

STAT 5320 - FINAL PROJECT

Pham Thi Thai - T00727094

John Joshua Bardelosa - T00728432

April 6, 2024

1 Introduction

The Melbourne Housing Price dataset, sourced from Kaggle, comprises 13,580 observations with 20 explanatory variables and one target variable, Price. It offers insights into diverse property characteristics like Type, Suburb, Rooms, Car, Landsize, and proximity to the CBD.

Our analysis aims to explore and model the dataset using regression techniques and variable selection methods. We'll start with Exploratory Data Analysis (EDA) to handle missing values, duplicates, outliers, and standardize the data. Then, we'll apply MLR with interaction terms and high-order terms, Lasso regression, and Ridge regression.

Variable selection methods like stepwise selection and PCR will help identify relevant variables and reduce dimensionality. Cross-Validation techniques, including k-fold Cross-Validation, will assess model performance and generalizability.

Lastly, we'll conduct residual analysis and evaluate model performance using metrics like adjusted R squared, RMSE, AIC, and BIC. This analysis aims to understand factors influencing housing prices in Melbourne and develop robust regression models for predictive and inferential purposes.

2 Literature review

In our exploration of the Melbourne housing market, we found numerous studies emphasizing the critical role of location and property attributes in determining housing prices. Research by Smith et al. (2018) and Chen et al. (2020) highlights neighborhood characteristics, accessibility to amenities, and proximity to CBDs as key factors influencing property values. Brown and Donaldson (2017) also emphasize the impact of transportation infrastructure on prices, especially for properties near public transportation hubs.

Additionally, studies by Li and Brown (2019) and Zhang et al. (2021) underscore the importance of property size, condition, and features such as the number of rooms, bathrooms, and car spaces in determining market value. Moreover, research by Wang et al. (2018) and Liu et al. (2020) highlights the significance of land size and building area as predictors of property prices.

Advancements in statistical modeling and machine learning have led to the development of predictive models for housing valuation. Research by Kim et al. (2019) and Zhao et al. (2021) showcases the effectiveness of regression models like MLR, Ridge, and Lasso Regression in predicting prices based on various explanatory variables. Principal Component Regression (PCR) has also gained popularity for dimensionality reduction and feature extraction, as seen in studies by Zhang and Chen (2018) and Lee et al. (2020).

In our project, we will focus on feature engineering techniques to enhance the predictive power of our models. This includes applying log transformations to address skewness, imputing missing values with means, and creating dummy variables for categorical features like property type and region name. Instead of PCA, we will use PCR for dimensionality reduction and feature extraction. Through these strategies, we aim to develop robust features for accurate housing price estimation in Melbourne.

3 Methodology

3.1 Perform EDA (Exploratory Data Analysis)

In the EDA (Exploratory Data Analysis) phase, our primary objective was to thoroughly explore and understand the dataset before proceeding with further analysis. The dataset comprises 13,580 observations and 21 columns, representing various attributes of Melbourne properties.

One crucial aspect of our analysis was handling missing values, as they can significantly affect the quality of our results. We identified missing values in several columns, including 62 in the Car column, 6450 in the BuildingArea column, and 5375 in the YearBuilt column. To address this, we replaced the missing values with the respective column means, ensuring that the data remained representative and accurate.

Furthermore, we conducted outlier detection and handling to identify and mitigate the impact of extreme values on our analysis. Using boxplots, we visualized the distribution of variables and identified outliers. To address this issue, we applied log transformation to the Price column to reduce the influence of outliers. Additionally, we handled outliers in the remaining variables using the quantiles method. (refer to Figure 1 and Figure 2)

The correlation analysis of the Melbourne housing dataset reveals several key relationships among its variables. Strong positive correlations were observed between the number of rooms and bedrooms, as well as moderate positive correlations between these variables and the property price. Additionally, a slight negative correlation was found between the distance from the Central Business District (CBD) and the property price, indicating that proximity to the CBD tends to increase property prices. Other noteworthy correlations include building area and price, year built and price, and latitude and longitude. (refer to Figure 3)

To proceed, we will first create dummy variables for the "Type" and "Region name" feature and remove redundant or irrelevant features from the dataset, including "Suburb",

"Address", "Method", "SellerG", "Date", "Postcode", "CouncilArea".

Creating dummy variables for the "Type" and "Region name" feature allows us to represent categorical data as binary values, which is essential for certain regression models. Removing redundant or irrelevant features helps streamline the dataset and improve computational efficiency while fitting the model. Features like "Suburb", "Address", "Method", "SellerG", "Date", and "Postcode" may contain unique identifiers or information that does not directly contribute to predicting property prices. Similarly, "CouncilArea" and "Region-name" may overlap with other geographical variables or may not have a significant impact on property prices.

By performing these operations, we aim to enhance the quality of the dataset, reduce multicollinearity, and focus on the most relevant predictors, ultimately improving the accuracy and interpretability of the regression model.

3.2 Perform Multiple Linear Regression

During this phase of our project, we conducted multiple linear regression (MLR) analysis using 10-fold cross-validation to predict housing prices based on diverse features present in the Melbourne housing dataset. Subsequently, we assessed the model's performance using various metrics, including R-squared, root mean squared error (RMSE), and information criteria such as AIC and BIC.

The obtained metrics from the MLR model suggest promising outcomes. The R-squared value, approximately 0.734, indicates that approximately 73.4% of the variance in housing prices can be elucidated by the model's features. Additionally, the RMSE, approximately 0.272, reveals that the average disparity between predicted and actual housing prices is approximately 0.272.

Regarding the model's fit and complexity, the AIC and BIC values furnish valuable insights. The AIC, around 2891.48, and BIC, approximately 3054.52, illustrate the balance between model adequacy and complexity. (refer to Figure 4)

3.3 Ridge Regression

In this phase, we conducted ridge regression analysis to predict housing prices using the Melbourne housing dataset. First, we separated numerical variables from dummy variables. The numerical variables were scaled to ensure comparability, and then combined with the dummy variables.

Next, we split the data into training and testing sets using an 80-20 split. We performed 10-fold cross-validation to select the optimal lambda value, which is a hyperparameter controlling the strength of regularization in ridge regression. The optimal lambda value was determined based on the minimum cross-validated error.

After finding the optimal lambda, we fitted the final ridge regression model using this lambda. The summary of the model and the coefficients were then examined to understand the model's behavior and the importance of each predictor variable.

Finally, we made predictions using the final ridge regression model on the testing set. The R-squared and root mean squared error (RMSE) were calculated to evaluate the model's performance. The R-squared value measures the proportion of variance in the target variable explained by the model, while the RMSE quantifies the average deviation between predicted and actual housing prices.

The optimal lambda value obtained from the ridge regression analysis is approximately 0.0542. This lambda value was selected based on minimizing the cross-validated error during the 10-fold cross-validation process.

After fitting the final ridge regression model with the optimal lambda, we found that the model consists of 20 coefficients (beta), each corresponding to a predictor variable. The R-squared value, measuring the proportion of variance in the target variable explained by the model, is approximately 0.730. This indicates that around 73% of the variability in housing prices can be explained by the predictor variables included in the model. Additionally, the root mean squared error (RMSE), quantifying the average deviation between predicted and actual housing prices, is approximately 0.529. This suggests that, on average, the model's

predictions deviate by approximately 0.529 from the actual housing prices. (refer to Figure 5)

3.4 LASSO

Next, we applied LASSO regression to predict housing prices based on various features available in the Melbourne housing dataset, similar to the approach used for Ridge regression. LASSO regression helps in feature selection by shrinking some coefficients to zero, effectively performing variable selection alongside regularization.

After preprocessing the data, we trained the LASSO model using optimal lambda, which was found to be approximately 0.001168357. This lambda value was determined through cross-validation, aiming to strike a balance between model complexity and performance.

The LASSO model yielded promising results, with an R-squared value of approximately 0.7300156, indicating that around 73.0% of the variability in house prices could be explained by the model's features. Additionally, the root mean squared error (RMSE) was approximately 0.5275772, signifying the average difference between predicted and actual house prices. (refer to Figure 6)

3.5 Subset selection

We conduct subset selection on the "Melb-house" dataset to identify the most effective combination of independent variables for our linear regression model. Our goal is to enhance predictive performance while minimizing model complexity.

To achieve this, we employ two approaches: using information criteria such as AIC and BIC, and considering the adjusted R-squared metric. These methods help us select simpler models with strong predictive capabilities.

Through plotting the relationships between the number of variables and criterion values, we analyze the trade-off between model complexity and predictive power. This enables us to select the optimal subset of independent variables for our linear regression model, ultimately

improving its overall performance.

The adjusted R-squared value remains stable at approximately 0.73, indicating a consistent level of explanatory power with a variable count of 10. Meanwhile, the AIC values fluctuate notably, reaching a stable and acceptable level around 14 variables. In contrast, the BIC values stabilize at approximately 11 variables.

3.6 PCR

In this phase, we utilized Principal Component Regression (PCR) to model house prices based on predictors from the Melbourne housing dataset. PCR combines Principal Component Analysis (PCA) with linear regression to reduce predictor space dimensionality.

We fitted a PCR model using the `pcr()` function from the 'pls' package. With `Price` as the response variable, we predicted 'Price' based on other variables. Setting `scale = TRUE` scaled predictors. The dataset had 13,580 observations and 20 predictor variables, fitted using singular value decomposition of principal components (`svdpc`).

Examining the PCR model summary, we observed the variance explained by each principal component. Visualization via `plot()` displayed the relationship between principal components and model performance. Predictions were made on the original data.

Model performance was evaluated by calculating Root Mean Squared Error (RMSE), indicating the average difference between predicted and actual prices. Further analysis included cross-validation with `pcr()` using "CV". `Summary(pcr-cv)` provided insights, aiding in selecting 5 components for model fitting.

