. Each data point contains $3$ features: feature $A$, feature $B$ and feature $C$. The true particle energy $E$ is also reported. Albert wants to design a linear regression algorithm to predict the energy $E$ of each particle, given a combination of one or more of feature $A$, $B$, and $C$. As the first step, Albert calculated the correlation coefficients among $A$, $B$, $C$ and $E$. He wrote it down in the following table, where each cell of the table represents the correlaton of two terms:

|   | $A$ | $B$ | $C$ | $E$ |
|---|------|------|------|------|
| $A$ | 1 | -0.99 | 0.13 | 0.8 |
| $B$ | -0.99 | 1 | 0.25 | -0.95 |
| $C$ | 0.13 | 0.25 | 1 | 0.72 |
| $E$ | 0.8 | -0.95 | 0.72 | 1 |

Albert wants to add another feature to his linear regression in part (a) to further boost the model's performance. (i.e. $y = w_0 + w_1 x + + w_2 x_2$) Which feature should he choose as $x_2$ to make additional improvements?

○ $A$

○ $B$

● $C$

Yosen

Albert collected $400$ data points from a radiation detector. Each data point contains $3$ features: feature $A$, feature $B$ and feature $C$. The true particle energy $E$ is also reported. Albert wants to design a linear regression algorithm to predict the energy $E$ of each particle, given a combination of one or more of feature $A$, $B$, and $C$. As the first step, Albert calculated the correlation coefficients among $A$, $B$, $C$ and $E$. He wrote it down in the following table, where each cell of the table represents the correlaton of two terms:

|   | $A$ | $B$ | $C$ | $E$ |
|---|------|------|------|------|
| $A$ | 1 | -0.99 | 0.13 | 0.8 |
| $B$ | -0.99 | 1 | 0.25 | -0.95 |
| $C$ | 0.13 | 0.25 | 1 | 0.72 |
| $E$ | 0.8 | -0.95 | 0.72 | 1 |

Albert further refines his algorithm by fitting a prediction rule of the form:

$$H(A,B,C) = w_0 + w_1 \cdot \boxed{A \cdot C}^{b_1} + w_2 \cdot \boxed{B^{C-7}}^{b_2}$$

Given this prediction rule, what are the dimensions of the design matrix $X$?

$n \times (d+1)$ ✓

$$\begin{bmatrix} & & \\ & & \\ & & \end{bmatrix}_{r \times c}$$

So, what are $r$ and $c$ in $r$ rows $\times$ $c$ columns?

$400$ ✓   $d+1 = 3$ ✓

$$\begin{bmatrix} 1 & b_1^{(1)} & b_1^{(2)} \\ \vdots & b_2^{(1)} & b_2^{(2)} \\ \vdots & \vdots & \vdots \\ 1 & b_{400}^{(1)} & b_{400}^{(2)} \end{bmatrix}$$

Yosen

# Questions?

# Thanks for Coming!

Good luck everyone!