

A Machine Learning Framework for House Price Estimation

Adebayosoye Awonaike, Seyed Ali Ghorashi, Rawad Hammad

School of Architecture, Arts, Computing and Engineering
University of East London
London, United Kingdom
emails: {U1440387, s.a.ghorashi, r.hammad} @uel.ac.uk

Abstract House prices estimation has been the focus of both commercial and academic researches with various approaches being explored. Depending on the location, size, age, time and other factors, the value of houses may vary. This paper presents a modularized, process oriented, data enabled and machine learning based framework, designed to help the decision makers within the housing ecosystem to have more realistic estimation of the house prices. The development of the framework leverages the Design Science Research Methodology (DSRM) and the HM Land Registry Price Paid Data is ingested into the framework as the base transactions data. 1.1 million London based transaction records between January 2011 and December 2020 have been exploited for model design and evaluation. The proposed framework also leverages a range of neighborhood data including the location of rail stations, supermarkets and bus stops to explore the possible impact on house prices. Five machine learning algorithms have been exploited and three evaluation metrics have been presented and with a focus on RMSE. Results show that an increase in the variety of parameters enables improved accuracy which ultimately will enable decision making. The potential for future work based on this paper can explore the impact of the introduction of other groups of data on the accuracy of machine learning models designed for the estimation of house prices.

Keywords — *House price estimation, Machine learning, Neighborhood data*

1.0 Introduction

Residential housing has been a significant need of human being for a long time. With the average income in the United Kingdom currently around £38,600 for people in full-time jobs and £13,803 for people in part-time jobs [1] the estimated earnings after tax based on the gov.uk tax service calculator is £29,889 and £13,027 respectively. Based on these figures, a 40-year working career for either of these working groups will generate an estimated £1,195,592 and £521,113 respectively of lifetime after-tax income in today's money. Furthermore, with average residential property in the United Kingdom is valued at £231,885 [2] a significant part of our working life is committed to creating the wealth required to own a home. Therefore, though owning a home may probably not be for everyone, everyone needs a place to live and this helps to highlight the range of possible stakeholders in the housing market to include, renters, landlords, investors, developers, housing associations and even local government. To all these stakeholders, being able to estimate the value of residential housing is essential to making critical decisions and ultimately their behavior as player in the housing market.

Hong, Choi and Kim [3] used elapsed year, floor area, floor level of the property and heating system structural factors to determine their impact on price of a house. The outcome showed that elapsed year has a negative correlation with price while area has a positive impact on the price of house. Similarly, Ferlan, Bastic and Psunder [4] observed the same findings on structural factors of a house in Slovenia. An assertion was made that impact of floor level on apartment depends on context i.e., a floor level is a disadvantage if the apartment block has no elevator and but then an incentive if there is an elevator.

With housing price being accepted on a larger scale as a research interest as well as a business interest's economic indicator, researchers in [5] proposed a fine-grained model for price predictions. The study describes that how the proposed machine learning model will help to manage existing challenges with property pricing. This fine-grained housing price forecasting model used economic such as GDP, mortgage deposit ratio and social features such as population.

For the estimation of house prices in Arlington, North Virginia, a benchmark of Random Forest Machine Learning algorithm with linear regression model was created [6]. It was observed that Random forest algorithm is able to capture hidden non-linear relations among various features of a house and ultimately give a better house price estimation. Therefore, the resultant model can be used to predict future real estate prices. In the model, the researchers included influencing house prices factors such as zip code, location of the house, year the house was built, house price, and lot size. A total of 27649 data points were collected from Arlington country, Virginia, USA in 2015. All the data were of a single-family house. Random forest algorithm performed better in terms of R-square and RMSE.

While there are quite a few academic and commercial machine learning based research exist on the subject of house price estimation, there has not been enough focus on the design of a robust framework that continues to learn through batch data ingestion.

