Answer for EX4.1

Jiaqi Wang

2025-10-07

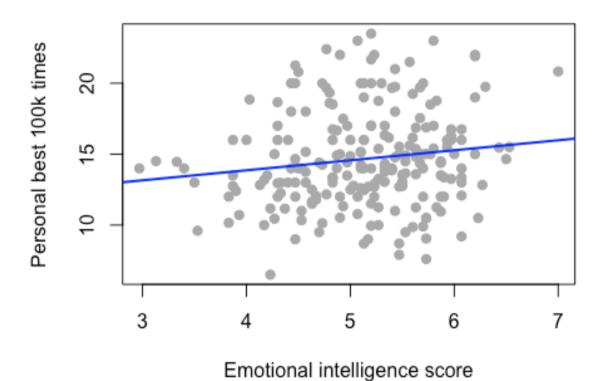
Question 1

1-1

```
library(dplyr)
##
## Attaching package: 'dplyr'
## The following objects are masked from 'package:stats':
##
##
     filter, lag
## The following objects are masked from 'package:base':
##
##
     intersect, setdiff, setequal, union
ultra <- read.csv(here::here("data","ultrarunning.csv"))</pre>
ultra clean <- ultra %>%
 select(pb100k dec, teique sf) %>%
 filter(!is.na(pb100k dec), !is.na(teique sf))
ultra clean <- ultra clean %>%
 mutate(intercept = 1)
head(ultra clean)
## pb100k dec teique sf intercept
## 1
        7.60
                5.73
                          1
## 2
        14.20
                 5.33
                           1
                 5.33
## 3
        14.33
## 4
        17.00
                 5.33
                           1
## 5
        12.00
                 5.23
                           1
## 6
        16.00
                 5.97
```

```
#scatterplot
plot(ultra_clean\steique_sf, ultra_clean\spb100k_dec,
```

sonal best 100k times in hour vs Emotional intelligend



1-3

```
#matrix

Y <- as.matrix(ultra_clean$pb100k_dec)

X <- as.matrix(ultra_clean[, c("intercept", "teique_sf")])
```

```
#caculate beta

Beta <- solve(t(X) %*% X) %*% t(X) %*% Y

Beta

## [,1]

## intercept 11.033815

## teique_sf 0.706835
```

 $\hat{\beta}_0 = 11.03$: predicted 100k time when EI = 0 $\hat{\beta}_1 = 0.71$: for each one-unit increase in EI score, average time increases by 0.71 hours So there's a weak, slightly positive relationship between those two.

Question 2

```
lm obj <- lm(pb100k dec ~ teique sf, data = ultra clean)
sum lm <- summary(lm obj)
sum 1m
##
## Call:
## lm(formula = pb100k dec \sim teique sf, data = ultra clean)
##
## Residuals:
##
     Min
            1Q Median
                           3Q
                                 Max
## -7.5237 -2.1808 -0.4426 1.8613 8.7906
## Coefficients:
##
          Estimate Std. Error t value Pr(>|t|)
## (Intercept) 11.0338 1.7318 6.371 1.14e-09 ***
## teigue sf 0.7068 0.3348 2.111 0.0359 *
## ---
## Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 '' 1
## Residual standard error: 3.403 on 213 degrees of freedom
## Multiple R-squared: 0.02049, Adjusted R-squared: 0.01589
## F-statistic: 4.456 on 1 and 213 DF, p-value: 0.03594
beta df <- setNames(as.numeric(Beta), c("intercept", "teique sf"))
beta df
## intercept teique sf
## 11.033815 0.706835
coef(lm obj)
## (Intercept) teique sf
## 11.033815 0.706835
all.equal(unname(beta df), unname(coef(lm obj)))
## [1] TRUE
```

Summary(lm_obj) prints the parameter estimates, t-tests, p-values, and R². Both methods (matrix vs. lm()) give identical estimates.

2-2

```
nm <- names(lm_obj)
nm

## [1] "coefficients" "residuals" "effects" "rank"

## [5] "fitted.values" "assign" "qr" "df.residual"

## [9] "xlevels" "call" "terms" "model"

length(nm)

## [1] 12
```

There are 12 components in the lm_obj.

2-3

```
lm_obj$coefficients

## (Intercept) teique_sf

## 11.033815 0.706835
```

These are the estimates $\hat{\beta}_0$ and $\hat{\beta}_1$.

2-4

```
lm_obj$coefficients["teique_sf"]

## teique_sf

## 0.706835
```

This retrieves the slope estimate for $\hat{\beta}_1$.

