MATH 484: HW 1

**Problem 1** The purity of oxygen (Y) produced by a fractional distillation process is thought to be related to the percentage of hydrocarbons (X) in the main condensor of the processing unit. Twenty samples were measured and shown in the attached data sheet.

1) Fit a simple linear regression model to the data, provide the liner regression line equation  $\hat{Y} = \hat{\beta}_0 + \hat{\beta}_1 X$ .

$$\vec{x} = 91.818$$
 $\vec{y} = 1.1825$ 

$$\hat{\beta}_{1} = \frac{2(x_{1} - \overline{x})(y_{1} - \overline{y})}{2(x_{1} - \overline{x})^{2}} \longrightarrow = 1.064975$$

$$= 11.80103$$

$$\hat{\beta}_0 = \hat{y} - \hat{\beta}_1 \overline{\chi}$$

$$= 91.818 - 11.80103 \times 1.1825$$

$$= 77.86328$$

$$\hat{Y} = \hat{\beta}_0 + \hat{\beta}_1 X$$

$$\hat{Y} = 77.863 + 11.801 X$$

2) Test the hypothesis  $H_0: \beta_1 = 0$  against the alternative  $H_A: \beta_1 \neq 0$ , and conclude if there is significant linear relationship between the purity of oxygen and the percentage of hydrocarbons.

$$5SE = \frac{2(\gamma_{1} - \hat{\gamma})^{2}}{n-2} = \frac{232.8344}{20-2}$$

$$5.e. (\hat{\beta}_{1}) = \sqrt{\frac{\hat{\delta}^{2}}{2(\chi_{1} - \overline{\chi})^{2}}} = \sqrt{12.93524}/1.044975 = 3.4851$$

$$\begin{cases} H_{0}: \beta_{1} = 0 \\ H_{a}: \beta_{1} \neq 0 \end{cases} \qquad t.test: T = \frac{\hat{\beta}_{1} - \beta_{1}^{0}}{5.e.(\hat{\beta}_{1})} = \frac{11.801 - 0}{3.4851} = 3.386$$

$$t_{a/2, h-2} = t_{0.025, 18} = 2.10$$

$$1 > t_{a/2, h-2} \rightarrow reject H_{0}$$

3) Calculate the coefficient of determination,  $r^2$ .

$$r^{2} = 1 - \frac{55E}{55T} \rightarrow 55E = 232.834$$
  
 $r^{2} = 0.389$ 

4) Find a 95% confidence interval on the slope.

95% C.1.: 
$$\hat{\beta}$$
,  $\pm$   $t_{0.05/2, n-2} \times S.E.(\hat{\beta},)$ 

11.801  $\pm$  2.10  $\times$  3.485  $\longrightarrow$  [4.4791, 19.123]

5) Find a 95% confidence interval on the mean purity when the hydrocarbon percentage is 1.05.

$$\hat{Y}_{0} = \hat{\beta}_{0} + \hat{\beta}_{1} \times 1.05$$

$$= 77.863 + 11.801 \times 1.05$$

$$= 90.2541$$

959. C.1. for 
$$E(\hat{Y}_{0})$$
:
$$\hat{Y}_{0} \pm t_{\alpha/2, \, N-2} \times S.E.(\hat{Y}_{0}) = \sqrt{\frac{5SE}{n-2}} \times \left(\frac{1}{n} + \frac{(x_{0} - \bar{x})^{2}}{\xi(x_{1} - x_{0})^{2}}\right)$$

$$= \sqrt{\frac{5SE}{n-2}} \times \left(\frac{1}{n} + \frac{(x_{0} - \bar{x})^{2}}{\xi(x_{1} - x_{0})^{2}}\right)$$

$$= \sqrt{\frac{232.934}{20.2}} \left(\frac{1}{20} + \frac{(1.05 - 1.1825)^{2}}{(\xi(x_{1} - 1.05)^{2})^{2}}\right)$$