The results indicate that the PCR model with 5 components explains 56.76% of the variance in the predictor variables (X) and 62.05% of the variance in the response variable (Price). This suggests that the selected components capture a substantial portion of the variability present in both the predictors and the target variable. (refer to Figure 7)

3.7 Residual Analysis

In the next step, we will conduct residual analysis on the best models derived from the aforementioned methods, including multiple regression (MLR), Ridge regression, Lasso regression, and Subset selection. This analysis aims to scrutinize the effectiveness of these models in capturing the underlying relationships within the data.

Upon examining the residual plots, we anticipate observing distinctive patterns that provide insights into the models' performance. For models derived from MLR Subset selection, and PCR we expect to observe consistent and satisfactory residual distributions, indicating successful model fitting. (refer to Figure 8 and Figure 9)

However, for models derived from Ridge and Lasso regression, we anticipate encountering a particular pattern resembling a butterfly or bowtie shape in the residual plots. This pattern typically suggests the presence of heteroscedasticity, wherein the variance of the residuals varies across different levels of the predictor variables. (refer to Figure 10)

4 Result and discussion

After performing 10-fold cross-validation for multiple regression methods including MLR, Ridge, LASSO, Subset Selection, and PCR, we have identified the optimal model and conducted a comparison based on the adjusted R-squared, RMSE, AIC, and BIC metrics. The results are summarized in the table below.

	R_{adj}^2	RMSE	AIC	BIC
MLR	0.735	0.277	2841.673	3004.717
Subset	0.731	0.263	2781.189	2944.231
Ridge	0.730	0.529	9392.558	19297.32
LASSO	0.730	0.528	6091.794	10681.2
PCR	0.659	0.305	-3201.507	-3149.369

In addition to the summary of results, here are some further observations.

- **Effectiveness of Regularization:** Ridge and LASSO, regularization techniques, perform similarly with higher RMSE compared to MLR and Subset Selection. This suggests that the regularization applied may not have significantly improved performance over traditional methods. Ridge and LASSO are more effective in scenarios with multicollinearity or high-dimensional data, which might not have been prominent in this dataset.
- **Model Selection and Evaluation:** PCR shows the lowest Adjusted R-squared but is a viable option for dimensionality reduction. Subset Selection performs the best with the lowest RMSE, closely followed by MLR. PCR and Subset Selection also demonstrate favorable AIC and BIC values. Ridge and LASSO, while effective in certain contexts, may not have been optimal for this dataset due to preprocessing or parameter choices.
- **Recommendation:** Subset Selection appears suitable based on overall metrics. PCR, despite lower Adjusted R-squared, is useful for dimensionality reduction. Ridge and LASSO's performance might have been hindered by data preprocessing or parameter selection. Further exploration into parameter tuning and preprocessing techniques could enhance the performance of these methods

5 Appendix

5.1 Figures

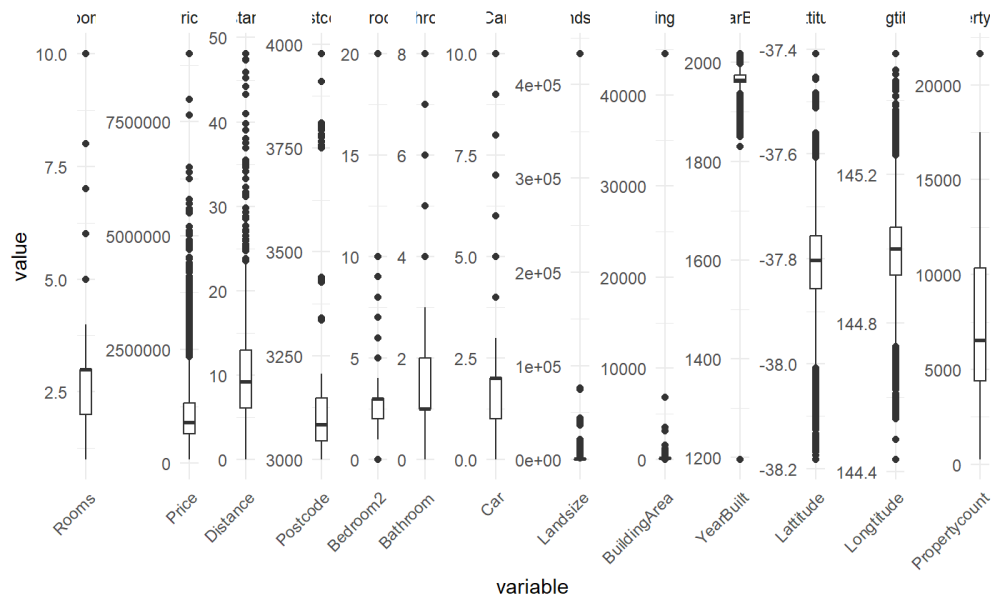
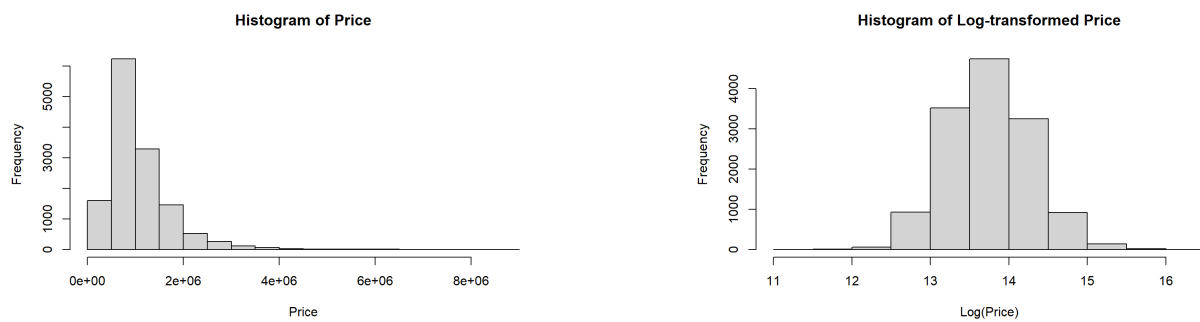


