

Lecture 13: Midterm Review

EECS 398: Practical Data Science, Winter 2025

practicaldsc.org • github.com/practicaldsc/wn25 •  See latest announcements [here on Ed](#)

Agenda

- We'll work through the first 8 questions of the Fall 2024 Final Exam:
study.practicaldsc.org/fa24-final.
- I'll post these annotated slides after lecture.
- The solutions + recording for yesterday's review session are also posted.

	Type	Brand	Name	Price	Rating	Num Ingredients	Sensitive
0	Eye cream	PERRICONE MD	PRE:EMPT SERIES™ Brightening Eye Cream	55	4.2	33	1
1	Cleanser	CLINIQUE	Pep-Start 2-in-1 Exfoliating Cleanser	190	3.1	36	0
2	Eye cream	PETER THOMAS ROTH	FIRMx™ 360 Eye Renewal	75	5.0	42	0
3	Treatment	KIEHL'S SINCE 1851	Clearly Corrective™ Dark Spot Solution	50	4.5	24	1
4	Cleanser	PETER THOMAS ROTH	Irish Moor Mud Purifying Cleanser Gel	38	3.6	23	0

one product

$np.count_nonzero(s) / s.shape[0]$

Problem 1

An expensive product is one that costs at least \$100.

Problem 1.1

Write an expression that evaluates to the proportion of products in `skin` that are expensive.

$(skin["Price"] \geq 100).mean()$

equivalent:

$skin[skin["Price"] \geq 100].shape[0] / skin.shape[0]$

False
True
False
True
False
...
) .mean()

(155, 7)
↑ ↑
[0] [1]

Fall 2024 Final Exam

	Type	Brand	Name	Price	Rating	Num Ingredients	Sensitive
0	Eye cream	PERRICONE MD	PRE:EMPT SERIES™ Brightening Eye Cream	55	4.2	33	1
1	Cleanser	CLINIQUE	Pep-Start 2-in-1 Exfoliating Cleanser	19	3.1	36	0
2	Eye cream	PETER THOMAS ROTH	FIRMx™ 360 Eye Renewal	75	5.0	42	0
3	Treatment	KIEHL'S SINCE 1851	Clearly Corrective™ Dark Spot Solution	50	4.5	24	1
4	Cleanser	PETER THOMAS ROTH	Irish Moor Mud Purifying Cleanser Gel	38	3.6	23	0

Problem 1.2

Fill in the blanks so that the expression below evaluates to the number of brands that sell **fewer than 5** expensive products.

`skin.groupby(__(i)__).__(ii)__(__(iii)__)["Brand"].nunique()`

(i):

- ☒ "Brand"
- ☐ "Name"
- ☐ "Price"
- ☐ ["Brand", "Price"]

(ii):

- ☐ agg
- ☐ count
- ☒ filter
- ☐ value_counts

(iii): (Free response)

group condition

function

lambda df: $(df["Price"] \geq 100).sum() < 5$

returns a single True
or a single False

Tangent: loc df[]

df.loc [^{which rows}
_{do I want?} , ^{which columns}
_{do I want?}]

e.g. to get the "Price" column out of skin,

skin.loc [: , "Price"]

skin["Price"]

e.g. to get the expensive products

skin.loc [skin["Price"] >= 100 , ["Price", "Sensitive"])

Problem 3

loc: labels, iloc: integer positions

Consider the Series `small_prices` and `vc`, both of which are defined below.

```
small_prices = pd.Series([36, 36, 18, 100, 18, 36, 1, 1, 1, 36])
```

Handwritten red annotations below the list in the code above:
- under 36, - under 36, + under 18, * under 100, + under 18, - under 36, . . . under 1, 1, 1, - under 36

```
vc = small_prices.value_counts().sort_values(ascending=False)
```

In each of the parts below, select the value that the provided expression evaluates to. If the expression errors, select "Error".

