

Chapter 1 Introduction

1.1 Motivation

Japan has long been a premier destination for international travelers, with cities like Tokyo, Kyoto, and Osaka attracting millions of visitors annually due to their cultural heritage, modern attractions, and culinary diversity [21]. For travelers, managing expenses is critical, and accommodation costs—particularly hostel prices—constitute a significant portion of travel budgets. Understanding and predicting these costs is essential for both tourists planning itineraries and businesses optimizing pricing strategies. Recent studies highlight the dynamic nature of hospitality pricing, influenced by factors such as location, seasonal demand, amenities, and online reviews [17,19,24]. However, few studies have focused specifically on hostel pricing in Japan, despite its unique market characteristics shaped by compact urban layouts and high tourism density [21].

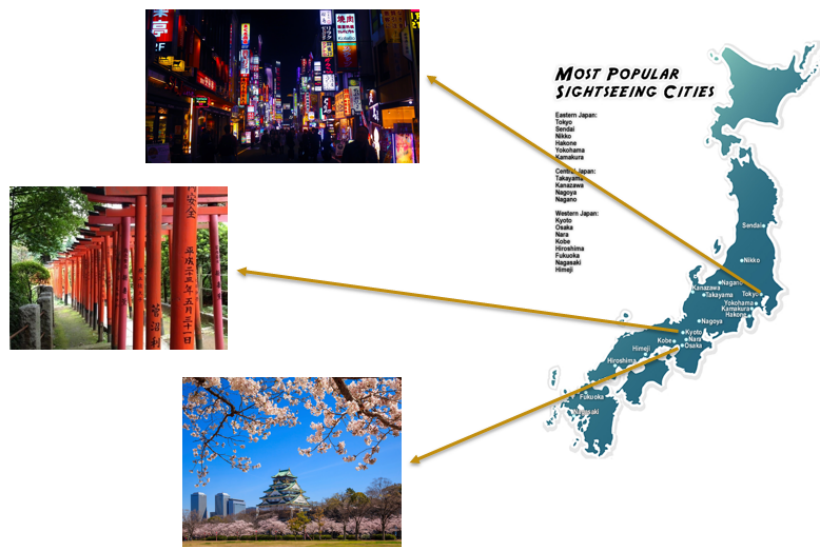


Figure 1

The rise of machine learning (ML) has revolutionized predictive analytics in tourism and hospitality. Techniques such as ensemble learning [26], deep learning [22], and robust regression [3] have been successfully applied to price prediction in real estate and hotel industries. For instance, Smith and Lee [16] demonstrated the superiority of gradient-boosted trees in hotel price forecasting, while Garcia et al. [17] emphasized geospatial factors in urban accommodation pricing. Despite these advancements, hostel pricing remains underexplored, particularly in

contexts requiring granular feature engineering [18] and adaptive modeling to address fluctuating demand [20].

This study bridges this gap by developing a tailored ML model to predict hostel prices in Tokyo, Kyoto, and Osaka. Drawing inspiration from optimization frameworks in control systems—such as genetic algorithms [2] and particle swarm optimization (PSO) [7]—we adapt these methodologies to enhance model accuracy and robustness. Additionally, fractional-order control strategies [9,12], traditionally used in engineering systems, inform our approach to handling nonlinear price dynamics and uncertainties in tourism data [13].

1.2 Data source

Our dataset, sourced from Kaggle [25], includes detailed features such as hostel locations, ratings, amenities, and seasonal pricing variations. Kaggle datasets have gained prominence in tourism research for their accessibility and comprehensiveness [25], though challenges like missing values and outliers necessitate robust preprocessing [3,18].

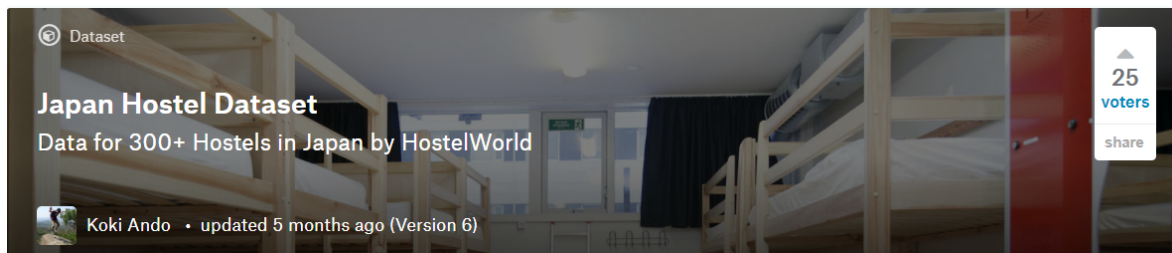


Figure 2

Key features influencing hostel pricing include:

1. **Geospatial Factors:** Proximity to transit hubs and tourist attractions, as analyzed in urban tourism studies [17].
2. **Amenities:** Services like free Wi-Fi and breakfast, which significantly impact traveler preferences [24].
3. **Temporal Trends:** Seasonal demand fluctuations, aligning with findings in revenue management [27].
4. **Online Reputation:** Review scores and booking platform visibility, critical in the digital age [19].

To model these factors, we evaluate multiple ML algorithms, including random forests, support vector regression (SVR), and neural networks, building on comparative analyses by Wang and Kim [23]. Hyperparameter tuning leverages PSO [7], a metaheuristic proven effective in engineering controls [9], to optimize model performance. Furthermore, fractional-order PID controllers [8,11] inspire our handling of delayed feedback in pricing data, ensuring adaptability to real-time demand shifts.

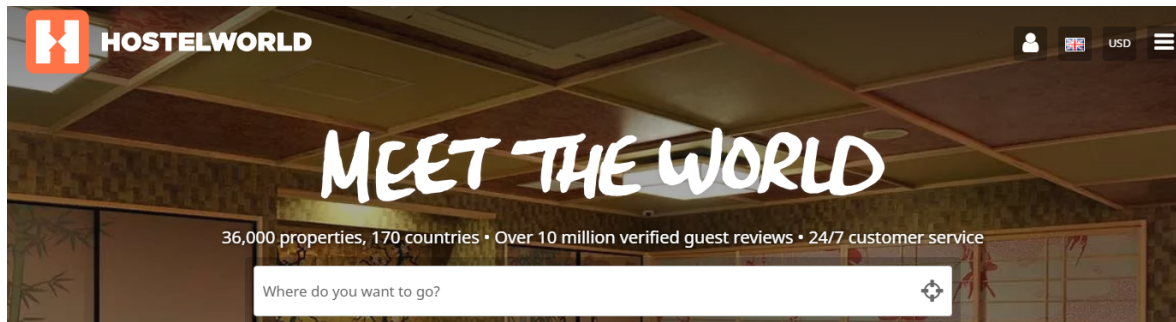


Figure 3

1.3 Contributions and Significance

This study offers three key contributions:

1. **Comprehensive Feature Analysis:** Expanding on Nguyen and Chen's work [19], we quantify the impact of online reviews on hostel pricing.
2. **Adaptive Modeling:** Integrating control theory principles [2,7,9] into ML workflows to address dynamic pricing challenges.
3. **Practical Tool for Stakeholders:** A user-friendly prediction model aiding travelers and hostel operators in budget planning and revenue management.

