

Recommended Exercise 7 in Statistical Linear Models, Spring 2021

alexaoH

14 mars, 2021

Problem 1 Inference about a new observation in multiple linear regression

a)

Since $\hat{\beta} = (X^T X)^{-1} X^T \mathbf{Y}$, we know that $\hat{\beta} \sim N(\beta, \sigma^2 (X^T X)^{-1})$. Since the Gaussian distribution is closed under linear transformations, this means that $\mathbf{x}_0^T \hat{\beta} \sim N(\mathbf{x}_0^T \beta, \sigma^2 \mathbf{x}_0^T (X^T X)^{-1} \mathbf{x}_0)$, which is univariate. Now, since $E \mathbf{x}_0^T \hat{\beta} = \mathbf{x}_0^T \beta$, $\mathbf{x}_0^T \hat{\beta}$ is an unbiased estimator of $EY_0 = \mathbf{x}_0^T \beta$, with distribution as given above.

b)

Generally, we know that if $Z \sim N(0, 1)$ and $U \sim \chi_q^2$, where Z and U are independent, then $\frac{Z}{\sqrt{U/q}} \sim t_q$. In this case, we can set $Z = \frac{\mathbf{x}_0^T \hat{\beta} - \mathbf{x}_0^T \beta}{\sqrt{\sigma^2 \mathbf{x}_0^T (X^T X)^{-1} \mathbf{x}_0}}$ and $U = \frac{1}{\sigma^2} \text{SSE} = \frac{(n-p)\hat{\sigma}^2}{\sigma^2} \sim \chi_{n-p}^2$. Moreover, Z and U are independent, because $\hat{\beta}$ and $\hat{\sigma}^2$ are independent. Hence, the estimator

$$T = \frac{Z}{\sqrt{U/q}} = \frac{\frac{\mathbf{x}_0^T \hat{\beta} - \mathbf{x}_0^T \beta}{\sqrt{\sigma^2 \mathbf{x}_0^T (X^T X)^{-1} \mathbf{x}_0}}}{\sqrt{\frac{(n-p)\hat{\sigma}^2}{\sigma^2} / (n-p)}} = \frac{\mathbf{x}_0^T \hat{\beta} - \mathbf{x}_0^T \beta}{\sqrt{\hat{\sigma}^2 \mathbf{x}_0^T (X^T X)^{-1} \mathbf{x}_0}} \sim t_{n-p}.$$

This means that

$$\begin{aligned} -t &\leq \frac{\mathbf{x}_0^T \hat{\beta} - \mathbf{x}_0^T \beta}{\sqrt{\hat{\sigma}^2 \mathbf{x}_0^T (X^T X)^{-1} \mathbf{x}_0}} \leq t, \\ \iff \mathbf{x}_0^T \hat{\beta} - t\hat{\sigma} \sqrt{\mathbf{x}_0^T (X^T X)^{-1} \mathbf{x}_0} &\leq \mathbf{x}_0^T \beta \leq \mathbf{x}_0^T \hat{\beta} + t\hat{\sigma} \sqrt{\mathbf{x}_0^T (X^T X)^{-1} \mathbf{x}_0}, \\ \iff \mathbf{x}_0^T \hat{\beta} - t\hat{\sigma} \sqrt{\mathbf{x}_0^T (X^T X)^{-1} \mathbf{x}_0} &\leq EY_0 \leq \mathbf{x}_0^T \hat{\beta} + t\hat{\sigma} \sqrt{\mathbf{x}_0^T (X^T X)^{-1} \mathbf{x}_0}. \end{aligned}$$

Hence,

$$\begin{aligned} 1 - \alpha &= P \left(-t_{\alpha/2} \leq \frac{\mathbf{x}_0^T \hat{\beta} - \mathbf{x}_0^T \beta}{\sqrt{\hat{\sigma}^2 \mathbf{x}_0^T (X^T X)^{-1} \mathbf{x}_0}} \leq t_{\alpha/2} \right) \\ &= P \left(\mathbf{x}_0^T \hat{\beta} - t_{\alpha/2} \hat{\sigma} \sqrt{\mathbf{x}_0^T (X^T X)^{-1} \mathbf{x}_0} \leq EY_0 \leq \mathbf{x}_0^T \hat{\beta} + t_{\alpha/2} \hat{\sigma} \sqrt{\mathbf{x}_0^T (X^T X)^{-1} \mathbf{x}_0} \right), \end{aligned}$$

which means that $\mathbf{x}_0^T \hat{\beta} \pm t_{\alpha/2} \hat{\sigma} \sqrt{\mathbf{x}_0^T (X^T X)^{-1} \mathbf{x}_0}$ is a $100(1 - \alpha)\%$ confidence interval for EY_0 .