☐ 0

Handwritten red "✓" next to the radio button

Handwritten red "36" next to the radio button

Handwritten red "4" next to the radio button

`vc.iloc[0]`

Handwritten blue "4" next to the code

☐ 1

`vc.loc[0]`

Handwritten blue "Error" next to the code

☐ 2

Handwritten red "1" next to the radio button

Handwritten red "3" next to the radio button

`vc.index[0]`

Handwritten blue "36" next to the code

☐ 3

☐ 4

Handwritten red "18" next to the radio button

Handwritten red "2" next to the radio button

`vc.iloc[1]`

Handwritten blue "3" next to the code

☐ 18

Handwritten red "100" next to the radio button

Handwritten red "1" next to the radio button

`vc.loc[1]`

Handwritten blue "3" next to the code

☐ 36

`vc.index[1]`

Handwritten blue "1" next to the code

☐ 100

☐ Error

☐ None of these

Problem 4

Consider the DataFrames `type_pivot`, `clinique`, `fresh`, and `boscia`, defined below.

```
type_pivot = skin.pivot_table(index="Type",
                               columns="Brand",
                               values="Sensitive",
                               aggfunc=lambda s: s.shape[0] + 1)
```

```
clinique = skin[skin["Brand"] == "clinique"]
fresh = skin[skin["Brand"] == "fresh"]
boscia = skin[skin["Brand"] == "BOSCIA"]
```

Three columns of `type_pivot` are shown below in their entirety.

Brand	CLINIQUE	FRESH	BOSCIA
Type			
Cleanser	5	NaN	2
Eye cream	3	NaN	2
Face Mask	2	4	4
Moisturizer	2	3	NaN
Sun protect	1	NaN	NaN

In each of the parts below, give your answer as an integer.

Problem 4.1

How many rows are in the following DataFrame?

```
clinique.merge(fresh, on="Type", how="inner")
```

110

Problem 4.2

How many rows are in the following DataFrame?

```
clinique.merge(fresh, on="Type", how="outer")
        .merge(boscia, on="Type", how="outer")
```

Handwritten calculation for Problem 4.2:

Intermediate result (from `clinique.merge(fresh, on="Type", how="outer")`):

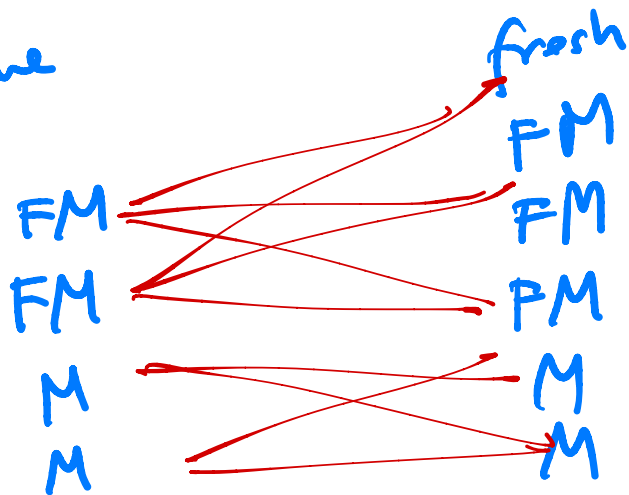
Type	Count
Cleanser	5
Eye cream	3
Face Mask	6
Moisturizer	4
Sun protect	1

Then merge with `boscia` (how="outer"):

Type	Count
Cleanser	5
Eye cream	3
Face Mask	6
Moisturizer	4
Sun protect	1
Face Mask (from boscia)	4

Total count: 5 + 3 + 6 + 4 + 1 + 4 = 23

clinique



$$\underbrace{FM}_{2 \times 3} + \underbrace{M}_{2 \times 2}$$

$$= 6 + 4$$

$$= \boxed{10}$$

Problem 5

Consider a sample of 60 skincare products. The name of one product from the sample is given below:

"our **drops** cream is the best **drops drops** for eye **drops drops** proven formula..."

The total number of terms in the product name above is unknown, but we know that the term **drops** only appears in the name 5 times.

Suppose the TF-IDF of **drops** in the product name above is $\frac{2}{3}$. Which of the following statements are **NOT possible**, assuming we use a base-2 logarithm? Select all that apply.

- ☐ All 60 product names contain the term **drops**, including the one above.
- ☐ 14 **other** product names contain the term **drops**, in addition to the one above.
- ☐ None of the 59 **other** product names contain the term **drops**.
- ☐ There are 15 terms in the product name above **in total**.
- ☐ There are 25 terms in the product name above **in total**.

see study site
solutions
for more
detail

Let T be # of terms in
product name above,
 n be # of documents
with "drops", including
document above

TF-IDF ("drops", doc above)

$$= \frac{5}{T} \cdot \log_2 \left(\frac{60}{n} \right) = \frac{2}{3}$$

the problem boils down to checking if the T, n in the question satisfy this eq'n.