In this paper we proposed a framework by focusing on the design of a modularized, process-oriented, data-enabled, machine learning-powered framework. This framework exploits publicly available data and enriches the HM Land Registry’s price paid data by geocoding it (making it spatially enabled) and blending it with neighborhood features such as distance from nearest rail stations, bus stops and supermarkets. Leveraging the capabilities of five different machine learning algorithms; LightGBM for accuracy, efficiency and low computational cost [15], Random Forest for strong performance, ability to handle categorical data with multiple levels, and working adequately with missing data [16], XGBoost for scalability [17] and Hybrid Regression and Stacked Generalization being ensemble of algorithms [18], [19] and [20]. This paper further takes a cumulative multi-parameter layering approach to explore the impact of groups of parameters shown in figure 1 on the accuracy of these machine learning algorithms.

The rest of this paper is organized as follows. Section 2.0 focuses on the methodology deployed including data collection In Section 3.0 modelling and results are discussed while section 4.0 explains model optimization and evaluation. Finally, we conclude in section 5.0.

2.0 Methodology

A variety of research methods that have been explored in existing researches including experimental [10], comparative study [11] and systematic sampling [12]. Since this paper aims to explore a feature layering approach for the design of a multi-parameter house prices estimation framework, the *Design Science Research Methodology* is explored [13]. As shown in figure 2, its process includes (i) problem identification and motivation, (ii) definition of the objectives for a solution, (iii) design and development, (iv) demonstration, (v) evaluation, and (vi) communication. The first three will be discussed in this section while remaining three will be spread across section 3.0 to 5.0.

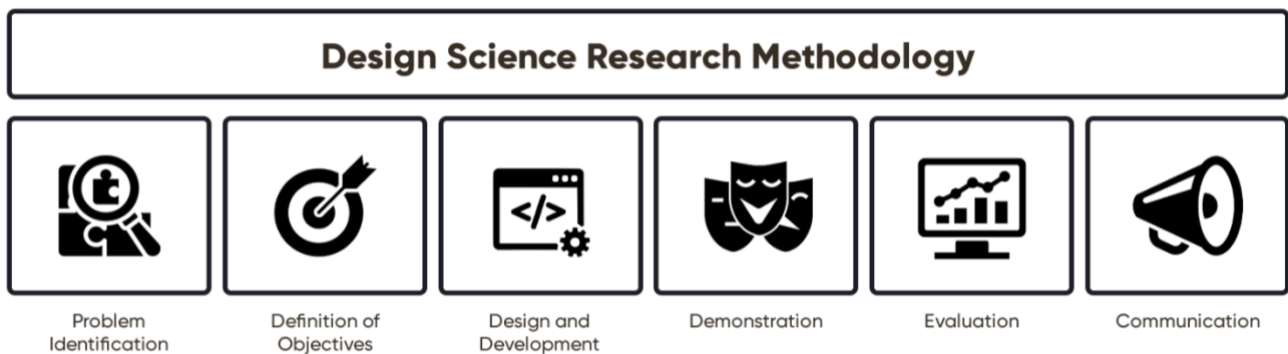


Fig. 2.1. Framework Design Methodology

2.1 Problem Identification

A number of machine learning driven and multi data enabled frameworks already exist for the estimation of house prices however, there is a lack of research to understand the potential impact a cumulative parameter layering on the selection of machine learning algorithms and their accuracy.

2.2 Definition of Objectives

The objectives of this paper include

- Develop a modularized framework to ensure ease of development
- Use an ensemble of machine learning techniques and identify which suites best to estimate house prices based on available data

2.3 Design and Development

The design and development explore data collection, the creation of the modules and pipeline to ensure ease of development and reusability.

2.3.1 Data Collection

All the datasets exploited in this paper are publicly available. These are (i) price paid data, (ii) Office of National Statistics – National Statistics Postcode Lookup (ons nspl) product, (iii) supermarkets, (iv) bus stops and (v) rail stations. These make up the first and second tiers of data exploited as shown in figure 2.2.

