

UNIVERSITY OF  
COPENHAGEN



A PHYSICIST'S APPROACH TO MACHINE LEARNING  
Understanding The Basic Bricks

CHRISTIAN MICHELSSEN  
Master's Thesis (Cand. Scient.)  
January 3<sup>rd</sup>, 2020

Supervised by  
Troels Petersen



UNIVERSITY OF COPENHAGEN

Copyright © 2020  
Christian Michelsen

[HTTPS://GITHUB.COM/CHRISTIANMICHelsen](https://github.com/CHRISTIANMICHelsen)

This thesis was inspired by the works of Edward R. Tufte and is based on the Tufte-L<sup>A</sup>T<sub>E</sub>X package.

*First printing, January 2020*

## *Abstract*

Here will be a decent abstract at some point<sup>TM</sup>.

*(This page is intentionally left blank).*

# *Contents*

<i>Abstract</i>	iii
<i>Table of Contents</i>	v
<i>Foreword</i>	xi
1 <i>Introduction</i>	1
<i>Part I</i>	3
2 <i>Machine Learning</i>	5
2.1 <i>Statistical Learning Theory</i> . . . . .	5
2.2 <i>Supervised Learning</i> . . . . .	6
2.3 <i>Generalization Bound</i> . . . . .	7
2.3.1 <i>Generalization Bound for Infinite Hypotheses</i> . . . . .	9
2.4 <i>Avoiding overfitting</i> . . . . .	10
2.4.1 <i>Model Regularization</i> . . . . .	10
2.4.2 <i>Cross Validation</i> . . . . .	12
2.4.3 <i>Early Stopping</i> . . . . .	13
2.5 <i>Loss functions</i> . . . . .	14
2.5.1 <i>Evaluation Function</i> . . . . .	16
2.6 <i>Decision Trees</i> . . . . .	16
2.6.1 <i>Ensembles of Decision Trees</i> . . . . .	17
2.7 <i>Hyperparameter Optimization</i> . . . . .	19
2.7.1 <i>Grid Search</i> . . . . .	20
2.7.2 <i>Random Search</i> . . . . .	20
2.7.3 <i>Bayesian Optimization</i> . . . . .	21
2.8 <i>Feature Importance</i> . . . . .	23

3	<i>Danish Housing Prices</i>	27
3.1	<i>Data Preparation and Exploratory Data Analysis</i> . . . . .	28
3.1.1	<i>Correlations</i> . . . . .	30
3.1.2	<i>Validity of Input Variables</i> . . . . .	32
3.1.3	<i>Selection Criteria</i> . . . . .	33
3.2	<i>Feature Augmentation</i> . . . . .	33
3.2.1	<i>Time-Dependent Price Index</i> . . . . .	34
3.3	<i>Evaluation Function</i> . . . . .	35
3.4	<i>Initial Hyperparameter Optimization</i> . . . . .	36
3.5	<i>Hyperparameter Optimization</i> . . . . .	38
3.6	<i>Results</i> . . . . .	40
3.7	<i>Model Inspection</i> . . . . .	43
3.8	<i>Multiple Models</i> . . . . .	45
3.9	<i>Discussion</i> . . . . .	47
3.10	<i>Conclusion</i> . . . . .	49
<i>Part II</i>		51
4	<i>Particle Physics and LEP</i>	53
4.1	<i>The Standard Model</i> . . . . .	53
4.2	<i>Quark Hadronization</i> . . . . .	54
4.3	<i>The ALEPH Detector and LEP</i> . . . . .	56
4.4	<i>Jet Clustering</i> . . . . .	58
4.5	<i>The Variables</i> . . . . .	58
5	<i>Quark Gluon Analysis</i>	63
5.1	<i>Data Preprocessing</i> . . . . .	63
5.2	<i>Exploratory Data Analysis</i> . . . . .	64
5.2.1	<i>Dimensionality Reduction</i> . . . . .	66
5.2.2	<i>Correlations</i> . . . . .	67
5.3	<i>Loss and Evaluation Function</i> . . . . .	67
5.4	<i>b</i> - <i>Tagging Analysis</i> . . . . .	68
5.4.1	<i>b</i> - <i>Tagging Hyperparameter Optimization</i> . . . . .	69
5.4.2	<i>b</i> - <i>Tagging Results</i> . . . . .	70
5.4.3	<i>b</i> - <i>Tagging Model Inspection</i> . . . . .	71
5.5	<i>b</i> - <i>Tagging Efficiency</i> . . . . .	73
5.6	<i>g</i> - <i>Tagging Analysis</i> . . . . .	75
5.6.1	<i>Permutation Invariance</i> . . . . .	75
5.6.2	<i>Truncated Uniform PDF</i> . . . . .	76