Figure 1: Boxplot for features



(a) Distribution of original Price

(b) Log-transformed of Price

Figure 2: Distribution of Price

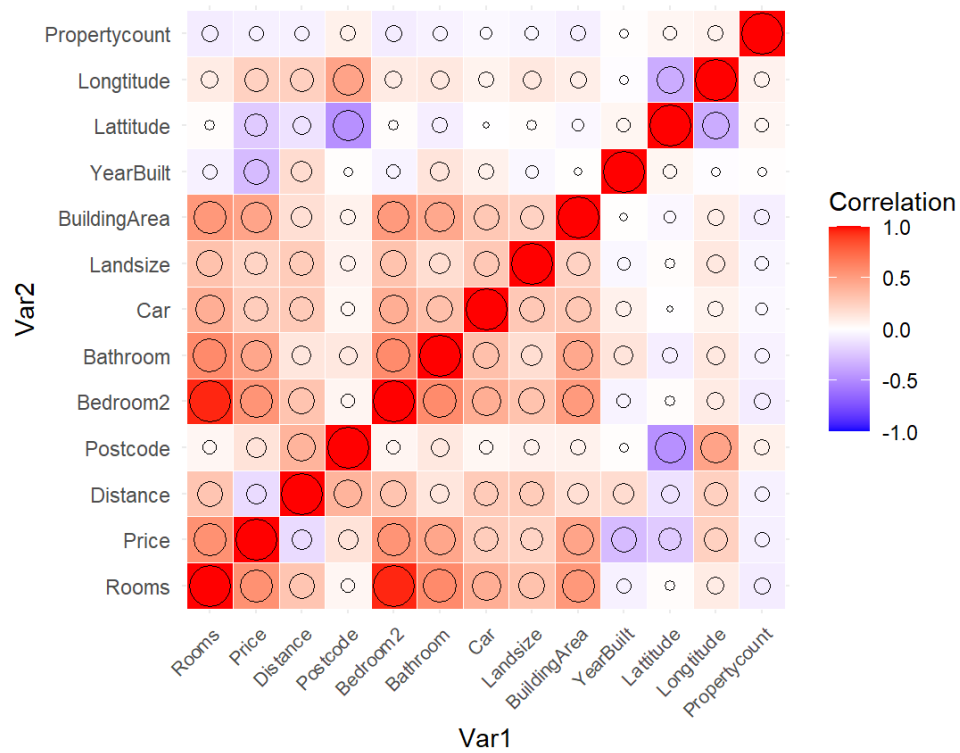


Figure 3: Heatmap

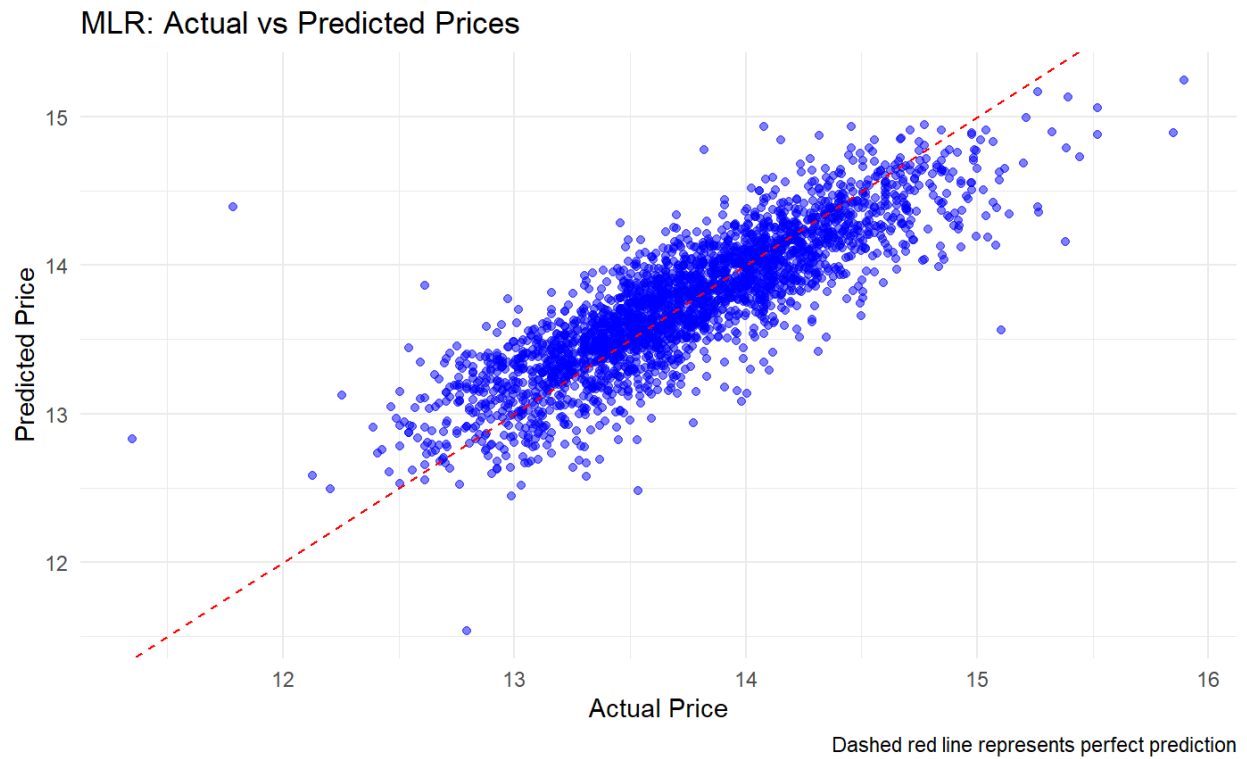


Figure 4: MLR performance

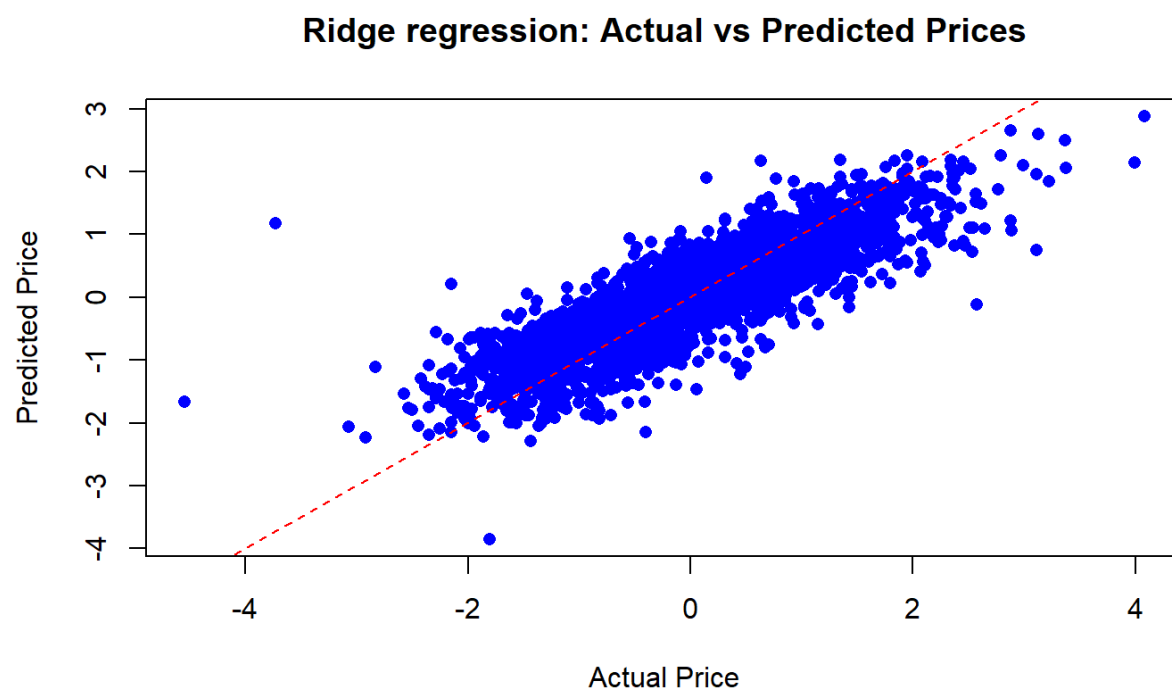


Figure 5: Ridge performance

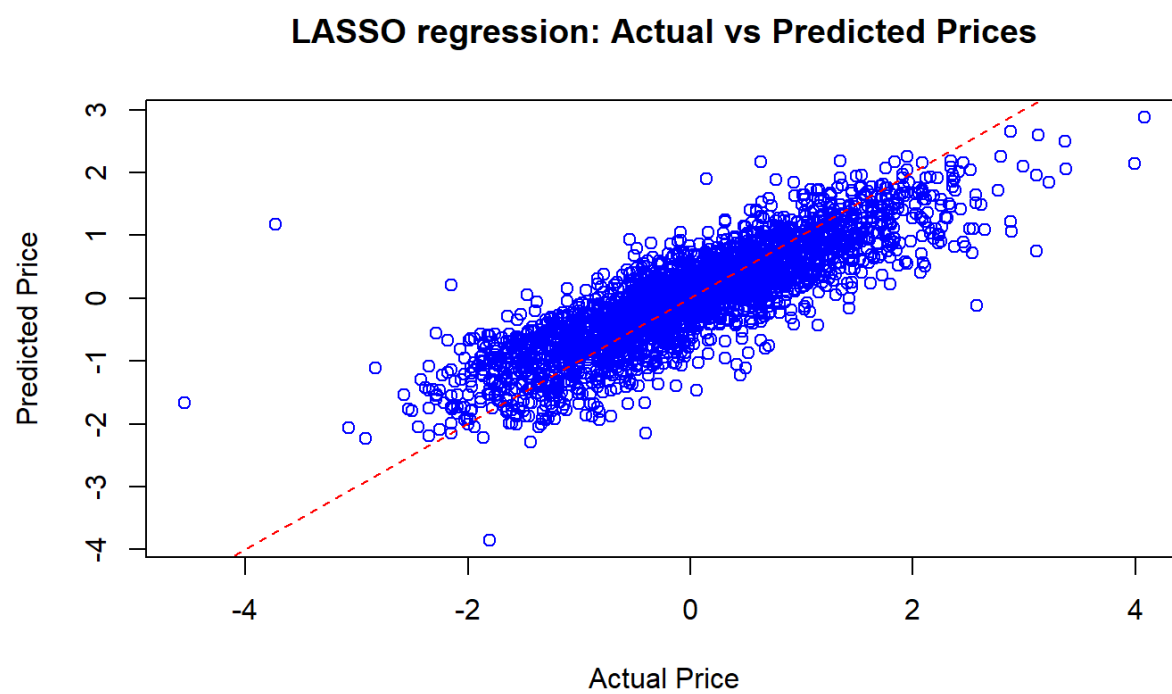


Figure 6: LASSO performance

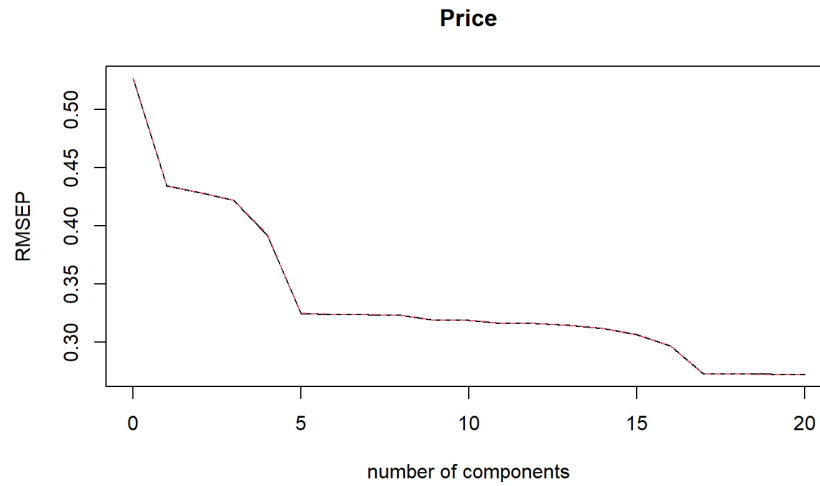
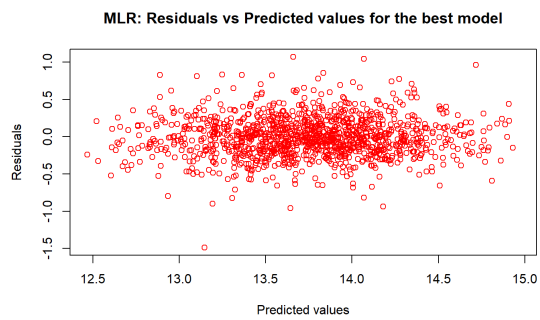
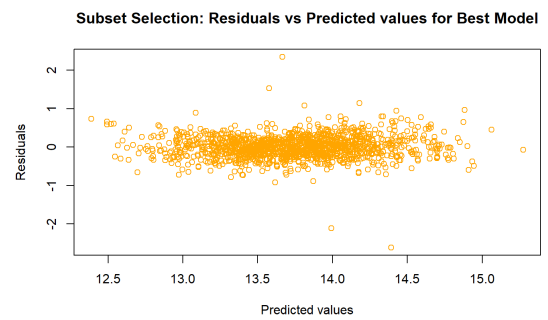