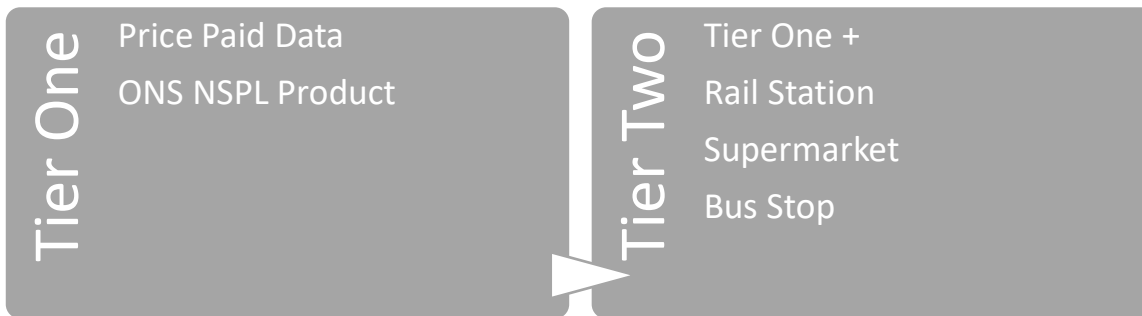


Fig. 2.2. Tiers for cumulative multi-parameter data enabled framework

Tier 1

The HM Land Registry Prices Paid Data is described as “the official house price dataset in England” [7], published by HM Land Registry. The single download file for this dataset contains over 25 million records with 16 variables representing information on all property sales transactions in England and Wales from 1 January 1995 to date [8]. In this paper we have exploited 1.1 million records which represents all the London based transactions.

The ONS produce two main postcode products. These are (i) ONS Postcode Directory (ONSPD) and (ii) National Statistics Postcode Lookup (NSPL).

These products are widely used by a range of customers including central and local government, commercial organizations and academia [9]. The ONSPL has been used in this paper to create a geocoded version of the HM Land Registry price paid data.

Tier 2

The “GB Rail Stations dataset” is a list of all the train stations in the United Kingdom. It is made up of 2,569 records and 39 variables including the volume of passengers travelling through each station either as the start or end of their journeys or even as an interchange is captured. This paper takes a scientific approach to explore the possible impact of the distance to the nearest train stations on house prices.

Retailers are also opening up in multiple location and in different formats, thereby providing customers with choice. With so many new stores, it becomes relatively easier to know where the competition is and consequentially, the new markets being targeted by retailers. Therefore, this paper also explores the impact of the distance to supermarkets on house prices using United Kingdom Supermarket “Retail Points”, a dataset of supermarkets. This Geolytix data has 16,991 records and 17 variables.

“The bus stop dataset” is published to the National Public Transport Data Repository by the Department for Transport. The data in this repository is available from October 2004 to October 2011. However, it is now static and superseded by the Traveline National Dataset.

Table 2.1. Publicly available datasets collected and exploited in machine learning models

High-Level Profile					
Tier	Dataset	No of Records	No of Variables	Format	Source
Tier 1	Price Paid Data	1,096,000	16	.csv	gov.uk
	ONS NSPL Product	2,661,131	41	.csv	ons.gov.uk
Tier 2	Rail Stations	2,569	39	.csv	doogal.co.uk
	Supermarkets	16,991	17	.csv	geolytics.com
	Bus Stops	406,873	25	.csv	data.gov.uk

2.3.2 Modular Programming

To introduce a robust solution that ensures ongoing and future development is done more quickly, a modular approach has been followed in this paper. A module approach has been followed to develop this framework. This module approach aims at applying software engineering principles including low coupling and high cohesiveness to produce a maintainable software tool that is scalable in the future.

Figure 3.3 shows the modules that define this framework and they include (i) **assets** – holding all raw data to be ingested as detailed in figure 1, (ii) **data ingestion** – comprising of functions designed to ingest all research data from the asset module and also extract the specific records required as baseline, (iii) **data processing** – this is made up of the functions designed for data cleansing and initial exploratory data analysis, (iv) **features engineering** – this caters for the engineering of features in datasets across both tiers (v) **model building** – this is where the machine learning algorithms are exploited on the training data, test data and validation data (vi) **params** – this module serves as the store for multiple dictionaries created to hold groups of data that belong to the same tier and (vii) **utils** – this comprises of classes and functions designed to handle processing tasks like the *label encoding* of categorical features and *feature union* of numerical features and (viii) **Main** is the primary environment where the functions and classes across other modules are called into action.

