

High-Dimensional Data Analysis

Homework#1

Student name: Lin Ting-Yu

1.Full in all values

Simple regression									
	$\widehat{\beta}_0$	$\widehat{\beta}_1$	$se(\widehat{\beta}_0)$	$se(\widehat{\beta}_1)$	t_0	t_1	Rss	P_0	P_1
TV(x_1)	7.033	0.048	0.458	0.003	15.360	17.668	2102.531	≈ 0	≈ 0
radio(x_2)	9.312	0.202	0.563	0.020	16.542	9.921	3618.48	≈ 0	≈ 0
newspaper(x_3)	12.351	0.055	0.621	0.017	19.876	3.3	5134.805	≈ 0	0.001148196
CC	2.529	0.860	0.345	0.025	7.335	34.321	779.557	≈ 0	≈ 0
H_0	$H_0: \beta_{TV} = \beta_{radio} = \beta_{newspaper} = 0$								

Multiple regression				
	$\widehat{\beta}_i$	$se(\widehat{\beta}_i)$	t_0	p
Intercept	2.938889369	0.311908	9.422	≈ 0
TV(x_{i1})	0.045764645	0.001395	32.809	≈ 0
Radio(x_{i2})	0.188530017	0.008611	21.893	≈ 0
Newspaper(x_{i3})	-0.001037493	0.005871	-0.177	0.86
Rss	556.8253			
F	F=570.2707			
H_0	$H_0: \beta_{TV} = \beta_{radio} = \beta_{newspaper} = 0$			

2. Use mathematical formulas to explain how to calculate all the values.

For simple linear regression

TV:

$$\bar{x}_1 = \frac{\sum_{i=1}^{200} x_i}{200} = 147.0425, \bar{y}_{sales} = \frac{\sum_{i=1}^{200} y_i}{200} = 14.0225$$

$$S_{x_1x_1} = \sum_{i=1}^{200} (x_i - \bar{x}_1)^2 = 1466819, S_{x_1y} = \sum_{i=1}^{200} (x_i - \bar{x}_1)(y_i - \bar{y}_{sales}) = 69727.65$$

$$\Rightarrow \hat{\beta}_1 = \frac{S_{x_1y}}{S_{x_1x_1}} = \frac{69727.65}{1466819} = 0.04753664, \quad \hat{\beta}_0 = \bar{y}_{sales} - \hat{\beta}_1 \times \bar{x}_1 = 7.032594$$

$$RSS = \sum_{i=1}^{200} (y_i - \hat{y}_i)^2 = \sum_{i=1}^{200} (y_i - \hat{\beta}_0 - \hat{\beta}_1 \times x_i)^2 = 2102.531$$

$$\hat{\sigma}^2 = \frac{1}{198} \sum_{i=1}^{200} (y_i - \hat{\beta}_0 - \hat{\beta}_1 \times x_i)^2 = 10.61884 \Rightarrow \hat{\sigma} = 3.258656$$

$$\Rightarrow se(\hat{\beta}_0) = \hat{\sigma} \sqrt{\frac{1}{n} + \frac{\bar{x}_1^2}{S_{x_1x_1}}} = 0.4578429, se(\hat{\beta}_1) = \sqrt{\frac{\hat{\sigma}^2}{S_{x_1x_1}}} = 0.002690607$$

$$\Rightarrow t - \text{statistic } t_0 = \frac{\hat{\beta}_0}{se(\hat{\beta}_0)} = \frac{7.032594}{0.002690607} = 15.36028, p - \text{value} = Pr(|t_{df=198}| > |t|) \approx 0$$

$$t - \text{statistic } t_1 = \frac{\hat{\beta}_1}{se(\hat{\beta}_1)} = \frac{0.04753664}{0.002690607} = 17.66763, p - \text{value} = Pr(|t_{df=198}| > |t|) \approx 0$$

radio:

$$\bar{x}_2 = \frac{\sum_{i=1}^{200} x_i}{200} = 23.264, \bar{y}_{sales} = \frac{\sum_{i=1}^{200} y_i}{200} = 14.0225$$

$$S_{x_2x_2} = \sum_{i=1}^{200} (x_i - \bar{x}_2)^2 = 43.86512, S_{x_2y} = \sum_{i=1}^{200} (x_i - \bar{x}_2)(y_i - \bar{y}_{sales}) = 8882.502$$

$$\Rightarrow \hat{\beta}_1 = \frac{S_{x_2y}}{S_{x_2x_2}} = \frac{8882.502}{43.86512} = 0.2024958, \quad \hat{\beta}_0 = \bar{y}_{sales} - \hat{\beta}_1 \times \bar{x}_2 = 9.311638$$