Problem 6

Suppose `soup` is a BeautifulSoup object representing the homepage of a Sephora competitor.

Furthermore, suppose `prods`, defined below, is a list of strings containing the name of every product on the site.

```
prods = [row.get("prod") for row in soup.find_all("row", class_="thing")]
```

Given that `prods[1]` evaluates to `"Cleansifier"`, which of the following options describes the source code of the site?

- Option 1:

```
<row class="thing">prod: Facial Treatment Essence</row>
<row class="thing">prod: Cleansifier</row>
<row class="thing">prod: Self Tan Dry Oil SPF 50</row>
...
```

- Option 2:

```
<row class="thing" prod="Facial Treatment Essence"></row>
<row class="thing" prod="Cleansifier"></row>
<row class="thing" prod="Self Tan Dry Oil SPF 50"></row>
...
```

- Option 3:

```
<row prod="thing" class="Facial Treatment Essence"></row>
<row prod="thing" class="Cleansifier"></row>
<row prod="thing" class="Self Tan Dry Oil SPF 50"></row>
...
```

- Option 4:

```
<row class="thing">prod="Facial Treatment Essence"</row>
<row class="thing">prod="Cleansifier"</row>
<row class="thing">prod="Self Tan Dry Oil SPF 50"</row>
...
```

needs to be within
<row prod="...">

Problem 7

Consider a dataset of n values, y_1, y_2, \dots, y_n , all of which are positive. We want to fit a constant model, $H(x) = h$, to the data.

Let h_p^* be the optimal constant prediction that minimizes average degree- p loss, $R_p(h)$, defined below.

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

For example, h_2^* is the optimal constant prediction that minimizes $R_2(h) = \frac{1}{n} \sum_{i=1}^n |y_i - h|^2$.

In each of the parts below, determine the value of the quantity provided. By "the data", we are referring to y_1, y_2, \dots, y_n .