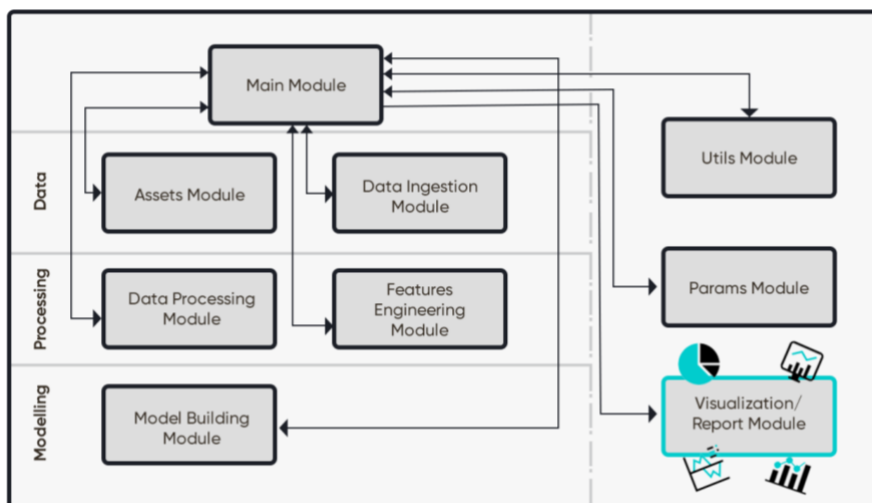


Fig. 2.3. A view of all the Modules in the Framework

3.0 Modelling and results

Baseline models were created using five modelling techniques but using default parameters (i.e., no tuning at this point). These are (i) Light GBM, (ii) Random Forest, (iii) XGBoost, (iv) Hybrid Regression and (v) Stacked Generalization. These regression algorithms have been selected because of the speed of learning, the handling of overfitting to improve accuracy and high flexibility.

Light GBM (Light Gradient Boosting Machine) is a machine learning algorithm based on decision trees. It has wide application in real-life such as ranking, classification and tasks based on machine learning. Its development focuses on performance and scalability. It's advantages include sparse optimization, early stopping, parallel training, bagging,

regularization as well as multiple loss functions. Exclusive Feature Bundling (EFB) and Gradient-based one-side sampling (GOSS) are the two powerful techniques used by LightGBM to improve accuracy, efficiency and memory consumption as well as speed [15].

Random Forest is a machine learning (ML) algorithm based on weak learner decision trees [16]. The ensembling manner of the decision trees eliminates instability problems as well as overcome the high variance of decision trees. Since decision trees are generated by random sampling method, hence the name is random forest.

XGBoost exists as an open-source package for tree boosting and can be described as a machine learning system that is scalable. Its scalability feature allows the users to quickly define their objectives.

This algorithm is also very fast as it performs parallel computation, it accepts a wide array of inputs, it has sparsity, customization and acceptable performance [17].

Hybrid Regression is an ad-hoc and user-defined ML method. In this regard, hybrid regression entails combining two or more ML methods to develop a unique ML method. The combined outcome of these ML methods is far better than the results of each ML method by itself. For instance, [18] tested houses Prices in Beijing with a model consisting of 33.33% Random Forest, 33.3% LightGBM and 33.3% XGBoost.

The hybrid model achieved a better result of RMSE of 0.14969 far better than each algorithm run solely. Similarly, [19] realized a better overall result of a hybrid model consisting 65% Lasso and 35% Gradient Boosting.

Stacked Generalization was introduced by Wolpert [20] and it's a Python based package. The main idea behind this method is to use predictions of the previous models as features for the present model. Stacked Generalization uses K-fold cross-validation to avoid overfitting. For instance, [18] used 2-level stacking architecture, the first stack comprised on Random Forest and LightGBM while the second stack comprised on XGBoost to predict the house prices. [18] They also noted that the combined results are not as impressive as the Hybrid Regression

Table 3.1 show the results for all models using the default parameters for each algorithm on only the geocoded HM Land Registry Price Paid Data while table 3.2 show the results for all models using the default parameters for each algorithm after a layer of neighborhood data (shown in figure 2.2) have been introduced. In this paper, the introduction of the tier 2 datasets is described as cumulative multi-parameter layering.

