

PSTAT131HW02

Yifei Zhang

2022-04-10

Question 1

Your goal is to predict abalone age, which is calculated as the number of rings plus 1.5. Notice there currently is no `age` variable in the data set. Add `age` to the data set.

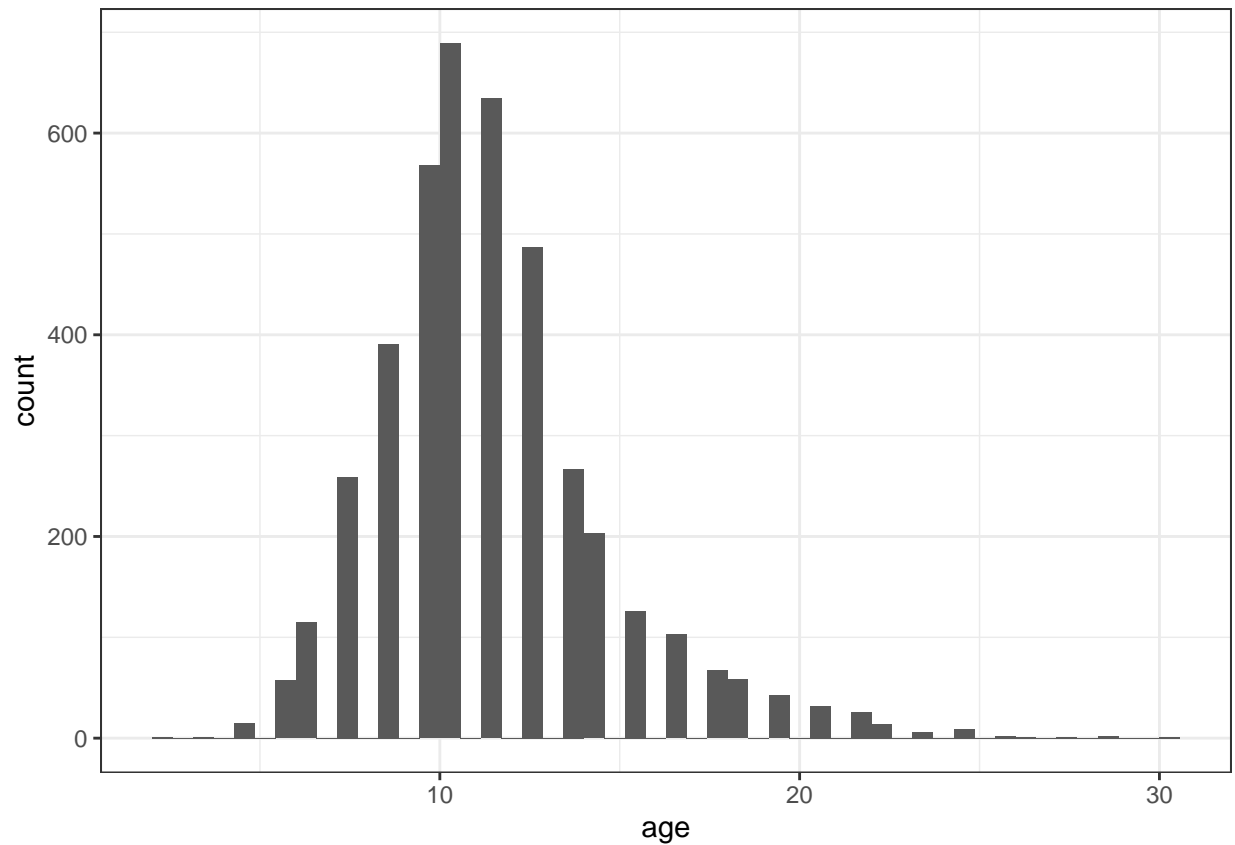
```
abalone <- read.csv("abalone.csv")

new_abalone <- abalone %>%
  mutate(age = rings + 1.5)
```

Assess and describe the distribution of `age`.

As we access the distribution of the new variable `age`, we can see it looks relatively normal, but is skewed a bit to the right, centering around 10.

```
new_abalone %>%
  ggplot(aes(x = age)) +
  geom_histogram(bins = 50) +
  theme_bw()
```



Question 2

Split the abalone data into a training set and a testing set. Use stratified sampling. You should decide on appropriate percentages for splitting the data.