$$RSS = \sum_{i=1}^{200} (y_i - \hat{y}_i)^2 = \sum_{i=1}^{200} (y_i - \hat{\beta}_0 - \hat{\beta}_1 \times x_i)^2 = 3618.48$$

$$\hat{\sigma}^2 = \frac{1}{198} \sum_{i=1}^{200} (y_i - \hat{\beta}_0 - \hat{\beta}_1 \times x_i)^2 = 18.27515 \Rightarrow \hat{\sigma} = 4.274944$$

$$\Rightarrow se(\hat{\beta}_0) = \hat{\sigma} \sqrt{\frac{1}{200} + \frac{\bar{x}_2^2}{S_{x_2x_2}}} = 0.5629005, se(\hat{\beta}_1) = \sqrt{\frac{\hat{\sigma}^2}{S_{x_2x_2}}} = 0.02041131$$

$$\Rightarrow t - \text{statistic } t_0 = \frac{\hat{\beta}_0}{se(\hat{\beta}_0)} = \frac{9.311638}{0.5629005} = 16.54225, p - \text{value} = Pr(|t_{df=198}| > |t|) \approx 0$$

$$t - \text{statistic } t_1 = \frac{\hat{\beta}_1}{se(\hat{\beta}_1)} = \frac{0.2024958}{0.02041131} = 9.920765, p - \text{value} = Pr(|t_{df=198}| > |t|) \approx 0$$

newspaper:

$$\bar{x}_3 = \frac{\sum_{i=1}^{200} x_i}{200} = 30.554, \bar{y}_{sales} = \frac{\sum_{i=1}^{200} y_i}{200} = 14.0225$$

$$S_{x_3x_3} = \sum_{i=1}^{200} (x_i - \bar{x}_3)^2 = 94387.36, S_{x_3y} = \sum_{i=1}^{200} (x_i - \bar{x}_3)(y_i - \bar{y}_{sales}) = 5162.337$$

$$\Rightarrow \hat{\beta}_1 = \frac{S_{x_3y}}{S_{x_3x_3}} = \frac{5162.337}{94387.36} = 0.0546931, \hat{\beta}_0 = \bar{y}_{sales} - \hat{\beta}_1 \times \bar{x}_3 = 12.35141$$

$$RSS = \sum (y_i - \hat{y}_i)^2 = \sum (y_i - \hat{\beta}_0 - \hat{\beta}_1 \times x_i)^2 = 5134.805$$

$$\hat{\sigma}^2 = \frac{1}{198} \sum (y_i - \hat{\beta}_0 - \hat{\beta}_1 \times x_i)^2 = 25.93336 \Rightarrow \hat{\sigma} = 5.09248$$

$$\Rightarrow se(\hat{\beta}_0) = \hat{\sigma} \sqrt{\frac{1}{200} + \frac{\bar{x}_3^2}{S_{x_3x_3}}} = 0.6214202, se(\hat{\beta}_1) = \sqrt{\frac{\hat{\sigma}^2}{S_{x_3x_3}}} = 0.01657572$$

$$\Rightarrow t - \text{statistic } t_0 = \frac{\hat{\beta}_0}{se(\hat{\beta}_0)} = \frac{12.35141}{0.6214202} = 19.8761, p - \text{value} = Pr(|t_{df=198}| > |t|) \approx 0$$

$$t - \text{statistic } t_1 = \frac{\hat{\beta}_1}{se(\hat{\beta}_1)} = \frac{0.0546931}{0.01657572} = 3.299591, p - \text{value} = Pr(|t_{df=198}| > |t|)$$

$$= 0.001148196$$

CC:

$$\bar{CC} = \frac{\sum_{i=1}^{200} CC_i}{200} = 13.37186, \bar{y}_{sales} = \frac{\sum_{i=1}^{200} y_i}{200} = 14.0225$$

$$S_{CC} = \sum_{i=1}^{200} (CC_i - \bar{CC})^2 = 6277.578, S_{CCy} = \sum_{i=1}^{200} (CC_i - \bar{CC})(y_i - \bar{y}_{sales}) = 5395.632$$

$$\Rightarrow \hat{\beta}_1 = \frac{S_{CCy}}{S_{CC}} = \frac{5395.632}{6277.578} = 0.8595085, \hat{\beta}_0 = \bar{y}_{sales} - \hat{\beta}_1 \times \bar{CC} = 2.529271$$

$$RSS = \sum_{i=1}^{200} (y_i - \hat{y}_i)^2 = \sum_{i=1}^{200} (y_i - \hat{\beta}_0 - \hat{\beta}_1 \times CC_i)^2 = 779.5574$$