Table 3.1. Modelling results using default parameters (Tier 1)

Model	RMSE		MAE		R-Square	
	Train	Test	Train	Test	Train	Test
Light GBM	2399452.94	3520071.29	340137.82	359972.67	0.2249	0.0878
Random Forest	1417480.98	3702246.49	215498.63	370818.23	0.7295	-0.0090
XGBoost	2224870.41	3705299.63	353982.24	389103.92	0.3336	-0.0106
Hybrid Regression	1943280.22	3552304.53	280043.09	347270.39	0.4916	0.0711
Stacked Generalization	2605422.57	3705613.80	374086.52	401373.52	0.0861	-0.0108

Table 3.2. Modelling results using default parameters (Tier 2)

Model	RMSE		MAE		R-Square	
	Train	Test	Train	Test	Train	Test
Light GBM	2248053.95	3477427.88	325454.33	348603.42	0.3196	0.1098
Random Forest	945781.32	3492294.16	107019.28	289294.10	0.8795	0.1021
XGBoost	2224870.41	3705299.63	353982.24	389103.92	0.3336	-0.0106
Hybrid Regression	1640064.48	3416892.70	236584.54	309849.29	0.6379	0.1405
Stacked Generalization	2345842.61	3631184.76	348189.47	375124.47	0.2591	0.0293

RMSE, Root Mean Square Error is a common metric used to measure the error of a model predicting quantitative data [21]. It estimates the standard deviation of an observed value from the model prediction. According to [21], the observed value is equal to the sum of the predicted value and predictably distributed random noise with mean zero. If the noise is negligible as estimated by RMSE, the model is assessed as good at predicting the observed data. However, if RMSE is large, it means that the model is not accounting for important features in the data

Table 3.3. A comparison of modelling results for Tier 1 and Tier 2 with a focus on RMSE

Model	RMSE		% Improvement
	Tier_1_Default	Tier_2_Default	
	Test	Test	
Light GBM	3520071.2880	3477427.8808	1.2%
Random Forest	3702246.4881	3492294.1565	5.7%
XGBoost	3705299.6315	3705299.6315	0.0%
Hybrid Regression	3552304.5328	3416892.6971	3.8%
Stacked Generalisation	3705613.8017	3631184.7601	2.0%

$$RMSE = \sqrt{\frac{\sum_{i=1}^n (\hat{y}_i - y_i)^2}{n}} \quad [1]$$

Where: y_1, y_2, \dots, y_n are predicted values; y_1, y_2, \dots, y_n are observed values; n is the number of observations

A reduction in RMSE values for all models in Tier 2 compared with Tier 1 shows an improvement in the accuracy of the models based on the introduction of neighborhood data. As shown in figure 2.2. the introduction of neighborhood data is a cumulative layering of neighborhood datasets on the geocoded price paid data. Random Forest model showed the most improvement while XGBoost had no change.

4.0 Model Optimization and Evaluation

Model optimization in machine learning is, one of the most challenging aspects of the implementation of ML solutions. There is immense attention given to deep learning theories and machine learning to achieve the optimization of models. There exist two types of parameters in models of machine learning; model parameters - which possess the ability to be initiated and consequently updated through data learning and hyper-parameters - which have to be set before the training of the machine learning model since they are associated with the configuration of the machine learning model [22].

Hyperparameters were calculated using the ‘Bayesian Optimization’ method and applied to three of the baseline models i.e. (i) Light GBM, (ii) Random Forest and (iii) XGBoost. The hyperparameters used are shown in *table 4.1*. According to [23], there are four common methods explored for the hyperparameter optimization of machine learning models and Bayesian model-based Optimization is assessed as most efficient. The others include random search, grid search and manual.