☐ The standard deviation of the data

☐ The variance of the data

☐ The mean of the data

☐ The median of the data

☐ The midrange of the data, $\frac{y_{\min} + y_{\max}}{2}$

☐ The mode of the data

☐ None of the above

h_0^* minimizes $\frac{1}{n} \sum_{i=1}^n |y_i - h|^0$

$\min \frac{1}{n} \sum_{i=1}^n |y_i - h|^1$

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

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

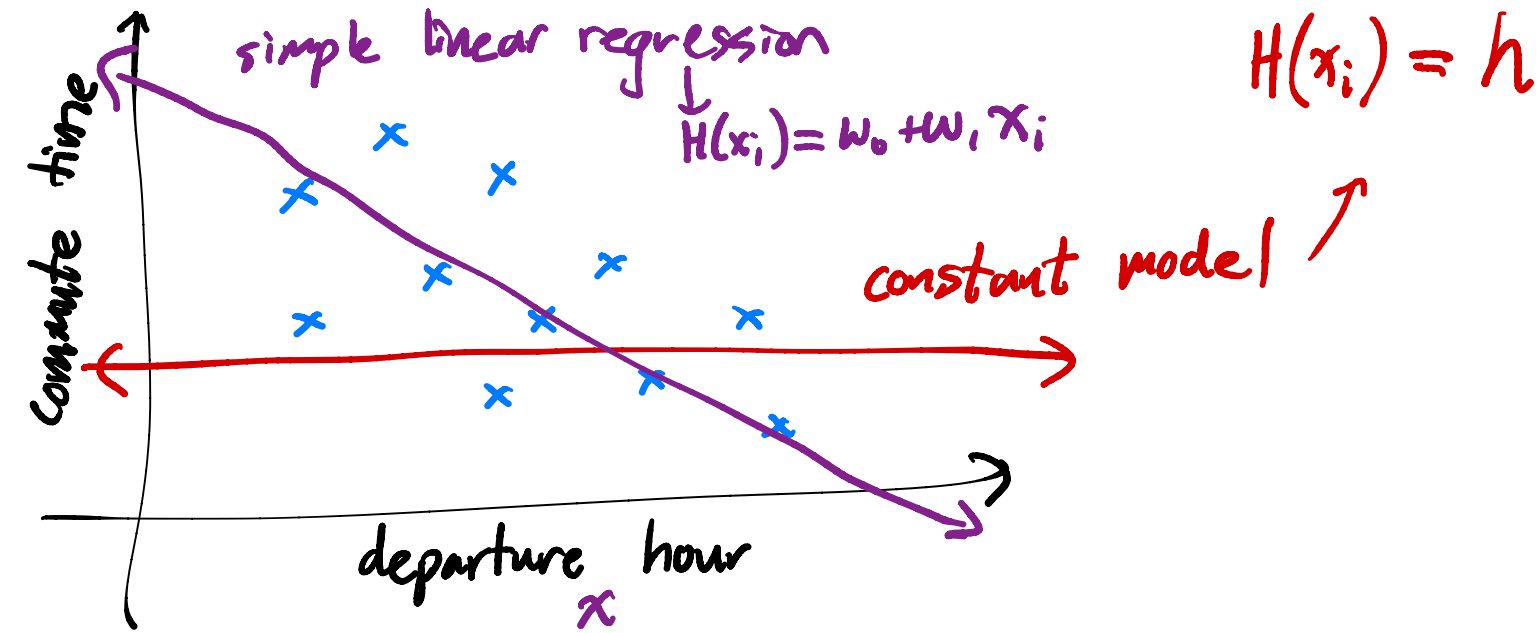
h_0^* : mode
 h_1^* : median

minimizes absolute loss

$R_1(h_1^*)$ none

h_2^* mean

$R_2(h_2^*)$ variance
mean



loss functions are how we find which
constant line or SLR line to use,
i.e. how we find parameters h or w_0, w_1 .

constant model: ignore departure hours, x_i

$$H(x_i) = h$$

how our predictions
are made

Loss function: how wrong an individual
prediction is

- squared loss: $(\text{actual} - \text{predicted})^2$
- absolute loss: $|\text{actual} - \text{predicted}|$

Given: $y_1, y_2, \dots, y_n \rightarrow$ e.g. $y_1 = 10$ $y_2 = 50$ $y_3 = 55$

average loss

$$\frac{1}{3} \left((10-h)^2 + (50-h)^2 + (55-h)^2 \right)$$

find the h^* that
minimizes average loss!!!

Chosen:

- constant model, $H(x_i) = h$
- squared loss, (actual - predicted)²

average squared loss

$$R_2(h) = \frac{1}{n} \left[(y_1 - h)^2 + (y_2 - h)^2 + \dots + (y_n - h)^2 \right]$$

empirical Risk (aka average loss)

Goal: we want best predictions

=> to do that, we choose the parameter that minimizes $R_2(h)$

$$\Rightarrow h^* = \text{Mean}(y_1, y_2, \dots, y_n)$$

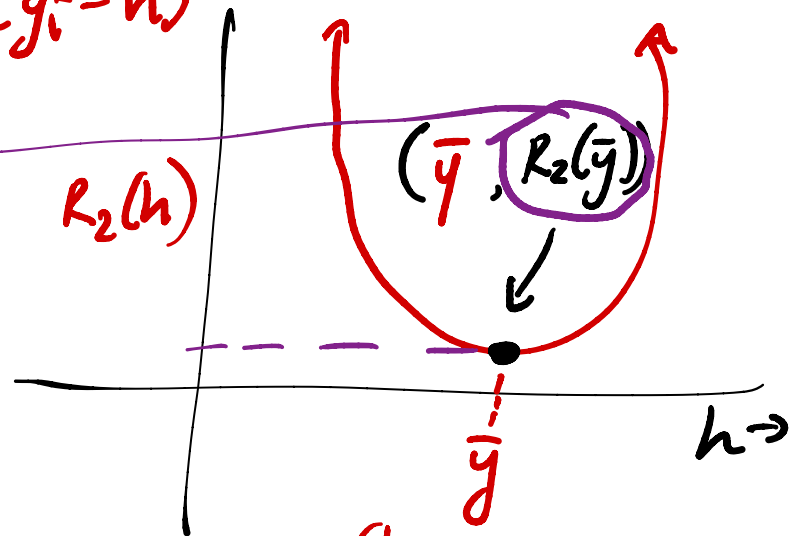
- if we choose squared loss : $h^* = \text{Mean}$
- if we choose abs loss : $h^* = \text{Median}$
- ⋮

