

DADM  
Aknil Patil

Q1 Problem1

(b)

⇒ We have following equations,

$$R^2 = 1 - \frac{RSS}{TSS}$$

$$= 1 - \frac{\sum_i (y_i - \hat{y}_i)^2}{\sum_i y_i^2}$$

With  $\hat{y}_i = \hat{\beta}_1 x_i$  we may write

$$R^2 = 1 - \frac{\sum_i (y_i - \sum_j x_j y_j / \sum_j x_j^2 x_i)^2}{\sum_i y_i^2}$$

$$= \frac{\sum_i y_i^2 - \left( \sum_i y_i^2 - 2 \sum_i y_i \left( \sum_j x_j y_j / \sum_j x_j^2 \right) x_i + \sum_i \left( \sum_j x_j y_j / \sum_j x_j^2 \right)^2 x_i^2 \right)}{\sum_i y_i^2}$$

$$\therefore R^2 = \frac{2 \left( \sum_i x_i y_i \right)^2 / \sum_i x_i^2 - \left( \sum_i x_i y_i \right)^2 / \sum_i x_i^2}{\sum_i y_i^2}$$

$$= \frac{\left( \sum_i x_i y_i \right)^2}{\sum_i x_i^2 \sum_i y_i^2} = \text{Corr}(X, Y)^2$$

Q4. b) Case Study -

⇒ Hawkins Model's Overview

By splitting the dataset, we have

Training set:

$$\text{No Breast cancer} = 70/100 * 900 = 630$$

$$\text{Breast cancer} = 70/100 * 100 = 70$$

Test set:

$$\text{No Breast cancer} = 0.3 * 900 = 270$$

$$\text{Breast cancer} = 0.3 * 100 = 30$$

In order to get what went wrong with the prediction model of Hawkins let's consider a single case of one woman.

Say,

we have an assumption  $H$ , that woman is suffering from Breast pain and not Breast cancer but another assumption  $H_0$  says woman was suffering from Breast cancer.

if Assumption  $H_0$  is true (positive) - Breast cancer.  
else Assumption  $H_0$  is false (negative) - No Breast cancer.

The table below represents what happens if this other assumption Ho is true or not.

|   | Null Hypothesis (Ho)<br>valid: Breast cancer                  | Null Hypo (Ho) is<br>invalid: NO Breast cancer                                   |
|---|---|--|
| Accept Ho:<br>women has Breast<br>cancer.             | TP - Women has<br>Breast cancer.                              | FP - False Alarm,<br>Women has temporary<br>Breast pain but<br>No Breast cancer. |
| Reject Ho:<br>women doesn't<br>have Breast<br>cancer. | FN - Women is healthy<br>but she's dying<br>of Breast cancer. | TN - Women<br>does not have<br>Breast cancer.                                    |

Where,

TP - True Positive, FP - False Positive  
FN - False Negative, TN = True Negative.

### - Hawkins's Model Prediction Results:-

After training the model with 70% of dataset Hawkins scientist tested model with 30% to evaluate the model for its accuracy. Their model got 270 predictions right out of 300.

$$\text{Hawkins accuracy} = 270/300 = 0.9$$

~~the~~ What went wrong then,

By re-evaluating the model we get,

False Negative = 30

True Positive = 0

True Negative = 270

False Positive = 0

Building a confusion matrix with above data,

|                      | Actual-cancer | Actual-Not cancer | Total |
|----------------------|---------------|-------------------|-------|
| Predicted - cancer   | TP=0          | FP=0              | 0     |
| Predicted-Not cancer | FN=30         | TN=270            | 300   |
| Total                | 30            | 270               | 300   |

In summary,

Hawkins model correctly classified 270 women who do not have breast cancer as 'No Breast Cancer' while it incorrectly classified 30 women who have breast cancer as 'No Breast Cancer'.

Here, the model has conveniently classified all the test data as "No Breast Cancer". Accuracy is 90%.

But,

None of the 'Breast cancer' data is correctly labeled. We can ~~to~~ evaluate this by checking precision & recall.

By observing the original results we can indicate that people who have breast cancer the model isn't precise model.

We can prove this by,

$$\text{Precision} = \frac{TP}{TP + FP}$$

$$\text{Precision} = \frac{0}{0+0} = 0$$

$$\text{Precision in \%} = 0\%$$

$$\text{Recall} = \frac{TP}{TP + FN}$$

$$\text{Recall} = \frac{0}{0+30} = 0$$

$$\text{Recall in \%} = 0\%$$

What this means is that, this model will always classify any data passed into it as 'No Breast cancer'.

This explains why some patients were where showing symptoms of Breast cancer.

Houk's model did not work.

Q1a) →

As per table,  
The Null hypothesis indicates that advertising budgets of TV, Newspaper or radio do not have any effect on sales.  
More precisely,

$$H_0^{(1)}: \beta_1 = 0$$

$$H_0^{(2)}: \beta_2 = 0$$

$$H_0^{(3)}: \beta_3 = 0$$

The corresponding p-values are highly significant for 'TV' & 'radio' & not significant for 'Newspaper'.

So, we reject  $H_0^{(1)}$  &  $H_0^{(2)}$  & not  $H_0^{(3)}$ .

We can conclude that Newspaper advertising budget do not affect sales.

Q4:  
a)

$$\hat{\beta}_0 = -6, \hat{\beta}_1 = 0.03, \hat{\beta}_2 = 1$$

Substituting  $\beta$  values in eq of predicted probabilities;

$$\hat{p}(x) = \frac{e^{-6+0.03x_1+x_2}}{1 + (e^{-6+0.03x_1+x_2})}$$

$$= 0.3775$$



Q4.

ii)

eq of predicted probabilities tells us that,

$$z = \frac{e^{-6+0.05x_1+x_2}}{1+e^{-6+0.05x_1+x_2}}$$

$$z = 0.5$$

which is equivalent to,

$$e^{-6+0.05x_1+x_2} = 1$$

Taking log on both sides,

$$x_1 = \frac{2.5}{0.05} = 50.$$

Student should study 50 hours.