Table 4.1. Hyperparameters calculated using Bayesian Optimization method

	Tier 1	Tier 2
Light GBM	{'colsample_bytree': 0.7139109116423488, 'learning_rate': 0.1133802923850488, 'max_depth': 1, 'min_child_samples': 300.0, 'n_estimators': 2, 'num_leaves': 400.0, 'reg_alpha': 2.697524711103592, 'reg_lambda': 7.792201841119942, 'subsample': 0.850664221897626, 'subsample_for_bin': 26000.0}	{'colsample_bytree': 0.9298462716109409, 'learning_rate': 0.10075805121869513, 'max_depth': 0, 'min_child_samples': 300.0, 'n_estimators': 3, 'num_leaves': 600.0, 'reg_alpha': 2.671148237480234, 'reg_lambda': 0.31242913248603976, 'subsample': 0.8930209167711038, 'subsample_for_bin': 30000.0}
Random Forest	{'max_depth': 16.80565381273467, 'min_samples_leaf': 4.548341736478191, 'n_estimators': 302.9997495928546}	{'max_depth': 19.633732543557223, 'min_samples_leaf': 4.900530193117771, 'n_estimators': 366.02562034116534}
XGBoost	{'colsample_bytree': 0.62, 'gamma': 0.32, 'learning_rate': 0.02, 'max_depth': 2, 'min_child_weight': 10.0, 'n_estimators': 4, 'subsample': 0.54}	{'colsample_bytree': 0.73, 'gamma': 0.31, 'learning_rate': 0.09, 'max_depth': 2, 'min_child_weight': 8.0, 'n_estimators': 2, 'subsample': 0.72}

Table 4.2. Modelling results using optimized parameters (Tier 1)

Model	RMSE		MAE		R-Square	
	Train	Test	Train	Test	Train	Test
Light GBM	2706799.10	3668646.78	453807.88	472872.95	0.0136798	0.00922325
Random Forest	2222173.82	3523322.99	320058.47	358772.42	0.3352445	0.08616261
XGBoost	2597208.98	3563082.33	409654.74	425546.49	0.0919293	0.06542163

Table 4.3. Modelling results using optimized parameters (Tier 2)

Model	RMSE		MAE		R-Square	
	Train	Test	Train	Test	Train	Test
Light GBM	2639464.50	3611587.57	397709.05	415669.82	0.0621411	0.03980308
Random Forest	1980820.53	3415804.71	238764.94	299770.67	0.4718028	0.14108522
XGBoost	2762658.11	3714580.93	534439.47	555323.00	-0.0274487	-0.01574258

Table 4.2 and Table 4.3 show the model results for three optimized models using three different metrics. However, with a focus on RMSE, table 4.4 shows a consistent improvement for Tier 2 compared to Tier 1. This also confirms that the introduction of neighborhood data leads to improved performance of the Random Forest model.

Table 4.4. A comparison of modelling optimization results for Tier 1 and Tier 2 with a focus on RMSE

Model	RMSE		% Improvement
	Tier_1_Optimised Test	Tier_2_Optimised Test	
Light GBM	3668646.7761	3611587.5656	1.6%
Random Forest	3523322.9870	3415804.7102	3.1%
XGBoost	3563082.3306	3714580.9281	-4.3%

