

BIOST 546 HW 1

My-Anh Doan

2023-01-09

Q1. In this problem, we will make use of the data set `Medical_Cost_2.RData`.

- (a) Load the data set with the command `load` and check if there are missing data.
- (b) If any, remove the missing data using the command `na.omit`.

```
# load data to environment
load("./dataset/Medical_Cost_2.RData")
str(df)
## 'data.frame': 1338 obs. of 7 variables:
## $ age : int 19 18 28 33 32 31 46 37 37 60 ...
## $ sex : Factor w/ 2 levels "female","male": 1 2 2 2 2 1 1 1 2 1 ...
## $ bmi : num 27.9 33.8 33 22.7 28.9 ...
## $ children: int 0 1 3 0 0 0 1 3 2 0 ...
## $ smoker : Factor w/ 2 levels "no","yes": 2 1 1 1 1 1 1 1 1 1 ...
## $ region : Factor w/ 4 levels "northeast","northwest",...: 4 3 3 2 2 3 3 2 1 2 ...
## $ charges : num 16885 1726 4449 21984 3867 ...

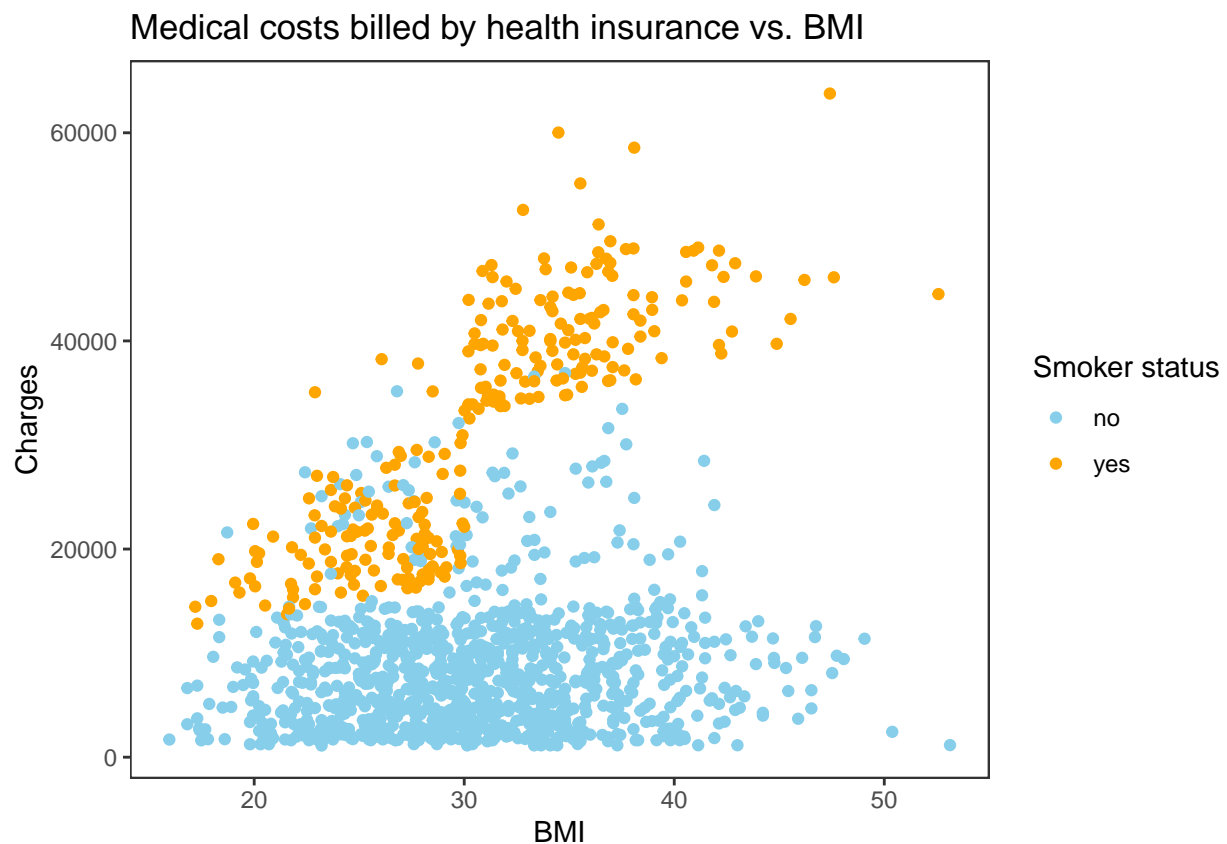
# check for missing data in `df`
# if returned value is non-zero, there are missing data
sum(is.na(df) == TRUE)
## [1] 60

# remove the missing data from `df`
df <- na.omit(df)
str(df)
## 'data.frame': 1278 obs. of 7 variables:
## $ age : int 19 18 28 33 32 31 46 37 37 60 ...
## $ sex : Factor w/ 2 levels "female","male": 1 2 2 2 2 1 1 1 2 1 ...
## $ bmi : num 27.9 33.8 33 22.7 28.9 ...
## $ children: int 0 1 3 0 0 0 1 3 2 0 ...
## $ smoker : Factor w/ 2 levels "no","yes": 2 1 1 1 1 1 1 1 1 1 ...
## $ region : Factor w/ 4 levels "northeast","northwest",...: 4 3 3 2 2 3 3 2 1 2 ...
## $ charges : num 16885 1726 4449 21984 3867 ...
## - attr(*, "na.action")= 'omit' Named int [1:60] 22 37 39 40 84 111 129 193 248 252 ...
## ..- attr(*, "names")= chr [1:60] "22" "37" "39" "40" ...
```