c)

A similar process to the one used in b) is used. For a future observation Y_0 we can derive that $Y_0 - \mathbf{x}_0^T \hat{\beta} \sim N(0, \sigma^2(1 + \mathbf{x}_0^T(X^T X)^{-1}\mathbf{x}_0))$, since $Y_0 = \mathbf{x}_0^T \beta + \varepsilon_0$ and $\mathbf{x}_0^T \hat{\beta}$ is an unbiased estimator of EY_0 . Setting $Z = \frac{Y_0 - \mathbf{x}_0^T \hat{\beta}}{\sigma \sqrt{1 + \mathbf{x}_0^T(X^T X)^{-1}\mathbf{x}_0}}$ and set U as in b). Solving the same inequalities as in b) then gives

$$1 - \alpha = P\left(\mathbf{x}_0^T \hat{\beta} - t_{\alpha/2} \hat{\sigma} \sqrt{1 + \mathbf{x}_0^T(X^T X)^{-1}\mathbf{x}_0} \leq Y_0 \leq \mathbf{x}_0^T \hat{\beta} + t_{\alpha/2} \hat{\sigma} \sqrt{1 + \mathbf{x}_0^T(X^T X)^{-1}\mathbf{x}_0}\right),$$

which means that $\mathbf{x}_0^T \hat{\beta} \pm t_{\alpha/2} \hat{\sigma} \sqrt{1 + \mathbf{x}_0^T(X^T X)^{-1}\mathbf{x}_0}$ is a $100(1 - \alpha)\%$ prediction interval for Y_0 . Why does Y_0 independent of ε give independence from $\mathbf{x}_0^T \hat{\beta}$?

d)

```
acidrain <- read.table("https://www.math.ntnu.no/emner/TMA4267/2018v/acidrain.txt", header=TRUE)
fit <- lm(y~., data = acidrain)

# Confidence interval: Done manually via the equations found in previous tasks.
X <- model.matrix(fit)
x0 <- c(1, 3, 50, 1, 50, 2, 1, 0)
n <- dim(X)[1]
p <- dim(X)[2]
quantile <- qt(0.025, n-p, lower.tail = FALSE)
prediction <- t(x0) %*% fit$coefficients
upper.conf <- prediction + quantile * summary(fit)$sigma * sqrt(t(x0) %*% solve(t(X) %*% X) %*% x0)
lower.conf <- prediction - quantile * summary(fit)$sigma * sqrt(t(x0) %*% solve(t(X) %*% X) %*% x0)
prediction

#>           [,1]
#> [1,] 5.531684
lower.conf

#>           [,1]
#> [1,] 5.446329
upper.conf

#>           [,1]
#> [1,] 5.617039

# Confidence interval: Done in R.
newdata <- data.frame(x1=3, x2=50, x3=1, x4=50, x5=2, x6=1, x7=0)
predict(fit, newdata, level=.95, interval="confidence")

#>           fit           lwr           upr
#> 1 5.531684 5.446329 5.617039

# Prediction interval: Done manually via the equations found in previous tasks.
upper.pred <- prediction + quantile * summary(fit)$sigma * sqrt(1 + t(x0) %*% solve(t(X) %*% X) %*% x0)
lower.pred <- prediction - quantile * summary(fit)$sigma * sqrt(1 + t(x0) %*% solve(t(X) %*% X) %*% x0)
prediction

#>           [,1]
#> [1,] 5.531684
```

```

lower.pred

#>           [,1]
#> [1,] 5.272555
upper.pred

#>           [,1]
#> [1,] 5.790813
# Prediction interval: Done in R.
predict(fit, newdata, level=.95, interval="prediction")

#>      fit      lwr      upr
#> 1 5.531684 5.272555 5.790813

```

e)

Problem 2 Plant stress

Problem 3 Multiple testing with plant stress

```

pvalues <- scan("https://www.math.ntnu.no/emner/TMA4267/2018v/damagePvalues.txt")
m <- length(pvalues)

```

a)

The family-wise error rate (FWER) is the probability of one or more false positive findings $P(V > 0)$, where V is the number of false positive findings among the m tests. The false discovery rate (FDR) gives the expected proportion of false positive results among the m hypotheses tested. A false positive is the same as a type I error, a case where the null hypothesis is rejected despite it being true. In the case of multiple linear regression this means that it is concluded that there is a linear relationship between the predictors and the response, despite the relationship not existing.

b)

Control FWER at level $\alpha = 0.05$ with the Bonferroni method. This gives $\alpha_{loc} = \frac{0.05}{m} = 5 \times 10^{-6}$. In this case one would reject `length(pvalues[pvalues<5e-6]) = 19` null-hypotheses in the data. Requirements when using Bonferroni's method are