

STA130 W11

Monday, November 25, 2019 2:12 PM

- Correlation (r) to quantify linear association between two numerical variables
- Simple linear regression model Strength and direction

Response Intercept, slope, random error term

$$y_i = \beta_0 + \beta_1 x_i + \epsilon_i$$

numerical

Numerical predictor

- Fitted regression line (estimated regression coefficients from $\text{lm}()$ or closed form expressions)

Estimated regression coefficients from sample data

$$\hat{y}_i = \hat{\beta}_0 + \hat{\beta}_1 x_i \quad (x_1, y_1), \dots, (x_n, y_n)$$

- Coefficient of determination R^2

$$R^2 = \frac{\text{"Explained" Variation}}{\text{Total Variation}} = \frac{\sum_{i=1}^n (\hat{y}_i - \bar{y})^2}{\sum_{i=1}^n (y_i - \bar{y})^2}$$

- RMSE to assess predictive performance of a linear regression model for prediction

$$RMSE = \sqrt{\frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2} = \sqrt{\frac{1}{n} \sum_{i=1}^n e_i^2} \quad 2 / 40$$

R^2	$RMSE$
$1 - \frac{\sum_{i=1}^n (y_i - \hat{y}_i)^2}{\sum_{i=1}^n (y_i - \bar{y})^2}$	$\sqrt{\frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2}$
Range: 0 - 1	Range: 0 to ∞
Large value \Rightarrow good fit	Small value \Rightarrow good fit
No units	Same units as y
Tells us proportion of variation in y which is explained by the model	Used to assess prediction error (with training or test data)

* R^2 only makes sense to calculate for the data used to build the model

*RMSE: compare RMSE for training and testing data

Hypothesis tests to determine if there is a **real** linear association between two variables

Linear regression with a categorical predictor

Multiple linear regression models (i.e. regression with more than one predictor)

- Model fitting Comparing predictive models

Predicting Body Mass Index (BMI) from waist circumference

Body measurements for 507 physically active adults

```
body <- read.csv("bodydat.csv")
body %>% select(age, weight, height, gender, waist) %>% glimpse()
```

```
## Observations: 507
## Variables: 5
## $ age      <int> 21, 23, 28, 23, 22, 21, 26, 27, 23, 21, 23, 22, 20, 26,...
## $ weight   <dbl> 65.6, 71.8, 80.7, 72.6, 78.8, 74.8, 86.4, 78.4, 62.0, 8...
## $ height   <dbl> 174.0, 175.3, 193.5, 186.5, 187.2, 181.5, 184.0, 184.5,...
## $ gender   <fct> Male, Male, Male, Male, Male, Male, Male, Male, M...
## $ waist    <dbl> 71.5, 79.0, 83.2, 77.8, 80.0, 82.5, 82.0, 76.8, 68.5, 7...
```

- lengths in cm
- weight in kg

$$BMI = \frac{\text{weight (in kg)}}{[\text{height (in meters)}]^2}$$

```
body <- body %>% mutate(BMI = weight / (height/100)^2)
```

```
summary(lm(BMI ~ waist, data=body))$coefficients
##           Estimate Std. Error t value Pr(>|t|)
## (Intercept) 4.6566078 0.585776319  7.949464 1.233300e-14
## waist       0.2443089 0.007532967 32.431966 1.495616e-125
summary(lm(BMI ~ waist, data=body))$rsquared
## [1] 0.6756234
```

Approximately 68% of the values of BMI can be predicted from waist circumference

Wrong - 68% of the variable in BMI is explained by **own** model

Q1 - What would the regression line look like if there was no association between waist circumference and BMI?
Horizontal line with slope of 0 (i.e. $\beta_1\text{-hat} = 0$)

Q2 - How can we determine whether the association we observed is "real" or just due to sampling variability?
Significant test(hypothesis test!)

Inference for simple linear regression

Model:

$$y = \beta_0 + \beta_1 x_i + \epsilon_i$$

Population level parameter

where y_i and x_i are the BMI and waist circumference of individual i , respectively, and ϵ_i is the difference between individual i 's response and the mean BMI for people with the same waist circumference as them.

Based on this model, how can we write hypotheses to test if there is an association between waist circumference and BMI in physically-active adults?