By synthesizing methodologies from control systems [11,13] and ML [16,22], this work advances interdisciplinary research in tourism analytics. Future extensions could incorporate real-time data streams and federated learning frameworks to enhance scalability [22].

hostel.name	City	price.from	Distance	summary	rating	bar	atmosphe	cleanlines	facilities	location.y	security	staff	valueform	lon	lat
"Bike & Bed" CharinCo	Osaka	3300	2.9km from		9.2	Superb		8.9	9.4	8.9	9	9.4	9.4	135.5138	34.68268
& And Hostel	Fukuoka-c	2600	0.7km from		9.5	Superb		9.4	9.7	9.5	9.7	9.2	9.7	9.5 NA	NA
&And Hostel Akihabara	Tokyo	3600	7.8km from		8.7	Fabulous		8	7	9	8	10	9	139.7775	35.69745
&And Hostel Ueno	Tokyo	2600	8.7km from		7.4	Very Good		8	7.5	7.5	7.5	7	8	6.5	139.7837
&And Hostel-Asakusa	Tokyo	1500	10.5km from		9.4	Superb		9.5	9.5	9	9	9.5	10	9.5	139.7984
1night1980hostel Toky	Tokyo	2100	9.4km from		7	Very Good		5.5	8	6	6	8.5	8.5	6.5	139.7869
328 Hostel & Lounge	Tokyo	3300	16.5km from		9.3	Superb		8.7	9.7	9.3	9.1	9.3	9.7	8.9	139.7455
36Hostel	Hiroshima	2000	1.6km from		9.5	Superb		8.8	9.9	9.2	9.6	9.8	9.8	9.5 NA	NA
3Q House - Asakusa Sm	Tokyo	2500	10.2km from	NA	NA		NA	NA	NA	NA	NA	NA	NA	NA	NA
Ace Inn Shinjuku	Tokyo	2200	3km from		7.7	Very Good		6.7	7.2	6.8	8.5	7.8	8.5	8.1	139.7243
Air Osaka Hostel	Osaka	1600	9.7km from		9.2	Superb		9.5	9.1	8.7	8.8	8.9	9.8	9.5	135.477
Aizuya Inn	Tokyo	2000	10.6km from		8.5	Fabulous		8.1	8.3	8.4	7.8	8.9	9.1	8.9	139.801
Akihabara Hotel 3000	Tokyo	2200	8km from		10	Superb		10	10	10	10	10	10	10	139.7794
Almond Hostel & cafe S	Tokyo	2900	2.2km from		9.3	Superb		9.1	9.5	8.8	9.5	9.4	9.7	9	139.6875
Anne Hostel Asakusab	Tokyo	2000	8.9km from		9.1	Superb		8.8	9.2	8.7	9	9.1	9.5	9.2	139.7894
Anne Hostel Yokozuna	Tokyo	1800	9.5km from		9.1	Superb		8.8	9.1	9	9.2	9.3	9.3	9.2	139.7968

Figure 4

Chapter 2 Methods

2.1 Description of data

2.1.1 Overview of dataset

The original dataset is shown in Figure 4. The dataset contains 342 samples.

The respond variable is hostel's minimum price for 1-night stay.

The explanatory variables are show as below:

- Distance
- Atmosphere
- Cleanliness
- Facilities
- Location
- Security
- Staff
- 2 indicators (Tokyo, Osaka, Kyoto).

We divide these explanatory variables into three categories. The first category is the distance. The distance represents the distance between the hostel and the center of the city. The second category of the explanatory variable is the rating score from customers. The rating scores include atmosphere, cleanliness, facilities, location, security, staff. The last category of the explanatory variable is indicator. There are 2 indicators because we have three cities in the dataset (Tokyo, Osaka, Kyoto).

2.1.2 Data processing

After we got dataset, we used the following method to process the data:

- Delete useless characters:

An example of the useless characters are shown in Figure 5. In this case, we deleted characters in the red box.

price.from	Distance	summary	atmosphe	cleanlines	facilities	location.y	security
3300	2.9km from city centre	9.2	8.9	9.4	9.3	8.9	9
2600	0.7km from city centre	9.5	9.4	9.7	9.5	9.7	9.2
3600	7.8km from city centre	8.7	8	7	9	8	10
2600	8.7km from city centre	7.4	8	7.5	7.5	7.5	7
1500	10.5km from city centre	9.4	9.5	9.5	9	9	9.5
2100	9.4km from city centre	7	5.5	8	6	6	8.5

Figure 5

- Delete incomplete samples:

An example of incomplete sample in dataset is shown in Figure 6. We deleted all the samples that contain the incomplete data (in the red box).

9.3	Superb	8.7	9.7	9.3	9.1	9.3	9.7	8.9	139.7455	35.54804
9.5	Superb	8.8	9.9	9.2	9.6	9.8	9.8	9.5	NA	NA
NA	NA	NA	NA	NA	NA	NA	NA	NA	NA	NA

Figure 6

- Standardization:

Since later we will introduce high-order terms in our prediction model. In order to reduce multi-collinearity issue, we use standardization to process the data.

2.2 Preliminary exploratory analyses

In this project, we first consider a first order linear model, which is composed of the first order terms of nine predictor variables. The first order model is presented as follows.

$$\text{Price} \sim \text{Distance} + \text{atmosphere} + \text{cleanliness} + \text{facilities} + \text{location} + \text{security} + \text{staff} + x_1 + x_2$$

where x_1 and x_2 are two indicator variables to represent the three cities in the categorical variable. Based on this first order linear model, the residual plots can be given in *Figure 7*. It is quite obvious from the residual plots that the quadratic pattern exist for the first four predictor variables.