- (c) We decide to focus on the outcome variable `charges` (individual medical costs billed by health insurance) and the predictors `bmi` (body mass index) and `smoker` (whether the subjects is a smoker or not). Make a scatter plot with `bmi` on the x-axis, `charges` on the y-axis, and with the color of each dot representing whether the subject is a smoker or not.

```
library(ggplot2)

q1c_plot <- ggplot(df, aes(x = bmi, y = charges, color = smoker)) +
  geom_point() +
  scale_color_manual(values = c("no" = "sky blue",
                                "yes" = "orange")) +
  labs(title = "Medical costs billed by health insurance vs. BMI",
        x = "BMI",
        y = "Charges",
        color = "Smoker status") +
  theme_bw() +
  theme(panel.grid.minor = element_blank(),
        panel.grid.major = element_blank())
q1c_plot
```



(d) Fit a least-squares linear model, with intercept, in order to predict:

- `charges` using `bmi` as the only predictor;
- `charges` using `bmi` and `smoker` as predictors;
- `charges` using `bmi` and `smoker` as in the previous model; but allowing for an interaction term between the variables `bmi` and `smoker`

For each of the three models:

- Present your results in the form of a table where you report the estimated regression coefficients and their interpretations (be careful with the dummy variables);
- Report the 95% confidence interval for the coefficient of the variable `bmi`, and provide a sentence explaining the meaning of this confidence interval;
- Draw the regression line(s) of the model on the scatter plot produced in point (c);
- Report the (training set) mean squared error of the model;
- Predict the medical costs billed by the health insurance company to a smoker with a `bmi` that is 29 and 31.5;
