

## Linear Regression

# Homework 2

Code ▾

PSTAT 131/231

## Linear Regression

For this lab, we will be working with a data set from the UCI (University of California, Irvine) Machine Learning repository (see website here (<http://archive.ics.uci.edu/ml/datasets/Abalone>)). The full data set consists of 4,177 observations of abalone in Tasmania. (Fun fact: Tasmania (<https://en.wikipedia.org/wiki/Tasmania>) supplies about 25% of the yearly world abalone harvest.)



Fig 1. Inside of an abalone shell.

The age of an abalone is typically determined by cutting the shell open and counting the number of rings with a microscope. The purpose of this data set is to determine whether abalone age (**number of rings + 1.5**) can be accurately predicted using other, easier-to-obtain information about the abalone.

The full abalone data set is located in the `\data` subdirectory. Read it into *R* using `read_csv()`. Take a moment to read through the codebook (`abalone_codebook.txt`) and familiarize yourself with the variable definitions.

Make sure you load the `tidyverse` and `tidymodels`!

### 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.

Assess and describe the distribution of `age`.

Hide

```
library(tidyverse)
library(ggplot2)
library(tidymodels)
library(corrplot)
library(ggthemes)
tidymodels_prefer()
```

Hide

```
abalone_data <- read_csv('abalone.csv')
head(abalone_data)
```

```
##   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
## 1          0.150   15
## 2          0.070    7
## 3          0.210    9
## 4          0.155   10
## 5          0.055    7
## 6          0.120    8
```

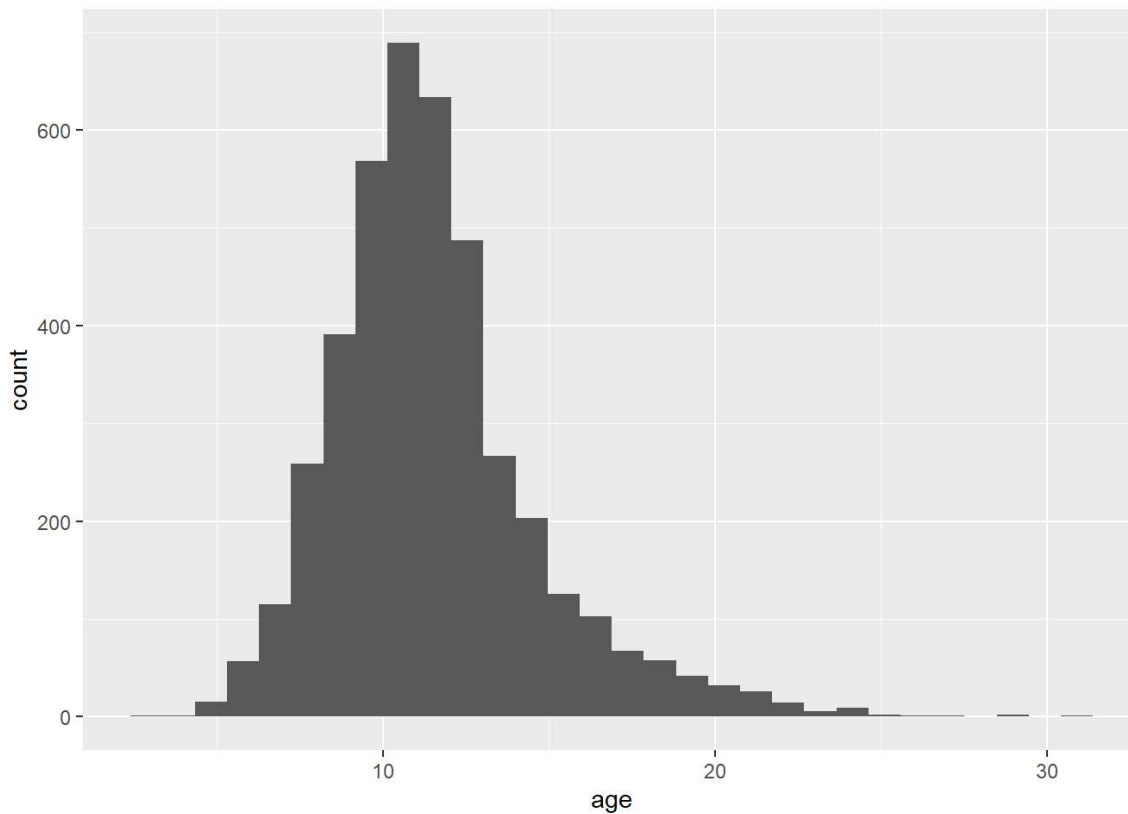
Hide

```
abalone <- abalone_data %>%
  mutate(abalone_data, age=rings+1.5)
head(abalone)
```

```
##   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
```