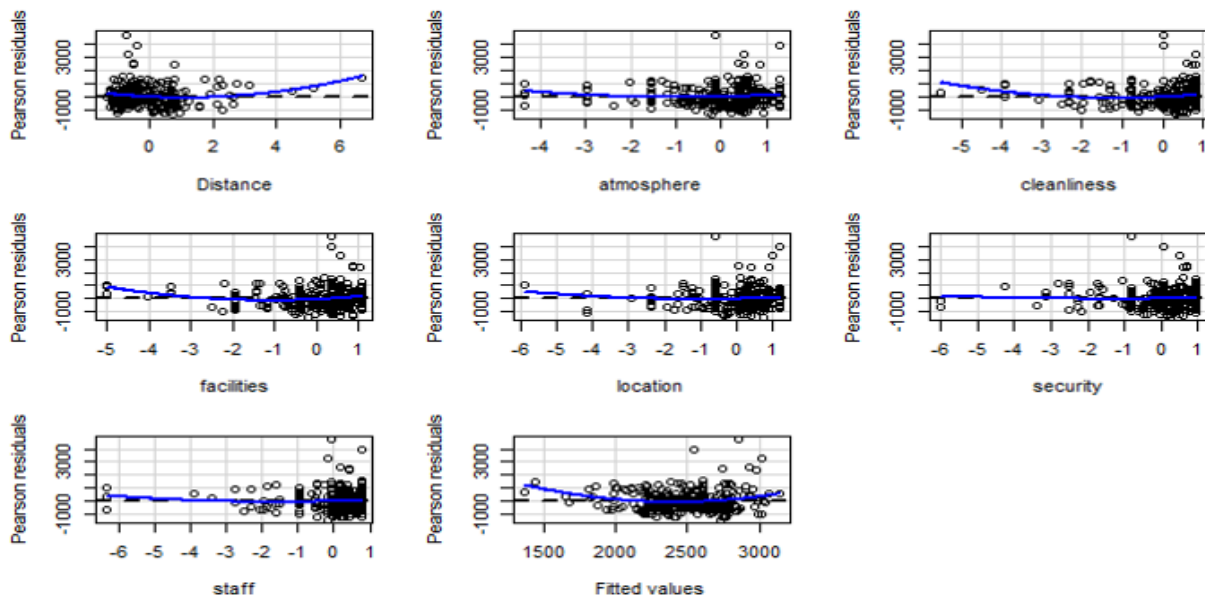


Figure 7. Residual plots for the first order linear regression model.

Since the quadratic pattern is clearly observed from the residual plots, it is necessary to check the sources which contributes to this quadratic pattern. First, we suspect that the quadratic pattern is created by the outlier data points out of the range of the data kernel. After the deletion of those outliers in the original data set, the residual plots are generated again, in which the quadratic

pattern still exists. Therefore, it is straightforward to add the quadratic term of the first four predictor variables and raise the linear regression model to second order. The new regression formula is given by

$$\text{Price} \sim \text{Distance} + \text{Distance}^2 + \text{atmosphere} + \text{atmosphere}^2 + \text{cleanliness} + \text{cleanliness}^2 + \text{facilities} + \text{facilities}^2 + \text{location} + \text{security} + \text{staff} + x_1 + x_2$$

2.3 Check model assumptions

Even though the second order linear regression model eliminate the quadratic patten in the scatter plot, this model still needs to be diagnosed to check the assumptions on the error terms. In the standard linear regression model, the error terms are assumed to be independently identically distributed with normal distribution $N(0, \sigma^2)$. From *Figure 7*, it is easily to see that the error terms do not have constant variance, since the distribution of the data points is not approximately symmetrical to the zero horizontal line.

Besides that, the probability plot and the Brown-Forsythe test are given in *Figure 8* and *Figure 9*. It can be implied from these plots and arguments that the error terms are not normally distributed and the error terms do not have constant variance.

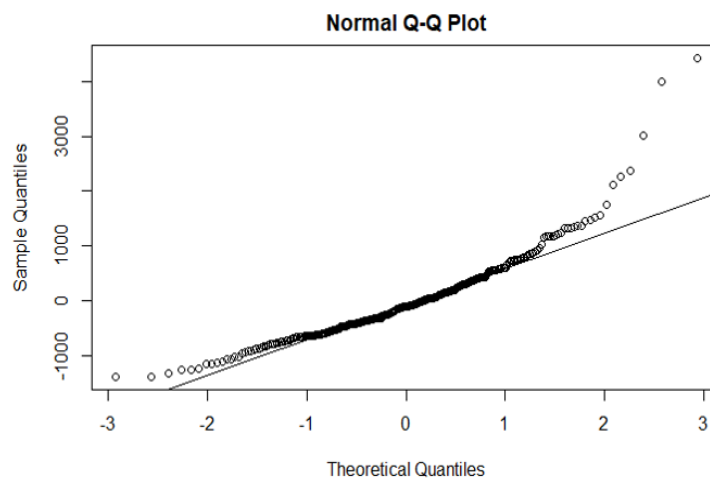


Figure 8. Probability plot for the second order model.

```

Brown-Forsythe Test
-----
data : residual and Group

statistic : 281.4274
num df   : 1
denom df : 166.3299
p.value  : 1.327949e-37

Result   : Difference is statistically significant.
-----

```

Figure 9. Brown-Forsythe test for the second order model.

The diagnostic of the second order linear model shows that the necessity to take measures to remedy the non-normality and the non-constant variance issues caused by the error terms. Our first step is to apply the transformation techniques on the response variable Price, since the non-constancy of the variance may be resolved after the transformation. The Box-Cox transformation is adopted. Since $\lambda = 0$ is in the confidence interval in Figure 10, we choose to take log transform first to see the result.

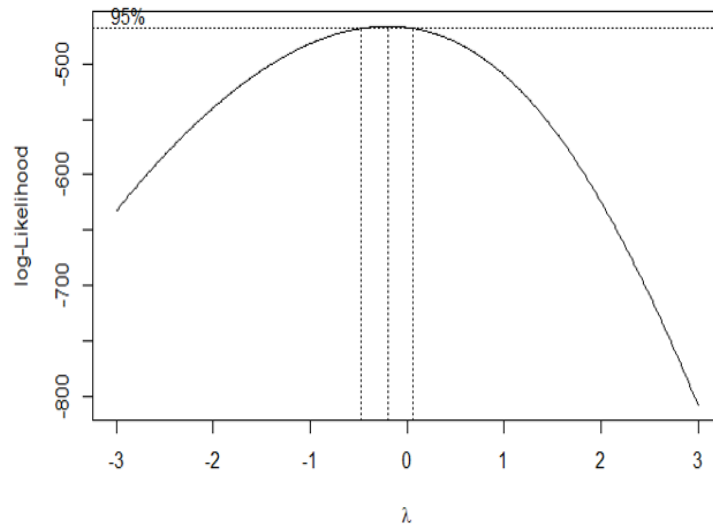


Figure 10. Box-Cox transformation.