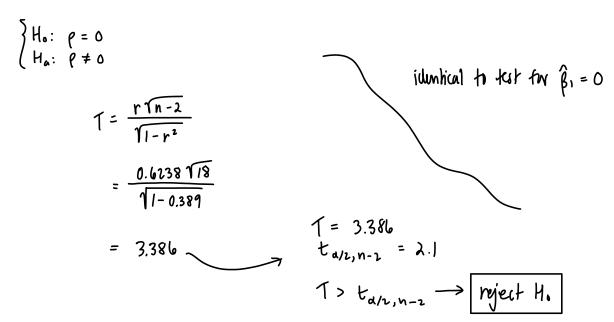
$$= \sqrt{0.807}$$

$$= 0.8984$$

6) What is the correlation coefficient between Y and X?

$$r = \sqrt{r^2} = \sqrt{0.389} = 0.6238$$

7) Test the hypothesis:  $H_0: \rho=0$  against  $H_A: \rho\neq 0$  using a t-test based on the correlation coefficient computed from the previous step.



**Problem 2** Consider the simple linear regression model  $Y = \beta_0 + \beta_1 X + \epsilon$ . Assume that  $E(\epsilon) = 0$  and  $Var(\epsilon) = \sigma^2$ , and  $\epsilon_i$  are independent of each other. Show that  $Cov(\bar{Y}, \hat{\beta}_1) = 0$ .

$$COV(\bar{y}, \hat{k}_1) = ?$$

$$\hat{\beta}_{i} = \frac{2(\chi_{i} - \bar{\chi})(\gamma_{i} - \bar{\gamma})}{2(\chi_{i} - \bar{\chi})^{2}}$$

$$= \frac{2(\chi_{i} - \bar{\chi})(\gamma_{i} - \bar{\gamma})}{2(\chi_{i} - \bar{\chi})^{2}}$$

$$= \frac{2(\chi_{i} + (\chi_{i} - \bar{\chi}))}{2(\chi_{i} - \bar{\chi})^{2}}$$

$$= \frac{2(\chi_{i} + (\chi_{i} - \bar{\chi}))}{2(\chi_{i} - \bar{\chi})^{2}}$$

$$= \frac{2(\chi_{i} + (\chi_{i} - \bar{\chi}))}{2(\chi_{i} - \bar{\chi})^{2}}$$

$$= 2(\chi_{i} + (\chi_{i} - \bar{\chi})) - (\bar{\gamma} + \bar{\chi} - \bar{\gamma} + \bar{\chi} - \bar{\chi})$$

$$= 2(\chi_{i} + (\chi_{i} - \bar{\chi})) - (\bar{\gamma} + \bar{\chi} - \bar{\chi} - \bar{\chi} - \bar{\chi})$$

$$= 2(\chi_{i} + (\chi_{i} - \bar{\chi})) - (\bar{\gamma} + \bar{\chi} - \bar{\chi}$$

$$\overline{\gamma} = \frac{\xi \gamma_i}{n}$$

$$COV(\bar{Y}_{1}, \hat{k}_{1}) = COV(\frac{\xi Y_{i}}{n}, \xi k_{i}Y_{i})$$

$$= \sum_{i=1}^{n} \frac{1}{n} \cdot k_{i} \cdot var(Y_{i}) + \sum_{i=1}^{n} \frac{1}{j^{2}} \cdot \frac{1}{n} \cdot k_{i} \cdot cov(X_{i}, Y_{j})$$

$$= \sum_{i=1}^{n} \frac{1}{n} \cdot k_{i} \cdot var(Y_{i})$$

$$= \frac{\sigma^{2}}{n} \sum_{i=1}^{n} k_{i}$$

$$= \frac{\sigma^{2}}{n} \sum_{i=1}^{n} k_{i}$$

$$= \frac{\sigma^{2}}{n} \sum_{i=1}^{n} (X_{i} - \bar{X})^{2}$$

$$= 0$$

**Problem 3** Show that the sample correlation coefficient r between X and Y is a value between -1 and 1, i.e.,  $-1 \le r \le 1$ . Use two different methods to show it, with one using only the data, and the other one in the context of simple linear regression.

