

House Price Prediction with Confidence: Empirical Results from the Norwegian Market

Anders Hjort

Department of Mathematics, University of Oslo and Eiendomsverdi AS

Abstract

Automated Valuation Models are statistical models used by banks and other financial institutions to estimate the price of a dwelling, typically motivated by financial risk management purposes. The preferred choice of model for this task is often tree based machine learning models such as gradient boosted trees or random forest, where uncertainty quantification is a major challenge. In this empirical contribution, we compare split conformal inference, conformalized quantile regression and Mondrian conformalized quantile regression on data from the Norwegian housing market, and use random forest as a point prediction. The data consists of $N = 29\,993$ transactions from Oslo (Norway) from the time period 2018-2019. The results indicate that the methods using conformalized quantile regression create narrower confidence regions than split conformal inference.

1 Introduction

Automated Valuation Models. Automated Valuation Models (AVMs) are models used by banks or other financial institutions to get an assessment of the estimated value of a dwelling. The most common model for this is a hedonic model that estimates the price based on the dwellings' attributes. While this historically has been a linear regression model (Bailey et al. 1963), many recent studies indicate that machine learning models such as gradient boosted trees or random forest often have better prediction accuracy (Sing et al. 2021, Kim et al. 2021, Hjort et al. 2022).

Conformal prediction. Conformal prediction (CP) is a model-agnostic framework for uncertainty quantification. The distribution of the absolute residuals $|y_i - \hat{y}_i|$ is used to form confidence regions for new and unobserved instances. One extension of the original CP framework is to construct a *Mondrian* CP (Shafer et al. 2007), where the feature space is split into a set of non-overlapping categories and confidence regions are created separately in each category. Another recently proposed extension by Romano et al. 2019 is *conformalized quantile regression* (CQR), that combines the idea of conformal prediction intervals with the quantile regression framework.

Bellotti 2016 and Bellotti et al. 2021 are both using conformal inference to create confidence bands for AVMs with applications to the UK housing market. The literature is otherwise quite sparse in terms of studies of uncertainty quantification in AVMs.

2 Preliminary results

This paper studies a novel data set of $N = 29\,993$ transactions from Oslo, the capital of Norway, from the two year period 2018 – 2019. The mean sale price in the data is 4.7 million kroner. We use random forest as a point prediction with 500 trees, each with a max depth of 10. The sale price is predicted based on a total of $p = 13$ covariates such as size of the dwelling (measured in m^2), the number of bedrooms, the coordinates of the dwelling and neighborhood characteristics. We then experiment with different ways of creating confidence bands around the predictions. Normalized split CP; CQR and a Mondrian CQR set up where we utilize the 15 different city districts in the data set and create confidence regions with CQR in each city district.

Table 1: The results of various CP methods applied to the data set of transactions from Oslo. The interval sizes are given in million Norwegian kroner (NOK). 1 NOK \approx 0.1 USD per July 2022.

Method	Coverage (%)	Mean interval size	Median interval size
Split CP	89.54	1.85	1.61
CQR	90.25	1.79	1.23
Mondrian CQR	90.40	1.85	1.25

The Root Mean Squared Error (RMSE) of the random forest point prediction is 11.9%. A comparison of the CP methods used to create confidence regions valid at $\alpha = 0.1$ on the test set can be seen in `tab:results`. The split CP method use $\sigma(\cdot) = \exp\{\gamma \cdot \hat{\mu}(x_i)\}$, where $\hat{\mu}(x_i)$ is a GAM model fitted on the residuals from the training set and γ is a hyper parameter. This follows the notation from Bellotti et al. 2021. The CQR methods use the `quantregForest` package to create the quantiles.

The results indicate that CQR methods create narrower confidence regions (measured in Norwegian kroner) than normalized split CP, especially when considering the median interval size.

3 Future work

A major challenge in house price prediction tasks is the temporal dimension; house prices change over time as a result of market movements. This also affects the task of creating confidence regions, as we also would expect these to change with time. A natural direction for future research is to utilize the growing literature on conformal inference for time series, for instance by building on the works of Xu et al. 2021. Another enticing option is to create confidence bands that account for covariate shifts, as outlined in Tibshirani et al. 2019. This is particularly relevant to our application, as different segments of the housing market tend to have very different characteristics. For instance, the dwellings in City A might be quite different from the dwellings in City B. It will be very useful to investigate how we can achieve valid and efficient confidence regions by accounting for this shift in distribution for some or all of the characteristics of a dwelling.

3.1 Acknowledgements

Thanks to Eiendomsverdi AS for providing the data. I am also grateful to Johan Pensar, Ida Scheel and Dag Einar Sommervoll for fruitful discussions about the topic.

References

- Bailey, Martin J., Richard F. Muth, and Hugh O. Nourse (1963). “A regression method for real estate price index construction”. In: *Journal of the American Statistical Association* 58.304, pp. 933–942.
- Bellotti, Anthony (2016). “Reliable region predictions for automated valuation models”. In: *Annals of Mathematics and Artificial Intelligence* 81, pp. 71–84.
- Bellotti, Anthony and Zhe Lim (2021). “Normalized nonconformity measures for automated valuation models”. In: *Expert Systems with Applications* 180, pp. 115–165.
- Hjort, Anders, Johan Pensar, Ida Scheel, and Dag Einar Sommervoll (2022). “House price prediction with gradient boosted trees under different loss functions”. In: *Journal of Property Research* 0.0, pp. 1–27. DOI: 10.1080/09599916.2022.2070525.
- Kim, Jungsun, Jaewoong Won, Hyeongsun Kim, and Joonghyeok Heo (Sept. 2021). “Machine-Learning-Based Prediction of Land Prices in Seoul, South Korea”. In: *Sustainability* 13.23, pp. 202–211. DOI: 10.3905/jpm.2017.43.6.202.
- Romano, Yaniv, Evan Patterson, and Emmanuel Candes (2019). “Conformalized Quantile Regression”. In: *Advances in Neural Information Processing Systems*. Ed. by H. Wallach, H. Larochelle, A. Beygelzimer, F. d'Alché-Buc, E. Fox, and R. Garnett. Vol. 32. Curran Associates, Inc. URL: <https://proceedings.neurips.cc/paper/2019/file/5103c3584b063c431bd1268e9b5e76fb-Paper.pdf>.
- Shafer, Glenn and Vladimir Vovk (July 2007). “A tutorial on conformal prediction”. In: *Journal of Machine Learning Research* 9. DOI: 10.1145/1390681.1390693.

- Sing, Tien Foo, Jesse Jingye Yang, and Shi Ming Yu (Sept. 2021). “Boosted Tree Ensembles for Artificial Intelligence Based Automated Valuation Models (AI-AVM)”. In: *The Journal of Real Estate Finance and Economics* 43, pp. 202–211. DOI: 10.3905/jpm.2017.43.6.202.
- Tibshirani, Ryan J, Rina Foygel Barber, Emmanuel Candes, and Aaditya Ramdas (2019). In: *Advances in Neural Information Processing Systems*.
- Xu, Chen and Yao Xie (2021). “Conformal prediction interval for dynamic time-series”. In: *Proceedings of the 38th International Conference on Machine Learning*. Ed. by Marina Meila and Tong Zhang. Vol. 139. Proceedings of Machine Learning Research. PMLR, pp. 11559–11569.