STATS 7022 - Data Science

Practical Exam

Prithviraj Banerjee Trimester 1, 2023

Overview

This practical exam is an opportunity for you to demonstrate your ability to analyse data in R using the tidymodels package. You will be analysing the diamonds dataset from the ggplot2 package to predict the price of a diamond.

Instructions 1. Download the file DSPG_Practical_Exam.Rmd from MyUni.

2. Within this .Rmd file, use the function start_exam() to obtain the data, your specified method, and your specified response variable.

3. Perform a complete analysis of the data. Your analysis should include at least the following: a. exploratory data analysis

b. data cleaning and pre-processing c. fitting a model using your specified method d. hyperparameter tuning e. assessing the fit of your model f. choosing a single model for prediction g. using your model to make predictions, using an example of your choice 4. You may use this file as a template for your analysis.

5. From the time the exam opens on MyUni, you have 24 hours to complete and submit the analysis.

Submission

You MUST submit your code as a .R or .Rmd file. Your marker may run your code. You may lose marks if your code does not produce the output

You MUST submit your analysis at the Practical Exam submission portal on MyUni.

contained in your analysis. No other file submissions will be accepted. Start the Exam

Enter your student ID in the code block below, then run the code block to start the exam. You MUST enter your ID as an integer. For example, if your student number is a1234567, then please enter your ID as ID <- 1234567. You MUST enter your own student ID. Using a different ID may result in a loss of marks.

library(rsample)

Importing the necessary libraries

You MUST submit your analysis as a .pdf file. You may export your analysis directly to .pdf, or export as an .html file and convert to .pdf.

#install.packages("parsnip") library(parsnip) #install.packages("tune")

library(ranger) library(tune) library(dplyr)

library(workflows) library(doParallel)

#install.packages("workflows")

Loading the data and the method ID <- 1753181 data <- start_exam(ID)</pre>

Please fit a random forest model to the data returned by this function, specifying price as the response varia ble. Analysis Inspect the data head(data)

table

< | db|>

56.1

55.0

57.0

8.06

58.0

60.0

depth

<dpl>

61.7

62.2

61.4

61.1

62.1

62.3

color

<ord>

F

F

F

J

D

D

Z

<qpl>

2.89

3.23

3.25

4.18

4.10

4.14

p25

p0

sd

p75

p100

5.01

95.00

79.00

31.80

10.74

p50

X

<dpl>

4.65

5.22

5.28

6.83

6.62

6.59

clarity

<ord>

VS1

VS1

SI1

SI2

SI2

SI2

939 0.39 Ideal 0.53 1908 Ideal 1340 0.56 Ideal 4824 1.21 Good

1.08

1.13

carat

<dpl>

6 rows

5138

5408

price

<int>

We get an idea about the data. The data contains 43152 rows and 9 predictors with price as the response variable. **EDA**

cut

<ord>

Premium

Very Good

Display tables skimr::skim_without_charts(data)	
Name	data
Number of rows	43152
Number of columns	9
Column type frequency:	
factor	3
numeric	6
Group variables	None

1 TRUE cut 0

complete_rate

n_missing

clarity

Variable type: factor

Variable type: numeric

skim_variable

n_missing n_unique top_counts skim_variable complete_rate ordered 5 Ide: 17252, Pre: 11031, Ver: 9703, Goo: 3859 0 1 TRUE color 7 G: 9104, E: 7844, F: 7589, H: 6617 8 SI1: 10461, VS2: 9753, SI2: 7351, VS1: 6574 0 1 TRUE

mean

0 3933.10 3979.33 326.0 953.00 2418.50 5334.25 18823.00 price 1 0 carat 1 0.80 0.47 0.2 0.40 0.70 1.04 table 0 1 57.45 2.24 44.0 56.00 57.00 59.00 0 depth 1 61.75 43.0 61.00 61.80 62.50 1.43 0 4.03 Z 1 3.54 0.71 0.0 2.91 3.53 Χ 0 1 5.73 1.12 0.0 4.72 5.70 6.54