## method 1:

$$r^{2} = 1 - \frac{55e}{557}$$

$$= 1 - \frac{2(4i - \hat{Y})^{2}}{2(7i - \bar{Y})^{2}}$$

$$= 1 - \frac{237.834}{381.147}$$

$$= 0.389$$

$$r = \pm 0.624$$

It will be the same sign as  $\hat{\beta}$ .

β, is positive

.. r = 0.624

Satisfies condition -1≤ r≤1 √√

## method 2:

$$SST = SSE + SSR = 2(y_i - \overline{y})^2$$

$$\Gamma^{2} = 1 - \frac{SSE}{SST}$$

$$= 1 - \frac{SSE}{SSE + SSR}$$

SSE will be 0 if y: = ŷ; for all i

Otherwise, SSE will be greakr than or equal to zero, since
it is a squared value.

55k is also a squared value — greater than or equal to zero

$$\frac{SSE}{SSE + SSR}$$
 is a fraction between 0 and 1, inclusive

$$\int min \left(1 - \frac{SSE}{SSE + SSR}\right) = 1 - 1 = 0$$

$$\int max \left(1 - \frac{SSE}{SSE + SSR}\right) = 1 - 0 = 1$$

$$r^2 \le 1$$
 and  $-1 \le r \le 1$ 

**Problem 4** Consider the simple linear regression model  $Y = \beta_0 + \beta_1 X + \epsilon$ . Assume that  $\epsilon_i \overset{i.i.d.}{\sim} N(0, \sigma^2)$ . Find the variance of  $\hat{Y}_h$ , where  $\hat{Y}_h = \hat{\beta}_0 + \hat{\beta}_1 x_h$ .

$$\hat{\beta}_{1} = \frac{\mathcal{E}(X; -\overline{X})(Y; -\overline{Y})}{\mathcal{E}(X; -\overline{X})^{2}} = \frac{\mathcal{E}Y; *(X; -\overline{X})}{\mathcal{E}(X; -\overline{X})^{2}}$$

$$\hat{\beta}_{0} = \overline{Y} - \hat{\beta}_{1}\overline{X}$$

$$\begin{split} \text{VAY} \Big( \widehat{\boldsymbol{y}}_h \Big) &= & \text{VAY} \Big( \widehat{\boldsymbol{\beta}}_0 + \widehat{\boldsymbol{\beta}}_1 \boldsymbol{X}_h \Big) \\ &= & \text{VAY} \Big( \overline{\boldsymbol{\gamma}} - \widehat{\boldsymbol{\beta}}_1 \overline{\boldsymbol{X}} + \widehat{\boldsymbol{\beta}}_1 \boldsymbol{X}_h \Big) \\ &= & \text{VAY} \Big( \frac{\underline{\zeta} \boldsymbol{y}_i}{n} + \widehat{\boldsymbol{\beta}}_1 (\overline{\boldsymbol{X}} - \boldsymbol{X}_h) \Big) \\ &= & \frac{\underline{\zeta} (\boldsymbol{X}_i - \overline{\boldsymbol{X}})^2}{\left[ \underline{\zeta} (\boldsymbol{X}_i - \overline{\boldsymbol{X}})^2 \right]^2} \, \text{VAY} \Big( \boldsymbol{y}_i \Big) \\ &= & \frac{1}{n} \, \sigma^2 + \left( \overline{\boldsymbol{X}} - \boldsymbol{X}_h \right)^2 \cdot \frac{1}{\underline{\zeta} (\boldsymbol{X}_i - \overline{\boldsymbol{X}})^2} \, \sigma^2 \\ &= & \frac{1}{n} \, \sigma^2 + \left( \overline{\boldsymbol{X}} - \boldsymbol{X}_h \right)^2 \cdot \frac{1}{\underline{\zeta} (\boldsymbol{X}_i - \overline{\boldsymbol{X}})^2} \, \sigma^2 \end{split}$$

**2.9** Let Y and X denote the labor force participation rate of women in 1972 and 1968, respectively, in each of 19 cities in the United States. The regression output for this data set is shown in Table 2.10. It was also found that SSR = 0.0358 and SSE = 0.0544. Suppose that the model  $Y = \beta_0 + \beta_1 X + \varepsilon$  satisfies the usual regression assumptions.