- Compute the predicted difference in charges between a smoker with `bmi` 31.5 and one with `bmi` 29. Do the same for non-smokers. Comment on the results

```
# functions for answering modeling and prediction questions

# returns a table containing model coefficient estimates, standard errors, t-values,
# p-values, and upper and lower limits of 95% CIs
results_table <- function(linear_model) {
  coeffs <- cbind(as.data.frame(summary(linear_model)$coefficients),
                  confint(linear_model))
  colnames(coeffs) <- c("Estimate", "Std. Error", "t-value", "p-value", "95% CI LL", "95% CI UL")
  coeffs <- coeffs %>%
    mutate(Estimate = round(Estimate, 2),
           `Std. Error` = round(`Std. Error`, 2),
           `t-value` = round(`t-value`, 2),
           `p-value` = format(`p-value`, scientific = TRUE, digits = 3),
           `95% CI LL` = round(`95% CI LL`, 2),
           `95% CI UL` = round(`95% CI UL`, 2))
  coeffs
}

# calculates the mean squared error of the training model
calc_mse <- function(linear_model) {
  x <- mean(linear_model$residuals^2)
  print(paste("The training set mean squared error of the model is:", round(x, 0)))
}

# calculates the predicted charges for smokers and non-smokers at different BMIs
pred_costs <- function(linear_model, smoker, bmi) {
  round(predict(linear_model, data.frame(smoker, bmi)), 2)
}

# calculates the difference in charges at different BMIs
diff_costs <- function(linear_model, smoker, bmi) {
  round(max(pred_costs(linear_model, smoker, bmi)) - min(pred_costs(linear_model, smoker, bmi)), 2)
}
```

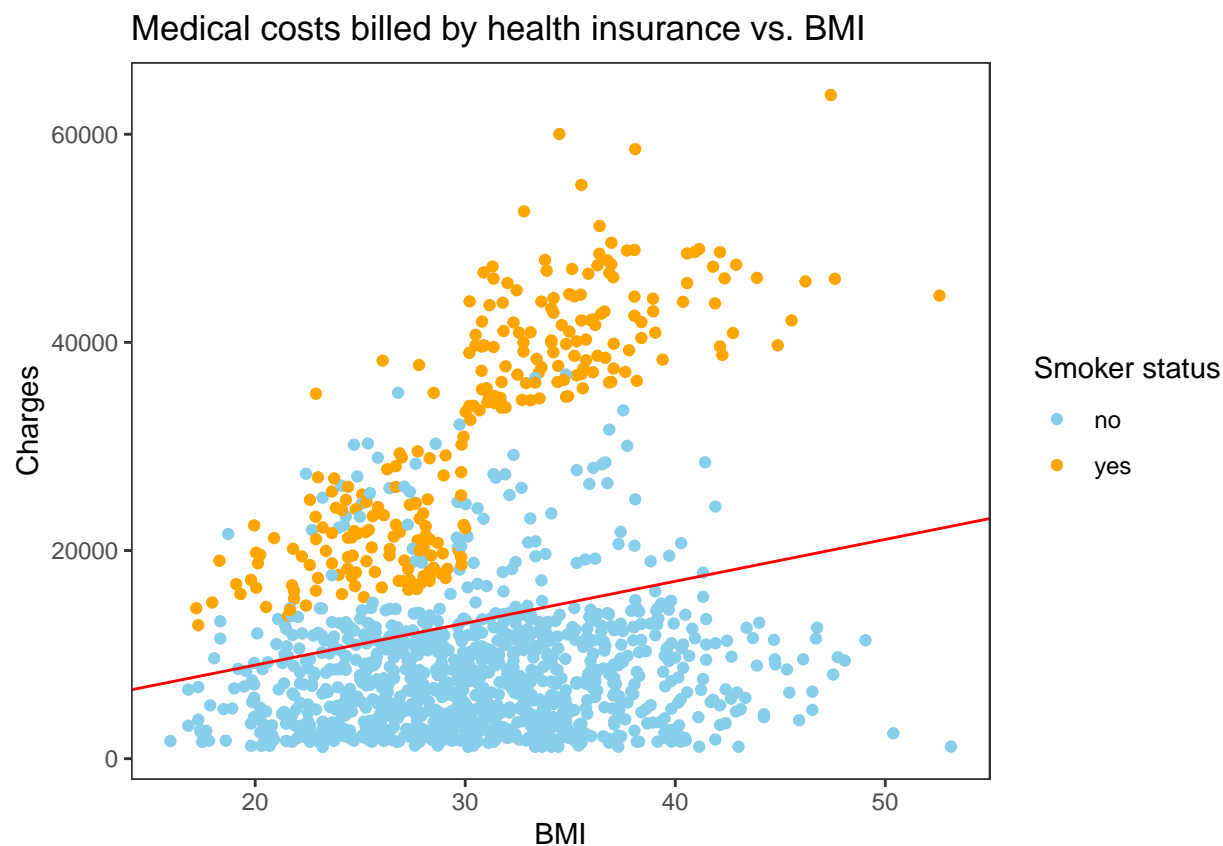
General linear regression models can be defined as:

$$y = \beta_0 + \beta_1 \times x_1 + \dots + \beta_p \times x_p$$

Model 1: charges ~ bmi

```
# linear model using `bmi` as the only predictor
charges_lm1 <- lm(charges ~ bmi, data = df)
kable(results_table(charges_lm1))

q1c_plot +
  geom_abline(slope = coef(charges_lm1)[2],
             intercept = coef(charges_lm1)[1],
             col = "red")
```



	Estimate	Std. Error	t-value	p-value	95% CI LL	95% CI UL
(Intercept)	938.44	1682.00	0.56	5.77e-01	-2361.35	4238.23
bmi	402.65	53.85	7.48	1.40e-13	297.01	508.28

The general model for linear regression models with only one predictor is defined as:

$$y = \beta_0 + \beta_1 \times x_1$$

Thus, the model above is given as:

$$charges = \beta_0 + \beta_1 \times bmi,$$

where β_0 (the intercept) is equal to 938.44 and β_1 (the slope) is equal to 402.65. The coefficient of variable `bmi` was estimated to be 402.65 with a 95% confidence interval of [297.01, 508.28]. In other words, there

is a 95% chance that the interval (297.01, 508.28) contains the true value of the coefficient of variable `bmi`. Additionally, for every unit increase in `bmi`, there will be an average increase in medical `charges` of between \$297.01 and \$508.28.

```
calc_mse(charges_lm1)
## [1] "The training set mean squared error of the model is: 138358366"

# predict the medical costs billed by the health insurance company to a smoker
# with a `bmi` that is 29 and 31.5 and compute the difference in charges.
pred_costs(charges_lm1, smoker = "yes", bmi = c(29, 31.5))
##           1           2
## 12615.22 13621.84
diff_costs(charges_lm1, smoker = "yes", bmi = c(29, 31.5))
## [1] 1006.62

# Do the same for non-smokers. Comment on the results
pred_costs(charges_lm1, smoker = "no", bmi = c(29, 31.5))
##           1           2
## 12615.22 13621.84
diff_costs(charges_lm1, smoker = "no", bmi = c(29, 31.5))
## [1] 1006.62
```

