

CS512 FUN Projects - Spring 2016

Yi Zhong yz614
Qingqiao Hu qh41
Qi Dong qd33

Computer Science Dep., Rutgers University
Piscataway, NJ, USA

Email: yz614@scarletmail.rutgers.edu qh41@scarletmail.rutgers.edu qd33@scarletmail.rutgers.edu

Abstract— Our project aims at building a house recommendation system so that users can receive feedback on house's value and suggestion selling price. The data are received from the website, and the predictions are based on regression unit, which provides data analysis process by converting all the attributions to a new variable.

I. PROJECT DESCRIPTION

Our project aims at construct a visual personalized recommendation system which can be used by different types of users in renting-related activities: house buyers, leasing agencies, and landlords. This system allows buyers to find a satisfactory house by evaluating the options, based on personal preferences that depend on both intrinsic attributes of the houses and their geographic relations to other places of interest. It also allows landlords to see the preferences and distributions of potential buyers so that they can adjust the pricing and improve the conditions of their houses to attract buyers more profitably. The main obstacles are automatically gathering the heterogeneous data related to each housing option, and allowing the user to adjust evaluation rules conveniently. We plan to mimic the real world thinking processes of the above-mentioned users and handle the problem step by step.

The project has four stages: Gathering, Design, Infrastructure Implementation, and User Interface.

A. Stage1 - The Requirement Gathering Stage.

A visual interface, such as a web page with a dynamic map and various options, that allows potential home seekers to filter known for-sale houses by their personal preferences; they can set the preferences on the map (for example, places or types of places that they want to be close to or far from) and on the tags and attributes of each house (such as areas, number of bedrooms, number of bathrooms). Landlords can use the same system to query the statistics of house seekers that are interested in their houses, and see if their preferences are satisfied by the house, and if not, figure out how to improve it to attract buyers.

- The general system description:
- Our system is very deliverable because our interface would be very simple to use. Users only need to input their requirements as some instructions on the interface. Our system aims at construct a visual personalized recommendation system which can be used by different types of

users in buying-related activities: House seeker, leasing agencies, and landlords. This system allows buyers to find a satisfactory house by evaluating the options, based on personal preferences that depend on both intrinsic attributes of the houses and their geographic relations to other places of interest. It also allows landlords to see the preferences and distributions of potential buyers so that they can adjust the pricing and improve the conditions of their houses to attract buyers more profitably. The main obstacles are automatically gathering the heterogeneous data related to each housing option, and allowing the user to adjust evaluation rules conveniently. We plan to mimic the real world thinking processes of the above-mentioned users and handle the problem step by step.

- The three types of users (grouped by their data access/update rights):
- User 1. House seekers: they want to rent a house or room in a certain area, satisfying some geographic preferences like being close to their workplace or any park, as well as functional preferences like having a number of bedrooms and garages. These preferences are personal and can not necessarily be inferred automatically, so user input is necessary.
- User 2. Landlords: they have houses or rooms to sale, and want to find out how to get the most profit, so they need to find out about the preferences of potential buyers; some of these preferences cannot be satisfied by changing the house itself (like location and number of rooms in most cases) but others like furniture and pricing can be adjusted for a better profit.
- User 3. Leasing agencies: they have more houses than typical landlords so can benefit more from knowing a wide range of renter preferences and improve their conditions accordingly.
- The user's interaction modes: User 1: Use keyboard and mouse to input requirements, and system feeds back houses' Informations according to the requirements.
- The real world scenarios:
 - Scenario1 description: Buy a non-commercial house.
 - System Data Input for Scenario1: Price, location, room area, number of bedrooms, number of bathrooms, parking lot size, and so on.