5.6.3	<i>g</i> -Tagging Hyperparameter Optimization . . . . .	76
5.6.4	<i>PermNet</i> . . . . .	78
5.6.5	1D Comparison of LGB and <i>PermNet</i> . . . . .	78
5.6.6	<i>g</i> -Tagging Results . . . . .	79
5.7	<i>g</i> -Tagging Efficiency . . . . .	81
5.8	Generalized Angularities in 3-Jet Events . . . . .	83
5.9	Gluon splitting . . . . .	85
5.9.1	Variables . . . . .	85
5.9.2	Efficiencies . . . . .	87
5.9.3	Closure Test . . . . .	88
5.9.4	Results . . . . .	90
5.10	Discussion . . . . .	92
5.11	Conclusion . . . . .	94
A	<i>Housing Prices Appendix</i>	95
A.1	<i>Input Distributions</i> . . . . .	95
A.2	<i>Correlations</i> . . . . .	110
A.3	<i>Validity Correlations</i> . . . . .	111
A.4	<i>Energy Rating Mapping</i> . . . . .	112
A.5	<i>Facebook Prophet for Houses</i> . . . . .	113
A.6	<i>Final Variable List</i> . . . . .	114
A.7	<i>Initial Hyperparameter Results Tables</i> . . . . .	116
A.8	<i>Initial Hyperparameter Results Parallel Coordinates</i> .	120
A.9	<i>Hyperparameter Results Parallel Coordinates</i> . . .	121
A.10	<i>Hyperparameter Optimization as a Function of Iterations</i> . . . . .	123
A.11	<i>Histogram of z Values for Apartments with the Tight Cut</i> . . . . .	125
A.12	<i>Multiple Models for Houses</i> . . . . .	125
A.13	<i>Performance Metrics for the Tight Selection</i> . . . .	126
B	<i>Quarks vs. Gluons Appendix</i>	127
B.1	<i>Fraction of Different Types of Jets</i> . . . . .	127
B.2	<i>Random Search PDFs for XGB</i> . . . . .	127
B.3	<i>Global SHAP Feature Importances for the g-Tagging Models in 3-Jet Events</i> . . . . .	127
B.4	<i>UMAP Parameters Comparison</i> . . . . .	128
B.5	<i>t-SNE Parameters Comparison</i> . . . . .	128
B.6	<i>HPO results for 3-jet b-Tagging</i> . . . . .	129
B.7	<i>HPO results for 4-jet b-Tagging</i> . . . . .	129
B.8	<i>b-Tag Scores in 3-Jet Events</i> . . . . .	130