Remember that you'll need to set a seed at the beginning of the document to reproduce your results.

```
set.seed(710)

abalone_split <- initial_split(new_abalone, prop = 0.80,
                               strata = age)

abalone_train <- training(abalone_split)

abalone_test <- testing(abalone_split)
abalone_train %>%
  head()
```

##	type	longest_shell	diameter	height	whole_weight	shucked_weight
## 5	I	0.330	0.255	0.080	0.2050	0.0895
## 6	I	0.425	0.300	0.095	0.3515	0.1410
## 17	I	0.355	0.280	0.085	0.2905	0.0950
## 19	M	0.365	0.295	0.080	0.2555	0.0970
## 36	M	0.465	0.355	0.105	0.4795	0.2270
## 43	I	0.240	0.175	0.045	0.0700	0.0315

```
##      viscera_weight shell_weight rings age
## 5          0.0395         0.055    7 8.5
## 6          0.0775         0.120    8 9.5
## 17         0.0395         0.115    7 8.5
## 19         0.0430         0.100    7 8.5
## 36         0.1240         0.125    8 9.5
## 43         0.0235         0.020    5 6.5
```

Question 3

Using the **training** data, create a recipe predicting the outcome variable, **age**, with all other predictor variables. Note that you should not include **rings** to predict **age**. Explain why you shouldn't use **rings** to predict **age**.

I should not use rings to predict age, because we have used rings to calculate age in step 1, and we won't be able to see how other predictors can predict by having rings as one of the predictors, since we already have $\text{age} = \text{rings} + 1.5$ function.

Steps for your recipe:

1. dummy code any categorical predictors

```
new_abalonegenderm <- ifelse(new_abalonetype == "M", 1, 0)
```

```
new_abalonegenderf <- ifelse(new_abalonetype == "F", 1, 0)
```

```
new_abalonegenderi <- ifelse(new_abalonetype == "I", 1, 0)
```

this don't seem to work with the latter steps

```
##      type longest_shell diameter height whole_weight
## "character" "numeric" "numeric" "numeric" "numeric"
## shucked_weight viscera_weight shell_weight rings age
## "numeric" "numeric" "numeric" "integer" "numeric"
```

2. create interactions between

- type and shucked_weight,
- longest_shell and diameter,
- shucked_weight and shell_weight

3. center all predictors, and

4. scale all predictors.

You'll need to investigate the **tidymodels** documentation to find the appropriate step functions to use.

Creating recipe

```
new_abalone %>%
  head()
```

```
##      type longest_shell diameter height whole_weight shucked_weight viscera_weight
## 1      M          0.455    0.365  0.095         0.5140         0.2245         0.1010
## 2      M          0.350    0.265  0.090         0.2255         0.0995         0.0485
```

```
## 3      F      0.530    0.420  0.135      0.6770      0.2565      0.1415
## 4      M      0.440    0.365  0.125      0.5160      0.2155      0.1140
## 5      I      0.330    0.255  0.080      0.2050      0.0895      0.0395
## 6      I      0.425    0.300  0.095      0.3515      0.1410      0.0775
##   shell_weight rings  age
## 1      0.150     15 16.5
## 2      0.070      7  8.5
## 3      0.210      9 10.5
## 4      0.155     10 11.5
## 5      0.055      7  8.5
## 6      0.120      8  9.5
```

```
abalone_recipe <- recipe(age ~ longest_shell + diameter + height +
                          whole_weight + shucked_weight +
                          viscera_weight + shell_weight + type,
                          data = new_abalone) %>%
  step_dummy(all_nominal_predictors()) %>%
  step_center(longest_shell, diameter, height, whole_weight,
              shucked_weight, viscera_weight, shell_weight
              ) %>%
  step_scale(longest_shell, diameter, height, whole_weight,
             shucked_weight, viscera_weight, shell_weight
             ) %>%
  step_interact(terms = ~ type : shucked_weight) %>%
  step_interact(terms = ~ longest_shell : diameter) %>%
  step_interact(terms = ~ shell_weight : shucked_weight)

abalone_recipe
```

```
## Data Recipe
##
## Inputs:
##
##      role #variables
##      outcome      1
##      predictor      8
##
## Operations:
##
## Dummy variables from all_nominal_predictors()
## Centering for longest_shell, diameter, height, ...
## Scaling for longest_shell, diameter, height, ...
## Interactions with type:shucked_weight
## Interactions with longest_shell:diameter
## Interactions with shell_weight:shucked_weight

%>%
prep(abalone_train)
```