**Table 2.10** Regression Output When Y is Regressed on X for Labor Force Participation Rate of Women

Variable	Coefficient	s.e.	t-Test	p-value
Constant	<b>6.</b> 0.203311	0.0976	2.08	0.0526
X	<b>B</b> <sub>1</sub> 0.656040	0.1961	3.35	< 0.0038
n = 19	$R^2 = 0.397$	$R_a^2 = 0.362$	$\hat{\sigma} = 0.0566$	df = 17

(a) Compute Var(Y) and Cor(Y, X).

$$Var(Y) = \theta^{2} \qquad cor(Y,X) = \sqrt{\frac{SSR}{SSR}}$$

$$= \frac{SSE}{n-2} = \sqrt{\frac{SSR}{SSR} + SSE}$$

$$= \frac{0.0544}{19-2} = \sqrt{\frac{0.0358}{0.0358 + 0.0549}}$$

$$= 0.0566^{2}$$

$$Var(Y) = 0.0032$$
  
 $cor(Y, X) = 0.63$ 

(b) Suppose that the participation rate of women in 1968 in a given city is 45%. What is the estimated participation rate of women in 1972 for the same city?

$$Y = 0.203311 + 0.65604 * 0.45$$
  
= 0.4985

Estimated participation: 49.85%

(c) Suppose further that the mean and variance of the participation rate of women in 1968 are 0.5 and 0.005, respectively. Construct the 95% confidence interval for the estimate in (b).

C.1. = 0.4985 ± 
$$t_{\alpha/2}$$
,  $n-2$  \*  $6\sqrt{\frac{1}{n} + \frac{(x_0 - \overline{x})^2}{(x_0 - \overline{x})^2}}$ 

$$= 0.4985 \pm 2.11 * 0.0566 \sqrt{\frac{1}{19} + \frac{(0.45 - 0.5)^2}{0.09}}$$

$$t_{0.05/2, |9-2} = 2.11$$

(d) Construct the 95% confidence interval for the slope of the true regression line,  $\beta_1$ .

Given: 
$$\hat{\beta}_1 = 0.203311$$
  
S.e.  $(\hat{\beta}_1) = 0.1961$  (given)

C.1. = 
$$\hat{\beta}_1 \pm t_{\alpha/2, h-2} * s.e.(\hat{\beta}_1)$$
  
= 0.65604 ± 2.11 \* 0.1961  $= [0.2423, 1.06981]$ 

(e) Test the hypothesis:  $H_0: \beta_1=1$  versus  $H_1: \beta_1>1$  at the 5% significance level.

$$T = \frac{\hat{\beta}_{1} - \beta_{1}^{\circ}}{5.\ell.(\hat{\beta}_{1})}$$
 Nject if  $T > t_{\alpha, n-2}$ 

$$= \frac{0.65604 - 1}{0.1961}$$

$$t_{\alpha, n-2} = t_{0.05, 17}$$

$$= 1.74$$
T is not greater than t, so we fail to reject Ho.

(f) If Y and X were reversed in the above regression, what would you expect  $\mathbb{R}^2$  to be?

$$r^{2} = \left(\operatorname{cor}(Y, X)\right)^{2}$$

$$= \frac{\left[2(X_{i} - \overline{X})(Y_{i} - \overline{Y})\right]^{2}}{2(X_{i} - \overline{X})^{2} \cdot 2(Y_{i} - \overline{Y})^{2}}$$

$$\vdots e_{i}, \left[\operatorname{cor}(Y, X)\right]^{2} = \left[\operatorname{cor}(X, Y)\right]^{2}$$

$$\vdots r^{2} \text{ will not change}$$