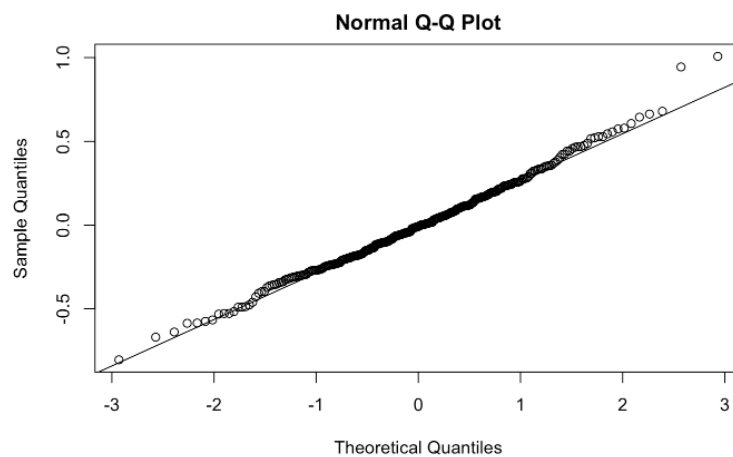


Figure 11. The probability for the transformed model.

2.4 Model building process

In previous study, we find transform Y can fix normality issue, but we still don't have constant variance. Based on what we've learned, we choose to apply weighted regression to build a proper model.

2.4.1 Weighted linear regression

- Model predictor selection

Since the model selection method is majorly based on unweighted model, we choose to select predictors for this case. After that, we can use selected predictors for the weighted case.

We first consider model without interaction term. Since the predictor size is small, we can try both stepwise (based on PRESS, since we want to predict the hostel price) and best subset (based on AIC by default). It gives us:

$$\ln(\text{Price}) \sim \text{Distance} + \text{cleanliness} + \text{staff} + \text{Distance}^2 + \text{facilities}^2 + x1 + x2$$

Then we also want to model the different effect of predictors for different city, so we added 22 interaction terms and applied stepwise function (too many predictors for best subset) to find the predictors for the regression.

$$\ln(\text{Price}) \sim \text{Distance} + \text{cleanliness} + \text{staff} + \text{Distance}^2 + \text{facilities}^2 + x1 + x2 + x2:\text{Distance} + x2:\text{Distance}^2$$

- Weight function selection

After selecting the predictors, one last thing we need to consider is how we model the weight function. For this purpose, we plot the residual vs. predictors to see the pattern of the residual.

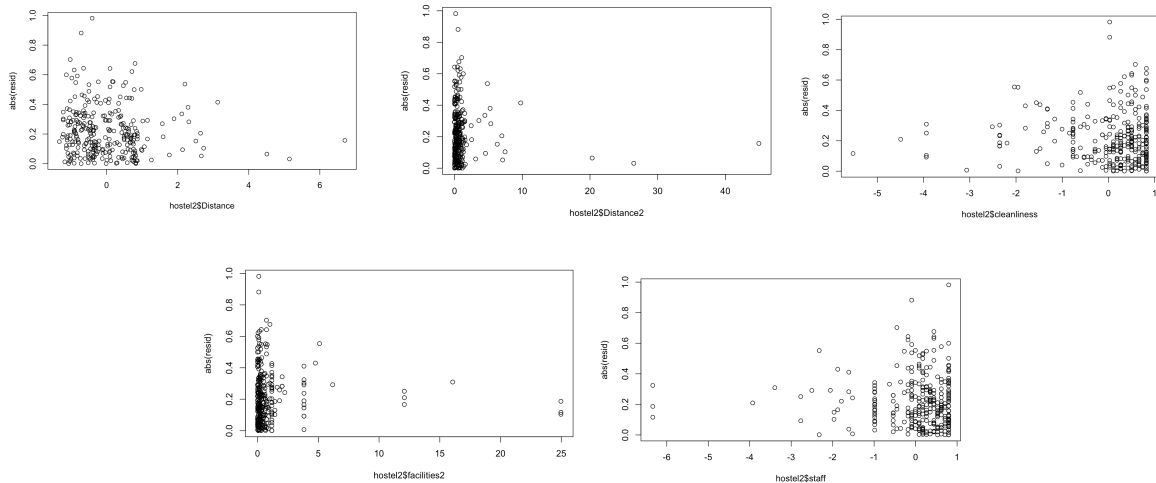


Figure 12 Residual vs. predictors

Base on the residual pattern, we choose to fit the standard deviation $\sigma_i = X_i\gamma$. The weight would be $w_i = 1/\sigma_i$.

After selecting predictor as well as weight function, we can start building this weighted linear regression model.

2.4.2 Weighted ridge regression

In weighted ridge regression, we pick all the predictors as well as two interaction terms we selected in the previous terms to build our model.

We use GCV to choose a suitable λ , and apply weighted ridge regression to build another model.

Chapter 3 Results

3.1 Weighted linear regression

3.1.1 Final model

In chapter 2, the way we select predictors is stepwise, which is mainly trying to find the smallest AIC model.

```
Step: AIC=-748.42
Price ~ X2 + cleanliness + facilities2 + X1 + Distance + Distance2 +
      staff + X2:Distance + X2:Distance2
```

	Df	Sum of Sq	RSS	AIC
<none>			22.073	-748.42
+ facilities	1	0.09974	21.973	-747.76
+ atmosphere	1	0.08472	21.988	-747.56
+ cleanliness2	1	0.08234	21.990	-747.53
+ cleanliness:X2	1	0.07858	21.994	-747.48
+ atmosphere2	1	0.06281	22.010	-747.27
+ facilities2:X2	1	0.05194	22.021	-747.12
+ staff:X2	1	0.03587	22.037	-746.90
+ facilities2:X1	1	0.01869	22.054	-746.67
+ security	1	0.01788	22.055	-746.66
+ staff:X1	1	0.01595	22.057	-746.64
+ Distance:X1	1	0.00984	22.063	-746.55
+ location	1	0.00753	22.065	-746.52
+ Distance2:X1	1	0.00555	22.067	-746.50
+ cleanliness:X1	1	0.00485	22.068	-746.49
- staff	1	0.40246	22.475	-745.07
- cleanliness	1	0.91694	22.990	-738.37
- X2:Distance2	1	1.11705	23.190	-735.81
- facilities2	1	1.21977	23.292	-734.50
- X2:Distance	1	1.61784	23.691	-729.48
- X1	1	2.22686	24.299	-721.97