min(data\$price) ## [1] 326

9000

6000

15000

5000

Cleaning

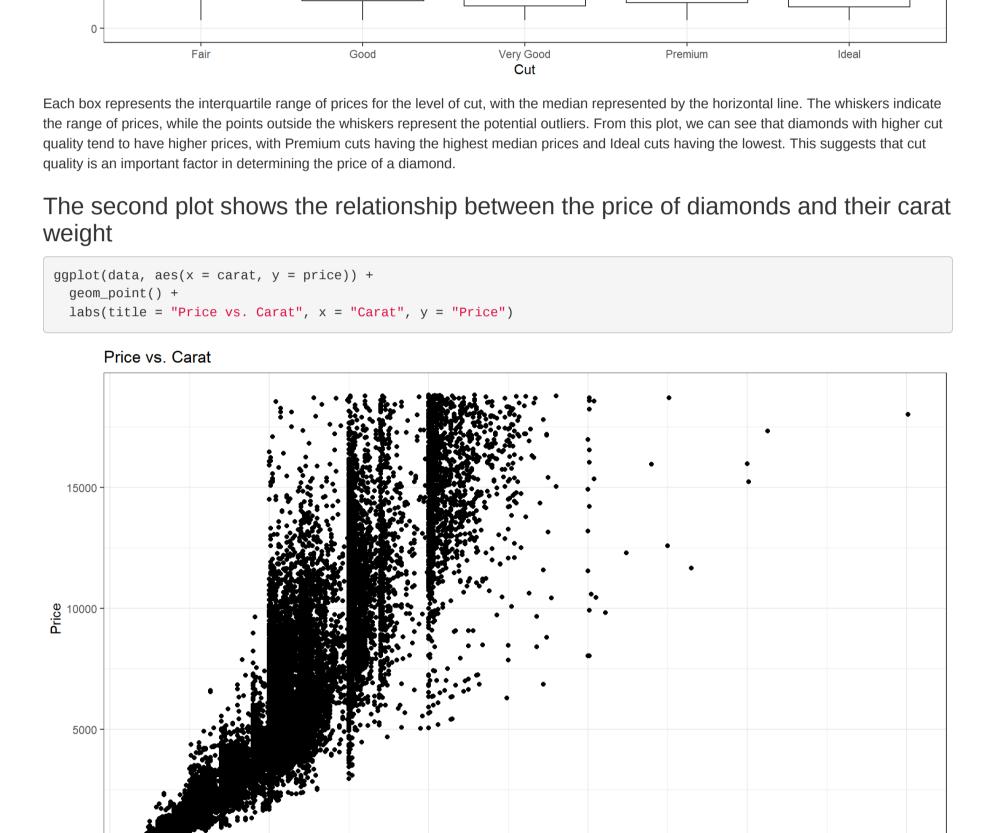
may drive up the price.

mean(data\$price)

[1] 3933.105

max(data\$price) ## [1] 18823 The sample mean of the response variable (price) is 3933.105 and the minimum value of price is 326 and the max value is 18823 from the table above. Visualising the response data %>% ggplot(aes(price)) + geom_histogram(col = "black", fill = "#9f4fff")

3000 5000 10000 15000 20000 price Since the response variable price is a quantitative continuous variable a histogram will give us a better understanding to visualize the data. From the histogram it is clearly visible that it is a right-skewed uni modal distribution. A box plot can be obtained between the response variable (price) and the predictor variable (cut) ggplot(data, aes(x = cut, y = price)) +geom_boxplot() + labs(title = "Price vs. Cut", x = "Cut", y = "Price") Price vs. Cut