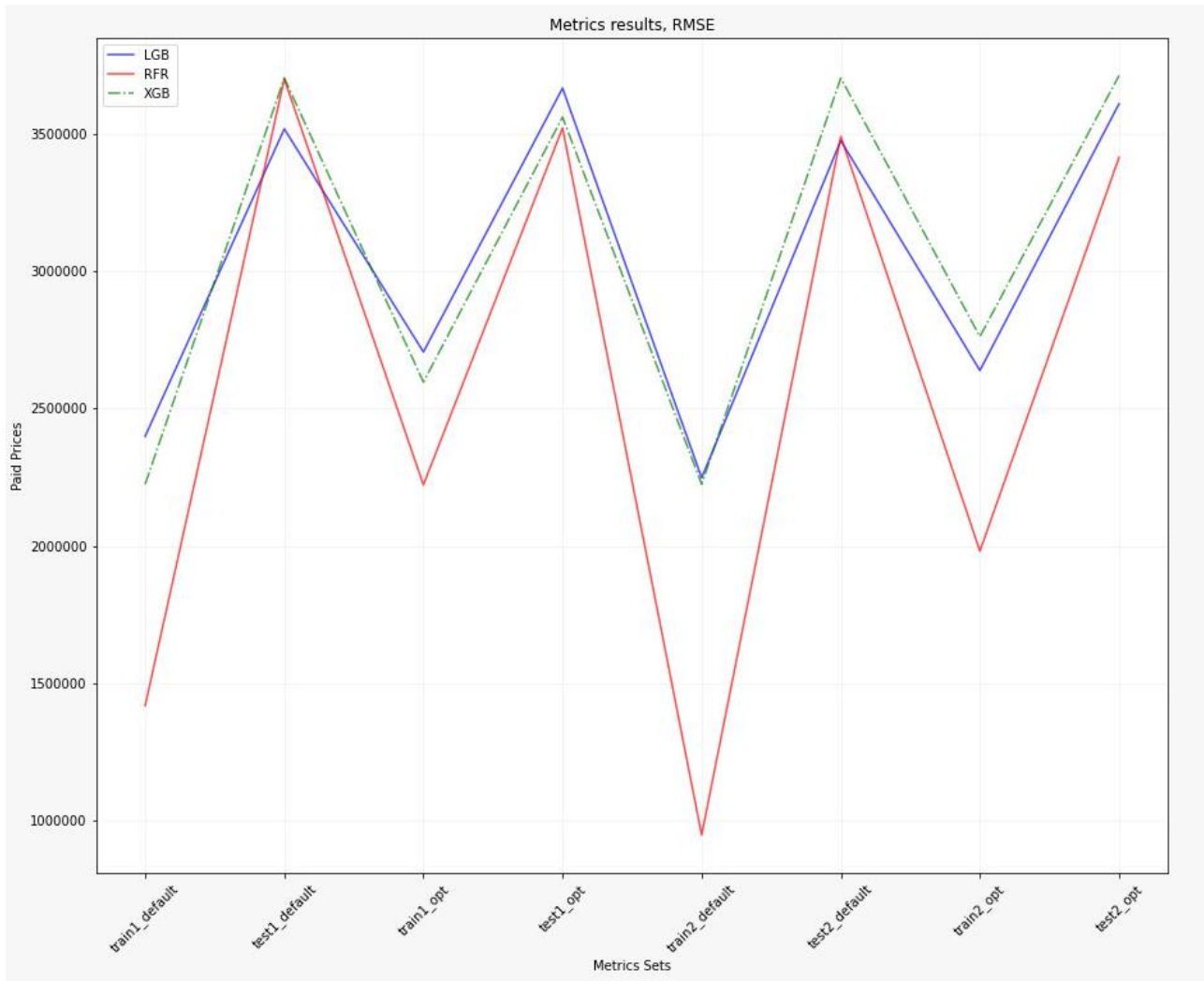


Fig. 4.2. Root Mean Square Error result for three models (both default and optimized parameters) across tiers 1 and 2 (both train and test data)

5.0 Discussion and Conclusion

This paper has presented a modularized, process-based, data driven and machine learning enabled framework for housing price estimation. It is comprised of eight modules leveraging the design science research methodology and five publicly available datasets categorized into two tiers. Five Machine Learning methods and techniques including LightGBM, Random Forest, XGBoost, Hybrid Regression and Stacked Generalization have been exploited and analyzed for optimal estimations.

Although these methods and techniques produced some intriguing results, they all have their advantages and disadvantages. For random forest model results using RMSE as an evaluation metric shows better accuracy in the performance of the model for tier 2 as compared with tier 1. This shows that the cumulative layering of neighborhood factors improved the accuracy of house price estimation.

With continuous improvement of the estimated house price, developers or builders will have data enabled insights on where to build and what kind of homes to build owing to the demography of potential buyers or investors. Investors and landlords will also have better insights to calculate their return on investment and new or serial tenants are able to identify geographic areas that best suit their personal preferences. Further work will be required to explore the potential impact of ingesting an additional layer of factors.

