



Predicting Valuation Prices of Danish Real Estate Property

- Forget about realtors

Group 40.

Exam numbers: 115, 144, 146, 204.

Afleveret den: 30/08/2019

Typeenheder:

Indhold

1	Introduction	3
2	Literature Review	3
3	Data Description & Ethics	3
3.1	Ethical Considerations in the Current Research Project	3
3.2	Data Scraping Process	4
3.3	Merging Data	5
3.4	Data Cleaning	5
3.5	Descriptive Statistics	5
3.5.1	The Dataset	5
3.5.2	Selecting Features	5
4	Methods	5
4.1	Supervised Machine Learning	5
4.2	Fitting the model	6
4.2.1	Regularization	6
4.2.2	Lasso	6
4.2.3	Ridge	7
4.2.4	Elastic-net	7
4.3	Optimizing the Hyperparameter	7
5	Analysis	8
6	Results	8
7	Discussion	8
7.1	Data critique	8
8	Conclusion	8
9	An Ethical Overview	8
10	Litterature	10

1 Introduction

2 Literature Review

3 Data Description & Ethics

3.1 Ethical Considerations in the Current Research Project

In the current paper the appropriate care and consideration has been given to the ethical concerns regarding the scraping, processing, and the presentation of the data. Drawing upon the European Comissions¹ ethical guidelines and principles for ethical conduct in social data science, the potential harm to users of Boliga's web-page were carefully considered. As a step to prevent the mosaic effect and in an effort to anonymize the scraped data only aggregated data will be presented in this paper, so that no single observation can be identified from the analysed data.

No informed consent has been obtained from the users of the site, prompting us to consider the consequences of the lack thereof, as informed consent is paramount to the proper, ethical conduct in social science. However, sometimes, as in this instance, informed consent can be logistically impossible to collect from all participants in the study. Salganik mentions that informed consent for everything is an ideal, but in practice impossible to obtain and researchers should instead strive to follow an alternative rule, that he describes as: "some form of consent for most things."² Adhering to this, more complex understanding of the practicality of informed consent, we chose to contact Boliga to inform them of our intent to scrape their website and use the data in an educational context. Boliga responded positively to our inquiry, which we took as informed consent from a third party on behalf of the participating users in our study.

In considering the legal ramifications of our research and to make sure we adhere to the seven principles of GDPR, and other appropriate legislation and legal contracts, we consulted the general guidelines introduced by the Consumer Data Research Center. In particular we noted that we are justified to collect and use the data on the lawful basis of legitimate interests. Furthermore, we consulted with Boliga's Terms and Conditions to avoid any legal ramifications and the appropriate contractual terms of interest can be seen in §10 Terms and Conditions. In order to comply with these terms we refrained from burdening their website's performance by implementing a time.sleep function, which causes each scraping iteration to pause for 0.5 seconds before commencing on scraping the next page (see appendix XX). Furthermore, Boliga prohibits the use of automated scrapers and bots, which we did not do, as our scraping was done in a specified time-frame and we did not automate the procedure to be done on multiple occasions.

¹European Commission (2018). *Ethics in Social Sciences and the Humanities*

²Salganik, Matthew J. (2019) *Bit by Bit Bit - Social research in the digital age*. p.281

3.2 Data Scraping Process

In the following paragraphs our scraping efforts will be described. The scrapers can be examined in the attached Jupyter Notebooks labelled **XX** and **XX**.

Boliga

Our data comes from Boliga.dk, the largest independent online web-portal for real-estate sales in Denmark and has access to unique features such as "liggetid", price-development and access to BBR - the Danish Building and Housing Register ³. Giving us unique insights into the pricing of real-estate in all 98 municipalities of Denmark. At the time of scraping 65,950 properties were for sale.

The scraping process was conducted on Friday the 23rd of August 2019. We had advised Boliga of our intent to collect data from their website via our scraper in an effort to identify ourselves and our intent⁴ (Shiab 2015). In order to scrape the data of interest we familiarized ourselves with the HTML-structure of Boliga. On the basis of these insights we constructed a code, which were able to scrape every page, containing information pertaining the currently listed real-estates on Boliga. The scraper requested all information available from each individual page, which surmounted to 65,950 URL requests.

For each URL, 34 features and the target variable (price) was collected. Table **XX** provides an overview of the features with a short description, and whether the feature has been dropped or saved for later usage. There are three reasons for a feature to be dropped:

1. The feature does not act with independent characteristics according to the our research,
2. The feature contains insufficient data,
3. The feature is poorly formatted and cannot efficiently be recreated.