<i>B.9 b-Tag Scores in 4-Jet Events . . . . .</i>	130
<i>B.10 ROC curve for 3-jet b-tagging . . . . .</i>	131
<i>B.11 ROC curve for 4-jet b-tagging . . . . .</i>	131
<i>B.12 Distribution of b-Tags in 3-Jet Events . . . . .</i>	132
<i>B.13 Distribution of b-Tags in 4-Jet Events . . . . .</i>	132
<i>B.14 Global Feature Importances for the LGB b-Tagging Algorithm on 3-Jet Events . . . . .</i>	133
<i>B.15 Global Feature Importances for the LGB b-Tagging Algorithm on 4-Jet Events . . . . .</i>	133
<i>B.16 b-tag Efficiency for b-Jets in the b-Signal Region for 3-jet events . . . . .</i>	134
<i>B.17 b-tag Efficiency for g-Jets in the g-Signal Region for 3-jet events . . . . .</i>	134
<i>B.18 b-tag Efficiency for b-Jets in the g-Signal Region for 3-jet events . . . . .</i>	135
<i>B.19 b-tag Efficiency for g-Jets in the b-Signal Region for 3-jet events . . . . .</i>	135
<i>B.20 Parallel Coordinate Plot of HPO Results for 3-Jet g-Tagging for Energy Ordered Jets . . . . .</i>	136
<i>B.21 Parallel Coordinate Plot of HPO Results for 3-Jet g-Tagging for Shuffled Jets . . . . .</i>	136
<i>B.22 Parallel Coordinate Plot of HPO Results for 4-Jet g-Tagging for Energy Ordered Jets . . . . .</i>	137
<i>B.23 Parallel Coordinate Plot of HPO Results for 4-Jet g-Tagging for Shuffled Jets . . . . .</i>	137
<i>B.24 PermNet Architecture . . . . .</i>	138
<i>B.25 1D LGB and PermNet Comparison for 4-Jet Events . .</i>	139
<i>B.26 1D LGB and PermNet Comparison for 3-Jet Events . .</i>	140
<i>B.27 g-Tag Scores in 4-Jet Events . . . . .</i>	141
<i>B.28 g-Tag Scores in 3-Jet Events . . . . .</i>	141
<i>B.29 ROC Curve for g-Tag in 4-Jet Events . . . . .</i>	142
<i>B.30 ROC Curve for g-Tag in 3-Jet Events . . . . .</i>	142
<i>B.31 Distribution of g-Tag Scores in 4-Jet Events for Signal and Background . . . . .</i>	143
<i>B.32 Distribution of g-Tag Scores in 3-Jet Events for Signal and Background . . . . .</i>	143
<i>B.33 Generalized Angularities for Charged Gluons Tracks .</i>	144
<i>B.34 Generalized Angularities for Neutral Gluons Clusters . . . . .</i>	146
<i>B.35 g-Tagging Efficiency for 4-Jet Events in MC . . . . .</i>	148

<i>B.36 Closure Plot for g-Tagging Model in 4-Jet Events . . . . .</i>	150
<i>B.37 Gluon Splitting Distribution Comparison in MC and Data for Phase Space Region A . . . . .</i>	153
<i>B.38 Gluon Splitting Distribution Comparison in MC and Data for Phase Space Region B . . . . .</i>	154
<i>B.39 Gluon Splitting Distribution Comparison in MC and Data for Phase Space Region C . . . . .</i>	155
<i>B.40 Gluon Splitting Distribution Comparison in MC and Data for Phase Space Region D . . . . .</i>	156
<i>List of Figures</i>	163
<i>List of Tables</i>	166
<i>Bibliography</i>	167

*(This page is intentionally left blank).*

# Part I

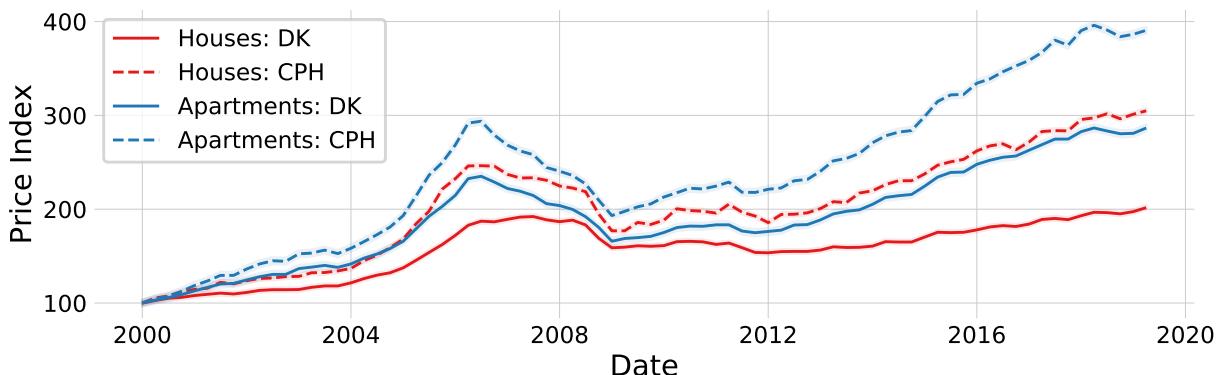
*(This page is intentionally left blank).*