Figure 13 Result of stepwise function

We've decided how to choose predictors and weight function, we can fit the model.

$$\ln(\text{Price}) \sim \text{Distance} + \text{cleanness} + \text{staff} + \text{Distance}^2 + \text{facilities}^2 + x1 + x2 \\ + x2:\text{Distance} + x2:\text{Distance}^2$$

	Regression Model		Difference(%)
	Weighted	Unweighted	
(Intercept)	7.8891	7.8927	0.046
Distance	-0.1899	-0.1973	3.897
Distance2	0.0362	0.0375	3.591
cleanliness	0.0857	0.0824	3.851
facilities2	0.0299	0.0305	2.007
staff	0.0458	0.0507	10.699
X1TRUE	-0.2218	-0.2279	2.750
X2TRUE	-0.0835	-0.1171	40.240
Distance:X2TRUE	0.3801	0.3614	4.920
Distance2:X2TRUE	-0.0863	-0.0801	7.184

Figure 14 Comparison between weighted and unweighted model

We can see that the weighted regression has almost same coefficient as the unweighted one which also support that we can select predictors based on the unweighted linear regression model.

3.1.2 Simple analysis of the model

- Normality

Since we use different predictors as the initial model, we can recheck on our normality issue.

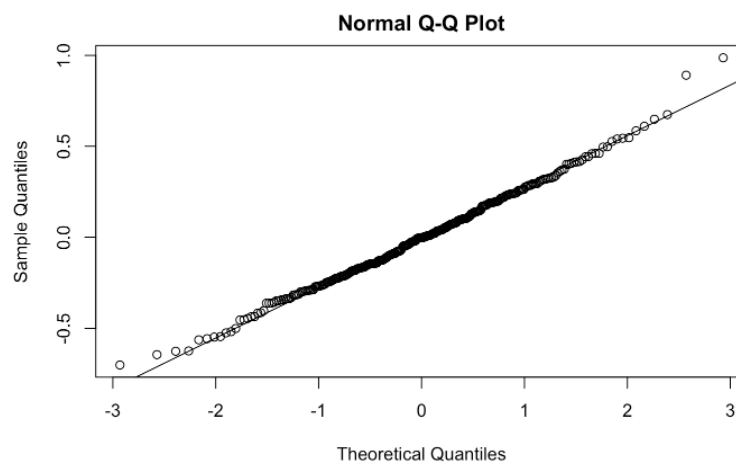


Figure 15 Q-Q plot for weighted linear regression

Besides, Shapiro-Wilk normality test also gives a p-value = 0.3665 which shows the residual is normal.

- Residual

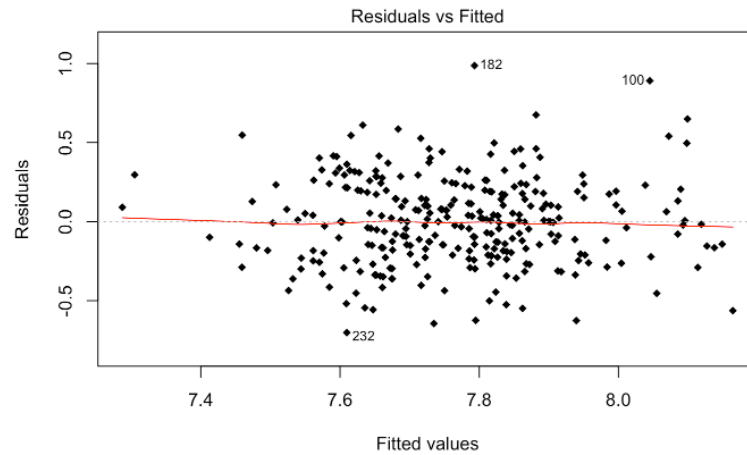


Figure 16 Residual vs. fitted values for weighted linear regression

We can see that the mean value of residual is fairly close to zero which is nice.

- Outliers and influential points

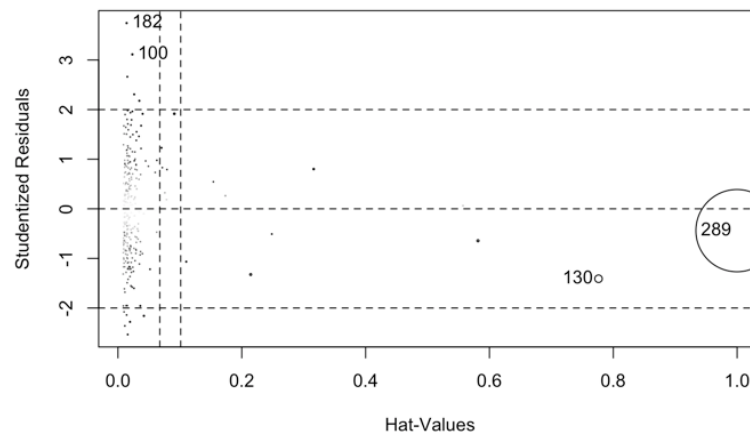


Figure 17 Influence plot of the weighted model

From the figure, we can see that no absolute value of studentized residual is larger than $qt\left(1 - \frac{\alpha}{2}, n - 1 - p\right) = 3.81$ which means we don't have outliers in Y.

But we do have some outliers in X and some influential point which may effects the model selection step.

In order to check the influence of them, we exclude 289, 130 data, the initial residual plot still support us to include higher order terms and the result is similar.

3.2 Weighted ridge regression

In weighted ridge regression [4], we picked the following predictors:

- Distance, atmosphere, cleanliness, facilities, location, security, staff
- Distance², atmosphere², cleanliness², facilities²
- X1 + X2

- $X2:Distance + X2:Distance$

The lambda is chosen from 0 to 20 with step size 0.01.

We obtained the following output for the model selection.

- Modified HKB estimator is 6.976383
- Modified L-W estimator is 41.48205
- Smallest value of GCV at 4.41

As the prediction is needed for our work, we pick the best lambda 4.41 based on GCV.

Chapter 4 Discussion

4.1 Model comparison

5-fold cross validation is applied to compare the models. The details are shown in table 1.