- Input Data Types for Scenario1: A list of the values and choices above.
- System Data Output for Scenario1: A list of rooms and their informations that accord with the requirements of users' input.
- Output Data Types for Scenario1: A table of different attributes of rooms.
- Scenario2 description: Buy a commercial house.
- System Data Input for Scenario2: Price, location, room area, number of bedrooms, number of bathrooms, parking lot size, and so on.
- Input Data Types for Scenario2: A table of the values or choices above.
- System Data Output for Scenario2: A list of s that accord with the requirements of users' input.
- Output Data Types for Scenario2: A table of different attributes of houses.
- The user's interaction modes: User 2. Landlords: Use keyboard and mouse to input their house's information, system feeds back a list of houses' informations that similar with the informations of user's input and suggested prices of the house according to our collected data and informations.
- The real world scenarios:
 - Scenario1 description: Sale a non-commercial house.
 - System Data Input for Scenario1: Price, location, room area, number of bedrooms, number of bathrooms, parking lot size, and so on.
 - Input Data Types for Scenario1: A table of the values or choices
 - System Data Output for Scenario1: A list of houses' informations that similar with the informations of user's input, and suggestions of the room price.
 - Output Data Types for Scenario1: A table of different attributes of rooms.
 - Scenario2 description: Sale a commercial house.
 - System Data Input for Scenario2: Price, location, room area, number of bedrooms, number of bathrooms, parking lot size, and so on.
 - Input Data Types for Scenario2: A table of the values or choices
 - System Data Output for Scenario2: A list of houses' informations that similar with the informations of user's input, and suggestions of the house price.
 - Output Data Types for Scenario2: A table of different attributes of houses and texts.

User 3. Leasing agencies: Use keyboard and mouse to input requirements(these requirements usually much less than User's). After analyse our collected data according to the requirements, system provide different suggestions and recommendations to help users gain more profit.

- The real world scenarios:
 - Scenario1 description: Find rooms' informations

- System Data Input for Scenario1: Location, room size, room numbers.
- Input Data Types for Scenario1: A list of the values and choices above
- System Data Output for Scenario1: A list of rooms and their informations that accord with the requirements of users' input and suggestions about how to gain more agency fee.
- Output Data Types for Scenario1: A table of different attributes of rooms and texts.
- Scenario2 description: Find houses' informations
- System Data Input for Scenario2: Location, house type, room numbers, garage numbers.
- Input Data Types for Scenario2: A table of the values or choices above
- System Data Output for Scenario2: A list of houses and their informations that accord with the requirements of users' input and suggestions about how to gain more agency fee.
- Output Data Types for Scenario2: A table of different attributes of houses and texts.
- Project Time line and Divison of Labor.
Timeline: 4/1, Finish Stage 2. 4/15 Finish Stage 3. 5/1 Finish the project and report
Yi Zhong: algorithm design, documentation, evaluation. Qi Dong: interface implementation. Qingqiao Hu: project report, testing and power point presentation.

B. Stage2 - The Design Stage.

- Short Textual Project Description. The recommendation system works along with users' inputs to give suggestions as a feedback to users about users' houses' prices. We build the system database using crawling, and deposit user's information to the database so that the recommendation system can calculate users' houses to new variables. The feedback includes prediction on house's price, recommended selling price. It is the same to other types of users. For example, landlords can decide when to sell the house for the maximum profit.
- Flow Diagram: see Fig 1
- High Level Pseudo Code System Description.
 - * Information Integration Module: Integrating Unit(website-info) We collect the houses' information of New Jersey from the website "weichert" using crawling to obtain the property information. Then we save the information to Excel files. Return Processed data
 - * Data Analysis Module: Data Analysis Unit(Processed data) Pass Processed data to Data Analysis Unit, then process it through Regression Unit. The process procedure includes to uniform the file style, to remove the redundant information from the original dataset, to transform the non-digital information to digital information. Deposit it to system database

Flow Diagram:

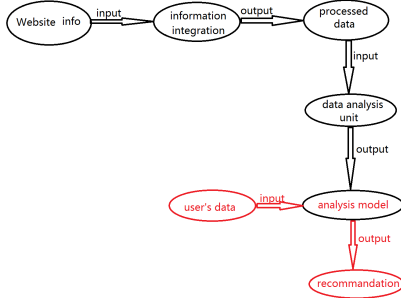


Fig. 1.

- * User's Preference input() Input User's preference
- * Regression Unit() Rate user's input house, then give the prediction using linear regression on house's price. Form a suggestion feedback to users
- Algorithms and Data Structures. Algorithms and Data Structures: First we use web crawling algorithm to stretch useful data from website. This algorithm works with website code and get different digital data from website. Secondly we process these data to table structure type which our data analysis unit could use. Finally, we use our model deal with digital data and output useful recommendations.
- As to the data structure, we use list to save houses' information from the website, process the data using string, use data frame to achieve regression and price prediction.
- The main algorithms used in our project are crawling to get the total information from the website, bubble sort which is used for sorting the house list according to the prices provided from the website and regression which is used for providing the standard of price prediction and providing similar houses from the system data-list.
- Flow Diagram Major Constraints.
 - * First integrity constraint: Houses data on the website must contain price, room numbers, year, area.
 - * Second integrity constraint: Input data to the data analysis unit must be table format with numbers.