Hvorlangterder.dk

The scraping of hvorlangterder.dk was achieved by writing a function that took in the GPS-coordinates gathered from Boliga and a list of points of interest. It then constructed an URL, which was then scraped and the json response was then returned as a dictionary of distances to the points of interest. The values from each key in the dictionary was then extracted as a new column in a Pandas DataFrame.

As Jupyter is bad for running long asynchronous tasks and an estimated running time of 18 hours this procedure was run in Visual Studio.

Social and economic factors

Social and economic factors on municipality level is collected from statistik.politi.dk and Danmarks Statistik respectively. These factors includes income, reported crime and level of

³www.boligagruppen.dk

⁴Shiab, Nael (2015) *On the Ethics of Web Scraping and Data Journalism*. Global Investigative Journalism Network.

completed education etc. These are transformed to into ratios, by taking the total population in a given municipality into account.

3.3 Merging Data

3.4 Data Cleaning

A few of the listed housing, have been for sale for over a decade. These are considered outliers, which are overly priced considering their characteristics. The outliers are excluded, because they would otherwise only contribute to a bias of overestimation. The total number of excluded housing are approximately 4000.

3.5 Descriptive Statistics

3.5.1 The Dataset

3.5.2 Selecting Features

Danmark					
	Sjælland		Ikke Sjælland		
Type/Region	Hovedstaden	Sjælland	Syddanmark	Midtjylland	Nordjylland
Villa	513,36	459,67	503,62	461,96	
Rækkehus	470,67	404,95	460,93	407,24	
Ejerlejlighed	482,30	463,94	455,87	449,54	
Fritidsbolig	513,36	459,67	503,62	461,96	
Kolonihave	470,67	404,95	460,93	407,24	
Andelsbolig	482,30	463,94	455,87	449,54	
Landejendom	513,36	459,67	503,62	461,96	
Helårsgrund	470,67	404,95	460,93	407,24	
Fritidsgrund	482,30	463,94	455,87	449,54	
Villalejlighed	513,36	459,67	503,62	461,96	
Lystejendom	470,67	404,95	460,93	407,24	

4 Methods

4.1 Supervised Machine Learning

The objective by applying Machine Learning is to train a model that can predict pricing of housing listed in near future. Two types of predictive models will be used for the sake of returning

different outputs: Regression-model to return the pricing as a continuous output. As well as a classification-model for predicting the interval which the price is included in.

4.2 Fitting the model

The potential problems of underfitting and overfitting should be assessed when fitting a model. A model is underfitted if it hardly captures the variation of the sample data. It is then said that the model has *high bias*. A model is overfitted, when it is overly sensitive to the idiosyncrasy of the sample data and captures the variation in too great detail. This problem often comes with the introduction of a sizeable number of features. Overfit models are said to have *high variance*⁵. In both cases, the model will generalize poorly. A key step in defining a decent model in machine learning is to find an optimal bias-variance-balance, by tuning the complexity of one's model. This is done through *regularization*.

4.2.1 Regularization

The use of a sizeable amount of features leads to a significant risk of over-fitting the prediction model. For the purpose of dealing with over-fitting, regularization needs to be carried out. Regularization can be done with different approaches, in this project three types are applied; Lasso, Ridge and Elastic net.

4.2.2 Lasso

Regularization by Lasso, will penalize complexity of the model by adjustment of parameter estimates. This penalty will make the model less complex and more appropriate for prediction.
6 7

Lasso minimize:

$$L_{Lasso}(\hat{\beta}) = \left(\sum_{i=1}^n (y_i - \hat{y}_i(\beta))^2 + \lambda \sum_{j=1}^p |\hat{\beta}_j| \right) \quad s.t \quad \lambda \geq 0$$

Another convenient attribute of the Lasso penalty is that some estimates are set equal to zero and thereby produce sparse models⁸. Lasso thereby performs the feature selection automatic.

⁵Rashka, Sebastian; Mirjalli, Vahid; *Python Machine Learning, Machine Learning and Deep Learning with Python, scikit-learn, and Tensorflow*. p.73

⁶Foster, Ian; Rayid Ghani Ron S. Jarmin, Frauke Kreuter, Julia Lane; *Big Data in Social Sciences, A Practical Guide to Methods and Tools* p. 173

⁷Rashka, Sebastian; Mirjalli, Vahid; *Python Machine Learning, Machine Learning and Deep Learning with Python, scikit-learn, and Tensorflow* p. 332

⁸Hal R. Varian. *Big data: New tricks for econometrics*. *Journal of Economic Perspectives*. p.19

4.2.3 Ridge

The ridge model add squared terms of the coefficient as penalty. In opposed to Lasso, Ridge do not force features to be omitted. Ridge minimizes:

$$L_{Ridge}(\hat{\beta}) = \left(\sum_{i=1}^n (y_i - \hat{y}_i(\beta))^2 + \lambda \sum_{j=1}^p \hat{\beta}_j^2 \right) \quad s.t \quad \lambda \geq 0$$

Like Lasso, the Ridge model add bias in exchange for a decrease in variance, resulting in a better prediction model.