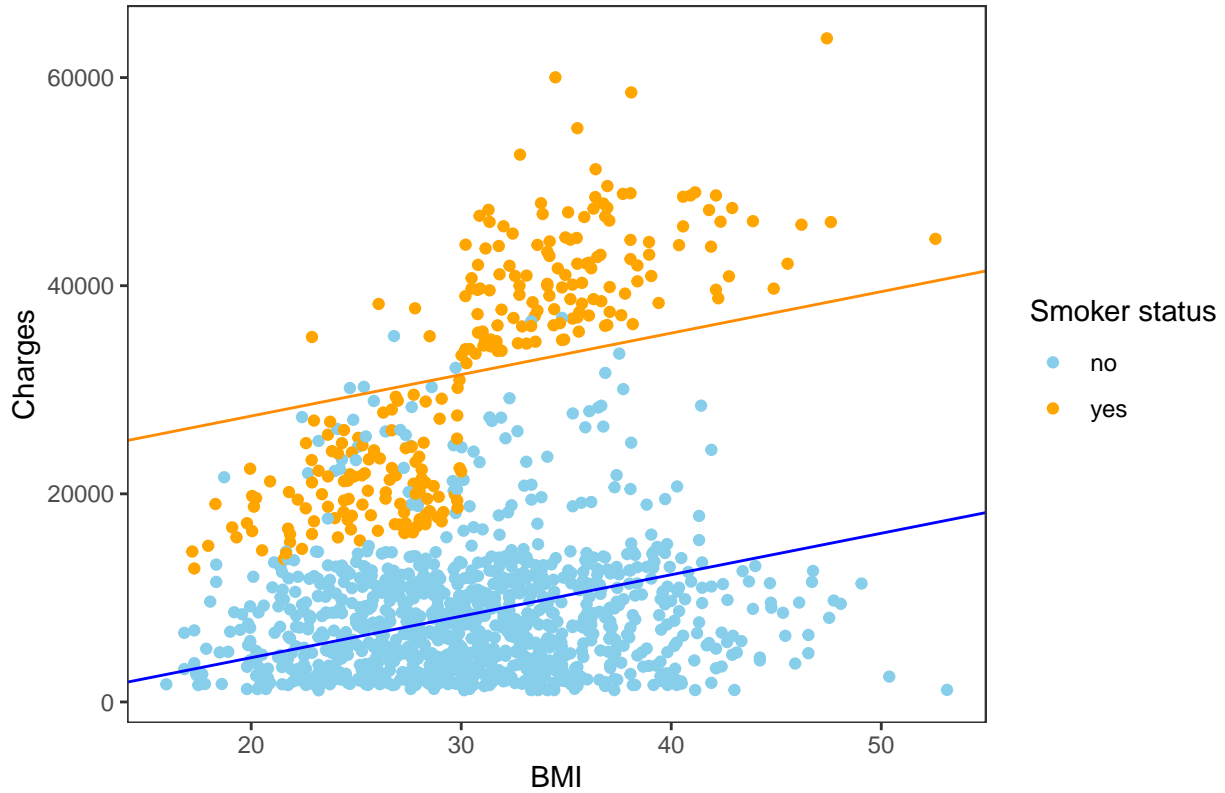
The predicted difference in charges between individuals with BMI 31.5 (\$13,621.84 for both smokers and non-smokers) and BMI 29 (\$12615.22 for both smokers and non-smokers) is the same between smokers and non-smokers, which is a difference of \$1006.62. This makes sense because the model is not dependent on smoking status and only uses BMI as the predictor (i.e. the model treats smokers and non-smokers as the same and only looks at BMI).

Model 2: Using `bmi` and `smoker` as predictors

```
# linear model using `bmi` and `smoker` as predictors
charges_lm2 <- lm(charges ~ bmi + smoker, data = df)
kable(results_table(charges_lm2))

q1c_plot +
  # smoker model
  geom_abline(slope = coef(charges_lm2)[2],
             intercept = coef(charges_lm2)[1] + coef(charges_lm2)[3],
             col = "dark orange") +
  # non-smoker model
  geom_abline(slope = coef(charges_lm2)[2],
             intercept = coef(charges_lm2)[1],
             col = "blue")
```

Medical costs billed by health insurance vs. BMI



	Estimate	Std. Error	t-value	p-value	95% CI LL	95% CI UL
(Intercept)	-3711.68	1018.78	-3.64	2.80e-04	-5710.35	-1713.02
bmi	398.47	32.46	12.27	8.12e-33	334.78	462.15
smokeryes	23218.77	491.04	47.28	1.02e-282	22255.43	24182.11

The general model for linear regression models with two predictors is defined as:

$$y = \beta_0 + \beta_1 \times x_1 + \beta_2 \times x_2$$

or in this case:

$$charges = \beta_0 + \beta_1 \times bmi + \beta_2 \times smoker$$

Thus the models are as follows:

- If the individual is a smoker: $charges = \beta_0 + \beta_2 + \beta_1 \times bmi$
- If the individual is a non-smoker: $charges = \beta_0 + \beta_1 \times bmi$

The smoker model is has a slope of β_1 and intercept of $\beta_0 + \beta_2$, whereas the non-smoker model has a slope of β_1 and intercept of β_0 , where $\beta_0 = -3711.68$, $\beta_1 = 398.47$, and $\beta_2 = 2.32 \times 10^4$. The coefficient of variable **bmi** was estimated to be 398.47 with a 95% confidence interval of (334.78, 462.15). For every unit increase in **bmi**, there will be an average increase in medical **charges** of between \$334.78 and \$462.15.

```

calc_mse(charges_lm2)
## [1] "The training set mean squared error of the model is: 50246296"

# predict the medical costs billed by the health insurance company to a smoker
# with a `bmi` that is 29 and 31.5 and compute the difference in charges.
pred_costs(charges_lm2, smoker = "yes", bmi = c(29, 31.5))
##          1          2
## 31062.62 32058.78
diff_costs(charges_lm2, smoker = "yes", bmi = c(29, 31.5))
## [1] 996.16

# Do the same for non-smokers. Comment on the results
pred_costs(charges_lm2, smoker = "no", bmi = c(29, 31.5))
##          1          2
## 7843.84 8840.01
diff_costs(charges_lm2, smoker = "no", bmi = c(29, 31.5))
## [1] 996.17

```

The predicted difference in charges between individuals with BMI 31.5 and BMI 29 is the same between smokers and non-smokers, which is \$996.17. This makes sense because the smoker and non-smoker models have the same slope of $\beta_1 = 398.47$, so the change in charges will be the same with each unit increase in BMI. However, the individual costs between smokers and non-smokers at the different BMIs are different. Smokers are charged \$31,062.62 and \$32,058.78 at BMI 29 and 31.5, respectively, while non-smokers will be billed \$7,843.84 and \$8,840.01 at BMI 29 and 31.5, respectively. This is a difference of about \$23,000, the approximate value of β_2 (which is the difference between the two models).

Model 3: Using bmi and smoker as predictors with an interaction term between bmi and smoker