### 3. Danish Housing Prices

*"Buy land, they're not making it anymore."*

— Mark Twain

HOUSING MARKETS have always been a playing field for economists, property speculators, real estate agents, and realtors. Personally, I am by no metric any of these, not even close to. Yet, the issue of estimating housing prices purely on a data driven basis was too interesting to ignore. Estimating housing prices is a classical economical discipline as seen in papers from the Danish National Bank developing a regional model of the Danish housing market [58] to an analysis of the financial crisis in 2008-2009 and its effects on the Copenhagen metropolitan area [73].



If one takes a look at the time development of the Danish housing market, the Danish governmental organization for statistics, Statistics Denmark, releases a price index [44] for both one-family houses (OFH) and owner-occupied apartments (OOA) quarterly, see Figure 3.1. Here it is easy to see the effect of the financial crisis around 2008, but also the steady increase in the housing market in both Copenhagen and the entire country since then. Housing in this context means both actual houses and privately owned apartments, and will be called residences in general in this project.

The goal of this subproject is not to predict any future collapse in the financial markets as we saw upwards of 10 years ago. Instead, it is to learn patterns in the price of houses in steady times. The goal is training a computer to automatically<sup>1</sup> be able to find these patterns and see if we can improve this model.

Figure 3.1: Prince index of Danish one-family houses and owner-occupied apartments where **houses** are shown in red and **apartments** in blue, where full lines are for the entirety of Denmark and dashed lines are only for Copenhagen. Errorbars (scaled up with a factor of 2) are shown as colored bands. The price index and its uncertainty is based on numbers from DST [44], however, rescaled to 100 in 2000 (instead of 2006 as it was in the data).

<sup>1</sup> In contrary to Hviid [58], Mulaic et al. [73] and others who base their models on macro-economic principles.

This chapter is organized as follows. The data will be introduced in [section 3.1](#) where also some initial plots are shown. The data set is feature augmented in [section 3.2](#) and the different evaluation functions will be discussed in [section 3.3](#). The choice of evaluation function will be based on an initial hyperparameter optimization process in [section 3.4](#). The model is fitted and optimized in [section 3.5](#) and the results from the final model are presented in [section 3.6](#). The forecasting performance of the model will be tested in XXX **TODO!**. To better understand the model, [section 3.7](#) deals with the issue of model inspection. The model is extended by combining multiple ML models into an ensemble, a process that is described in [section 3.8](#). Finally, the project concerning Danish housing prices will be discussed in [section 3.9](#) and concluded in [section 3.10](#).

### 3.1 Data Preparation and Exploratory Data Analysis

*“80 % of data science is cleaning the data and 20 % is complaining about cleaning the data.”*

---

— Anthony Goldbloom, Kaggle

The first part of any data science project is actually getting the data and being able to read it. This has been an iterative process that has improved over time. The last data transfer we got was September 3<sup>rd</sup> 2019 which consisted of a 522.4 MB CSV file with dimensions (711 212, 171). This section will go through the data cleaning process.

Before any further data analysis is performed, all of the data is loaded, except columns which only contain internal information for Boligsiden<sup>2</sup>. To get a better understanding of all the original input variables, histograms of the one-dimensional distributions of each variable are shown in Figures<sup>3</sup> [A.1](#) to [A.14](#).