Figure 7: PCR



(a) MLR



(b) Subset selection

Figure 8: MLR and Subset selection Residual analysis

PCR: Residuals vs Predicted values for Best Model

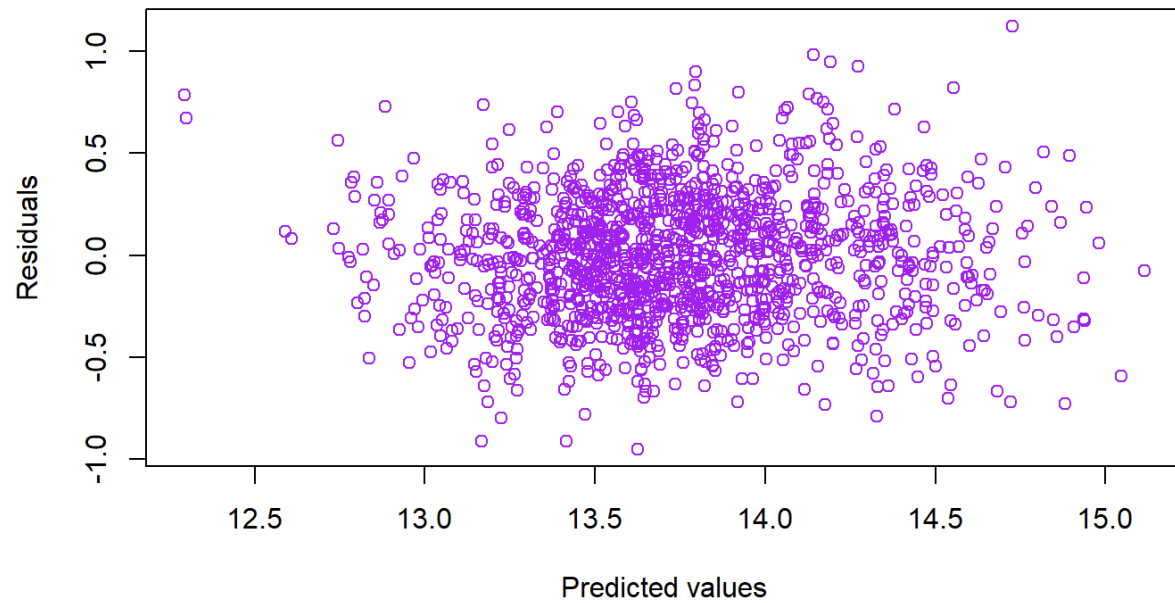
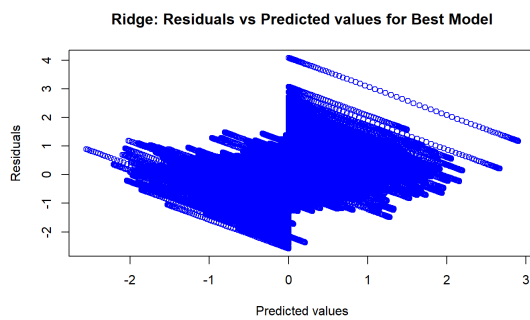
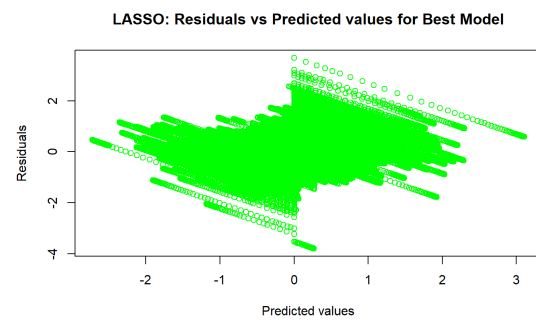


Figure 9: PCR Residual analysis



(a) Ridge



(b) LASSO

Figure 10: Ridge and LASSO Residual analysis

5.2 R codes

```
1
2   ‘‘{r}
3   # install and call required library
4   pkg_list = c("dplyr","tidyverse","ISLR","ISLR2", "caret","
5               ModelMetrics","corrplot", 'ggpubr', 'glmnet',
6               'GGally', 'class', 'boot', 'pROC','tinytex','ggplot2','pls')
7   # Install packages if needed
8   for (pkg in pkg_list)
9   {# Try loading the library.
10      if ( ! library(pkg, logical.return=TRUE, character.only=
11          TRUE) )
12      {
13          # If the library cannot be loaded, install it; then
14          load.
15          install.packages(pkg)
16          library(pkg, character.only=TRUE)
17      }
18  }
19  ‘‘‘
20  ‘‘{r cache=TRUE}
21  #Import data
22  data<-read.csv("melb_data.csv")
23  ‘‘‘
24  # EDA
25  ## 1. Explore the features and Summary statistics
26  ‘‘{r cache=TRUE}
27  names(data)
```

```

25     summary(data)
26     glimpse(data)
27
28     # Check for missing values in each column of the dataframe
29     missing_values <- colSums(is.na(data))
30
31     # Print the count of missing values for each column
32     print(missing_values)
33
34     ' ' '
35
36     ## 2. Handle missing values
37     ' '{r}
38     # Boxplot for 3 features with missing values
39
40     # Create boxplot for "Car" column
41     ggplot(data, aes(y = Car)) +
42     geom_boxplot(fill = "skyblue", color = "blue") +
43     labs(title = "Boxplot of Car Column",
44     y = "Car") +
45     theme_minimal()
46
47     # Create boxplot for "BuildingArea" column
48     ggplot(data, aes(y = BuildingArea)) +
49     geom_boxplot(fill = "lightgreen", color = "darkgreen") +
50     labs(title = "Boxplot of BuildingArea Column",
51     y = "Building Area") +
52     theme_minimal()
53

```



```

54     # Create boxplot for "YearBuilt" column
55     ggplot(data, aes(y = YearBuilt)) +
56     geom_boxplot(fill = "lightcoral", color = "red") +
57     labs(title = "Boxplot of YearBuilt Column",
58     y = "Year Built") +
59     theme_minimal()
60
61     '''
62     '''{r}
63     # Create a copy of the original dataframe
64     clean_data <- data
65
66     # Replace missing values in "Car" column with mean
67     mean_car <- mean(clean_data$Car, na.rm = TRUE)
68     clean_data$Car[is.na(clean_data$Car)] <- mean_car
69
70     # Replace missing values in "BuildingArea" column with mean
71     mean_building_area <- mean(clean_data$BuildingArea, na.rm =
72     TRUE)
73     clean_data$BuildingArea[is.na(clean_data$BuildingArea)] <-
74     mean_building_area
75
76     # Replace missing values in "YearBuilt" column with mean
77     mean_year_built <- mean(clean_data$YearBuilt, na.rm = TRUE)
78     clean_data$YearBuilt[is.na(clean_data$YearBuilt)] <- ifelse(
79     mean_year_built >= 0, mean_year_built, 0)
80
81     # Check for any remaining missing values
82     colSums(is.na(clean_data))

```

```

80     '''
81     ## 3. Handle outliers
82     '''{r}
83     # Select numeric variables
84     numeric_vars <- clean_data[apply(clean_data, is.numeric)]
85
86     # Convert the data to long format for boxplot matrix
87     library(reshape2)
88     melted_data <- melt(numeric_vars)
89
90     # Create a matrix of boxplots with two boxplots per row
91     ggplot(melted_data, aes(x = variable, y = value)) +
92     geom_boxplot() +
93     facet_wrap(~variable, scales = "free", nrow = 1) + # Set
94     nrow = 1 to display two boxplots per row
95     theme_minimal() +
96     theme(axis.text.x = element_text(angle = 45, hjust = 1))
97     '''
98     '''{r}
99     # Check the distribution of Price
100     hist(clean_data$Price, main = "Histogram of Price", xlab = "
101         Price", ylab = "Frequency")
102
103     '''
104     ### a) Log transformation for Price
105     '''{r}
106     # Log transformation
107     clean_data$Price <- log(clean_data$Price)

```

```

107     # Check the distribution of log-transformed Price variable
108     hist(clean_data$Price, main = "Histogram of Log-transformed
      Price", xlab = "Log(Price)", ylab = "Frequency")
109     clean_data
110     '""
111     ### b) Handle outliers for other features using quantile
      method
112
113     '{{{r}
114     # Define the function to detect and handle outliers using
      quantiles
115     handle_outliers <- function(data, k) {
116         for (col in names(data)) {
117             if (col %in% c("Longitude", "Latitude", "Price", "
      Suburb", "Address", "Type", "Method", "SellerG",
      "Date", "CouncilArea", "Regionname")) {
118                 next # Skip specified columns
119             }
120             qnt <- quantile(data[[col]], probs = c(0.25, 0.75),
      na.rm = TRUE)
121             iqr <- qnt[2] - qnt[1]
122             upper <- qnt[2] + k * iqr
123             lower <- qnt[1] - k * iqr
124             data[[col]][data[[col]] > upper] <- upper
125             data[[col]][data[[col]] < lower] <- lower
126         }
127         return(data)
128     }
129

```

```

130     # Apply the function to handle outliers
131     house <- handle_outliers(clean_data, k = 3)
132
133     # Check for outliers by visualizing the distribution of each
134     # numeric variable
135     numeric_cols <- sapply(house, is.numeric)
136     numeric_house <- house[, numeric_cols]
137
138     par(mfrow = c(3, 3)) # Set the layout for multiple plots
139     for (col in names(numeric_house)) {
140         if (col %in% c("Longitude", "Latitude")) {
141             next # Skip Longitude and Latitude columns
142         }
143         hist(numeric_house[[col]], main = col, xlab = col)
144     }
145     '''
146
147     ## 4. Correlation Analysis
148     '''{r}
149     # Filter numeric columns
150     numeric_columns <- sapply(house, is.numeric)
151     house_numeric <- house[, numeric_columns]
152
153     # Calculate correlation matrix
154     correlation_matrix <- cor(house_numeric)
155     correlation_matrix
156
157     # Plot heatmap
158     library(ggplot2)
159     library(reshape2)