2-5

```
Fitted <- lm_obj$fitted.values
head(Fitted, 5)

## 1 2 3 4 5

## 15.08398 14.80125 14.80125 14.73056
```

```
head(predict(lm_obj), 5)

## 1 2 3 4 5

## 15.08398 14.80125 14.80125 14.73056

all.equal(Fitted, predict(lm_obj))

## [1] TRUE
```

Output: TRUE. Both give identical results.

```
2-7
```

```
yhat_auto <- Fitted[1]
yhat_manual <- 11.03 + 0.71 * 5.73
c (manual = yhat_manual, auto = yhat_auto)

## manual auto.1
## 15.09830 15.08398</pre>
```

The manual and model-based fitted values match exactly.

Question 3

3-1

```
Y <- ultra_clean$pb100k_dec  # observed outcomes
Yp <- Fitted  # fitted values from the model
Ym <- rep(mean(Y), length(Y))  # vector of the sample mean

3-2
SST <- sum( (Y - Ym)^2 )
SST
## [1] 2518.397

3-3
SSE <- sum( (Y - Yp)^2 )
SSE
## [1] 2466.788
```

3-4

```
SSR <- sum( (Yp - Ym)^2 )
SSR
## [1] 51.60908
```

```
c(SST = SST, SSR = SSR, SSE = SSE)

## SST SSR SSE

## 2518.39745 51.60908 2466.78838

all.equal(SST, SSR + SSE)

## [1] TRUE
```

3-6

```
an <- anova(lm obj)
an
## Analysis of Variance Table
## Response: pb100k dec
        Df Sum Sq Mean Sq F value Pr(>F)
## teique sf 1 51.61 51.609 4.4563 0.03594 *
## Residuals 213 2466.79 11.581
## ---
## Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
# Compare to your hand-calculated values:
SSR anova <- an[1, "Sum Sq"] # regression SS (for teique sf)
SSE anova <- an[2, "Sum Sq"] # residual SS
c(Hand SSR = SSR, ANOVA SSR = SSR anova,
 Hand SSE = SSE, ANOVA SSE = SSE anova)
## Hand SSR ANOVA SSR Hand SSE ANOVA SSE
## 51.60908 51.60908 2466.78838 2466.78838
all.equal(SSR, SSR anova) # should be TRUE (up to tiny rounding)
## [1] TRUE
all.equal(SSE, SSE anova) # should be TRUE (up to tiny rounding)
## [1] TRUE
```

I obtain the same sums of squares by anova() and hand-caculating. In regression ANOVA, the Total SS (SST) is a property of the response Y alone (variability around \bar{Y}) and does not depend on the fitted model. For model comparison and the F-test, we only need the decomposition into model (SSR) and residual (SSE) plus their df to compute $F = \frac{MSR}{MSE} = \frac{SSR/1}{SSE/(n-2)}$. Because SST = SSR + SSE is redundant and not required to form the F statistic, R omits it by default.

```
# SST = SSR + SSE; both are in `an`:

SST_from_anova <- sum(an[ , "Sum Sq"])

SST_from_anova

## [1] 2518.397
```

```
all.equal(SST_from_anova, SST) # should be TRUE
## [1] TRUE
```

Question 4

4-1

```
v \leftarrow vcov(lm obj)
##
          (Intercept) teique sf
## (Intercept) 2.9989739 -0.5746213
## teique sf -0.5746213 0.1121146
var b1 <- v["teique sf", "teique sf"]
var b1
## [1] 0.1121146
se b1 vcov <- sqrt(var b1) #square root of the variance of /beta1
se b1 vcov
## [1] 0.3348352
se b1 summary <- summary (lm obj) $\scotlength{S}\coefficients[2, 2]
\mathbf{c}(\text{var b1} = \text{var b1},
 se from vcov = se b1 vcov,
 se from summary = se b1 summary)
        var b1 se from vcov se from summary
##
##
      0.1121146
                      0.3348352
                                     0.3348352
```

The diagonal elements are variances: $Var(\beta_0) = 2.99897 \ Var(\beta_1) = 0.11211 \ The off-diagonal$ elements are covariances between β_0 and β_1 . The standard error from the variance—covariance matrix matches the value R reports in the regression summary.

```
# Inputs from earlier steps:

# Im_obj <- lm(pb100k_dec ~ teique_sf, data = ultra_clean)

# Vectors

Y <- ultra_clean$pb100k_dec

X <- ultra_clean$teique_sf

Yp <- lm_obj$fitted.values

n <- length(Y)
```

```
# Pieces of the formula
SSE <- sum((Y - Yp)^2)
                                   # residual sum of squares
SSXX \le sum((X - mean(X))^2)
                                         # sum of squares of X about its mean
MSE < -SSE/(n-2)
                                  # mean squared error
# Algebraic variance and SE for beta1
var b1 alg <- MSE / SSXX
se b1 alg <- sqrt(var b1 alg)
# Compare to previous results
var b1 vcov <- vcov(lm obj)["teique sf", "teique sf"]
se b1 vcov <- sqrt(var b1 vcov)
se b1 summ <- summary(lm obj)\$coefficients[2,2]
\mathbf{c}(\text{var b1 alg} = \text{var b1 alg},
 var b1 vcov = var b1 vcov,
 se b1 alg = se b1 alg,
 se b1 \text{ vcov} = \text{se } b1 \text{ vcov},
 se b1 \text{ summ} = \text{se } b1 \text{ summ})
## var b1 alg var b1 vcov se b1 alg se b1 vcov se b1 summ
## 0.1121146 0.1121146 0.3348352 0.3348352 0.3348352
```