$$\hat{\sigma}^2 = \frac{1}{198} \sum_{i=1}^{200} (y_i - \hat{\beta}_0 - \hat{\beta}_1 \times CC_i)^2 = 3.9837158 \Rightarrow \hat{\sigma} = 1.984227$$

$$\Rightarrow se(\hat{\beta}_0) = \hat{\sigma} \sqrt{\frac{1}{200} + \frac{\bar{CC}^2}{S_{CC}}} = 0.3447996, se(\hat{\beta}_1) = \sqrt{\frac{\hat{\sigma}^2}{S_{CC}}} = 0.02504352$$

$$\Rightarrow t - \text{statistic } t_0 = \frac{\hat{\beta}_0}{se(\hat{\beta}_0)} = \frac{2.529271}{0.3447996} = 7.335481, p - \text{value} = Pr(|t_{df=198}| > |t|)$$

$$\approx 5.549561e - 12$$

$$t - \text{statistic } t_1 = \frac{\hat{\beta}_1}{se(\hat{\beta}_1)} = \frac{0.8595085}{0.02504352} = 34.32059, p - \text{value} = Pr(|t_{df=198}| > |t|) \approx 0$$

For multiple linear regression

$$\hat{\beta} = (X^T X)^{-1} X^T Y, \text{ where } \hat{\beta} = \begin{bmatrix} \hat{\beta}_0 \\ \hat{\beta}_1 \\ \hat{\beta}_2 \\ \hat{\beta}_3 \end{bmatrix}_{4 \times 1}, X = \begin{bmatrix} 1 & x_{1,1} & x_{1,2} & x_{1,3} \\ 1 & x_{2,1} & x_{2,2} & x_{2,3} \\ \vdots & \vdots & \vdots & \vdots \\ 1 & x_{200,1} & x_{200,2} & x_{200,3} \end{bmatrix}_{200 \times 4}, Y = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_{200} \end{bmatrix}_{200 \times 1}$$

$$\hat{\beta}_0 = 2.938889369, \hat{\beta}_1 = 0.045764645, \hat{\beta}_2 = 0.188530017, \hat{\beta}_3 = -0.001037493$$

Let $\{(x_i, y_i) | i = 1, 2, \dots, 200\}$, where $x_i = (x_{i1}, x_{i2}, x_{i3}), y_i$

$$= \beta_0 + \beta_1 x_{i1} + \beta_2 x_{i2} + \beta_3 x_{i3} + \varepsilon_i, \varepsilon_i \sim N(0, \sigma^2)$$

$$\hat{Y} = X(X^T X)^{-1} X^T Y = X \hat{\beta} \Rightarrow SSR_{eg} = (y - X \hat{\beta})^T (y - X \hat{\beta}) = 556.8253$$

$$SST = \sum_{i=1}^{200} (y_i - \bar{y})^2 = 5417.149,$$

To test $H_0: \beta_1 = \beta_2 = \beta_3 = 0$ with $F = \frac{(SST - SSR_{eg})/3}{SSR_{eg}/196} = \frac{(5417.149 - 556.8253)/3}{556.8253/196} = 570.2707$ with

$$p\text{-value} = Pr(F_{(3,196)} > F)$$

$$se(\hat{\beta}_j) = \sqrt{\frac{SSR_{eg}}{196} C_{jj}}, \text{ where } C = (X^T X)^{-1}$$

$$se(\hat{\beta}_0) = \sqrt{\frac{SSR_{eg}}{196} C_{11}} = \sqrt{2.840945 \times 0.0342445} = 0.3119082$$

$$se(\hat{\beta}_1) = \sqrt{\frac{SSR_{eg}}{196} C_{22}} = \sqrt{2.840945 \times 6.848908 \times 10^{-7}} = 0.001394897$$

$$se(\hat{\beta}_2) = \sqrt{\frac{SSR_{eg}}{196} C_{33}} = \sqrt{2.840945 \times 2.610165 \times 10^{-5}} = 0.008611234$$

$$se(\hat{\beta}_3) = \sqrt{\frac{SSR_{eg}}{196} C_{44}} = \sqrt{2.840945 \times 1.213285 \times 10^{-5}} = 0.00587101$$

To test $H_0: \beta_0 = \beta_1 = \beta_2 = \beta_3 = 0$ with $t_i = \frac{\hat{\beta}_i}{\sqrt{\frac{SSR_{eg}}{196} C_{jj}}}$

$$t_0 = \frac{2.938889369}{0.3119082} = 9.422288, p\text{-value} = Pr(|t_{df=196}| > |t|) \approx 0$$

$$t_1 = \frac{0.045764645}{0.001394897} = 32.80862, p\text{-value} = Pr(|t_{df=196}| > |t|) \approx 0$$

$$t_2 = \frac{0.188530017}{0.008611234} = 21.8935, p\text{-value} = Pr(|t_{df=196}| > |t|) \approx 0$$

$$t_3 = \frac{-0.001037493}{0.00587101} = -0.1767146, p\text{-value} = Pr(|t_{df=196}| > |t|) = 0.8599132$$

3.Explain how to test H_0 :There is no relationship between sales and 3 budgets by compound covariate.State your conclusions.