4.2.4 Elastic-net

The Elastic net is a combination of Lasso and Ridge. Elastic net minimizes the function:

$$L_{elasticnet}(\hat{\beta}) = \frac{\sum_{i=1}^n (y_i - \hat{y}_i(\beta))^2}{2n} + \lambda \left(\frac{1-\alpha}{2} \sum_{j=1}^p \hat{\beta}_j^2 + \alpha \sum_{j=1}^p |\hat{\beta}_j| \right) \quad \text{Note : Anderledes Rashka!}$$
$$0 \leq \alpha \leq 1 \quad \wedge \quad \lambda \geq 0$$

Where α denotes the relative mixing between Lasso and Ridge.

4.3 Optimizing the Hyperparameter

To minimize the mean squared errors of our Lasso and Ridge regression we performed k-fold cross validation to optimize the hyperparameter λ . We split the data into a test set and a development set, consisting of respectively 20% and 80% of the total observations. Subsequently, we use k-fold cross-validation to randomly split the development set into k folds, where k-1 folds are used to train the model. The remaining fold is used to validate the model's generalizability by calculating the mean squared errors of the trained model's prediction of the left-out fold⁹. This process is repeated k times and each time a new fold is left out for validation. Since we are working with a relatively large dataset we chose to split our data into 5 folds, and computed the average MSE for the 5 iterations. By using the k-fold cross-validation method we relieve ourselves of the concern that the estimation of our model's performance is simply due to a lucky or unlucky split of the data.

We performed this procedure for 12 different values of λ spanning between $10^{(-4)}$ and 10^4 . We chose the value of λ which yields the smallest average MSE over the 5 folds. We both calculated the optimal hyperparameters for a Ridge regression model, a Lasso regression model and an Elastic Net regression model. **The following table** shows the performance of the different models, when trained with their optimal hyperparameter and predicting the test data.

⁹Rashka, Sebastian; Mirjalili, Vahid; *Python Machine Learning, Machine Learning and Deep Learning with Python, scikit-learn, and Tensorflow*. p.191

5 Analysis

6 Results

7 Discussion

7.1 Data critique

Another interesting prediction which could have been obtained using the same methods, would have been to estimate expected selling prices. This could have been done simply by scraping data on sold housing instead. This could be of more value for private agents, whose main interest should be the selling price of their housing.

In a prediction-model like this it is near impossible to evade some form of omitted variable bias. A significant amount of potential important factors can not be acquired. For an example the view from the listed housing will for sure have great impact of the valuation price. Another factor of interest could have been a evaluation of the condition of the housing, unfortunately the "statement of property"red. (tilstandrapport) are not publicly accessible.

The prediction of the cooperative housing valuation are subject to significant bias, since the cooperative housing that enters the marked through a realtor often would be those with critical amount of undesirable characteristics.

8 Conclusion

9 An Ethical Overview

The ethical principles of social research are anchored in the fundamental human rights, which are broadly formulated in the UN Declaration of Human Rights. Additional policies and declarations that codify principles of research ethics and the ethical treatment of research participants include the Nuremberg Code, the Helsinki Declaration, the Belmont Report, and the Menlo Report^{10 11}. These codes and addendums originate mostly in the biomedical field, though they encompass the central principles applied to all human research, which have led some academics to call for a Hippocratic oath for data scientists to safeguard against powerful new technologies under development in laboratories and tech firms^{12 13}. This discussion is nothing new however, as a tentative reformulation of the Hippocratic oath was introduced by Karl Popper¹⁴, wherein

¹⁰Salganik, Matthew J. (2019) *Bit by Bit - Social research in the digital age*.

¹¹European Commission (2018). *Ethics in Social Sciences and the Humanities*

¹²Rotblat, Joseph (1999) *A Hippocratic Oath for Scientists*

¹³Sample, Ian (2019, Fri 16 Aug 2019) *and tech specialists need Hippocratic oath, says academic*.

¹⁴Popper, Karl (1969) *The Moral Responsibility of the Scientist*

he stressed the importance of professional responsibility, a critical mind, and an overriding loyalty towards the betterment of mankind.