$H_0 = \beta_1 = 0$ (there is no linear regression between x and y)

$H_a = \beta_1 \neq 0$

The inference results given by `lm()` are based on the **Students' t distribution** (a continuous probability distribution) so rely on some assumptions:

A1: Linear pattern between x and y

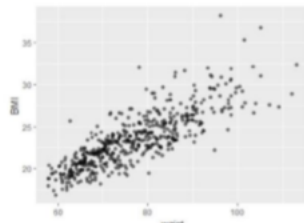
A2: Constant variance in y for all values of x

A3: Independent observations

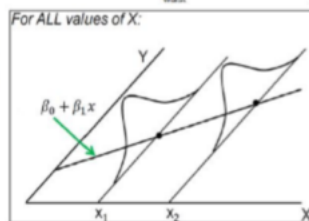
A4: Residuals follow a normal (a bell-shaped continuous probability distribution)

We checked that these conditions are reasonable for these data so it is valid to make inferences based on the p-values generated by R for this regression.

You'll learn more about strategies to check these assumptions and to deal with violations in future statistics courses



Unequal variance might look like



If one or more of the assumptions above are not reasonable, then the inference results given by `lm()` may not be valid, although the model may still be used for prediction.

```
summary(lm(BMI ~ waist, data=body))$coefficients
##           Estimate Std. Error t value Pr(>|t|)
## (Intercept) 4.6566078 0.585776319  7.949464 1.233300e-14 (beta0)
## waist       0.2443089 0.007532967 32.431966 1.495616e-125 (beta1)
```

R gives us p-values for hypothesis tests of the form

$H_0: \beta = 0$ $H_a: \beta \neq 0$

Do we have evidence against the null hypothesis that the slope is different from 0 (i.e. the null hypothesis that there is no association between waist circumference and BMI)? If yes, how strong is this evidence?

P-value : $1.49 \times 10^{-125} \approx 0$

Very strong evidence against the null hypothesis of no linear association between waist and BMI

What other factors might also affect BMI?

There are many other possible predictors in our dataset

Let's see if the distribution of BMI varies for men and women?

Regression with a categorical predictor

$$\text{Model: } y_i = \beta_0 + \beta_1 x_i + \epsilon_i$$

where y_i is the BMI for individual i , ϵ_i is the random error term for individual i , and

$$x_i = I(\text{individual } i \text{ is male}) = \begin{cases} 1 & \text{if individual } i \text{ is male} \\ 0 & \text{if individual } i \text{ is female} \end{cases}$$

Gender is a categorical predictor with levels "male" and "female", but we need to convert these levels to numbers

- We encode categorical predictors as **indicator variables** (also called **dummy variables**)
- We need to pick a **baseline value** (i.e. the level corresponding to $x = 0$); here **female** is the baseline

In the R output, the *genderMale* indicates that "Male" is **not** the baseline level: the level which is missing from the output is the baseline!

```
summary(lm(BMI ~ gender, data=body))$coefficients
##      Estimate Std. Error t value Pr(>|t|)
## (Intercept) 22.27793  0.1886217 118.109055 0.000000e+00
## genderMale  2.43330  0.2702385  9.004269 4.435794e-18
```

Male is not the base line here

The fitted regression line is $\hat{y} = 22.3 + 2.4 x_i$, where \hat{y} is the BMI and $x_i = I(\text{individual } i \text{ is male})$.

** x is either 0 or 1, so this really only gives us 2 predictions, one for men and one for women

Men: $\hat{y} = \beta_0 + \beta_1(1) + \epsilon$

Women: $\hat{y} = \beta_0 + \beta_1(0) + \epsilon$

$H_0: \beta_1 = 0$

$H_a: \beta_1 \neq 0$

->very strong evidence against the H_0

What other method (which we covered in this course) could we also use to investigate whether there is a difference in mean BMI for men and women?

- Randomization test, by shuffling the group labels (M F) to sim. Values under H_0

Model 1:

$$\underset{\text{BMI}}{y_i} = \beta_0 + \underset{\text{waist}}{\beta_1 x_{1i}} + \underset{I(i \text{ is male})}{\beta_2 x_{2i}} + \epsilon_i$$

where y_i and x_{1i} are the BMI and waist circumference for individual i and

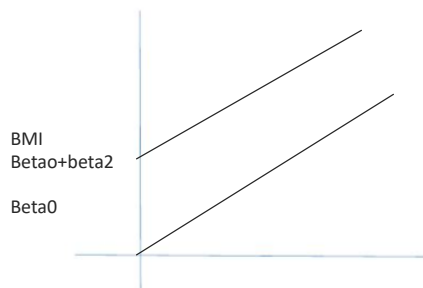
$x_{2i} = I(\text{individual } i \text{ is male})$.