variance, standard deviation

$$R_2(h) = \frac{1}{n} \sum_{i=1}^n (y_i - h)^2$$

$$R_2(\bar{y}) = \frac{1}{n} \sum_{i=1}^n (y_i - \bar{y})^2$$

variance of y_1, y_2, \dots, y_n



(because mean
minimizes $R_2(h)$)

$$R_2(\bar{y}) = \text{variance of } y_1, y_2, \dots, y_n$$

$$\sqrt{R_2(\bar{y})} = \text{standard deviation of } y_1, y_2, \dots, y_n$$

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

$$R_0(h) = \frac{1}{n} \sum_{i=1}^n |y_i - h|^0$$

$$x^0 = 1$$

0^0 undefined, but assume $0^0 = 0$.

either 0 or 1

$[1, 1, 2, 2, 2, 3]$

$$h=0 \rightarrow \text{all wrong} \rightarrow \frac{1}{6} \sum_{i=1}^6 (1) = 1$$

$$h^*=1 \rightarrow 2 \text{ right, } 4 \text{ wrong} \rightarrow \frac{1}{6} \cdot (4) = \frac{4}{6}$$

$$h^*=2 \rightarrow 3 \text{ right, } 3 \text{ wrong} \rightarrow \frac{3}{6}$$

$$h^*=3 \rightarrow \frac{5}{6}$$

Now, suppose we want to find the optimal constant prediction, h_U^* , using the "Ultra" loss function, defined below.

y_i are all positive!

$$L_U(y_i, h) = y_i (y_i - h)^2$$

$$L_2(y_i, h) = (y_i - h)^2$$

$$(5) (5 - h)^2$$

Problem 7.6

To find h_U^* , suppose we minimize average Ultra loss (with no regularization). How does h_U^* compare to the mean of the data, M ?

☐ $h_U^* > M$

☒ $h_U^* \geq M$

☐ $h_U^* = M$

☐ $h_U^* \leq M$

☐ $h_U^* < M$

L_U big when y_i is big

L_U small when y_i small (close to 0)

Big penalty when y_i is a large number

if all y_i are the same, $h_U^* = M$

e.g.

2, 3, 70

$$R_2(h) = \frac{1}{3} \left[(2-h)^2 + (3-h)^2 + (70-h)^2 \right]$$

$$R_n(h) = \frac{1}{3} \left[2(2-h)^2 + 3(3-h)^2 + \underbrace{70(70-h)^2}_{\text{to make this small}_2} \right]$$

we make $(70-h)$
small, i.e. bring
 h closer to 70

Now, to find the optimal constant prediction, we will instead minimize **regularized** average Ultra loss, $R_\lambda(h)$, where λ is a non-negative regularization hyperparameter:

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

It can be shown that $\frac{\partial R_\lambda(h)}{\partial h}$, the derivative of $R_\lambda(h)$ with respect to h , is:

objective function

$$\frac{\partial R_\lambda(h)}{\partial h} = -2 \left(\frac{1}{n} \sum_{i=1}^n y_i (y_i - h) - \lambda h \right) = 0$$

not in the sum!

Problem 7.7

Find h^* , the constant prediction that minimizes $R_\lambda(h)$. Show your work, and put a box around your final answer, which should be an **expression in terms of** y_i , n , and/or λ .

$$\frac{1}{n} \sum_{i=1}^n y_i (y_i - h) - \lambda h = 0$$

$$\frac{1}{n} \sum (y_i^2 - h y_i) - \lambda h = 0$$

$$\frac{1}{n} \left(\sum y_i^2 - \sum h y_i \right) - \lambda h = 0$$

not in the sum!

$$\frac{1}{n} \sum_{i=1}^n y_i (y_i - h) - \lambda h = 0$$

$$\frac{1}{n} \sum (y_i^2 - h y_i) - \lambda h = 0$$

$$\frac{1}{n} (\sum y_i^2 - \sum h y_i) - \lambda h = 0$$

$$\frac{1}{n} (\sum y_i^2 - h \sum y_i) - \lambda h = 0$$

$$\frac{1}{n} \sum y_i^2 - \underbrace{h}_{\frac{\sum y_i}{n}} - \underbrace{\lambda h}_{\lambda h} = 0$$