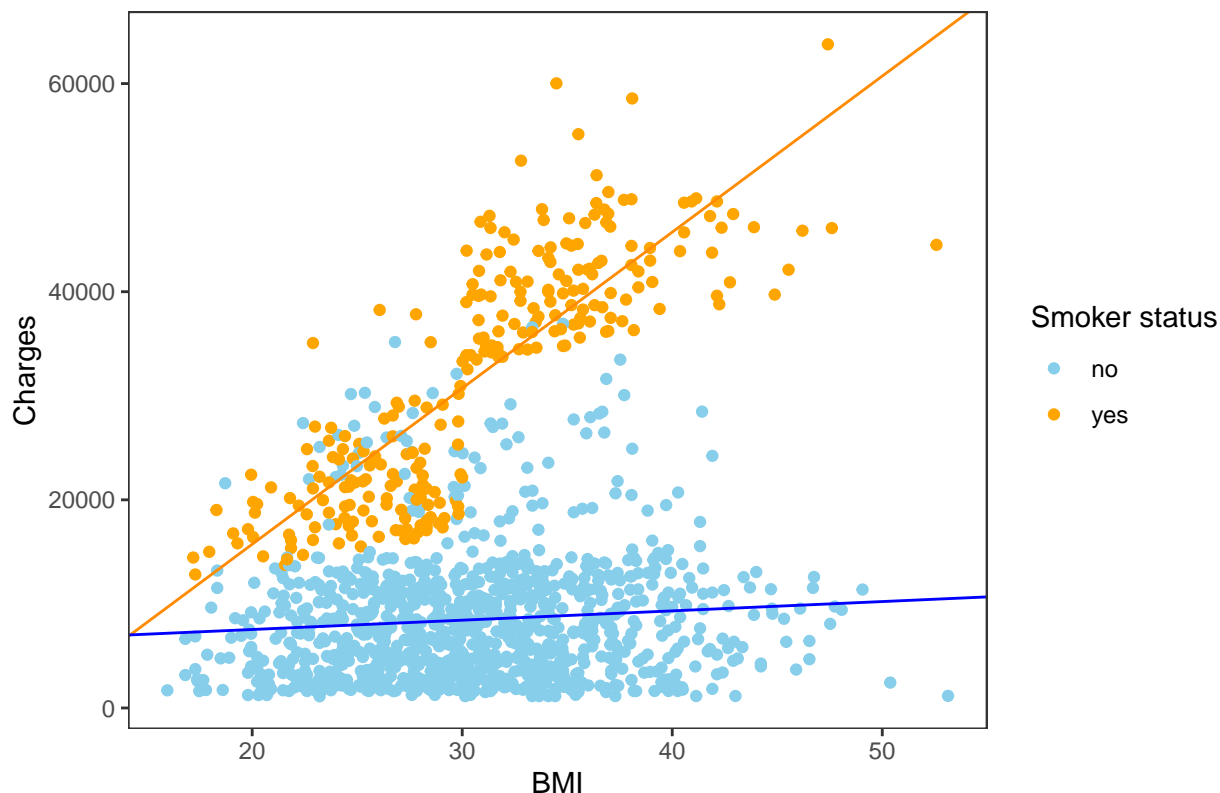
```

# linear model using `bmi` and `smoker` as predictors but also allowing for an
# interaction term between predictors `bmi` and `smoker`
charges_lm3 <- lm(charges ~ bmi + smoker + bmi*smoker , data = df)
kable(results_table(charges_lm3))

q1c_plot +
  # smoker model
  geom_abline(slope = coef(charges_lm3)[2] + coef(charges_lm3)[4],
             intercept = coef(charges_lm3)[1] + coef(charges_lm3)[3],
             col = "dark orange") +
  # non-smoker model
  geom_abline(slope = coef(charges_lm3)[2],
             intercept = coef(charges_lm3)[1],
             col = "blue")

```

Medical costs billed by health insurance vs. BMI



	Estimate	Std. Error	t-value	p-value	95% CI LL	95% CI UL
(Intercept)	5750.97	991.45	5.80	8.33e-09	3805.92	7696.03
bmi	89.47	31.76	2.82	4.92e-03	27.17	151.78
smokeryes	-20008.39	2122.67	-9.43	1.94e-20	-24172.70	-15844.07
bmi:smokeryes	1410.05	67.84	20.78	7.42e-83	1276.96	1543.15

The general model with two predictors and an interaction term between the two predictors is defined as:

$$y = \beta_0 + \beta_1 \times x_1 + \beta_2 \times x_2 + \beta_3 \times x_1 x_2$$

or in this case:

$$charges = \beta_0 + \beta_1 \times bmi + \beta_2 \times smoker + \beta_3 \times bmi \times smoker$$

Thus the models are as follows:

- If the individual is a smoker: $charges = \beta_0 + \beta_2 + bmi \times (\beta_1 + \beta_3)$
- If the individual is a non-smoker: $charges = \beta_0 + \beta_1 \times bmi$