REFERENCES

- [1] Office of National Statistics, Employee earnings in the UK: 2020, viewed 30 May 2021, <<https://www.ons.gov.uk/employmentandlabourmarket/peopleinwork/earningsandworkinghours/bulletins/annualsurveyofhoursandearnings/2020>>
- [2] HM Land Registry 2020, UK House Price Index for March 2020, GOV.UK, viewed 05 January 2021, <<https://www.gov.uk/government/news/uk-house-price-index-for-march-2020>>
- [3] Hong, J., Choi, H. and Kim, W.S., 2020. A house price valuation based on the random forest approach: the mass appraisal of residential property in south korea. *International Journal of Strategic Property Management*, 24(3), pp.140-152.
- [4] Ferlan, N., Bastic, M. and Psunder, I., 2017. Influential factors on the market value of residential properties. *Engineering Economics*, 28(2), pp.135-144.
- [5] Ge, C. et al., 2019. An Integrated Model for Urban Subregion House Price Forecasting: a Multi-source Data Perspective. *2019 IEEE International Conference on Data Mining (ICDM)*, pp. 1054-1059.
- [6] Wang, C. and Wu, H., 2018. A new machine learning approach to house price estimation. *New Trends in Mathematical Sciences*, 6(4), pp.165-171.
- [7] Chi B. et al, 2019. Creating a new dataset to analyse house prices in England. Available at <https://discovery.ucl.ac.uk/id/eprint/10082766/>. [Accessed 29 September 2021]
- [8] GOV.UK, Statistical Dataset Price Paid Data, viewed 30 May 2021, <https://www.gov.uk/government/statistical-data-sets/price-paid-data-downloads>
- [9] Office of National Statistics, Postcode products: 2021, viewed 30 May 2021, <<https://www.ons.gov.uk/methodology/geography/geographicalproducts/postcodeproducts>>
- [10] Park B and Bae J.K., 2014. Using machine learning algorithms for housing price prediction: The case of Fairfax County, Virginia housing data. [online] *Expert Systems with Application*. Available at <https://doi.org/10.1016/j.eswa.2014.11.040> [Accessed 30 March 2021]
- [11] Madhuri, C. R., Anuradha, G. & Pujitha, M. V., 2019. House Price Prediction Using Regression Techniques: A Comparative Study.. *2019 International Conference on Smart Structures and Systems (ICSSS)*.
- [12] Rico-Juan, J.R and Taltavull de La Paz, P., 2021. Machine learning with explainability or spatial hedonics tools? An analysis of the asking prices in the housing market in Alicante, Spain. *Expert Systems with Applications*
- [13] Hammad R.K.M., (2018). *A Hybrid E-Learning Framework: Process-Based, Semantically-Enriched and Service-Oriented*. PhD Thesis. University of West England, Bristol.
- [14] Scikit-Learn, sklearn.pipeline.Pipeline, accessed 2 October 2021, <<https://scikit-learn.org/stable/modules/generated/sklearn.pipeline.Pipeline.html>>
- [15] Ke, G., Meng, Q., Finley, T., Wang, T., Chen, W., Ma, W., Ye, Q. and Liu, T.Y., 2017. Lightgbm: A highly efficient gradient boosting decision tree. *Advances in neural information processing systems*, 30, pp.3146-3154.
- [16] Wang, C. and Wu, H., 2018. A new machine learning approach to house price estimation. *New Trends in Mathematical Sciences*, 6(4), pp.165-171.
- [17] Chen, T. and Guestrin, C., 2016a, August. Xgboost: A scalable tree boosting system. In *Proceedings of the 22nd acm sigkdd international conference on knowledge discovery and data mining* (pp. 785-794).
- [18] Truong, Q., Nguyen, M., Dang, H. and Mei, B., 2020. Housing price prediction via improved machine learning techniques. *Procedia Computer Science*, 174, pp.433-442.
- [19] Lu, S., Li, Z., Qin, Z., Yang, X. and Goh, R.S.M., 2017, December. A hybrid regression technique for house prices prediction. In *2017 IEEE International Conference on Industrial Engineering and Engineering Management (IEEM)* (pp. 319-323). IEEE.
- [20] Wolpert, D.H., 1992. Stacked generalization. *Neural networks*, 5(2), pp.241-259.
- [21] Moody, J. (2019). What does RSME really mean? [online] viewed 11 September 2021, <<https://towardsdatascience.com/what-does-rmse-really-mean-806b65f2e48e>>
- [22] Zhang, A. (2012). Evaluating Machine Learning Models. [online] O'Reilly Online Learning. Available at: <https://www.oreilly.com/library/view/evaluating-machine-learning/9781492048756/ch04.html> [Accessed 28 Jun. 2021].
- [23] Koehrsen, W. (2018). A Conceptual Explanation of Bayesian Hyperparameter Optimization for Machine Learning. [online] viewed 06 June 2021, <[A Conceptual Explanation of Bayesian Hyperparameter Optimization for Machine Learning | by Will Koehrsen | Towards Data Science](#)>

