

Cheng King Fung Michael

SID: 2009 2275

1. a) log ~~trans~~ transformation. From scatter plot 1, we can see that the log-log scatter plot have a linear ~~with~~ pattern ~~with~~ while the original does not.

1b) i

~~$e^{-13.363} = 2.457 \times 10^{-6}$~~
 $e^{-13.363} = 1.572 \times 10^{-6}$

$e^{-7.801} = 3.839 \times 10^{-4}$



We are 95% confident that the 'medv' is in ~~between~~ between the range of 1.572×10^{-6} to 3.839×10^{-4}

1b) ii The ~~need~~ linearity, equal variance and the

normality. From the residual plot,

we can see that there is a funnel in the error (violates equal var). More

the normal QQ plot have ~~skewness~~

right skew (normality ~~is~~ and linearity)

1 c)

I agree with him as the two models both have the same amount of predictors and response

I don't agree with him as the response of model 1 is 'crim' while model 2 is $\log(\text{crim})$, so we cannot compare directly

1 c) ii

$$e^{-0.91993} \approx 0.3985$$

$$e^{3.765798} \approx 43.198$$

The 95% prediction interval is (0.3985, 43.198)

2 a)

There are 7 dummies in the full model

2 b) i

Dallas and RH

b) ii No, as the coefficient of 'Type_of_locationurban' is only 0.0165 ≈ 0 , so it makes no significant difference

2b) 771

$$y = 0.1935 + 0.0156 X_1 - 0.0017 X_2 + 0.0035 X_3 \\ - 0.0141 X_4 + 0.0071 X_5 + 3.326 \times 10^{-5} X_6 + \\ -0.0435 + 0.0724 + 0.1481 + 0.1481 \times (-0.0033) X_5$$

$$y = 0.3705 + 0.0156 X_1 - 0.0017 X_2 + 0.0035 X_3 - 0.0141 X_4 \\ + 6.61127 \times 10^{-3} X_5 + 3.326 \times 10^{-5} X_6$$

2b) TV

The β_0 changes from 0.1935 to

0.3705

while $\beta_{\text{nx.husorr}}$ changes from 0.0071 to 6.61127×10^{-3}

3a) Administration - 13 It has the highest p value

(~~sig~~ > 0.05)
(p-value)

~~let $R^2 = 0.948$~~

b) ~~$y = 50120 + 0.8057 \times R\&D$~~

$$y = 50120 + 0.8057 \times R\&D - 0.0268 \times \text{administration} +$$

$$0.0272 \times \text{marketing}$$

$$H_0 = F < 4$$

c) $F = 296$ $p\text{-value} = 4.53 \times 10^{-30}$

\therefore Reject H_0

d) $0.602 > 0.05$

Reject H_0

e) $RMSE = \sqrt{(1 - R_{adj}^2) S_y^2}$

$$= \sqrt{(1 - 0.948)(40306)^2}$$

$$= 9191.22$$

Cheng King Fung Michael



SID: 2069 2275

3 f)

$$Y = 50120 + 75000 (0.8057) + 150000 (-0.0268) + 200000 (0.0272)$$

$$Y = \cancel{61821.5} \quad 111967.5$$

expected profit is \$ 111967.5 the ~~perf~~ performance of the company is very good. It outperforms our model.

3 g) No, as the VIF of all ~~int~~ predictors < 10



4a)

Best for 1 predictors = Age

Best for 2 predictors = (Age, HP)

Best for 3 predictors = (Age, KM, HP)

b)

1	(Age)
2	(Age, HP)
3	(Age, KM, HP)

c)

$$1 + p \left(\frac{p+1}{2} \right)$$

$$= 1 + 3 \frac{(3+1)}{2} = 7$$

d)

1	4979.9561	(Age)
2	4965.8036	(Age, HP)
3	4967.1074	(Age, KM, HP)

5a)

No, because classification method requires a response (Y), but customer segmentation in this case does not

5b)

Label = 2

c) Yes, we can see that the p.t decrease dramatically after $k=3$

d)

Label = 0 (window shopper)

Label = 1 (loyal customer)

Label = 2 (normal customer)

It is because ~~that~~ label 0 spent the less ~~and~~ money and the most in time ^(it's window shopper) while label 1 spent a ~~lot~~ most money and a ~~big~~ large period of time, so it is loyal, while the leftover is normal customer (label 2)

Cheng King Fung Michael

[Signature]

SID: 2009 2275

5e) ~~Yes, as he can use hierarchical clustering to visualize the behaviours by a dendrogram~~ Yes, as he can use hierarchical clustering to visualize the behaviours by a dendrogram

6. 3rd $92.8 + 0.507(x)$
 4th $91.701 + 0.522(x)$

$$3^{rd} \text{ MSE} = \sum_{i=1}^n (y_i - \hat{y}_i)^2$$

$$= 1.227773$$

$$4^{th} \text{ MSE} = 0.806012$$

4-fold
 CV
 MSE $= 0.6241 + 0.690561 + 1.227773$
 $+ 0.806012$

$$= 3.348$$

3rd

\hat{y}_1	$= 123.22$
\hat{y}_2	$= 111.052$
\hat{y}	$= 2110.038$
\hat{y}	$= 120.635$

4th

\hat{y}_1	123.021
\hat{y}_2	112.493
\hat{y}_3	109.449
\hat{y}_4	120.411

Name: (Chang Hong Feng
Michael)

Student ID: 20692275

Question 8

7. a)

0

$$b) \text{ TPR} = \frac{(\text{Real and Predict Y})}{(\text{Real Y})} = \frac{5}{9} = 0.556$$

$$c) \text{ FPR} = \frac{(\text{Predicted } 0)}{(\text{Real } 0)} = \frac{0}{5} = 0$$

$$d) \frac{12}{28} = 35.7\%$$

e) No, it is only achieved when

$$TP + FN = FP + TN$$

8 a)

$$\frac{1}{1 + e^{-(12.2871 + 3.5312(3.2) + 0.2149(9))}}$$

≈ 1

$$c) L = 12.2871 + 3.5312 \times (GA + 0.2149 \times H_{\text{uv}})$$

d) positive

17.9

2012275
hierarchical
a dendrogram

Cheng long Feng Michael

2012275

Lyf.

e)

It's 90.625 % chance will give be correct
~~high score a pass to the~~ about the
true positive (predict passed and really pass)

f)

KS measure γ

Around ~~0.63~~ 0.67

For

~~TPR and FPR~~