Matthew Salganik offer four principles deduced from the Belmont and Menlo Report that should guide the ethical deliberations of the researcher: 1) the respect for persons, that is individuals should be treated as autonomous and if circumstances require it individuals should be entitled to additional protections. 2) Beneficence stresses the importance of doing no harm and to maximize the possible benefits and minimizing any potential harms. 3) The principle of justice touches upon the importance of the distribution of burdens and benefits of the social scientist's research. This principle stress that is should not be a single stratum of society that bears the costs of the research while another stratum benefits. 4) The fourth and final principle is the respect for law and public interest, according to Salganik, the principle consists of two distinct elements, that is compliance to relevant laws and legal contracts and transparency-based accountability. It is worth noting that Popper's tentative Hippocratic oath mirrors the first three principles put forth by Salganik, stressing the importance of the ethical conduct of the researcher.

The need for rigid ethical standards within computational social science was made apparent by the Cambridge Analytica Scandal that broke on the 17th of March 2018. Where Steve Bannon could reveal that between 2013 and 2015 Cambridge Analytica had exploited a loophole in Facebook's API which allowed the company to harvest profile data from 87 million Facebook users, without the user's permission and use the harvested data to construct a massive targeted marketing database based on the user's likes and interests¹⁵. Other examples of misuse of data acquired from Facebook include the Harvard-run experiment, where students' data was used to create new knowledge about how social networks form and how these networks and their actors' behavior co-evolve, and the emotional contagion experiment from 2012, where approximately 700,000 users were involved in a research experiment to examine the extent to which a person's emotions are affected by the emotions of the people they interact with (Salganik. 2019).

From this it should be evident that clear ethical guidelines are required in order to protect the user's privacy from tech-savvy companies. To this end, the European Parliament introduced the General Data Protection Regulation (GDPR). With its seven overarching principles the GDPR seeks to formalize the procedures involved in the data processing and storing of sensitive and private information (CDRC¹⁶). These principles include and expand upon the principles found in the Belmont and Menlo Report. The most important consideration, however, must be that even a dataset comprising tens of thousands of observations involve human beings who must be protected from adverse side-effects of the social research. There is considerable evidence that point to the fact that even in anonymized data sets it can be possible to backtrack an individual's identity. A researcher must therefore be mindful of the mosaic effect if the dataset combines large amount of data from various sources (European Commission 2018).

¹⁵Vox.com (2018) *The Cambridge Analytica Facebook scandal* [Online]

¹⁶Consumer Data Research Center, UK [CDRC] (2018) *The General Data Protection Regulation & Social Science Research* [Online]:

10 Literature

Foster, Ian; Rayid Ghani Ron S. Jarmin, Frauke Kreuter, Julia Lane; *Big Data in Social Sciences, A Practical Guide to Methods and Tools*, CRS Press 2017

Rashka, Sebastian; Mirjalli, Vahid; *Python Machine Learning, Machine Learning and Deep Learning with Python, scikit-learn, and Tensorflow*, 2nd edition, Packt Publishing 2017.

Hal R. Varian. *Big data: New tricks for econometrics*. *Journal of Economic Perspectives*, 28(2):3–28, 2014.

Consumer Data Research Center, UK [CDRC] (2018) *The General Data Protection Regulation & Social Science Research* [Online]: ec.europa.eu/research/participants/data/ref/h2020/other/hi/h2020_ethics-soc-science-humanities_en.pdf [Accessed on 25/08/2019]

European Commission (2018). *Ethics in Social Sciences and the Humanities* [Online]: https://ec.europa.eu/research/participants/data/ref/h2020/other/hi/h2020_ethics-soc-science-humanities_en.pdf [Accessed on 25/08/2019]

Popper, Karl (1969) *The Moral Responsibility of the Scientist*. Encounter, March 1969, pp. 52-56 [Online]: <http://www.unz.com/print/Encounter-1969mar-00052> [Accessed on 25/08/2019]

Rotblat, Joseph (1999) *A Hippocratic Oath for Scientists*. Science, November 1999: Vol. 286, Issue 5444, pp. 1475 [Online]: <https://science.sciencemag.org/content/286/5444/1475.full> [Accessed on 25/08/2019] DOI: 10.1126/science.286.5444.1475

Salganik, Matthew J. (2019) *Bit by Bit Bit - Social research in the digital age*. Princeton, NJ: Princeton University Press.

Sample, Ian (2019, Fri 16 Aug 2019) *and tech specialists need Hippocratic oath, says academic*. The Guardian [Online]: [Accessed on 25/08/2019]

Vox.com (2018) *The Cambridge Analytica Facebook scandal* [Online]: [Accessed on 26/08/2019]

Boligagruppen.dk (2019) *Om boliga grupper* [Online]: www.boligagruppen.dk: [Accessed on 27/08/2019]

Shiab, Nael (2015) *On the Ethics of Web Scraping and Data Journalism*. Global Investigative Journalism Network. [Online]: gijn.org/2015/08/12/on-the-ethics-of-web-scraping-and-data-journalism/: [Accessed on 27/08/2019]