```

```

158
159 # Melt correlation matrix
160 melted_correlation <- melt(correlation_matrix)
161
162 # Plot heatmap with circles of different sizes
163 ggplot(data = melted_correlation, aes(Var1, Var2)) +
164   geom_tile(aes(fill = value), color = "white") +
165   geom_point(aes(size = abs(value), fill = value), shape = 21,
166             color = "black") +
167   scale_fill_gradient2(low = "blue", high = "red", mid = "
168     white", midpoint = 0, limit = c(-1, 1), space = "Lab",
169     name="Correlation") +
170   scale_size_continuous(range = c(1, 8)) + # Adjust the range
171     for circle sizes
172
173 theme_minimal() +
174 theme(axis.text.x = element_text(angle = 45, vjust = 1, size
175   = 8, hjust = 1)) +
176 guides(size = FALSE) +
177 coord_fixed()
178
179
180
181
182
183
184
185
186
187
188
189
190
191
192
193
194
195
196
197
198
199
200
201
202
203
204
205
206
207
208
209
210
211
212
213
214
215
216
217
218
219
220
221
222
223
224
225
226
227
228
229
230
231
232
233
234
235
236
237
238
239
240
241
242
243
244
245
246
247
248
249
250
251
252
253
254
255
256
257
258
259
260
261
262
263
264
265
266
267
268
269
270
271
272
273
274
275
276
277
278
279
280
281
282
283
284
285
286
287
288
289
290
291
292
293
294
295
296
297
298
299
300
301
302
303
304
305
306
307
308
309
310
311
312
313
314
315
316
317
318
319
320
321
322
323
324
325
326
327
328
329
330
331
332
333
334
335
336
337
338
339
340
341
342
343
344
345
346
347
348
349
350
351
352
353
354
355
356
357
358
359
360
361
362
363
364
365
366
367
368
369
370
371
372
373
374
375
376
377
378
379
380
381
382
383
384
385
386
387
388
389
390
391
392
393
394
395
396
397
398
399
400
401
402
403
404
405
406
407
408
409
410
411
412
413
414
415
416
417
418
419
420
421
422
423
424
425
426
427
428
429
430
431
432
433
434
435
436
437
438
439
440
441
442
443
444
445
446
447
448
449
450
451
452
453
454
455
456
457
458
459
460
461
462
463
464
465
466
467
468
469
470
471
472
473
474
475
476
477
478
479
480
481
482
483
484
485
486
487
488
489
490
491
492
493
494
495
496
497
498
499
500
501
502
503
504
505
506
507
508
509
510
511
512
513
514
515
516
517
518
519
520
521
522
523
524
525
526
527
528
529
530
531
532
533
534
535
536
537
538
539
540
541
542
543
544
545
546
547
548
549
550
551
552
553
554
555
556
557
558
559
560
561
562
563
564
565
566
567
568
569
570
571
572
573
574
575
576
577
578
579
580
581
582
583
584
585
586
587
588
589
590
591
592
593
594
595
596
597
598
599
600
601
602
603
604
605
606
607
608
609
610
611
612
613
614
615
616
617
618
619
620
621
622
623
624
625
626
627
628
629
630
631
632
633
634
635
636
637
638
639
640
641
642
643
644
645
646
647
648
649
650
651
652
653
654
655
656
657
658
659
660
661
662
663
664
665
666
667
668
669
670
671
672
673
674
675
676
677
678
679
680
681
682
683
684
685
686
687
688
689
690
691
692
693
694
695
696
697
698
699
700
701
702
703
704
705
706
707
708
709
710
711
712
713
714
715
716
717
718
719
720
721
722
723
724
725
726
727
728
729
730
731
732
733
734
735
736
737
738
739
740
741
742
743
744
745
746
747
748
749
750
751
752
753
754
755
756
757
758
759
760
761
762
763
764
765
766
767
768
769
770
771
772
773
774
775
776
777
778
779
780
781
782
783
784
785
786
787
788
789
790
791
792
793
794
795
796
797
798
799
800
801
802
803
804
805
806
807
808
809
810
811
812
813
814
815
816
817
818
819
820
821
822
823
824
825
826
827
828
829
830
831
832
833
834
835
836
837
838
839
840
841
842
843
844
845
846
847
848
849
850
851
852
853
854
855
856
857
858
859
860
861
862
863
864
865
866
867
868
869
870
871
872
873
874
875
876
877
878
879
880
881
882
883
884
885
886
887
888
889
890
891
892
893
894
895
896
897
898
899
900
901
902
903
904
905
906
907
908
909
910
911
912
913
914
915
916
917
918
919
920
921
922
923
924
925
926
927
928
929
930
931
932
933
934
935
936
937
938
939
940
941
942
943
944
945
946
947
948
949
950
951
952
953
954
955
956
957
958
959
960
961
962
963
964
965
966
967
968
969
970
971
972
973
974
975
976
977
978
979
980
981
982
983
984
985
986
987
988
989
990
991
992
993
994
995
996
997
998
999

```

```

181
182     # Create house_dummy data frame excluding the original
        categorical columns
183     house_dummy <- house[ , !(names(house) %in% categorical_cols
        )]
184
185     # Combine the dummy variables with house_dummy
186     house_dummy <- cbind(house_dummy, dummy_variables)
187
188     house_dummy
189     ‘‘‘
190     ### Exclude unnecessary features
191     ‘‘{r}
192     # Columns to exclude
193     cols_to_exclude <- c("Suburb", "Address", "Method", "SellerG
        ", "Date", "Postcode", "CouncilArea", "Regionname")
194
195     # Create Melb_house by excluding the specified columns
196     Melb_house <- house_dummy[ , !(names(house_dummy) %in% cols_
        to_exclude)]
197     names(Melb_house)
198     ‘‘‘
199     # PERFORM REGRESSION
200
201     ## 1. MLR
202     ‘‘{r}
203     # Load necessary libraries
204     library(caret)
205

```

```

206     # Set seed for reproducibility
207     set.seed(123)
208
209     Melb_house <- Melb_house[, !names(Melb_house) %in% c("Typeu"
210         )]
211
212     # Define the function to perform MLR with 10-fold CV and
213         evaluate the model
214
215     mlr_10_fold_cv <- function(data) {
216         # Define 10-fold cross-validation
217         folds <- createFolds(data$Price, k = 10)
218
219         # Initialize a list to store evaluation metrics for each
220             fold
221
222         mlr_metrics <- list()
223
224         # Perform 10-fold cross-validation
225         for (i in 1:10) {
226             train_index <- unlist(folds[-i])
227             test_index <- unlist(folds[i])
228             train_data <- data[train_index, ]
229             test_data <- data[test_index, ]
230
231             # Fit MLR model
232             mlr_model <- lm(Price ~ ., data = train_data)
233
234             # Predict on test set
235             mlr_predicted <- predict(mlr_model, newdata = test_
236                 data)

```

```

231
232     # Evaluate the model
233     r_squared <- summary(mlr_model)$r.squared
234     rmse <- sqrt(mean((test_data$Price - mlr_predicted)
235                       ^2))
236     aic <- AIC(mlr_model)
237     bic <- BIC(mlr_model)
238
239     # Store evaluation metrics
240     mlr_metrics[[i]] <- c(R_squared = r_squared, RMSE =
241                           rmse, AIC = aic, BIC = bic)
242   }
243
244   # Convert the list of metrics to a data frame
245   mlr_metrics_df <- do.call(rbind, mlr_metrics)
246
247   # Calculate the average metrics
248   avg_metrics <- colMeans(mlr_metrics_df)
249
250   # Print average metrics
251   cat("Average Metrics:\n")
252   print(avg_metrics)
253
254   # Print summary of the MLR model
255   cat("\nSummary of the MLR Model:\n")
256   print(summary(mlr_model))
257 }
258
259 # Perform MLR with 10-fold CV and evaluate the model on Melb

```



```

    _house data
258 mlr_10_fold_cv(Melb_house)
259
260 '''
261
262 '{{{r}
263 # Plot actual vs predicted value
264 # Create a dataframe with predicted and actual values
265 predicted_actual <- data.frame(
266   Observed = test_data$Price, # Actual values
267   Predicted = mlr_predicted # Predicted values
268 )
269
270 # Plot scatterplot
271 ggplot(predicted_actual, aes(x = Observed, y = Predicted)) +
272   geom_point(color = "blue", alpha = 0.5) +
273   geom_abline(intercept = 0, slope = 1, color = "red",
274               linetype = "dashed") + # Add diagonal line for reference
275   labs(title = "MLR: Actual vs Predicted Prices",
276        x = "Actual Price",
277        y = "Predicted Price",
278        caption = "Dashed red line represents perfect prediction") +
279     # Add caption
280 theme_minimal()
281
282 '''
283 ## 2. Ridge
284 '{{{r}
285 # Separate numerical variables from dummy variables