$$\frac{1}{n} \sum y_i^2 = \left(\frac{\sum y_i}{n} + \lambda \right) h$$

$$h^* = \frac{\frac{1}{n} \sum y_i^2}{\frac{1}{n} \sum y_i + \lambda}$$

$$h^* = \frac{\sum y_i^2}{\sum y_i + n\lambda}$$

find answer

$$L(y_i, H(x_i)) = (\text{some loss function})$$

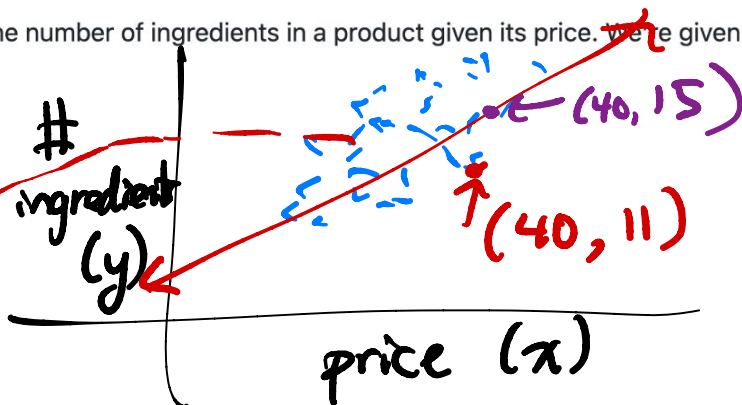
$$R(H) = \frac{1}{n} \sum_{i=1}^n L(y_i, \underbrace{H(x_i)}_{H(x_i)=h})$$

Problem 8

Suppose we want to fit a simple linear regression model (using squared loss) that predicts the number of ingredients in a product given its price. We're given that:

- The average cost of a product in our dataset is \$40, i.e. $\bar{x} = 40$.
- The average number of ingredients in a product in our dataset is 15, i.e. $\bar{y} = 15$.

The intercept and slope of the regression line are $w_0^* = 11$ and $w_1^* = \frac{1}{10}$, respectively.



Problem 8.1

Suppose Vectors' Veil (a skincare product) costs \$40 and has 11 ingredients. What is the squared loss of our model's predicted number of ingredients for Vectors' Veil? Give your answer as a **number**.

squared loss = (actual - predicted)² = (11 - 15)² = 16

$$H(x_i) = 11 + \frac{1}{10} x_i$$

$$H(40) = 11 + \frac{1}{10} \cdot 40 = 11 + 4 = 15$$

predicted # of ingredients

Problem 8

Suppose we want to fit a simple linear regression model (using squared loss) that predicts the number of ingredients in a product given its price. We're given that:

- The average cost of a product in our dataset is \$40, i.e. $\bar{x} = 40$.
- The average number of ingredients in a product in our dataset is 15, i.e. $\bar{y} = 15$.

The intercept and slope of the regression line are $w_0^* = 11$ and $w_1^* = \frac{1}{10}$, respectively.

Fact: the regression line

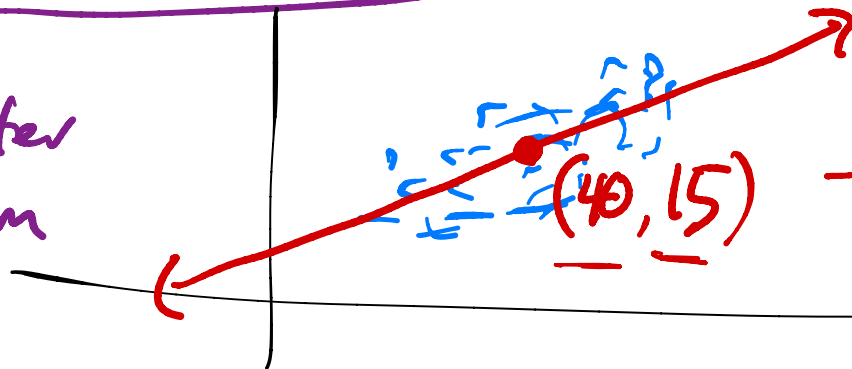
ALWAYS passes through the point (\bar{x}, \bar{y})

Problem 8.2

Is it possible to answer part (a) above **just** by knowing \bar{x} and \bar{y} , i.e. **without** knowing the values of w_0^* and w_1^* ?