Hide

```
ggplot(data=abalone, aes(age)) +
  geom_histogram()
```



## 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.*

[Hide](#)

```
set.seed(3435)

abalone_split <- initial_split(abalone, prop = 0.80,
                               strata = age)
abalone_train <- training(abalone_split)
abalone_test <- testing(abalone_split)

head(abalone_test)
```

##	type	longest_shell	diameter	height	whole_weight	shucked_weight
## 2	M	0.350	0.265	0.090	0.2255	0.0995
## 6	I	0.425	0.300	0.095	0.3515	0.1410
## 13	M	0.490	0.380	0.135	0.5415	0.2175
## 28	M	0.590	0.445	0.140	0.9310	0.3560
## 30	M	0.575	0.425	0.140	0.8635	0.3930
## 39	F	0.575	0.445	0.135	0.8830	0.3810

##	viscera_weight	shell_weight	rings	age
## 2	0.0485	0.07	7	8.5
## 6	0.0775	0.12	8	9.5
## 13	0.0950	0.19	11	12.5
## 28	0.2340	0.28	12	13.5
## 30	0.2270	0.20	11	12.5
## 39	0.2035	0.26	11	12.5

## 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`.

Steps for your recipe:

1. dummy code any categorical predictors
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.

Hide

```
abalone_train1 <- abalone_train %>% select(-rings)

abalone_recipe <- recipe(age ~ ., data = abalone_train1) %>%
  step_dummy(all_nominal_predictors())

abalone_recipe
```

```
## Recipe
##
## Inputs:
##
##      role #variables
## outcome      1
## predictor      8
##
## Operations:
##
## Dummy variables from all_nominal_predictors()
```

Hide

```
int_mod_1 <- abalone_recipe %>%
  step_interact(terms = ~ starts_with("type"):shucked_weight)

int_mod_2 <- int_mod_1 %>%
  step_interact(terms = ~ longest_shell:diameter)

int_mod_3 <- int_mod_2 %>%
  step_interact(terms = ~ shucked_weight:shell_weight)
int_mod_3 %>% prep() %>% bake(abalone_train1)
```

```
## # A tibble: 3,340 x 14
##   longest_shell diameter height whole_weight shucked_weight viscera_weight
##   <dbl>      <dbl> <dbl>      <dbl>      <dbl>      <dbl>
## 1      0.33      0.255 0.08      0.205      0.0895     0.0395
## 2      0.355     0.28  0.085     0.290      0.095     0.0395
## 3      0.365     0.295 0.08      0.256      0.097     0.043
## 4      0.465     0.355 0.105     0.480      0.227     0.124
## 5      0.45      0.355 0.105     0.522      0.237     0.116
## 6      0.24      0.175 0.045     0.07       0.0315    0.0235
## 7      0.205     0.15  0.055     0.042      0.0255    0.015
## 8      0.21      0.15  0.05      0.042      0.0175    0.0125
## 9      0.39      0.295 0.095     0.203      0.0875    0.045
## 10     0.46      0.375 0.12      0.460      0.178     0.11
## # ... with 3,330 more rows, and 8 more variables: shell_weight <dbl>,
## #   age <dbl>, type_I <dbl>, type_M <dbl>, type_I_x_shucked_weight <dbl>,
## #   type_M_x_shucked_weight <dbl>, longest_shell_x_diameter <dbl>,
## #   shucked_weight_x_shell_weight <dbl>
```

Hide

```
norm_trans <- int_mod_3 %>%
  step_normalize(all_predictors())
norm_trans %>% prep() %>% bake(abalone_train1)
```

```
## # A tibble: 3,340 x 14
##   longest_shell diameter height whole_weight shucked_weight viscera_weight
##   <dbl>      <dbl> <dbl>      <dbl>      <dbl>      <dbl>
## 1     -1.61    -1.54 -1.40      -1.27      -1.21      -1.28
## 2     -1.40    -1.29 -1.28      -1.10      -1.19      -1.28
## 3     -1.32    -1.14 -1.40      -1.17      -1.18      -1.25
## 4     -0.489   -0.532 -0.810     -0.711     -0.593     -0.514
## 5     -0.613   -0.532 -0.810     -0.623     -0.548     -0.582
## 6     -2.36    -2.35 -2.22      -1.54      -1.47      -1.43
## 7     -2.65    -2.60 -1.98      -1.60      -1.50      -1.51
## 8     -2.61    -2.60 -2.10      -1.60      -1.53      -1.53
## 9     -1.11    -1.14 -1.04      -1.27      -1.22      -1.23
## 10    -0.530    -0.330 -0.458     -0.749     -0.815     -0.641
## # ... with 3,330 more rows, and 8 more variables: shell_weight <dbl>,
## #   age <dbl>, type_I <dbl>, type_M <dbl>, type_I_x_shucked_weight <dbl>,
## #   type_M_x_shucked_weight <dbl>, longest_shell_x_diameter <dbl>,
## #   shucked_weight_x_shell_weight <dbl>
```