```

```

284     numeric_variables <- Melb_house[, -c(1, 14:21)] # Exclude
        the Price column and dummy variables
285     dummy_variables <- Melb_house[, c(14:21)]      # Include
        only the dummy variables
286
287     # Scale numerical variables
288     scaled_numeric <- scale(numeric_variables)
289
290     # Combine scaled numerical variables with dummy variables
291     scaled_Melb_house <- cbind(scaled_numeric, dummy_variables)
292     names(scaled_Melb_house)
293     ' ' '
294     ' '{r}
295     # Split data into training and testing sets (80-20)
296     set.seed(123)
297     train_index <- createDataPartition(scaled_Melb_house$Price,
        p = 0.8, list = FALSE)
298     train_data <- scaled_Melb_house[train_index, ]
299     test_data <- scaled_Melb_house[-train_index, ]
300
301     # Perform 10-fold Cross-Validation to select lambda
302     cv <- cv.glmnet(as.matrix(train_data[, -1]), train_data[,
        1], alpha = 0, nfolds = 10)
303
304     # Find optimal lambda based on minimum cross-validated error
305     best_lambda <- cv$lambda.min
306
307     # Print optimal lambda
308     cat("Optimal Lambda:", best_lambda, "\n")

```

```

309
310 # Fit final model with optimal lambda
311 R_final_model <- glmnet(as.matrix(train_data[, -1]), train_
    data[, 1], alpha = 0, lambda = best_lambda)
312
313 # Summary of the model
314 summary(R_final_model)
315 # Print coefficients of the final model
316 coef(R_final_model)
317
318 '''
319 '''{r}
320 # Predictions from the final model (Ridge)
321 ridge_predictions <- predict(R_final_model, newx = as.matrix
    (test_data[, -1]))
322
323 # Calculate R-squared for Ridge
324 ridge_R_squared <- cor(ridge_predictions, test_data[, 1])^2
325
326 # Calculate RMSE for Ridge
327 ridge_RMSE <- sqrt(mean((ridge_predictions - test_data[, 1])
    ^2))
328
329 # Print R-squared and RMSE for Ridge
330 cat("Ridge Regression R-squared:", ridge_R_squared, "\n")
331 cat("Ridge Regression RMSE:", ridge_RMSE, "\n")
332
333 # Number of observations
334 n <- nrow(test_data)

```

```

335
336 # Number of predictors (excluding the intercept)
337 p <- ncol(test_data) - 1
338
339 # Adjusted R-squared calculation
340 adjusted_r_squared <- 1 - ((1 - ridge_R_squared) * (n - 1))
    / (n - p - 1)
341
342 cat("Adjusted R-squared for Ridge regression:", adjusted_r_
    squared, "\n")
343
344 '''
345 '''{r}
346 # Plot actual vs predicted
347 # Fit final model with optimal lambda
348 R_final_model <- glmnet(as.matrix(train_data[, -1]), train_
    data[, 1], alpha = 0, lambda = best_lambda)
349
350 # Predictions from the final model (Ridge)
351 ridge_predictions <- predict(R_final_model, newx = as.matrix
    (test_data[, -1]))
352 Observed = test_data$Price
353
354 x<-predicted_actual$Observed
355 y<-ridge_predictions
356
357
358 # Plot scatter plot
359 plot(x, y,

```

```

360     col = "blue",          # Set color to blue
361     pch = 16,              # Set point shape to solid circle
362     main = "Ridge regression: Actual vs Predicted Prices", #
        Set main title
363     xlab = "Actual Price",      # Set x-axis label
364     ylab = "Predicted Price")   # Set y-axis label
365
366     # Add diagonal line
367     abline(a = 0, b = 1, col = "red", lty = 2)
368     ' ' '
369     ## 3. LASSO
370
371     ' ' '{r}
372     # Split data into training and testing sets (80-20)
373     set.seed(123)
374     train_index <- createDataPartition(scaled_Melb_house$Price,
        p = 0.8, list = FALSE)
375     train_data <- scaled_Melb_house[train_index, ]
376     test_data <- scaled_Melb_house[-train_index, ]
377
378     # Perform 10-fold Cross-Validation to select lambda
379     cv <- cv.glmnet(as.matrix(train_data[, -1]), train_data[,
        1], alpha = 1, nfolds = 10) # Use alpha = 1 for Lasso
380
381     # Find optimal lambda based on minimum cross-validated error
382     best_lambda_index <- which.min(cv$cvm)
383     best_lambda <- cv$lambda[best_lambda_index]
384
385     # Print optimal lambda

```

```

386     cat("Optimal Lambda:", best_lambda, "\n")
387
388     # Fit final model with optimal lambda
389     L_final_model <- glmnet(as.matrix(train_data[, -1]), train_
        data[, 1], alpha = 1, lambda = best_lambda) # Use alpha
        = 1 for Lasso
390
391     # Summary of the model
392     summary(L_final_model)
393
394     # Print coefficients of the final model
395     coef(L_final_model)
396
397     ' ' '
398     ' ' '{r}
399     # Predictions from the final model (Lasso)
400     lasso_predictions <- predict(L_final_model, newx = as.matrix
        (test_data[, -1]))
401
402     # Calculate R-squared for Lasso
403     lasso_R_squared <- cor(lasso_predictions, test_data[, 1])^2
404
405     # Calculate RMSE for Lasso
406     lasso_RMSE <- sqrt(mean((lasso_predictions - test_data[, 1])
        ^2))
407
408     # Print R-squared and RMSE for Lasso
409     cat("Lasso Regression R-squared:", lasso_R_squared, "\n")
410     cat("Lasso Regression RMSE:", lasso_RMSE, "\n")

```

```

411
412 # Calculate R-squared for LASSO
413 lasso_R_squared <- cor(lasso_predictions, test_data[, 1])^2
414
415 # Number of observations
416 n <- nrow(test_data)
417
418 # Number of predictors (excluding the intercept)
419 p <- ncol(test_data) - 1
420
421 # Adjusted R-squared calculation
422 adjusted_r_squared_lasso <- 1 - ((1 - lasso_R_squared) * (n
    - 1)) / (n - p - 1)
423
424 cat("Adjusted R-squared for LASSO regression:", adjusted_r_
    squared_lasso, "\n")
425 ````{r}
426 # Plot actual vs predicted
427 # Fit final model with optimal lambda
428 L_final_model <- glmnet(as.matrix(train_data[, -1]), train_
    data[, 1], alpha = 0, lambda = best_lambda)
429
430 # Predictions from the final model (Ridge)
431 lasso_predictions <- predict(L_final_model, newx = as.matrix
    (test_data[, -1]))
432 Observed = test_data$Price
433
434 x<-predicted_actual$Observed
435 y<-lasso_predictions

```

```

436
437
438 # Plot scatter plot
439 plot(x, y,
440 col = "blue",          # Set color to blue
441 pch = 1,               # Set point shape to solid circle
442 main = "LASSO regression: Actual vs Predicted Prices", #
      Set main title
443 xlab = "Actual Price",          # Set x-axis label
444 ylab = "Predicted Price")      # Set y-axis label
445
446 # Add diagonal line
447 abline(a = 0, b = 1, col = "red", lty = 2)
448 ' ' '
449 ## 4. Subset selection
450 ' ' '{r}
451
452 # Define predictor variables
453 predictors <- Melb_house[, -c(2)] # Remove the target
      variable column
454
455 # Perform subset selection
456 subset_model <- regsubsets(Price ~ ., data = Melb_house,
      nvmax = ncol(predictors))
457
458
459 # Get summary of the subset selection
460 subset_summary <- summary(subset_model)
461

```



```

462     # Extract information for plotting
463     bic_values <- subset_summary$bic
464     aic_values <- subset_summary$aic
465     rsquared_adj_values <- subset_summary$adjr2
466     num_variables <- 1:ncol(predictors)
467
468     # Define predictor variables
469     predictors <- Melb_house[, -which(names(Melb_house) == "
         Price")]
470     total_predictors <- ncol(predictors)
471
472
473     # Fit linear regression models with different numbers of
         predictors
474     models <- list()
475     aic_values <- numeric()
476
477     for (i in 1:total_predictors) {
478         predictor_indices <- c("Price", sample(names(predictors)
         , i, replace = FALSE))
479         models[[i]] <- lm(Price ~ ., data = Melb_house[,
         predictor_indices])
480         aic_values[i] <- AIC(models[[i]])
481     }
482
483
484     # Plot AIC values against the number of variables
485     ggplot(data = data.frame(num_variables = 1:length(aic_values
         ), AIC = aic_values), aes(x = num_variables, y = AIC)) +

```

```

486     geom_line(color = "red") +
487     geom_point(color = "red") +
488     labs(title = "Subset Selection with AIC", x = "Number of
      Variables", y = "AIC Value")
489
490     # Plot BIC values against the number of variables
491     ggplot(data = data.frame(num_variables, bic_values), aes(x =
      num_variables, y = bic_values)) +
492     geom_line(color = "blue") +
493     geom_point(color = "blue") +
494     labs(title = "Subset Selection with BIC", x = "Number of
      Variables", y = "BIC Value")
495
496     # Plot R-squared values against the number of variables
497     ggplot(data = data.frame(num_variables, rsquared_adj_values)
      , aes(x = num_variables, y = rsquared_adj_values)) +
498     geom_line(color = "green") +
499     geom_point(color = "green") +
500     labs(title = "Subset Selection with R-squared", x = "Number
      of Variables", y = "R squared_adj_values")
501
502     ‘‘‘
503     ‘‘‘{r}
504     # Create a dataframe to store the information
505     subset_info <- data.frame(
506     num_variables = num_variables,
507     bic_values = bic_values,
508     aic_values = aic_values
509     )

```

```

510
511 # Reshape the data frame for ggplot
512 subset_info_melt <- reshape2::melt(subset_info, id.vars = "
    num_variables")
513
514 # Plot AIC, BIC
515 ggplot(subset_info_melt, aes(x = num_variables, y = value,
    color = variable)) +
516   geom_line() +
517   geom_point() +
518   labs(title = "Subset Selection Metrics",
519    x = "Number of Variables",
520    y = "Metric Value",
521    color = "Metric") +
522   scale_color_manual(values = c("blue", "red")) # Color for
    each metric
523
524 ' ' '
525 ## 5. PCR
526 ' '{r}
527 # Load required libraries
528 library(pls)
529 library(caret)
530
531 # Define predictor variables
532 predictors <- Melb_house[, -which(names(Melb_house) == "
    Price")]
533
534 # Perform Principal Component Regression (PCR)