Fitted line for women

$\hat{y} = \hat{\beta}_0 + \hat{\beta}_1 x_i + \hat{\beta}_2(0)$

Fitted line for men

$\hat{y} = \hat{\beta}_0 + \hat{\beta}_1 x_i + \hat{\beta}_2(1)$
 $= (\hat{\beta}_0 + \hat{\beta}_2) + \hat{\beta}_1 x_i$



Fitted model -> for multiple linear regression

$(y \sim x_1 + x_2)$

```
parallel_lines <- lm(BMI ~ waist + gender, data = body)
summary(parallel_lines)$coefficients
##      Estimate Std. Error t value Pr(>|t|)
## (Intercept) 0.7454759 0.651078332  1.144987 2.527582e-01
## waist      0.3084726 0.009205015  33.511364 2.358772e-130
## genderMale -2.1104205 0.202610484 -10.416146 3.845116e-23
```

Equation:

$\hat{y} = 0.745 + 0.308 \times (\text{waist circumference}) - 2.11 \times I(\text{individual is male})$

Plotting parallel lines

The *augment* function (in the library *broom*) creates a data frame with predicted values (.fitted), residuals, etc ...

```
library(broom) augment(parallel_lines)
## # A tibble: 507 x 10
##   BMI waist gender .fitted .se.fit .resid .hat .sigma .cooks
```

```
## * <dbl> <dbl> <dbl> <dbl> <dbl> <dbl> <dbl> <dbl>
## 1 21.7 71.5 Male 20.7 0.161 0.976 0.00906 1.70 1.02e-3
...
```

Join up the fitted values to plot the parallel lines model

```
body %>% ggplot(aes(x=waist, y=BMI, color=gender)) +
  geom_point(alpha=0.5) +
  geom_line(data=augment(parallel_lines),
    aes(y=fitted, colour=gender), lwd=1.5)
# Line width
```

Scatterplot

To add the fitted regression line

Allowing for non-parallel lines

In Model 1, we assumed that the association between waist and BMI was the same for men and women (i.e. same slope)

Let's look to see if gender **modifies the relationship** between waist and BMI - to do this, we add a new independent variable to the model which is the product of *waist* and *gender*: this is an **interaction term**

Model 2:

Interaction term

$$y_i = \beta_0 + \beta_1 x_{1i} + \beta_2 x_{2i} + \beta_3 x_{1i} x_{2i} + \epsilon_i$$

(Note: The original image has a handwritten note "Same as before" with an arrow pointing to the first three terms of the equation.)

where y_i and x_{1i} are the BMI and waist circumference for individual i and

$x_{2i} = I(\text{individual } i \text{ is male})$.

To fit a linear model with an interaction term in R, use * instead of + between the predictors

```
summary(lm(BMI ~ waist * gender, data = body))$coefficients
```

##	Estimate	Std. Error	t value	Pr(> t)	B-hat0
## (Intercept)	-3.4908232	0.94161575	-3.707269	2.327883e-04	1
## waist	0.3691615	0.01341083	27.527117	1.985004e-102	2
## genderMale	6.1578438	1.38114262	4.458514	1.018556e-05	3
## waist:genderMale	-0.1083858	0.01792202	-6.047631	2.869562e-09	

-> coefficient for the interaction between waist cir. And gender

**row names are important!

Fitted line for women($x_2=0$)

$\hat{Y} = b_0\text{-hat} + b_1\text{-hat}x_1 + b_2\text{-hat}(0) + b_3\text{-hat}x_1(0)$

Fitted line for men

$\hat{Y} = b_0\text{-hat} + b_1\text{-hat}x_1 + b_2\text{-hat}(1) + b_3\text{-hat}x_1(1)$

$= (b_0\text{-hat} + b_2\text{-hat}) + (b_1\text{-hat} + b_3\text{-hat})x_1$

We can then plug in values of b-hat from the table of the top to make a prediction for a new observations