Question 4

Create and store a linear regression object using the "lm" engine.

```
lm_model <- linear_reg() %>%
  set_engine("lm")

# I ended up using what we used in the lab, but what is the difference?

fit <- lm(age ~ longest_shell + diameter + height +
          whole_weight + shucked_weight +
          viscera_weight + shell_weight + type,
          data = new_abalone)

summary(fit)

##
## Call:
## lm(formula = age ~ longest_shell + diameter + height + whole_weight +
##     shucked_weight + viscera_weight + shell_weight + type, data = new_abalone)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -10.4800  -1.3053  -0.3428   0.8600  13.9426
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)    5.39464    0.29157  18.502 < 2e-16 ***
## longest_shell  -0.45834    1.80912  -0.253   0.800
## diameter       11.07510    2.22728   4.972 6.88e-07 ***
## height         10.76154    1.53620   7.005 2.86e-12 ***
## whole_weight    8.97544    0.72540  12.373 < 2e-16 ***
## shucked_weight -19.78687    0.81735 -24.209 < 2e-16 ***
## viscera_weight -10.58183    1.29375  -8.179 3.76e-16 ***
## shell_weight    8.74181    1.12473   7.772 9.64e-15 ***
## typeI          -0.82488    0.10240  -8.056 1.02e-15 ***
## typeM           0.05772    0.08335   0.692   0.489
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 2.194 on 4167 degrees of freedom
## Multiple R-squared:  0.5379, Adjusted R-squared:  0.5369
## F-statistic: 538.9 on 9 and 4167 DF,  p-value: < 2.2e-16
```

```
summary(lm_model)
```

```
##           Length Class      Mode
## args      2      -none-    list
## eng_args  0      quosures list
## mode      1      -none-    character
## method    0      -none-    NULL
## engine     1      -none-    character
```

Question 5

Now:

1. set up an empty workflow,
2. add the model you created in Question 4, and
3. add the recipe that you created in Question 3.

```
lm_wflow <- workflow() %>%
  add_model(lm_model) %>%
  add_recipe(abalone_recipe)

lm_wflow
```

```
## == Workflow =====
## Preprocessor: Recipe
## Model: linear_reg()
##
## -- Preprocessor -----
## 6 Recipe Steps
##
## * step_dummy()
## * step_center()
## * step_scale()
## * step_interact()
## * step_interact()
## * step_interact()
##
## -- Model -----
## Linear Regression Model Specification (regression)
##
## Computational engine: lm
```

Question 6

Use your `fit()` object to predict the age of a hypothetical female abalone with `longest_shell = 0.50`, `diameter = 0.10`, `height = 0.30`, `whole_weight = 4`, `shucked_weight = 1`, `viscera_weight = 2`, `shell_weight = 1`.

```
lm_fit <- fit(lm_wflow, abalone_train)
```

```
## Warning: Interaction specification failed for: ~type:shucked_weight. No
## interactions will be created.
```

```
test <- data.frame(longest_shell = 0.50,
  diameter = 0.10,
  height = 0.30,
  whole_weight = 4,
  shucked_weight = 1,
  viscera_weight = 2,
  shell_weight = 1,
  type = "F")

lm_fit %>%
  extract_fit_parsnip() %>%
  tidy()
```

```
## # A tibble: 12 x 5
##   term                estimate std.error statistic  p.value
##   <chr>              <dbl>      <dbl>      <dbl>    <dbl>
## 1 (Intercept)        12.2        0.0875    139.      0
## 2 longest_shell     -0.838        0.260    -3.22  1.29e- 3
## 3 diameter           0.455        0.254     1.79  7.33e- 2
## 4 height             0.253        0.0694     3.65  2.69e- 4
## 5 whole_weight       4.72         0.396    11.9  4.51e-32
## 6 shucked_weight    -4.00         0.209   -19.2  1.23e-77
## 7 viscera_weight    -0.888        0.158    -5.61  2.23e- 8
## 8 shell_weight       1.69         0.182     9.24  4.23e-20
## 9 type_I            -0.806        0.115    -6.99  3.29e-12
##10 type_M             0.0220       0.0920     0.239 8.11e- 1
##11 longest_shell_x_diameter -0.447      0.0492    -9.08  1.78e-19
##12 shell_weight_x_shucked_weight -0.0615    0.0527    -1.17  2.43e- 1
```

```
predict(lm_fit, new_data = test)
```

```
## # A tibble: 1 x 1
##   .pred
##   <dbl>
## 1  24.1
```

```
lm_fit
```

```
## == Workflow [trained] =====
## Preprocessor: Recipe
## Model: linear_reg()
##
## -- Preprocessor -----
## 6 Recipe Steps
##
## * step_dummy()
## * step_center()
## * step_scale()
## * step_interact()
## * step_interact()
## * step_interact()
##
## -- Model -----
##
## Call:
## stats::lm(formula = ..y ~ ., data = data)
##
## Coefficients:
##              (Intercept)              longest_shell
##              12.17642                -0.83846
##              diameter                height
##              0.45498                 0.25308
##              whole_weight            shucked_weight
##              4.71785                -3.99668
##              viscera_weight          shell_weight
##              -0.88750                 1.68594
```

```
##           type_I           type_M
##      -0.80574      0.02196
## longest_shell_x_diameter shell_weight_x_shucked_weight
##      -0.44729      -0.06153
```