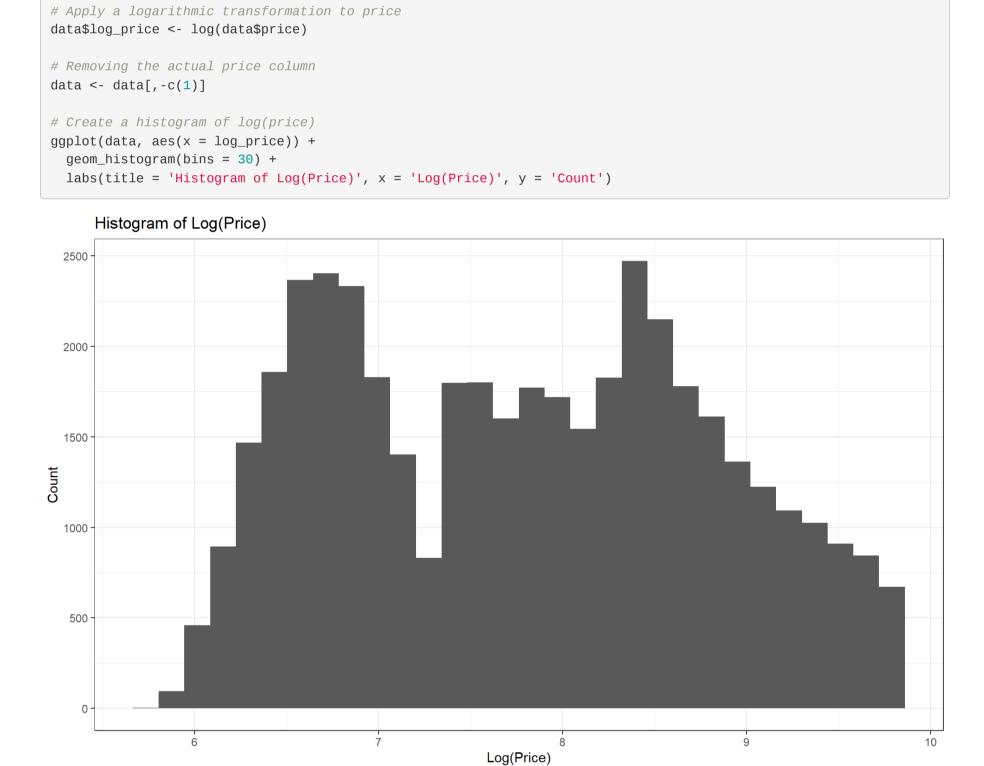
Carat

The x-axis represents the carat weight, while the y-axis represents the price. Each point on the plot represents a diamond, and the position of the point indicates its carat weight and price. From this plot, we can see that there is a strong positive co-relationship between the carat weight and

After looking at the distribution of price, it is clearly visible that the response variable

is rightly skewed. This suggests that there may be extreme outliers or values that

price of diamonds, as expected. This means that, in general, larger diamonds tend to be more expensive than smaller diamonds.



Therefore, a logarithmic transformation makes the distribution more symmetric and easier to analyse.

data\$cut <- factor(data\$cut, levels = c('Fair', 'Good', 'Very Good', 'Premium', 'Ideal'))</pre>

data\$clarity <- factor(data\$clarity, levels = c("I1", "SI2", "SI1", "VS2", "VS1", "VVS2", "VVS1", "IF"))</pre>

Convert the cut, color and clarity variable to a factor

There are no missing values in our data. So no imputation is required.

Splitting the data into test and training sets

data_split <- initial_split(data, strata = log_price)</pre>

we will partition our training set into 10 folds.

data_folds <- vfold_cv(data_train, v = 10, strata = log_price)</pre>

Convert ordinal to nominal

data\$color <- factor(data\$color)</pre>

Checking for null values

sum(is.na(data))

Pre-processing

data_train <- training(data_split)</pre> data_test <- testing(data_split)</pre>

set.seed(2022)

data_folds

-1.011596503

-0.969337947

-1.032725781

-0.990467225

-0.884820834

1-10 of 10,000 rows

 $min_n = tune(),$ trees = 10) %>%

set_mode("regression") %>%

set_engine("ranger")

Create the workflow

data_wf <- workflow() %>% add_recipe(data_recipe) %>%

add_model(data_model)

Preprocessor: Recipe ## Model: rand_forest()

== Workflow =

— Preprocessor ## 2 Recipe Steps

• step_center() ## • step_scale()

Main Arguments: mtry = tune() trees = 10 $min_n = tune()$

Computational engine: ranger

Creating the tuning grid

data_wf <- data_wf %>%

Preprocessor: Recipe ## Model: rand_forest()

data_wf

== Workflow =

— Preprocessor ## 2 Recipe Steps

• step_center() ## • step_scale()

Main Arguments: mtry = 5trees = 10 $min_n = 11$

Computational engine: ranger

--- Model -

0.990

0.985

0.980

0.975

0.970

.metric

##Prediction

Prediction for new observation

View the predicted log_price

prediction\$.pred

[1] 8.481783

Prediction for the first observation in test data

prediction <- predict(data_wf_fit, new_data = data_test[1,])</pre>

Prediction

Testing the model performance

data_fit %>% collect_metrics()

calculate the regular performance metrics

.estimator

data_fit <- data_wf %>% last_fit(split = data_split)

finalize_workflow(select_best(data_tune, metric = "rmse"))

Random Forest Model Specification (regression)

data_grid <- grid_regular(mtry(c(1,5)),</pre>

-- Model -

data_wf

Ideal

Very Good

Premium

Very Good

Very Good

Modelling and Tuning

data_model <- rand_forest(mtry = tune(),</pre>

Random Forest Model Specification (regression)