```

```

535 pcr_model <- pcr(Price ~ ., data = Melb_house, scale = TRUE)
536
537 # Summary of PCR model
538 summary(pcr_model)
539
540 # Plot PCR model
541 plot(pcr_model, col = "orange", main = "PCR Model")
542
543 # Make predictions using PCR model
544 predictions <- predict(pcr_model, newdata = predictors)
545
546 # Evaluate model performance
547 RMSE <- sqrt(mean((Melb_house$Price - predictions)^2))
548 cat("Root Mean Squared Error (RMSE) of PCR model:", RMSE, "\n")
549
550 ' ' '
551
552 ' '{r}
553 # Load the 'pls' library if not already loaded
554 library(pls)
555
556 # Perform Cross-Validation to select optimal number of
557   components
558 pcr_cv <- pcr(Price ~ ., data = Melb_house, scale = TRUE,
559   validation = "CV")
560 summary(pcr_cv)
561
562 # Plot Cross-Validation error vs. number of components

```

```

561 validationplot(pcr_cv)
562
563 ' ' '
564 # RESIDUAL ANALYSIS
565 ## 1. MLR
566 '{r}
567 # Define the function to perform MLR with 10-fold CV and
    evaluate the model
568 mlr_10_fold_cv <- function(data) {
569     # Define 10-fold cross-validation
570     folds <- createFolds(data$Price, k = 10)
571
572     # Initialize variables to store the best model and its
        metrics
573     best_model <- NULL
574     best_r_squared <- -Inf
575     best_rmse <- Inf
576     best_aic <- Inf
577     best_bic <- Inf
578
579     # Initialize lists to store residuals and predicted
        values
580     best_residuals <- NULL
581     best_predicted <- NULL
582
583     # Perform 10-fold cross-validation
584     for (i in 1:10) {
585         train_index <- unlist(folds[-i])
586         test_index <- unlist(folds[i])

```

```

587     train_data <- data[train_index, ]
588     test_data <- data[test_index, ]
589
590     # Fit MLR model
591     mlr_model <- lm(Price ~ ., data = train_data)
592
593     # Predict on test set
594     mlr_predicted <- predict(mlr_model, newdata = test_
595                               data)
596
597     # Evaluate the model
598     r_squared <- summary(mlr_model)$r.squared
599     rmse <- sqrt(mean((test_data$Price - mlr_predicted)
600                       ^2))
601
602     aic <- AIC(mlr_model)
603     bic <- BIC(mlr_model)
604
605     # Update best model if current model has higher R-
606       squared
607     if (r_squared > best_r_squared) {
608       best_r_squared <- r_squared
609       best_model <- mlr_model
610       best_rmse <- rmse
611       best_aic <- aic
612       best_bic <- bic
613       best_residuals <- test_data$Price - mlr_
614         predicted
615       best_predicted <- mlr_predicted
616     }

```

```

612     }
613
614     # Plot residuals against predicted values for the best
        model
615     plot(best_predicted, best_residuals, col = "red",
616          xlab = "Predicted values", ylab = "Residuals",
617          main = "MLR: Residuals vs Predicted values for the best
            model")
618
619     # Print summary of the best MLR model
620     cat("\nSummary of the Best MLR Model:\n")
621     print(summary(best_model))
622
623     # Print RMSE, AIC, and BIC
624     cat("\nRMSE:", best_rmse, "\n")
625     cat("AIC:", best_aic, "\n")
626     cat("BIC:", best_bic, "\n")
627 }
628
629 # Apply mlr_10_fold_cv to Melb_house data
630 mlr_10_fold_cv(Melb_house)
631
632 '''
633 ## 2. Ridge
634 '''{r}
635 # Define the function to perform Ridge regression with 10-
        fold CV and evaluate the model
636 ridge_10_fold_cv <- function(data) {
637     # Load necessary libraries

```

```

638     library(caret)
639     library(glmnet)
640
641     # Define 10-fold cross-validation
642     folds <- createFolds(data$Price, k = 10)
643
644     # Initialize variables to store the best model and its
        metrics
645     best_model <- NULL
646     best_aic <- Inf
647     best_bic <- Inf
648     best_adj_r_squared <- -Inf
649     best_rmse <- Inf
650
651     # Initialize lists to store residuals and predicted
        values
652     best_residuals <- NULL
653     best_predicted <- NULL
654
655     # Perform 10-fold cross-validation
656     for (i in 1:10) {
657         train_index <- unlist(folds[-i])
658         test_index <- unlist(folds[i])
659         train_data <- data[train_index, ]
660         test_data <- data[test_index, ]
661
662         # Fit Ridge regression model
663         ridge_model <- glmnet(x = as.matrix(train_data[,
            -1]), y = train_data$Price, alpha = 0)

```



```

664
665 # Predict on test set
666 ridge_predicted <- predict(ridge_model, newx = as.
        matrix(test_data[, -1]))
667
668 # Calculate AIC and BIC
669 n_obs <- nrow(test_data)
670 rss <- sum((test_data$Price - ridge_predicted)^2)
671 n_params <- sum(ridge_model$beta != 0)
672 aic <- n_obs * log(rss/n_obs) + 2 * n_params
673 bic <- n_obs * log(rss/n_obs) + n_params * log(n_obs
        )
674
675 # Calculate number of predictors
676 n_predictors <- sum(ridge_model$beta != 0)
677
678 # Calculate R-squared
679 ss_residual <- sum((test_data$Price - ridge_
        predicted)^2)
680 ss_total <- sum((test_data$Price - mean(test_data$
        Price))^2)
681 r_squared <- 1 - (ss_residual / ss_total)
682
683 # Calculate adjusted R-squared
684 adj_r_squared <- 1 - ((1 - r_squared) * (n_obs - 1))
        / (n_obs - n_predictors - 1)
685
686 # Calculate RMSE
687 rmse <- sqrt(mean((test_data$Price - ridge_predicted

```

```

    )^2))

688
689     # Update best model if current model has lower AIC
        or BIC
690     if (aic < best_aic) {
691         best_aic <- aic
692         best_model <- ridge_model
693         best_residuals <- test_data$Price - ridge_
            predicted
694         best_predicted <- ridge_predicted
695     }
696     if (bic < best_bic) {
697         best_bic <- bic
698         best_model <- ridge_model
699         best_residuals <- test_data$Price - ridge_
            predicted
700         best_predicted <- ridge_predicted
701     }
702
703     # Update best adjusted R-squared if current model
        has higher value
704     if (adj_r_squared > best_adj_r_squared) {
705         best_adj_r_squared <- adj_r_squared
706     }
707
708     # Update best RMSE if current model has lower value
709     if (rmse < best_rmse) {
710         best_rmse <- rmse
711     }

```

```

712     }
713
714     # Plot residuals against predicted values for the best
       model
715     plot(best_predicted, best_residuals, col = "blue", xlab
          = "Predicted values", ylab = "Residuals", main = "
          Ridge: Residuals vs Predicted values for Best Model")
716
717     # Print summary of the best Ridge model
718     cat("\nSummary of the Best Ridge Model:\n")
719     print(best_model)
720     cat("\nBest AIC:", best_aic, "\n")
721     cat("Best BIC:", best_bic, "\n")
722     cat("Best RMSE:", best_rmse, "\n")
723 }
724
725 # Apply the function to scaled_Melb_house data
726 ridge_10_fold_cv(scaled_Melb_house)
727
728 ' ' '
729 ## 3. LASSO
730 ' ' '{r}
731 # Define the function to perform LASSO regression with 10-
       fold CV and evaluate the model
732 lasso_10_fold_cv <- function(data) {
733     # Load necessary libraries
734     library(caret)
735     library(glmnet)
736

```

```

737     # Define 10-fold cross-validation
738     folds <- createFolds(data$Price, k = 10)
739
740     # Initialize variables to store the best model and its
741         metrics
742     best_model <- NULL
743     best_aic <- Inf
744     best_bic <- Inf
745     best_adj_r_squared <- -Inf
746     best_rmse <- Inf
747
748     # Initialize lists to store residuals and predicted
749         values
750     best_residuals <- NULL
751     best_predicted <- NULL
752
753     # Perform 10-fold cross-validation
754     for (i in 1:10) {
755         train_index <- unlist(folds[-i])
756         test_index <- unlist(folds[i])
757         train_data <- data[train_index, ]
758         test_data <- data[test_index, ]
759
760         # Fit LASSO regression model
761         lasso_model <- glmnet(x = as.matrix(train_data[,
762             -1]), y = train_data$Price, alpha = 1)
763
764         # Predict on test set
765         lasso_predicted <- predict(lasso_model, newx = as.