- ☐ Yes; the values of w_0^* and w_1^* don't impact the answer to part (a).
- ☒ No; the values of w_0^* and w_1^* are necessary to answer part (a).

proof after midterm

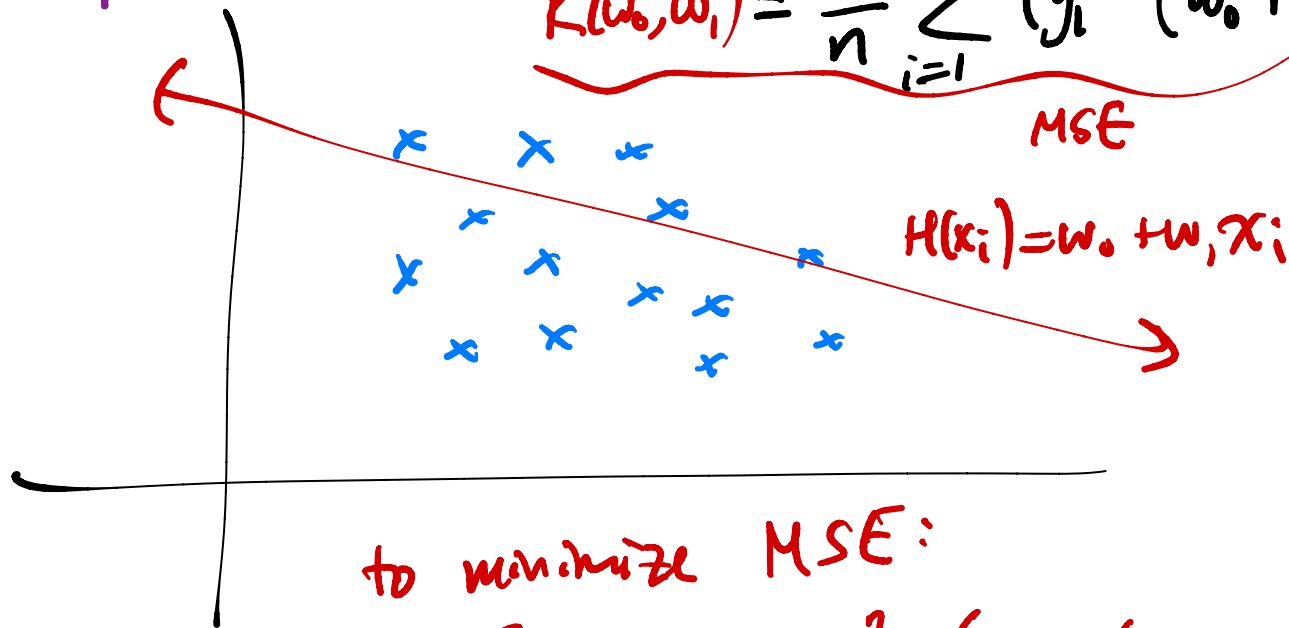


→ always true if w_0, w_1 minimize mean sq. error

Simple linear regression

$$R(w_0, w_1) = \frac{1}{n} \sum_{i=1}^n (y_i - (w_0 + w_1 x_i))^2$$

MSE



set to 0

to minimize MSE:

$$\frac{\partial R}{\partial w_0} = -\frac{2}{n} \sum_{i=1}^n (y_i - (w_0 + w_1 x_i)) = 0$$

solve for w_0^T, w_1^T

$$\frac{\partial R}{\partial w_1} = -\frac{2}{n} \sum_{i=1}^n (y_i - (w_0 + w_1 x_i)) x_i = 0$$

$$\underbrace{w_0^*}_{\text{intercept}} = \bar{y} - \underbrace{w_1^*}_{\text{slope}} \bar{x}$$

$$\underbrace{w_1^*}_{\text{slope}} = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sum_{i=1}^n (x_i - \bar{x})^2} = \frac{\sum (x_i - \bar{x}) y_i}{\sum (x_i - \bar{x}) x_i} = r \frac{\sigma_y \leftarrow \text{SD } y}{\sigma_x \leftarrow \text{SD } x}$$

correlation
coefficient

$$w_0^* = \bar{y} - w_1^* \bar{x}$$

$$w_1^* = r \frac{\sigma_y}{\sigma_x}$$

$$H(x_i) = w_0^* + w_1^* x_i$$

$$H(\bar{x}) = \vdots$$

$$= \bar{y}$$