C. Stage3 - The Implementation Stage.

The programming language we use for this project is Python 2.7, the programming environment is Windows 64-bit. The deliverables for this stage include the following items:

- Data Snippet: See Fig.2
- Regression and Prediction: See Fig.3
- Datacrawling: See Fig.4-6
- DataAnalysis: See Fig.7-8
- Demo and sample findings
 - * Data size: 1MB; Disk Resident: 10MB;

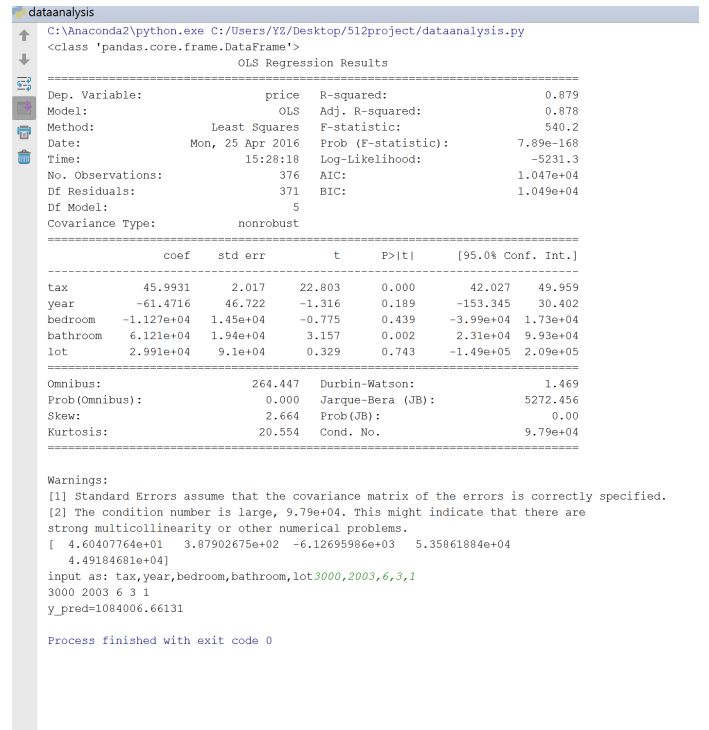


Fig. 2.

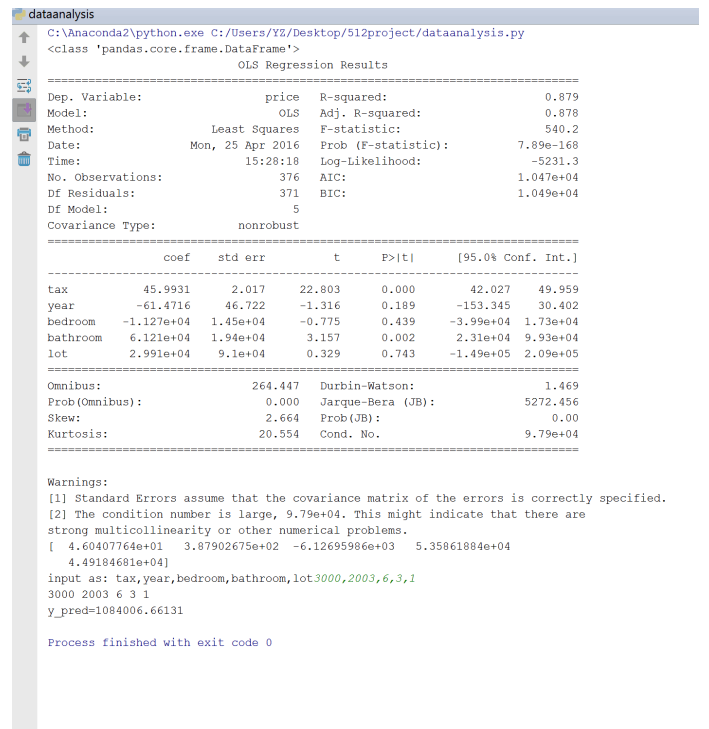


Fig. 3.