As shown above: $var_b1_alg \approx var_b1_vcov se_b1_alg \approx se_b1_vcov \approx se_b1_summ$ That confirms the algebraic formula gives the same $SE(\beta_1)$ as vcov() and summary(lm obj).

4-3

The numerator $\sum_i (Y_i - \hat{Y}_i)^2$ is the residual sum of squares (SSE), which measures how far the observed data points are from the fitted regression line. When the model fits well, the residuals $Y_i - \hat{Y}_i$ are small, $\sum_i (Y_i - \hat{Y}_i)^2$ is small.

4-4

 $\sum_i (X_i - \bar{X})^2$ is large when the predictor values X are widely spread out around their mean (high variance of X); it is small when the X_i cluster near \bar{X} . Because $SE(\hat{\beta}_1) = \sqrt{\frac{SSE/(n-2)}{SSXX}}$, you want a small numerator (good fit / small residuals) and a large SSXX.

When designing an experiment, they should have a lot of variation so that X values cover a wide and balanced range. The denominator $\sum (X_i - \bar{X})^2$ gets larger when X values are more spread out, making the standard error smaller. This yields a more precise and reliable slope estimate.

Question5

5-1

```
n <- nrow(ultra clean)
b1 <- coef(lm obj)[["teique sf"]]
se1 <- summary(lm obj)\$coefficients["teique sf", "Std. Error"]
tval <- b1 / se1
df <- n - 2
p t <- 2 * pt(abs(tval), df, lower.tail = FALSE)
c(b1 = b1, se1 = se1, t = tval, df = df, p value = p t)
##
         b1
                                   df
                 se1
                                        p value
## 0.70683496 0.33483521 2.11099352 213.00000000 0.03593904
summary(lm obj)$coefficients["teique sf", c("t value", "Pr(>|t|)")]
## t value Pr(>|t|)
## 2.11099352 0.03593904
```

The hand-caculated t matches the one from lm().

```
Y <- ultra clean$pb100k dec
Yp <- lm obj\sfitted.values
SSR \leftarrow sum((Yp - mean(Y))^2)
SSE \leftarrow sum((Y - Yp)^2)
MSR <- SSR / 1
MSE \leq -SSE/(n-2)
Fval <- MSR / MSE
p F \leftarrow pf(Fval, df1 = 1, df2 = n - 2, lower.tail = FALSE)
c(SSR = SSR, SSE = SSE, MSR = MSR, MSE = MSE, F = Fval, p value = p F)
        SSR
                  SSE
                                       MSE
                            MSR
                                                       p value
## 5.160908e+01 2.466788e+03 5.160908e+01 1.158117e+01 4.456294e+00 3.593904e-02
anova(lm obj)
## Analysis of Variance Table
## Response: pb100k dec
         Df Sum Sq Mean Sq F value Pr(>F)
## teique sf 1 51.61 51.609 4.4563 0.03594 *
```

```
## Residuals 213 2466.79 11.581
## ---
## Signif. codes: 0 '***' 0.001 '**' 0.05 '.' 0.1 ' ' 1
```

The hand-caculated F matches the one from anova().

5-3

```
tval <- summary(lm_obj)$coefficients["teique_sf", "t value"]

Fval <- anova(lm_obj)[1, "F value"]

c(t_value = tval,
    t_squared = tval^2,
    F_value = Fval)

## t_value t_squared F_value
## 2.110994 4.456294 4.456294
```

The F-statistic is the square of t-statistic.

5-4

At $\alpha = 0.05$, there is statistically significant evidence (p = 0.036) that emotional intelligence affects ultramarathon times. The estimated slope (0.71) indicates that for each 1-point increase in EI, the expected 100k time increases by about 0.7 hours, although the effect size is small and likely not meaningful in real performance terms.

5-5

Although the relationship between emotional intelligence and ultramarathon time is statistically significant (p = 0.036), the magnitude of the effect is very small. The estimated slope ($\hat{\beta}_1 = 0.71$) indicates that a one-point increase in the TEIQUE-SF score corresponds to an average increase of only about 0.7 hours (\approx 43 minutes) in the 100k finishing time. Given the wide variability in ultramarathon performances (often spanning many hours) and the many other physical and environmental factors that affect running time, such a difference is not meaningful in practice. Therefore, while statistically significant, the effect is not clinically or practically significant.