

1.The attached data set "data.csv" is a subset of the Ames House data set. The first row of this data set is the name of the variables, where "SalePrice" is the dependent variable (column 1) and "Lot_Area" is the independent variable (predictor; column 2). There are $n = 1598$ observations in total (rows 2-1599). Suppose we use a linear regression model to approximate the relationship between SalePrice and Lot_Area:

$$\text{SalePrice} = \beta_0 + \beta_1 \text{Lot_Area} + \epsilon$$

Suppose that the noise ϵ_i i.i.d. $\sim N(0, \sigma^2)$. Calculate the MLEs of $\beta_0, \beta_1, \sigma^2$

solve:

let $Y = \text{SalePrice}$ and $X = \text{Lot_Area}$

$$\begin{aligned} Y_i &\sim N(\beta_0 + \beta_1 X_i, \sigma^2) \\ \Rightarrow L(\beta_0, \beta_1, \sigma^2) &= P(Y_1) P(Y_2) \dots P(Y_n) = \frac{1}{(2\pi)^{\frac{n}{2}} \sigma^n} e^{-\frac{1}{2\sigma^2} \sum (Y_i - \beta_0 - \beta_1 X_i)^2} \\ \Rightarrow \ell(\beta_0, \beta_1, \sigma^2) &= \ln(L) = -n \ln(\sqrt{2\pi}\sigma) - \frac{1}{2\sigma^2} \sum (Y_i - \beta_0 - \beta_1 X_i)^2 \end{aligned}$$

To maximize the ℓ

$$\Rightarrow \begin{cases} \frac{\ell(\beta_0, \beta_1, \sigma^2)}{\partial \beta_0} = 0 \\ \frac{\ell(\beta_0, \beta_1, \sigma^2)}{\partial \beta_1} = 0 \\ \frac{\ell(\beta_0, \beta_1, \sigma^2)}{\partial \sigma^2} = 0 \end{cases} \Rightarrow \begin{cases} \hat{\beta}_0 = \bar{Y} - \hat{\beta}_1 \bar{X} = 1.330 \times 10^5 \\ \hat{\beta}_1 = \frac{\sum_i (X_i - \bar{X})(Y_i - \bar{Y})}{\sum_i (X_i - \bar{X})^2} = 4.688 \\ \hat{\sigma}^2 = \frac{\sum_i (Y_i - \hat{\beta}_0 - \hat{\beta}_1 X_i)^2}{n} = \frac{1596}{1598} \times 74930^2 \end{cases}$$

```

1 Call:
2 lm(formula = Y ~ X)
3
4 Residuals:
5      Min       1Q   Median       3Q      Max
6 -370781 -46560  -18455   33684  418123
7
8 Coefficients:
9             Estimate Std. Error t value Pr(>|t|)
10 (Intercept) 1.330e+05  3.795e+03   35.04  <2e-16 ***
11 X           4.688e+00  3.308e-01   14.17  <2e-16 ***
12 ---
13 Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
14
15 Residual standard error: 74930 on 1596 degrees of freedom
16 Multiple R-squared:  0.1118, Adjusted R-squared:  0.1112

```

```
17 | F-statistic: 200.8 on 1 and 1596 DF, p-value: < 2.2e-16
```

Code

```
1 | data=read.csv('/Users/liuchenghua/Downloads/data.csv',header = TRUE)
2 | Y=data['SalePrice']
3 | X=data['Lot_Area']
4 | Y<-unlist(Y)
5 | X<-unlist(X)
6 | fit=lm(Y~X)
7 | summary(fit)
```

2. Let $\mathbf{X} = (X_1, \dots, X_n)$ be a random sample of size $n = 15$ from a normal distribution $N(\mu, 0.25^2)$, and the observed values are:

2.9, 2.8, 3.0, 2.8, 3.1, 2.7, 2.3, 2.8, 2.4, 2.8, 2.6, 2.6, 3.1, 3.2, 2.9

Construct a 95% confidence interval for μ

We know that

$$\bar{X} \sim N\left(\mu, \frac{\sigma^2}{n}\right)$$

```
1 | > left
2 | [1] 2.673485
3 | > right
4 | [1] 2.926515
```

Code

```
1 | x=c(2.9,2.8,3.0,2.8,3.1,2.7,2.3,2.8,2.4,2.8,2.6,2.6,3.1,3.2,2.9)
2 | sigma=0.25
3 | alpha=0.05
4 | len=length(x)
5 | left=mean(x)-qnorm(1-alpha/2)*sigma/sqrt(len)
6 | right=mean(x)+qnorm(alpha/2)*sigma/sqrt(len)
```