The rings should not be used as predictor since the response age is obtained by adding every value of rings by 1.5. If we use the rings to predict age, we would definitely correctly predict all age with rings.

## Question 4

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

Hide

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

## 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.

Hide

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

## 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`.

Hide

```
lm_fit <- fit(lm_wflow, abalone_train1)
```

Hide

```
lm_fit %>%
  # This returns the parsnip object:
  extract_fit_parsnip() %>%
  # Now tidy the linear model object:
  tidy()
```

```
## # A tibble: 14 x 5
##   term                                estimate std.error statistic  p.value
##   <chr>                                <dbl>    <dbl>    <dbl>    <dbl>
## 1 (Intercept)                        11.4      0.0375   305.      0
## 2 longest_shell                       0.591     0.286     2.07    3.86e- 2
## 3 diameter                           2.06      0.313     6.61   4.59e-11
## 4 height                             0.236     0.0696     3.39   7.10e- 4
## 5 whole_weight                       4.29      0.387    11.1   4.66e-28
## 6 shucked_weight                     -4.06      0.250   -16.2   5.35e-57
## 7 viscera_weight                     -0.792     0.158    -5.00   6.12e- 7
## 8 shell_weight                       1.74      0.212     8.20   3.32e-16
## 9 type_I                             -0.942     0.117    -8.07   9.36e-16
##10 type_M                             -0.239     0.104    -2.29   2.21e- 2
##11 type_I_x_shucked_weight            0.525     0.0876     5.99   2.26e- 9
##12 type_M_x_shucked_weight            0.293     0.109     2.68   7.41e- 3
##13 longest_shell_x_diameter           -2.75      0.396    -6.95   4.32e-12
##14 shucked_weight_x_shell_weight     -0.00330    0.205    -0.0161 9.87e- 1
```

Hide

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

predict(lm_fit, new_data = new)
```

```
## # A tibble: 1 x 1
##   .pred
##   <dbl>
## 1  23.7
```

## 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.

Hide

```
library(yardstick)
```

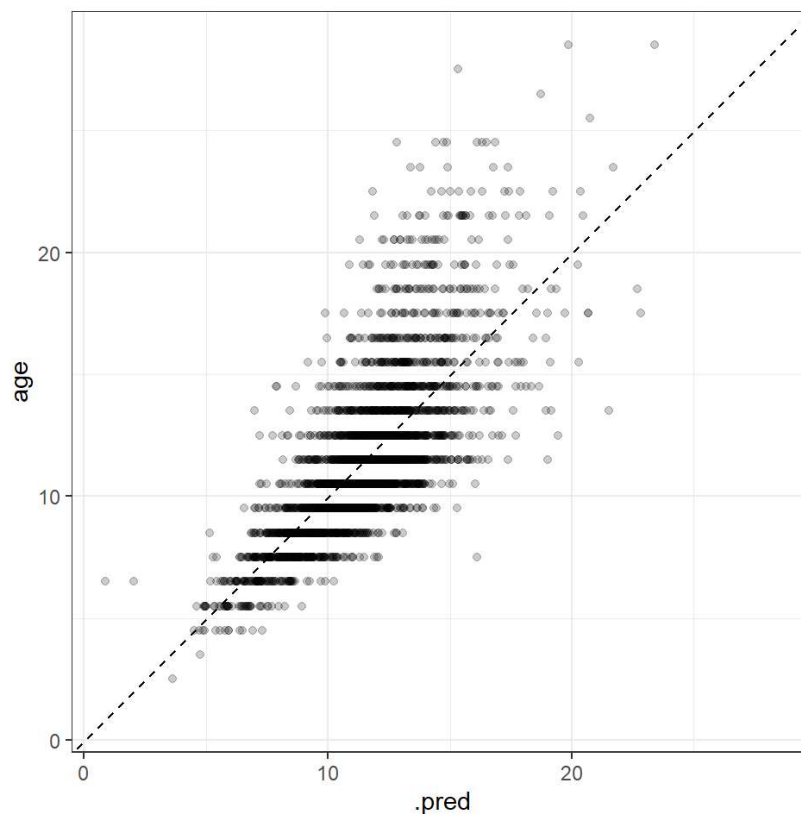
Hide

```
abalone_train_res <- predict(lm_fit, new_data = abalone_train1 %>% select(-age))
abalone_train_res <- bind_cols(abalone_train_res, abalone_train1 %>% select(age))
head(abalone_train_res)
```