Suppose that H_0 :There is no relationship between sales and 3 budgets by CC(compound covariate) then it is equivalent to test $H_0: \gamma_1 = 0$.

$$\text{sales} = \gamma_0 + \gamma_1 CC + \varepsilon, \varepsilon \sim N(0, \sigma^2)$$

$$\hat{\gamma}_0 = 2.529271, \hat{\gamma}_1 = 0.8595085$$

$$t - \text{statistic} = \frac{\hat{\gamma}_1}{se(\hat{\gamma}_1)} = \frac{0.8595085}{0.02504352} = 34.32059, p - \text{value} Pr(|t_{df=198}| > |34.32059|) \approx 0 < 0.05$$

\Rightarrow reject H_0 :There is no relationship between sales and 3 budgets ,there exists relationship bwtween sales and three budgets.

Appendix

R-code&their outputs

```
> data=read.table("C:\\Users\\teresa\\Desktop\\Advertising.csv",sep=";",header=T)
> t=data[1:200,2]
> r=data[1:200,3]
> n=data[1:200,4]
> s=data[1:200,5]
> mean(s)
[1] 14.0225
> ####simple#regression#TV
> mean(TV)
[1] 147.0425
> st=sum((t-mean(TV))^2)
```

```
> st
[1] 1466819
> sts=sum((t-mean(TV))*(s-mean(s)))
> sts
[1] 69727.65
> bt1tv=sum((t-mean(TV))*(s-mean(s)))/st
> bt1tv
[1] 0.04753664
> bt0tv=mean(s)-bt1tv*mean(TV)
> bt0tv
[1] 7.032594
> sigma2_hat_tv=(sum((s-bt0tv-bt1tv*t)^2))/198
> sigma2_hat_tv
[1] 10.61884
> sigma2_hat_tv^0.5
[1] 3.258656
> se_bt0tv=(sigma2_hat_tv^0.5)*((1/200+mean(TV)^2/st)^0.5)
> se_bt0tv
[1] 0.4578429
> se_bt1tv=(sigma2_hat_tv/st)^0.5
> se_bt1tv
[1] 0.002690607
> t0tv=bt0tv/se_bt0tv
> t0tv
[1] 15.36028
> t1tv=bt1tv/se_bt1tv
> t1tv
[1] 17.66763
```

```
> ptv1=(1-pt(t1tv,198))*2
> ptv1
[1] 0
> ptv0=(1-pt(t0tv,198))*2
> ptv0
[1] 0
> (1-pt(t1tv,198))
[1] 0
> #####simple#regression#radio
> mean(radio)
[1] 23.264
> sr=sum((r-mean(radio))^2)
> sr
[1] 43865.12
> srs=sum((r-mean(radio))*(s-mean(s)))
> srs
[1] 8882.502
> bt1ra=sum((r-mean(radio))*(s-mean(s)))/sr
> bt1ra
[1] 0.2024958
> bt0ra=mean(sales)-bt1ra*mean(radio)
> bt0ra
[1] 9.311638
> sigma2_hat_ra=(1/198)*(sum((s-bt0ra-bt1ra*r)^2))
> sigma2_hat_ra
[1] 18.27515
> sigma2_hat_ra^0.5
[1] 4.274944
```

```
> se_bt0ra=(sigma2_hat_ra^0.5)*((1/200+mean(radio)^2/sr)^0.5)
```

```
> se_bt0ra
```

```
[1] 0.5629005
```

```
> se_bt1ra=(sigma2_hat_ra/sr)^0.5
```

```
> se_bt1ra
```

```
[1] 0.02041131
```

```
> t0ra=bt0ra/se_bt0ra
```

```
> t0ra
```

```
[1] 16.54225
```

```
> t1ra=bt1ra/se_bt1ra
```

```
> t1ra
```

```
[1] 9.920765
```

```
> pra1=(1-pt(t1ra,198))*2
```

```
> pra1
```

```
[1] 0
```

```
> pra0=(1-pt(t0ra,198))*2
```

```
> pra0
```

```
[1] 0
```

```
> #####simple#regression#newspaper
```

```
> mean(newspaper)
```

```
[1] 30.554
```

```
> sn=sum((n-mean(newspaper))^2)
```

```
> sn
```

```
[1] 94387.36
```

```
> sns=sum((n-mean(newspaper))*(s-mean(s)))
```

```
> sns
```

```
[1] 5162.337
```

```
> bt1ne=sum((n-mean(newspaper))*(s-mean(s)))/sn
```



```
> bt1ne
[1] 0.0546931
> bt0ne=mean(sales)-bt1ne*mean(newspaper)
> bt0ne
[1] 12.35141
> sigma2_hat_ne=(1/198)*(sum((s-bt0ne-bt1ne*n)^2))
> sigma2_hat_ne
[1] 25.93336
> sigma2_hat_ne^0.5
[1] 5.09248
> se_bt0ne=(sigma2_hat_ne^0.5)*((1/200+mean(newspaper)^2/sn)^0.5)
> se_bt0ne
[1] 0.6214202
> se_bt1ne=(sigma2_hat_ne/sn)^0.5
> se_bt1ne
[1] 0.01657572
> t0ne=bt0ne/se_bt0ne
> t0ne
[1] 19.8761
> t1ne=bt1ne/se_bt1ne
> t1ne
[1] 3.299591
> pnews1=(1-pt(t1ne,198))*2
> pnews1
[1] 0.001148196
> pnews0=(1-pt(t0ne,198))*2
> pnews0
[1] 0
```

```

> #####simple#regression#cc
> cc=bt1tv*t+bt1ra*r+bt1ne*n
> mean(cc)
[1] 13.37186
> sccs=sum((cc-mean(cc))*(s-mean(s)))
> sccs
[1] 5395.632
> scc=sum((cc-mean(cc))^2)
> scc
[1] 6277.578
> bt1cc=sccs/scc
> bt1cc
[1] 0.8595085
> bt0cc=mean(s)-bt1cc*mean(cc)
> bt0cc
[1] 2.529271
> sigma2_hat_cc=sum((s-bt0cc-bt1cc*cc)^2)/198
> sigma2_hat_cc
[1] 3.937158
> sigma2_hat_cc^0.5
[1] 1.984227
> se_bt0cc=(sigma2_hat_cc^0.5)*(1/200+mean(cc)/scc^0.5)
> se_bt0cc
[1] 0.3447996
> se_bt1cc=(sigma2_hat_cc/scc)^0.5
> se_bt1cc
[1] 0.02504352
> t0cc=bt0cc/se_bt0cc