Fig. 4.

Fig. 5.

Fig. 6.

Fig. 7.

- * We collect houses' information from "weichert", calculate the prediction price through multiple linear regression. After collecting all the information from the website, we find out that the built-year of houses have no effects on houses' prices. And this is weird, because in common sense, the elder the house is, the lower the price will be, considering higher cost on maintenance, lower security standard. Maybe the website does not put the built-year into the calculation standard.

- Unzip the interface web page and related data into a directory, run the command "python -m SimpleHTTPServer 8000" (if using python 2) or "python -m http.server 8000"(if using python 3) in the directory,

```

dataanalysis.py | Data_CrawlingWrangling.py x
csv_cont = csv_to_list('house.csv')
# print('nooriginal CSV file:')
# print(csv_cont)
# print('noCSV sorted by column "price":')
convset_cells_to_float(csv_cont)
csv_sorted = sort_by_column(csv_cont, 'price')
# print(csv_sorted)
write_csv('housesorted.csv', csv_sorted)

data = pd.read_csv('housesorted.csv')
#sns.pairplot(data, x_vars=['tax','year','bedroom'], y_vars='price', size=7, aspect=0.8)
#plt.show()
data = data[['price','tax','year','bedroom','bathroom','lot']]
x_vars=['tax','year','bedroom','bathroom','lot']

X = data[x_vars]
print type(X)

y = data['price']

X_train,X_test, y_train, y_test = train_test_split(X, y, random_state=1)

linreg = LinearRegression()
model=linreg.fit(X, y)

linear_model = sm.OLS(y,X)
results = linear_model.fit()
print results.summary()
B0 = linreg.intercept_

B1,B2,B3,B4,B5 = linreg.coef_[0],linreg.coef_[1],linreg.coef_[2],linreg.coef_[3],linreg.coef_[4]
print linreg.coef_
x1,x2,x3,x4,x5 = input('input as: tax,year,bedroom,bathroom,lot')
print x1,x2,x3,x4,x5

y_pred = B1*x1 + B2*x2 + B3*x3 + B4*x4 + B5*x5
print 'y_pred:' + str(y_pred)

#print zip(x_vars, linreg.coef_)
#y_pred = linreg.predict(x_test)
#print y_pred
#x_input = (2000,1999,5,3,1)
#x_input = zip(x_vars,x_input)
#print x_input
#y_pred = linreg.predict(x_input)
#print y_pred
#

```

Fig. 8.

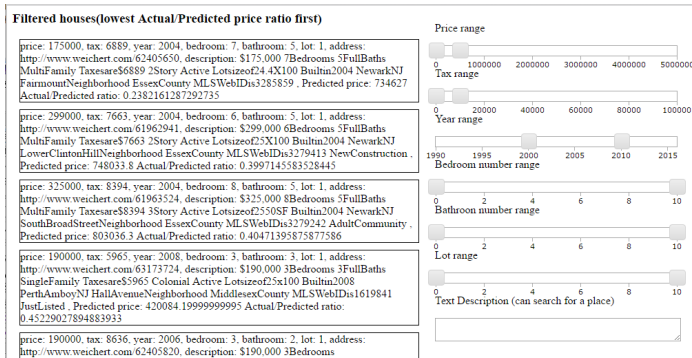


Fig. 9.

then use a web browser and navigate to localhost: 8000/ house data filter.html. The initial screen shot:

- Two different sample navigation user paths through the data exemplifying the different modes of interaction and the corresponding screen shots.
 1. The user wants to find cheap small houses, so slide the bedroom, bathroom, price and tax sliders to the left, then the web page will display the houses that seem like deals according to the actual price to predicted price ratio. The predicted price is calculated with the coefficients obtained from linear regression.
 2. The user wants to find any houses that are deals in the Edison region, so type the text "Edison" (case sensitive) into the Text description search box, and slide all sliders as wide apart as possible. The top houses displayed are the ones that seem to be deals

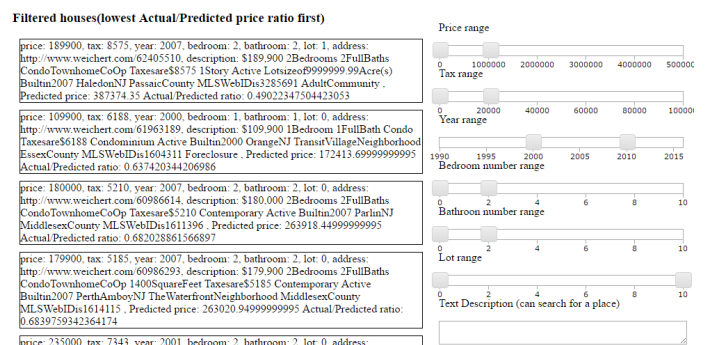


Fig. 10.