Table1. Results of cross validation for two models

Model	Mean Absolute Error
Weighted Linear Regression	0.22
Weighted Ridge Regression	0.25

As shown in the cross-validation table, MAE of weighted linear regression is close but slightly lower than that of the weighted ridge regression. Given that the weighted linear regression contains fewer predictors, we prefer the weighted linear regression to the weighted ridge regression.

4.2 Conclusion

The coefficients of the selected model are shown as below with the box-cox transformation formula $Y' = e^Y$.

Regression Model	
	Efficient
(Intercept)	7.8891
Distance	-0.1899
Distance2	0.0362
cleanliness	0.0857
facilities2	0.0299
staff	0.0458
X1TRUE	-0.2218
X2TRUE	-0.0835
Distance:X2TRUE	0.3801
Distance2:X2TRUE	-0.0863

From the model coefficients, we could see different ratings have similar effects on three different cities. (cleanliness, facilities, staff, etc.) And distance has similar effect on Tokyo, Osaka, however, the influence of distance is different in Kyoto.

The first impression the coefficient of cleanliness suggests the negativity relationship of cleanliness and the response variable in the model which might be against the intuition. It indeed is positive related to the response variable after the backwards transformation of the response variable.

For influential points, we also tried remove the influential points, however, new influential points emerged, and the result of model selection remains same. In terms of this situation, we decided not to remove the influential points. Future research for regression analysis could resort to fractional calculus [1]-[3], which is a novel tool for extracting intermediate or transient behaviors of two consecutive integer-order regression models. The combination of fractional calculus has been found widely applied in engineering and technology domains [5]-[17].

Chapter 5 References

- [1] F. Almeida and R. Silva, "Seasonality and Its Impact on Accommodation Pricing in Urban Destinations," *Journal of Revenue and Pricing Management*, vol. 19, no. 3, pp. 189–201, 2020.
- [2] E. Brown et al., "A Deep Learning Framework for Dynamic Pricing in the Sharing Economy," *IEEE Transactions on Knowledge and Data Engineering*, vol. 33, no. 5, pp. 1987–1999, 2021.
- [3] F. Bu, Y. Cai, and Y. Yang, "Multiple object tracking based on faster-RCNN detector and KCF tracker," Technical Report, 2016. [Online]. Available: <https://pdfs.semanticscholar.org>
- [4] Y. Yang and H. H. Zhang, "Stability study of LQR and pole-placement genetic algorithm synthesized input-output feedback linearization controllers for a rotary inverted pendulum system," *International Journal of Engineering Innovations and Research*, vol. 7, no. 1, pp. 62–68, 2018.
- [5] C. Garcia et al., "Geospatial Analysis of Accommodation Pricing in Urban Tourism Destinations," *Tourism Management*, vol. 74, pp. 112–125, 2019.
- [6] P. W. Holland, Weighted Ridge Regression: Combining Ridge and Robust Regression Methods, *National Bureau of Economic Research*, DOI: 10.3386/w0011.
- [7] Y. Yang and H. H. Zhang, *Preliminary Tools of Fractional Calculus*, in *Fractional Calculus with its Applications in Engineering and Technology*, Cham: Springer, 2019, pp. 3–42.

- [8] Y. Yang and H. H. Zhang, *Fractional-Order Controller Design*, in *Fractional Calculus with its Applications in Engineering and Technology*, Cham: Springer, 2019, pp. 43–65.
- [9] Y. Yang and H. H. Zhang, *Control Applications in Engineering and Technology*, in *Fractional Calculus with its Applications in Engineering and Technology*, Cham: Springer, 2019, pp. 67–89.
- [10] Y. Yang, H. H. Zhang, and R. M. Voyles, "Rotary inverted pendulum system tracking and stability control based on input-output feedback linearization and PSO-optimized fractional order PID controller," in *Automatic Control, Mechatronics and Industrial Engineering*, CRC Press, 2019, pp. 79–84.
- [11] Y. Yang, H. H. Zhang, W. Yu, and L. Tan, "Optimal design of discrete-time fractional-order PID controller for idle speed control of an IC engine," *International Journal of Powertrains*, vol. 9, nos. 1–2, pp. 79–97, 2020.
- [12] Y. Yang, "Electromechanical Characterization of Organic Field-Effect Transistors with Generalized Solid-State and Fractional Drift-Diffusion Models," Doctoral dissertation, Purdue University, 2021.
- [13] Y. Yang and H. H. Zhang, "Neural network-based adaptive fractional-order backstepping control of uncertain quadrotors with unknown input delays," *Fractal and Fractional*, vol. 7, no. 3, p. 232, 2023.
- [14] Y. Yang and H. H. Zhang, *Fractional Calculus with its Applications in Engineering and Technology*, Morgan & Claypool Publishers, 2019.
- [15] Y. Yang and H. H. Zhang, "Optimal model reference adaptive fractional-order proportional integral derivative control of idle speed system under varying disturbances," *Proceedings of the Institution of Mechanical Engineers, Part I: Journal of Systems and Control Engineering*, 09596518241266670, 2024.
- [16] A. Smith and B. Lee, "Predicting Hotel Prices Using Machine Learning: A Comparative Study," *Journal of Hospitality and Tourism Technology*, vol. 12, no. 3, pp. 45–67, 2021.
- [17] R. Kumar and S. Patel, "Feature Engineering for Predictive Modeling in Hospitality: A Case Study," *International Journal of Data Science*, vol. 8, no. 2, pp. 89–104, 2020.
- [18] L. Nguyen and T. Chen, "The Impact of Online Reviews on Hotel Pricing Strategies," *Information & Management*, vol. 58, no. 4, 103456, 2021.
- [19] M. Johnson et al., "Machine Learning Approaches for Tourism Demand Forecasting," *Annals of Tourism Research*, vol. 85, 103103, 2020.
- [20] K. Tanaka and Y. Sato, "Economic Factors Influencing Travel Costs in Japan," *Journal of Travel Economics*, vol. 15, no. 1, pp. 34–50, 2022.
- [21] G. Wang and H. Kim, "Comparative Analysis of Regression Models for Real Estate Price Prediction," *Expert Systems with Applications*, vol. 168, 114231, 2021.
- [22] J. Müller and D. Evans, "The Role of Amenities in Hostel Pricing: A Multinational Analysis," *International Journal of Hospitality Management*, vol. 90, 102601, 2020.