```
## # A tibble: 6 x 2
##   .pred age
##   <dbl> <dbl>
## 1  8.03  8.5
## 2  9.68  8.5
## 3 10.4   8.5
## 4 10.1   9.5
## 5 10.9   9.5
## 6  6.26  6.5
```

Hide

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



Hide

```
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         0.551
## 3 mae     standard         1.55
```

$R^2$ : 55.1% of variability of age can be explained by this linear regression model.

## Required for 231 Students

In lecture, we presented the general bias-variance tradeoff, which takes the form:

$$E[(y_0 - \hat{f}(x_0))^2] = \text{Var}(\hat{f}(x_0)) + [\text{Bias}(\hat{f}(x_0))]^2 + \text{Var}(\epsilon)$$

where the underlying model  $Y = f(X) + \epsilon$  satisfies the following:

- $\epsilon$  is a zero-mean random noise term and  $X$  is non-random (all randomness in  $Y$  comes from  $\epsilon$ );
- $(x_0, y_0)$  represents a test observation, independent of the training set, drawn from the same model;
- $\hat{f}(\cdot)$  is the estimate of  $f$  obtained from the training set.

### Question 8

Which term(s) in the bias-variance tradeoff above represent the reproducible error? Which term(s) represent the irreducible error?

reducible error:  $\text{Var}(\hat{f}(x_0))$  and  $[\text{Bias}(\hat{f}(x_0))]^2$  irreducible error:  $\text{Var}(\epsilon)$

### Question 9

Using the bias-variance tradeoff above, demonstrate that the expected test error is always at least as large as the irreducible error.

Answer: Even if we make the reducible error equal to zero, which means  $\text{Var}(\hat{f}(x_0)) = 0$  and  $[\text{Bias}(\hat{f}(x_0))]^2 = 0$ , the irreducible error still exists.

$$E[(y_0 - \hat{f}(x_0))^2] = \text{Var}(\hat{f}(x_0)) + [\text{Bias}(\hat{f}(x_0))]^2 + \text{Var}(\epsilon) = 0 + \text{Var}(\epsilon) = \text{Var}(\epsilon)$$

### Question 10

Prove the bias-variance tradeoff.

Hints:

- use the definition of  $\text{Bias}(\hat{f}(x_0)) = E[\hat{f}(x_0)] - f(x_0)$ ;
- reorganize terms in the expected test error by adding and subtracting  $E[\hat{f}(x_0)]$

Proof:

$$\begin{aligned} E[y_0 - \hat{f}(x_0)]^2 &= E[(f(x_0) - \epsilon - \hat{f}(x_0))^2] \\ &= E[(f^2(x_0) + \epsilon f(x_0) - f(x_0)\hat{f}(x_0) + \epsilon f(x_0) + \epsilon^2 - \epsilon \hat{f}(x_0) - f(x_0)\hat{f}(x_0) - \epsilon \hat{f}(x_0) + \hat{f}^2(x_0))] \\ &= E[(f(x_0) - \hat{f}(x_0))^2] + E[\epsilon] + 2E[(f(x_0) - \hat{f}(x_0))\epsilon] \\ &= E[(f(x_0) - \hat{f}(x_0))^2] + \text{Var}(\epsilon) \\ &= E[(f(x_0) + E[\hat{f}(x_0)] - E[\hat{f}(x_0)] - \hat{f}(x_0))^2] + \text{Var}(\epsilon) \\ &= E[(E[\hat{f}(x_0)] - f(x_0))^2] + E[(\hat{f}(x_0) - E[\hat{f}(x_0)])^2] - 2E[(f(x_0) - E[\hat{f}(x_0)])(E[\hat{f}(x_0)] - E[\hat{f}(x_0)])] + \text{Var}(\epsilon) \end{aligned}$$

After reorganizing the terms, get:

$$= (E[\hat{f}(x_0)] - f(x_0))^2 + E[(\hat{f}(x_0) - E[\hat{f}(x_0)])^2] + \text{Var}(\epsilon)$$

Applying the definition for bias and variance:

$$= (\text{bias}(\hat{f}(x_0)))^2 + \text{Var}(\hat{f}(x_0)) + \text{Var}(\epsilon)$$