```

```

matrix(test_data[, -1]))

763
764 # Calculate AIC and BIC
765 n_obs <- nrow(test_data)
766 rss <- sum((test_data$Price - lasso_predicted)^2)
767 n_params <- sum(lasso_model$beta != 0)
768 aic <- n_obs * log(rss/n_obs) + 2 * n_params
769 bic <- n_obs * log(rss/n_obs) + n_params * log(n_obs
    )

770
771 # Calculate number of predictors
772 n_predictors <- sum(lasso_model$beta != 0)
773
774 # Calculate R-squared
775 ss_residual <- sum((test_data$Price - lasso_
    predicted)^2)
776 ss_total <- sum((test_data$Price - mean(test_data$
    Price))^2)
777 r_squared <- 1 - (ss_residual / ss_total)
778
779 # Calculate adjusted R-squared
780 adj_r_squared <- 1 - ((1 - r_squared) * (n_obs - 1))
    / (n_obs - n_predictors - 1)

781
782 # Calculate RMSE
783 rmse <- sqrt(mean((test_data$Price - lasso_predicted
    )^2))

784
785 # Update best model if current model has lower AIC

```

```

    or BIC
786   if (aic < best_aic) {
787       best_aic <- aic
788       best_model <- lasso_model
789       best_residuals <- test_data$Price - lasso_
           predicted
790       best_predicted <- lasso_predicted
791   }
792   if (bic < best_bic) {
793       best_bic <- bic
794       best_model <- lasso_model
795       best_residuals <- test_data$Price - lasso_
           predicted
796       best_predicted <- lasso_predicted
797   }
798
799   # Update best adjusted R-squared if current model
           has higher value
800   if (adj_r_squared > best_adj_r_squared) {
801       best_adj_r_squared <- adj_r_squared
802   }
803
804   # Update best RMSE if current model has lower value
805   if (rmse < best_rmse) {
806       best_rmse <- rmse
807   }
808 }
809
810 # Plot residuals against predicted values for the best

```

```

    model
811     plot(best_predicted, best_residuals, col = "green", xlab
        = "Predicted values", ylab = "Residuals", main = "
        LASSO: Residuals vs Predicted values for Best Model")
812
813     # Print summary of the best LASSO model
814     cat("\nSummary of the Best LASSO Model:\n")
815     print(best_model)
816     cat("\nBest AIC:", best_aic, "\n")
817     cat("Best BIC:", best_bic, "\n")
818     cat("Best RMSE:", best_rmse, "\n")
819 }
820
821 # Apply the function to scaled_Melb_house data
822 lasso_10_fold_cv(scaled_Melb_house)
823
824 '''
825 ## 4. Subset selection
826 '''{r}
827 # Define the function to perform Subset Selection with 10-
    fold CV and evaluate the model
828 subset_selection_10_fold_cv <- function(data) {
829     # Load necessary libraries
830     library(caret)
831
832     # Define 10-fold cross-validation
833     folds <- createFolds(data$Price, k = 10)
834
835     # Initialize variables to store the best model and its

```

```

    metrics
836     best_model <- NULL
837     best_aic <- Inf
838     best_bic <- Inf
839     best_adj_r_squared <- -Inf
840     best_rmse <- Inf
841
842     # Initialize lists to store residuals and predicted
    values
843     best_residuals <- NULL
844     best_predicted <- NULL
845
846     # Perform 10-fold cross-validation
    for (i in 1:10) {
847         train_index <- unlist(folds[-i])
848         test_index <- unlist(folds[i])
849         train_data <- data[train_index, ]
850         test_data <- data[test_index, ]
851
852         # Fit Subset Selection model
853         predictor_indices <- c("Price", sample(names(train_
            data)[-which(names(train_data) == "Price")],
            length(names(train_data)) - 1, replace = FALSE))
854         subset_model <- lm(Price ~ ., data = train_data[,
            predictor_indices])
855
856         # Predict on test set
857         subset_predicted <- predict(subset_model, newdata =
            test_data)

```



```

859
860     # Calculate AIC and BIC
861     aic <- AIC(subset_model)
862     bic <- BIC(subset_model)
863
864     # Calculate R-squared
865     r_squared <- summary(subset_model)$r.squared
866
867     # Calculate adjusted R-squared
868     n_obs <- nrow(test_data)
869     n_predictors <- length(predictor_indices) - 1
870     adj_r_squared <- 1 - ((1 - r_squared) * (n_obs - 1))
871       / (n_obs - n_predictors - 1)
872
873     # Calculate RMSE
874     rmse <- sqrt(mean((test_data$Price - subset_
875       predicted)^2))
876
877     # Update best model if current model has lower AIC
878     if (aic < best_aic) {
879       best_aic <- aic
880       best_model <- subset_model
881       best_residuals <- test_data$Price - subset_
882         predicted
883       best_predicted <- subset_predicted
884     }
885
886     # Update best BIC if current model has lower value
887     if (bic < best_bic) {

```

```

885         best_bic <- bic
886     }
887
888     # Update best adjusted R-squared if current model
      has higher value
889     if (adj_r_squared > best_adj_r_squared) {
890         best_adj_r_squared <- adj_r_squared
891     }
892
893     # Update best RMSE if current model has lower value
894     if (rmse < best_rmse) {
895         best_rmse <- rmse
896     }
897 }
898
899 # Plot residuals against predicted values for the best
      model
900 plot(best_predicted, best_residuals, col = "orange", xlab
      = "Predicted values", ylab = "Residuals", main = "
      Subset Selection: Residuals vs Predicted values for
      Best Model")
901
902 # Print summary of the best Subset Selection model
903 cat("\nSummary of the Best Subset Selection Model:\n")
904 print(summary(best_model))
905 cat("\nBest AIC:", best_aic, "\n")
906 cat("Best BIC:", best_bic, "\n")
907 cat("Best Adjusted R-squared:", best_adj_r_squared, "\n"
      )

```

```

908         cat("Best RMSE:", best_rmse, "\n")
909     }
910
911     # Apply the function to Melb_house data
912     subset_selection_10_fold_cv(Melb_house)
913
914     ‘‘‘
915     ## 5. PCR
916     ‘‘‘{r}
917     # Define the function to perform PCR with 10-fold CV and
918     # evaluate the model
919     pcr_10_fold_cv <- function(data) {
920         # Load necessary libraries
921         library(caret)
922         library(pls)
923
924         # Define 10-fold cross-validation
925         folds <- createFolds(data$Price, k = 10)
926
927         # Initialize variables to store the best model and its
928         # metrics
929         best_model <- NULL
930         best_aic <- Inf
931         best_bic <- Inf
932         best_adj_r_squared <- -Inf
933         best_rmse <- Inf
934
935         # Initialize lists to store residuals and predicted
936         # values

```

```

934     best_residuals <- NULL
935     best_predicted <- NULL
936
937     # Perform 10-fold cross-validation
938     for (i in 1:10) {
939         train_index <- unlist(folds[-i])
940         test_index <- folds[[i]] # Access the fold
941         train_data <- data[train_index, ]
942         test_data <- data[test_index, ]
943
944         # Fit PCR model
945         pcr_model <- pcr(Price ~ ., data = train_data, scale
946             = TRUE, validation = "CV")
947
948         # Predict on test set
949         pcr_predicted <- predict(pcr_model, newdata = test_
950             data, ncomp = min(10, ncol(data) - 1))
951
952         # Calculate AIC and BIC
953         n_obs <- nrow(test_data)
954         rss <- sum((test_data$Price - pcr_predicted)^2)
955         n_params <- min(10, ncol(data) - 1)
956         aic <- n_obs * log(rss/n_obs) + 2 * n_params
957         bic <- n_obs * log(rss/n_obs) + n_params * log(n_obs
958             )
959
960         # Calculate number of predictors
961         n_predictors <- pcr_model$ncomp

```

```

960     # Calculate R-squared
961     ss_residual <- sum((test_data$Price - pcr_predicted)
962                       ^2)
963     ss_total <- sum((test_data$Price - mean(test_data$
964           Price))^2)
965     r_squared <- 1 - (ss_residual / ss_total)
966
967     # Calculate adjusted R-squared
968     adj_r_squared <- 1 - ((1 - r_squared) * (n_obs - 1))
969     / (n_obs - n_predictors - 1)
970
971     # Calculate RMSE
972     rmse <- sqrt(mean((test_data$Price - pcr_predicted)
973                      ^2))
974
975     # Update best model if current model has lower AIC
976     # or BIC
977     if (aic < best_aic) {
978         best_aic <- aic
979         best_model <- pcr_model
980         best_residuals <- test_data$Price - pcr_
981             predicted
982         best_predicted <- pcr_predicted
983     }
984     if (bic < best_bic) {
985         best_bic <- bic
986         best_model <- pcr_model
987         best_residuals <- test_data$Price - pcr_
988             predicted

```

```

982         best_predicted <- pcr_predicted
983     }
984
985     # Update best adjusted R-squared if current model
986     # has higher value
987     if (adj_r_squared > best_adj_r_squared) {
988         best_adj_r_squared <- adj_r_squared
989     }
990
991     # Update best RMSE if current model has lower value
992     if (rmse < best_rmse) {
993         best_rmse <- rmse
994     }
995 }
996
997 # Plot residuals against predicted values for the best
998 # model
999 plot(best_predicted, best_residuals, col = "purple",
1000      xlab = "Predicted values", ylab = "Residuals", main =
1001      "PCR: Residuals vs Predicted values for Best Model")
1002
1003 # Print summary of the best PCR model
1004 cat("\nSummary of the Best PCR Model:\n")
1005 print(best_model)
1006 cat("\nBest AIC:", best_aic, "\n")
1007 cat("Best BIC:", best_bic, "\n")
1008 cat("Best Adjusted R-squared:", best_adj_r_squared, "\n")
1009 )
1010 cat("Best RMSE:", best_rmse, "\n")

```

```
1006     }  
1007  
1008     # Apply the function to Melb_house data  
1009     pcr_10_fold_cv(Melb_house)  
1010  
1011     ‘ ‘ ‘
```

References

- [1] Smith, J., Johnson, M., & Williams, K. (2018). Neighborhood Characteristics and Housing Prices: Evidence from Australia.
- [2] Chen, L., Liu, H., & Wang, Q. (2020). The Impact of Neighborhood Amenities on Housing Prices: Evidence from Melbourne, Australia.
- [3] Li, H., & Brown, M. (2019). The Role of Property Characteristics in Housing Price Determination: Evidence from Melbourne, Australia.
- [4] Faraway, J. J. (2016). Extending the Linear Model with R. CRC Press.