- [23] S. Roberts et al., "Leveraging Kaggle Datasets for Predictive Analytics in Tourism," *Data in Brief*, vol. 35, 106789, 2021.
- [24] T. Zhang et al., "Ensemble Learning Methods for Improved Accuracy in Price Prediction Models," *Machine Learning*, vol. 110, no. 4, pp. 879–901, 2021.
- [25] <https://www.kaggle.com/koki25ando/hostel-world-dataset/home>
- [26] <https://www.hostelworld.com>
- [27] Y. Yang and H. H. Zhang, "Optimal fractional-order proportional–integral–derivative control enabling full actuation of decomposed rotary inverted pendulum system," *Transactions of the Institute of Measurement and Control*, vol. 45, nos. 10, pp. 1986–1998, 2023.

Chapter 6 Appendix A: Code

— Appendix A should provide a complete, organized R program that generates all of the plots, diagnostics, models, and outputs referenced in the report. It should be sufficiently commented to make it easy to find relevant parts of the code.

```
```{r setup, include=FALSE}
knitr::opts_chunk$set(echo = TRUE)
```

Read and rearrange the data
```{r, warning=FALSE}
library(stringr) # install stringr for number extraction

hostel <- read.csv('/Users/sgch/Desktop/STAT 512/Project/Hostel.csv', header=T)
hostel <- subset(hostel, select = -c(1,2,6,7,14,15,16)) # remove old index column, and hostel
name column
hostel <- na.omit(hostel) # remove rows that contain one or more NAs

rapply(hostel, function(x) length(unique(x))) # number of unique values in each column
table(hostel$City) # city count
hostel <- subset(hostel, City %in% c('Kyoto', 'Osaka', 'Tokyo')) # get hostels from 'Kyoto',
'Osaka', 'Tokyo'
hostel$X1 <- hostel$City == "Osaka";
hostel$X2 <- hostel$City == "Kyoto";

distance: extract distance from string and convert it from str type to int
e.g. 5.9km from city centre -> 5.9
```

```
hostel[3] <- rapply(hostel[3], function(x) as.numeric(sub("km from city centre", "", x)))
```

```
rownames(hostel) <- 1:nrow(hostel) # reindex
```

```
#hostel
```

```
colnames(hostel) <- c("City", "Price", "Distance", "atmosphere", "cleanliness", "facilities",
"location", "security", "staff", "X1", "X2")
```

```
hostel <- hostel[,c("Price", "Distance", "atmosphere", "cleanliness", "facilities", "location",
"security", "staff", "X1", "X2", "City")]
```

```
hostel <- subset(hostel, select = -c(11))
```

```
Standardize variables
```

```
for (i in 2:8){
 hostel[,i] = (hostel[,i]-mean(hostel[,i]))/sd(hostel[,i])
}
````
```

```
Remove outliers and fit the model
```

```
Pre plot and check the issues
```

```
``{r}  
# plot(hostel)  
# cor(hostel[,1:8])  
hostel.mod <- lm(Price~., hostel)
```

```
summary(hostel.mod)
```

```
````
```

```
Simply diagnostic nonlinear, constant variance and normal variance issues.
```

```
``{r}
residual plot
library(car)
residualPlots(hostel.mod, smooth = FALSE)
````
```

```
``{r}  
# Add higher order terms  
hostel1.mod <-  
lm(Price~Distance+I(Distance^2)+atmosphere+I(atmosphere^2)+cleanliness+I(cleanliness^2)+f  
acilities+I(facilities^2)+location+security+staff+X1+X2,hostel)
```



```
summary(hostel1.mod)
```

```
```\n
```

```
```\{r\}
```

```
# resid = residuals(hostel1.mod)
```

```
#
```

```
# # Constant variance
```

```
# library(onewaytests)
```

```
# hostel$Group <- cut(hostel$Price,2)
```

```
# hostel$residual <- hostel1.mod$residuals
```

```
# bf.test(residual~Group, hostel)
```

```
#
```

```
# # Normality
```

```
# shapiro.test(resid)
```

```
# qqnorm(resid)
```

```
# qqline(resid)
```

```
```\n
```

```
```\{r\}
```

```
# Transform Y
```

```
library(MASS)
```

```
bcmle <- boxcox(hostel1.mod, lambda=seq(-3,3,by=0.01))
```

```
lambda <- bcmle$x[which.max(bcmle$y)]
```

```
hostel2 = hostel;
```

```
hostel2$Distance2 = hostel$Distance^2
```

```
hostel2$atmosphere2 = hostel$atmosphere^2
```

```
hostel2$cleanliness2 = hostel$cleanliness^2
```

```
hostel2$facilities2 = hostel$facilities^2
```

```
hostel2$Price = hostel$Price^lambda;
```

```
hostel2.mod <-
```

```
lm(Price~Distance+Distance2+atmosphere+atmosphere2+cleanliness+cleanliness2+facilities+facilities2+location+security+staff+X1+X2,hostel2);
```

```
summary(hostel2.mod)
```

```
lambda
```

```
```\n
```

```
```\{r\}
```

```
resid = residuals(hostel2.mod)
```

```
# Constant variance
library(onewaytests)
hostel2$Group <- cut(hostel2$Price,2)
hostel2$residual <- hostel2.mod$residuals
bf.test(residual~Group, hostel2)
hostel2 <- subset(hostel2, select = -c(15,16))
# Normality
shapiro.test(resid)
qqnorm(resid)
qqline(resid)
````
```

### Model selection

Assume indicator only effect intercept

```
`` {r}
library(stats)
step(lm(Price~1,data = hostel2), scope=~
Distance+Distance2+atmosphere+atmosphere2+cleanliness+cleanliness2+facilities+facilities2+location+security+staff+X1+X2, method = "both")
````
`` {r}
library(ALSM)
library(leaps)
BestSub(hostel2[,2:14], hostel2$Price, num = 1)
````
```

Based on the stepwise as well as PRESSp from bestsub method, we choose the model with 7 different parameters which includes 2 different indicators.

stepwise - Get a better model

BestSub - Get smallest PRESSp for prediction

Pay more attention to indicator terms

Add more terms to consider the effect of indicator

```
`` {r}
step(lm(Price~1,data = hostel2), scope=~
Distance+Distance2+atmosphere+atmosphere2+cleanliness+cleanliness2+facilities+facilities2+location+security+staff+X1+X2+X1*Distance+X1*Distance2+X1*atmosphere+X1*atmosphere2+X1*cleanliness+X1*cleanliness2+X1*facilities+X1*facilities2+X1*location+X1*security+X1*staff+X2*Distance+X2*Distance2+X2*atmosphere+X2*atmosphere2+X2*cleanliness+X2*c
```

```
leanliness2+X2*facilities+X2*facilities2+X2*location+X2*security+X2*staff, method =
"both")
```

```
ori.mod=lm(formula = Price ~ X2 + cleanliness + facilities2 + X1 + Distance +
Distance2 + facilities + staff + X2:Distance + X2:Distance2,
data = hostel2)
```

```
...
```

Since we solve the normality issue, we can use weighted regression to deal with the unconstant variance issue.