-0.20174598

1.14001192

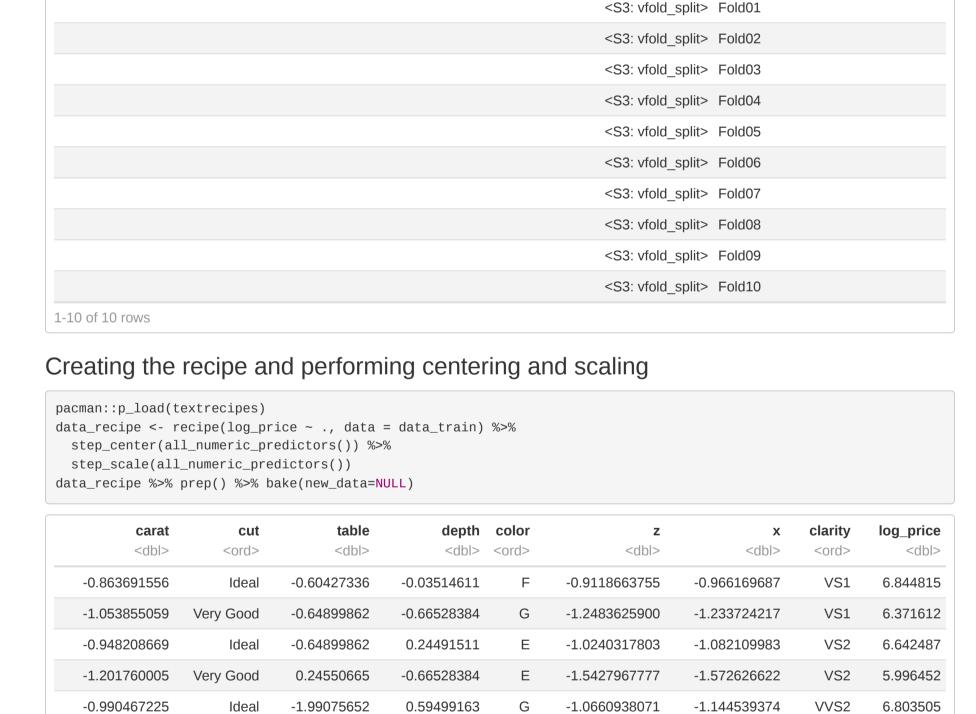
0.69275929

0.24550665

-0.20174598

hyperparameters we want to tune in our model:

[1] O



-0.24519202

0.10488450

0.31493041

-0.24519202

1.15511406

our response variable is a quantitative variable. We specify tune() for the

D

Η

G

We create a recipe object and perform pre-processing steps like centering and scaling to all the numeric predictors. We do this on the training data.

we specify a random forest model. Note that we set_mode() to regression because

-1.1642385364

-1.1081558339

-1.1782592120

-1.1081558339

-0.8698043487

-1.180213311

-1.171294827

-1.224805733

-1.135620889

-0.984006656

Previous 1 2 3 4

SI1

VS2

VVS2

SI1

SI2

6.378426

6.285998

6.810142

6.098074

6.553933

2

2

2

2

11

11

11

11

11

→ 30 **→** 40

5 6 ... 1000 Next

splits id <chr>

min_n(), levels = 5) data_grid mtry min_n <int> <int> 1 2 3 4 5 1 2 3 4 5 1-10 of 25 rows Previous 1 2 3 Next Tuning the model for optimal hyperparameters doParallel::registerDoParallel() data_tune <- tune_grid(data_wf,</pre> resamples = data_folds, grid = data_grid)

Choose the optimal hyperparameters based on tuning results

Plot the tuning model data_tune %>% autoplot() 0.175 0.150 0.125 Minimal Node Size

<chr> <dbl> <chr> 0.1017875 Preprocessor1_Model1 standard rmse standard 0.9898090 Preprocessor1_Model1 rsq 2 rows Visualising the model performance data_fit %>% collect_predictions() %>% ggplot(aes(log_price, .pred)) + geom_point() + geom_smooth(method = "lm") + geom_abline(intercept = 0, slope = 1)

.estimate .config

Randomly Selected Predictors

log_price Fitting the best model to the train data data_wf_fit <- fit(data_wf, data = data_train)</pre>