Question 7

Now you want to assess your model's performance. To do this, use the `yardstick` package:

1. Create a metric set that includes R^2 , RMSE (root mean squared error), and MAE (mean absolute error).
2. Use `predict()` and `bind_cols()` to create a tibble of your model's predicted values from the **training data** along with the actual observed ages (these are needed to assess your model's performance).
3. Finally, apply your metric set to the tibble, report the results, and interpret the R^2 value.

From what we got for R^2 which is roughly 0.55, it is not significant enough to show a strong correlation, it can only be considered relatively strong, but it is not significant enough to compare with our initial function which is $\text{age} = \text{rings} + 1.5$. Though the question didn't require us to plot a graph, but the scatter plot should be a clear visual representation that our model didn't do well.

```
abalone_train %>%
  head()
```

```
##   type longest_shell diameter height whole_weight shucked_weight
## 5    I      0.330      0.255  0.080      0.2050      0.0895
## 6    I      0.425      0.300  0.095      0.3515      0.1410
## 17   I      0.355      0.280  0.085      0.2905      0.0950
## 19   M      0.365      0.295  0.080      0.2555      0.0970
## 36   M      0.465      0.355  0.105      0.4795      0.2270
## 43   I      0.240      0.175  0.045      0.0700      0.0315
##   viscera_weight shell_weight rings age
## 5      0.0395      0.055      7 8.5
## 6      0.0775      0.120      8 9.5
## 17     0.0395      0.115      7 8.5
## 19     0.0430      0.100      7 8.5
## 36     0.1240      0.125      8 9.5
## 43     0.0235      0.020      5 6.5
```

```
abalone_train_res <- predict(lm_fit, new_data = abalone_train %>%
  select(-age, -rings))

abalone_train_res %>%
  head()
```

```
## # A tibble: 6 x 1
##   .pred
##   <dbl>
## 1  8.23
## 2  9.58
## 3  9.98
## 4 10.3
## 5 10.1
## 6  6.33
```



```

abalone_train_res <- bind_cols(abalone_train_res, abalone_train %>% select(age))

abalone_train_res %>%
  head()

```

```

## # A tibble: 6 x 2
##   .pred age
##   <dbl> <dbl>
## 1  8.23  8.5
## 2  9.58  9.5
## 3  9.98  8.5
## 4 10.3   8.5
## 5 10.1   9.5
## 6  6.33  6.5

```

```

rmse(abalone_train_res, truth = age, estimate = .pred)

```

```

## # A tibble: 1 x 3
##   .metric .estimator .estimate
##   <chr>   <chr>       <dbl>
## 1 rmse    standard      2.16

```

```

abalone_metrics <- metric_set(rmse, rsq, mae)

abalone_metrics(abalone_train_res, truth = age,
  estimate = .pred)

```

```

## # A tibble: 3 x 3
##   .metric .estimator .estimate
##   <chr>   <chr>       <dbl>
## 1 rmse    standard      2.16
## 2 rsq     standard      0.552
## 3 mae     standard      1.57

```

```

abalone_train_res %>%
  ggplot(aes(x = .pred, y = age)) +
  geom_point(alpha = 0.2) +
  geom_abline(lty = 2) +
  theme_bw() +
  coord_obs_pred()

```