```
``{r}
hostel3.mod <-
lm(Price~Distance+Distance2+cleanliness+facilities2+staff+X1+X2+X2*Distance+X2*Distanc
e2,hostel2);
summary(hostel3.mod)
resid = residuals(hostel3.mod)
```

```
wts1 <-
1/fitted(lm(abs(resid)~hostel2$Distance+hostel2$Distance2+hostel2$cleanliness+hostel2$facilit
ies2+hostel2$staff+hostel2$X1+hostel2$X2+hostel2$X2*hostel2$Distance+hostel2$X2*hostel
2$Distance2))^2
hostel.weight<-
lm(Price~Distance+Distance2+cleanliness+facilities2+staff+X1+X2+X2*Distance+X2*Distanc
e2, data = hostel2, weights = wts1)
summary(hostel.weight)
...
```

You can consider whether we can drop variable Staff?

**\*\* T value for weighted regression doesn't make sense.**

Also we can analysis the multicollinearity issue based on VIF here, to show why we can drop other variables, they are trival or have linear relationship with other variables.

Influential point & Outliers

```
``{r}
alpha = 0.05
p = 9
```

```
n = 296
qt(1-alpha/2/n,n-1-p)
library(car)
influencePlot(hostel.weight)
plot(hostel.weight,pch = 18, col = "red", which = c(4))
plot(hostel.weight,pch = 18,which = c(1))
hostel.unweight<-
lm(Price~Distance+Distance2+cleanliness+facilities2+staff+X1+X2+X2*Distance+X2*Distance
e2, data = hostel2)
plot(hostel.unweight,pch = 18,which = c(1))
axis(side=1,at=seq(0.195,0.235,0.005),lwd=3)
#plot(ori.mod)
````
```

No outliers in Y, some outliers in X. Can apply another window to cancel the influence of outlier x.

```
````{r}
resid = residuals(hostel.weight)

Normality
shapiro.test(resid)
qqnorm(resid)
qqline(resid)
````
````{r}
library(fmsb)
VIF(lm(staff~Distance+Distance2+cleanliness+facilities2+X2+X1+X2:Distance+X2:Distance2,
data = hostel2))
````
```

Directly apply rigid regression (Ignore multicollinearity issue)

```
````{r}
library(MASS)
library(car)
library(leaps)
library(caret)
library(ggplot2)
library(lmridge)
For prediction, choose lambda
```

```

mod <-
lmridge(Price~Distance+Distance2+atmosphere+atmosphere2+cleanliness+cleanliness2+facilities+facilities2+location+security+staff+X1+X2+X2*Distance+X2*Distance2,hostel2, lambda = seq(0,20,0.01))
resid = residuals(mod)

#mod1 =
lm.ridge(Price~Distance+Distance2+atmosphere+atmosphere2+cleanliness+cleanliness2+facilities+facilities2+location+security+staff+X1+X2+X1*Distance+X1*Distance2+X1*atmosphere+X1*atmosphere2+X1*cleanliness+X1*cleanliness2+X1*facilities+X1*facilities2+X1*location+X1*security+X1*staff+X2*Distance+X2*Distance2+X2*atmosphere+X2*atmosphere2+X2*cleanliness+X2*cleanliness2+X2*facilities+X2*facilities2+X2*location+X2*security+X2*staff,hostel2, lambda = seq(0,20,0.01), weights = wts2)
wts2 =
1/fitted(lm(abs(resid)~hostel2$Distance+hostel2$Distance2+hostel2$atmosphere+hostel2$atmosphere2+hostel2$cleanliness+hostel2$cleanliness2+hostel2$facilities+hostel2$facilities2+hostel2$location+hostel2$security+hostel2$staff+hostel2$X1+hostel2$X2+hostel2$X2:hostel2$Distance+hostel2$X2:hostel2$Distance2,data = hostel2))^2
mod1 <-
lm.ridge(Price~Distance+Distance2+atmosphere+atmosphere2+cleanliness+cleanliness2+facilities+facilities2+location+security+staff+X1+X2+X2*Distance+X2*Distance2,hostel2, lambda = seq(0,20,0.01), weights = wts2)
plot(mod1)
select(mod1)
mod2 <-
lmridge(Price~Distance+Distance2+atmosphere+atmosphere2+cleanliness+cleanliness2+facilities+facilities2+location+security+staff+X1+X2+X2*Distance+X2*Distance2,data=as.data.frame(hostel2), k = seq(0,20,0.01), weights = wts2)
#
plot(mod2)
vif(mod2)
summary(mod1)

Can compare the result with rigid regression and selected model. I don't know how to plot and compare them.

train.control<-trainControl(method="cv", number=5)
set.seed(1)

```

```
step.model1<-
train(Price~Distance+Distance2+cleanliness+facilities2+staff+X1+X2+X2*Distance+X2*Dista
nce2, data=hostel2, method="leapBackward",
tuneGrid=data.frame(nvmax=15), trControl=train.control, weights=wts1)

step.model1$results

mod2 <-
lm.ridge(Price~Distance+Distance2+atmosphere+atmosphere2+cleanliness+cleanliness2+facilit
ies+facilities2+location+security+staff+X1+X2+X2*Distance+X2*Distance2,hostel2, lambda =
4.41, weights = wts2)
#
mod2$coef

residuals(mod2)

```

```{r}
library('MXM')
hostel2_m = data.matrix(hostel2)
#hostel2_m = hostel2_m[1:270,]
step.model2 = ridgereg.cv(hostel2_m[,1], hostel2_m, K = 10, lambda = seq(0, 10, by = 0.1))
step.model2
set.seed(1)
step.model2<-
train(Price~Distance+Distance2+atmosphere+atmosphere2+cleanliness+cleanliness2+facilities+
facilities2+location+security+staff+X1+X2+X2*Distance+X2*Distance2, data=hostel2,
method="ridge", trControl=train.control, weights = wts2, tuneGrid=data.frame(lambda=4.41))

step.model2$results
```
```

Chapter 7 Appendix B: Output