```

```

> t0cc
[1] 7.335481
> t1cc=bt1cc/se_bt1cc
> t1cc
[1] 34.32059
> pcc1=(1-pt(t1cc,198))*2
> pcc1
[1] 0
> pcc0=(1-pt(t0cc,198))*2
> pcc0
[1] 5.549561e-12
> x=as.matrix(cbind(intercept=rep(1,200),data[,2:4]))
> y=as.matrix(data[,5])
> solve(t(x)%*%x)%*%t(x)%*%y
           [,1]
intercept  2.938889369
TV          0.045764645
radio       0.188530017
newspaper -0.001037493
> sst=sum((s-mean(s))^2)
> ssres=sum((s-2.938889369-0.045764645*t-0.188530017*r+0.001037493*n)^2)
> ssres
[1] 556.8253
> msr=(sst-ssres)/3
> mse=ssres/196
> F=msr/mse
> F
[1] 570.2707

```

```
> AA=solve(t(x)%**x)
> se0=(mse*AA[1,1])^0.5
> se0
[1] 0.3119082
> se1=(mse*AA[2,2])^0.5
> se1
[1] 0.001394897
> se2=(mse*AA[3,3])^0.5
> se2
[1] 0.008611234
> se3=(mse*AA[4,4])^0.5
> se3
[1] 0.00587101
> t0=2.938889369/se0
> t0
[1] 9.422288
> t1=0.045764645/se1
> t1
[1] 32.80862
> t2=0.188530017/se2
> t2
[1] 21.8935
> t3=-0.001037493/se3
> t3
[1] -0.1767146
> pp0=(1-pt(t0,198))*2
> pp0
[1] 0
```

```
> pp1=(1-pt(t1,198))*2
```

```
> pp1
```

```
[1] 0
```

```
> pp2=(1-pt(t2,198))*2
```

```
> pp2
```

```
[1] 0
```

```
> pp3=(1-pt(t3,198))*2
```

```
> pp3
```

```
[1] 0.8599132
```