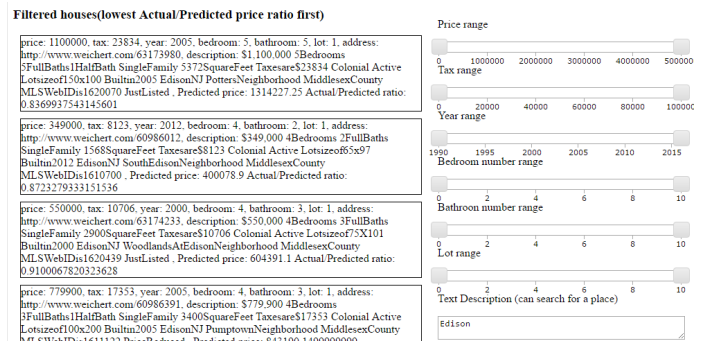


Fig. 11.

according to the linear regression.

- The error messages popping-up when users access and/or updates are denied (along with explanations and examples):
 - * The error message: There are no access denied error messages, because the data gathering/storing and the user interface are different parts in our application. The user interface has all access to data already scraped from the Web.
 - * The error message explanation (upon which violation it takes place): The houses display being empty means that the request cannot be met within the current data set.
 - * The error message example according to user(s) scenario(s): If the user inputs a place name incorrectly (for example, type "Edison" as "Eddison") then most likely there will be no results displayed. What's more, since the result is updated in real time for every character of text, the user can notice the mistake immediately when he types "Edd".
- The information messages or results that pop-up in response to user interface events.
 - * The information message: The "message" is the left side's houses display area being empty.
 - * The information message explanation and the corresponding event trigger Because the user starts typing "Edison" as "Eddison" in the text search box, the houses display becomes empty because

Filtered houses(lowest Actual/Predicted price ratio first)

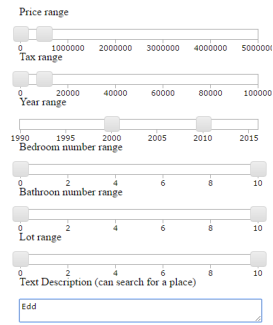


Fig. 12.

"Edd" doesn't match any data records.

- * The error message example in response to data range constraints and the corresponding user's scenario. If the user inputs conditions that cannot be met, the displayed houses area will immediately be empty, and the user can simply correct it and do not need to click on a pop-up. The interface showing no houses match a mistake in input:
- The interface mechanisms that activate different views.
- * The interface mechanism: There is only one view, and the filtering of information is controlled by sliders and search boxes.

II. PROJECT HIGHLIGHTS.

- * Our project aims at building a house-buying system for users like house-seekers, landlords, leasing agencies. The recommendation system provides predictions on house's prices according to all the preferences. The system database are built from the website "Weichert", and price prediction is fulfilled using multiple linear regression which ensures the accuracy of the system. Upon using the system interface, users can get the similar houses directly when they move the preference sliders, and it's simple to use.
- 1) Title: House buying system(Yi Zhong, Qingqiao Hu, Qi Dong)
 - 2) Project Goal: Building a recommendation system on house sale.
 - 3) Outline of the presentation:
 - 4) Description: Our project is based on the information from "Weichert", a website of house renting or buying, and we focus on the houses of New Jersey. We collect the house's information, and save them to our system database. We use linear regression to give prediction to the house price according to the preferences provided by users, and provide the concrete property information of houses for users to choose from.

- 5) Pictures are included in stage 4.
- 6) Project Stumbling Blocks: Information collection, Data processing, Linear Regression, UI Design
- 7) Future Extensions: Expanding the System Database to the US, even the world, Providing more preferences for users to choose from, Building a virtualization of the houses, so that users can get a clear view of similar houses according to the dynamic graph
- 8) References: Weichert