Plot of non-parallel lines

```
body %>% ggplot(aes(x=waist, y=BMI, color=gender)) +
  geom_point(alpha=0.5) + geom_smooth(method="lm", se=FALSE)
```

Different intercepts and slopes for men and women

Consider our four regression models

Model A:

$$y_i = \beta_0 + \beta_1 \times \text{waist}_i + \epsilon_i$$

Model B:

$$y_i = \beta_0 + \beta_1 \times I(\text{individual } i \text{ is male}) + \epsilon_i$$

Model C:

$$y_i = \beta_0 + \beta_1 \times \text{waist}_i + \beta_2 \times I(\text{individual } i \text{ is male}) + \epsilon_i$$

Model D:

$$y_i = \beta_0 + \beta_1 \times \text{waist}_i + \beta_2 \times I(i \text{ is male}) + \beta_3 \times \text{waist}_i \cdot I(i \text{ is male}) + \epsilon_i$$

Comparing our models: Which of our models performs best to predict BMI?

Strategy

1. Divide sample into training dataset and testing dataset
2. Fit a prediction model (today: a linear regression model) using the training data
3. Make predictions for the testing data and compare the predictions to the true responses

```
set.seed(1210);
```

```
n <- nrow(body)
```

```

training_indices <- sample(1:n, size=round(0.8*n))
train <- body[training_indices,]
y_train <- train$BMI;

# Testing dataset includes all observations NOT in the training data
test <- body[-training_indices,]
y_test <- test$BMI;

# Fit model to training data
modA_train <- lm(BMI ~ waist, data=train)
modB_train <- lm(BMI ~ gender, data=train)
modC_train <- lm(BMI ~ waist + gender, data=train)
modD_train <- lm(BMI ~ waist * gender, data=train)

# Make predictions for testing data using training model
yhat_modA_test <- predict(modA_train, newdata = test)
yhat_modB_test <- predict(modB_train, newdata = test)
yhat_modC_test <- predict(modC_train, newdata = test)
yhat_modD_test <- predict(modD_train, newdata = test)

# Make predictions for training data using training model
yhat_modA_train <- predict(modA_train, newdata = train)
yhat_modB_train <- predict(modB_train, newdata = train)
yhat_modC_train <- predict(modC_train, newdata = train)
yhat_modD_train <- predict(modD_train, newdata = train)

# Calculate RMSE for testing data
modA_test_RMSE <- sqrt(sum((y_test - yhat_modA_test)^2) / nrow(test))
modB_test_RMSE <- sqrt(sum((y_test - yhat_modB_test)^2) / nrow(test))
modC_test_RMSE <- sqrt(sum((y_test - yhat_modC_test)^2) / nrow(test))
modD_test_RMSE <- sqrt(sum((y_test - yhat_modD_test)^2) / nrow(test))

modA_train_RMSE <- sqrt(sum((y_train - yhat_modA_train)^2) / nrow(train))
modB_train_RMSE <- sqrt(sum((y_train - yhat_modB_train)^2) / nrow(train))
modC_train_RMSE <- sqrt(sum((y_train - yhat_modC_train)^2) / nrow(train))
modD_train_RMSE <- sqrt(sum((y_train - yhat_modD_train)^2) / nrow(train))

data_frame(Model = c("A", "B", "C", "D"),
  RMSE_testdata = c(modA_test_RMSE, modB_test_RMSE,
    modC_test_RMSE, modD_test_RMSE),
  RMSE_traindata = c(modA_train_RMSE, modB_train_RMSE,
    modC_train_RMSE, modD_train_RMSE),
  ratio_of_RMSEs = RMSE_traindata / RMSE_testdata)

## # A tibble: 4 x 4
##   Model RMSE_testdata RMSE_traindata ratio_of_RMSEs
##   <chr>      <dbl>      <dbl>      <dbl>
## 1 A         1.89        1.86        0.981
## 2 B         3.25        2.98        0.916
## 3 C         1.75        1.68        0.958
## 4 D         1.79        1.60        0.894

```

-> most useful

Which model makes the most accurate predictions, on average?
 Look at RMSE_testdata -> model C (D is close though)
 Is there evidence of overfitting?
 Not too much, but maybe a bit in model D