The smoker model has a slope of $\beta_1 + \beta_3$ and intercept of $\beta_0 + \beta_2$ while the non-smoker model has a slope of β_1 and intercept of β_0 , where $\beta_0 = 5750.97$, $\beta_1 = 89.47$, $\beta_2 = -2 \times 10^4$, and $\beta_3 = 1410.05$. The coefficient of variable **bmi** was estimated to be 89.47 with a 95% confidence interval of (27.17, 151.78). For every unit increase in **bmi**, there will be an average increase in medical **charges** of between \$27.17 and \$151.78.


```

calc_mse(charges_lm3)
## [1] "The training set mean squared error of the model is: 37522941"

# predict the medical costs billed by the health insurance company to a smoker
# with a `bmi` that is 29 and 31.5 and compute the difference in charges.
pred_costs(charges_lm3, smoker = "yes", bmi = c(29, 31.5))
##          1          2
## 29228.89 32977.71
diff_costs(charges_lm3, smoker = "yes", bmi = c(29, 31.5))
## [1] 3748.82

# Do the same for non-smokers. Comment on the results
pred_costs(charges_lm3, smoker = "no", bmi = c(29, 31.5))
##          1          2
## 8345.72 8569.40
diff_costs(charges_lm3, smoker = "no", bmi = c(29, 31.5))
## [1] 223.68

```

The predicted difference in charges between individuals with BMI 31.5 and BMI 29 is different between smokers and non-smokers. This makes sense because the models do not have the same slope so charges will not increase at the same rate with each unit increase in BMI. For smokers, the difference in charges at BMI 31.5 (\$32,977.71) and BMI 29 (\$29,228.89) is \$3,748.82. For non-smokers, the difference in charges at BMI 31.5 (\$8,569.40) and BMI 29 (\$8,345.72) is \$223.68. The change in charges has a smaller increase for every unit increase in BMI for non-smokers compared to smokers.

(e) Now define and add to the data set a new Boolean variable `smoker_bmi30p` that is `True` only if the subject is a smoker **and** has a `bmi` greater than 30. Use this newly defined variable, together with `bmi` and `smoker`, to fit the linear model represented in Figure 1 by carefully defining the interaction terms (allow each of the three straight lines to have their own intercept and slope, but use the command `lm` only once).

- Present your results in the form of one table where you report the estimated coefficients of the model.
- For each predictor, comment on whether you can reject the null hypothesis that there is no (linear) association between that predictor and `charges`, conditional on the other predictors in the model.
- Explain the interpretation of the non-significant variables in the model ($p > 0.05$) and explain how Figure 1 would change if we were to discard those variables, i.e., perform variable selection.
- According to this newly defined model, compute the predicted difference in `charges` between a smoker with `bmi` 31.5 and one with `bmi` 29. Do the same for non-smokers. Compare the analogous results in point (d) and comment on the results.

```

df <- df %>%
  mutate(smoker_bmi30p = ifelse(smoker == "yes" & bmi > 30, TRUE, FALSE))

charges_lm4 <- lm(charges ~ bmi + smoker + smoker_bmi30p + bmi:smoker + bmi:smoker_bmi30p, data = df)

#q1c_plot +
# smoker model
# geom_abline(slope = ,
#             intercept = ,
#             col = "dark orange") +
# non-smoker model

```

```
# geom_abline(slope = ,  
#             intercept = ,  
#             col = "blue")
```

Q2. This problem has to do with the notation of bias-variance trade-off. For (a) and (b), it's okay to submit hand-sketches: this is a conceptual exercise.

- (a) Make a plot, like the one we saw in class, with “flexibility” on the x-axis. Sketch the following curves: squared bias, variance, irreducible error, expected prediction error. Be sure to label each curve. Indicate which level of flexibility is *best*.
- (b) Make a plot with “flexibility” on the x-axis. Sketch curves corresponding to the training error and the test error. Be sure to label each curve. Indicate which level of flexibility is “best”.

Q3. This problem has to do with numerical explorations of the bias-variance trade-off phenomenon. You will generate simulated data, and will use these data to perform **linear regression**. Set the seed with `set.seed(0)` before you begin.

- (a) Use the `rnorm()` function to generate a predictor vector \mathbf{X} of length $n = 30$, and use `runif()` to generate a noise vector ϵ of length $n = 30$.
- (b) Generate a response vector Y of length $n = 30$ according to the model