Four particularly interesting distributions are shown in Figure [3.2](#): the date of the sale `SalgsDato` in subfigure [\(a\)](#), the type of residence `SagtypeNr` in subfigure [\(b\)](#), the longitude of the residence `GisX_WGS84` in subfigure [\(c\)](#), and the area of the residence `ArealBolig` in subfigure [\(d\)](#).

The distribution of the date of sale, Figure [3.2 \(a\)](#), is an interesting variable because it shows how Boligsiden has been collecting more and more data over time. Here 2007 and 2019 are clear outliers since their current database only contains sales from the end of 2007, and 2019 only contains data from the first eight months of the year. The `SagTypeNr` is a discrete code that Boligsiden uses to differentiate between different types of residences. The mapping between code and description is shown in Table [3.1](#).

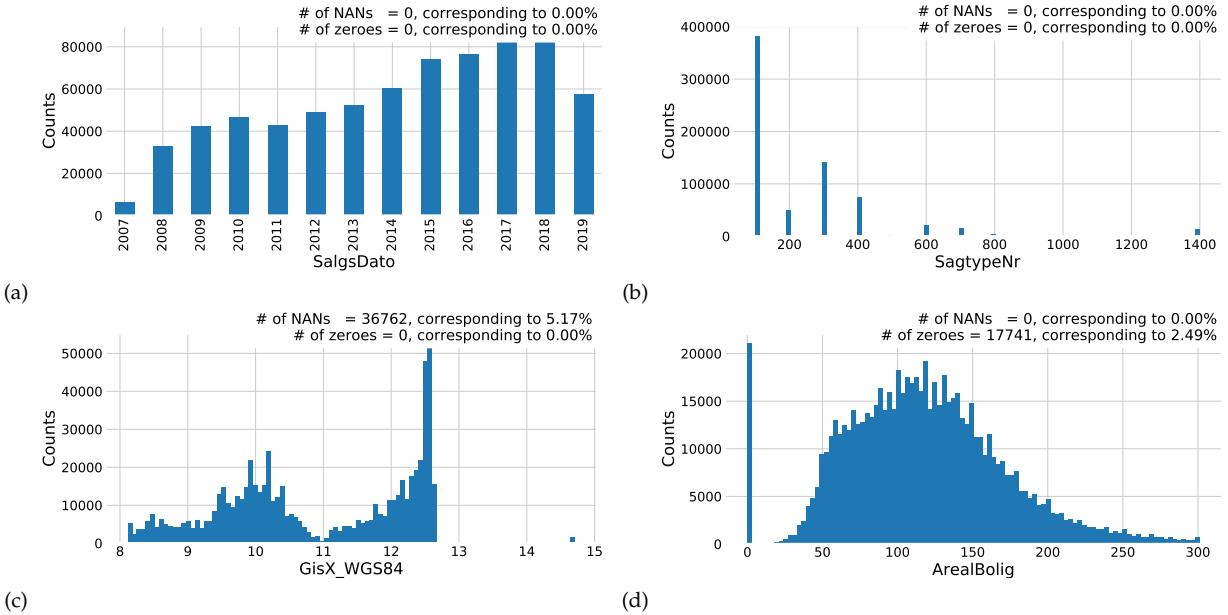
In this project only one-family houses (“Villas”) with code 100 and owner-occupied apartments (“Ejerlejlighed”) with code 300 are considered. In Figure [3.2 \(b\)](#) it is shown that these two types of residences are also the most frequent sales with close to 400 000 and 150 000 sales in total. The longitude distribution, Figure [3.2 \(c\)](#), is

<sup>2</sup> The variables `Sag_Kvhx`, `Bygning_GOP_Matrikelnr`, and `Enhed_GOP_BoligtypeKode`.

<sup>3</sup> Where the A indicates that the figures are in the the [appendix](#).

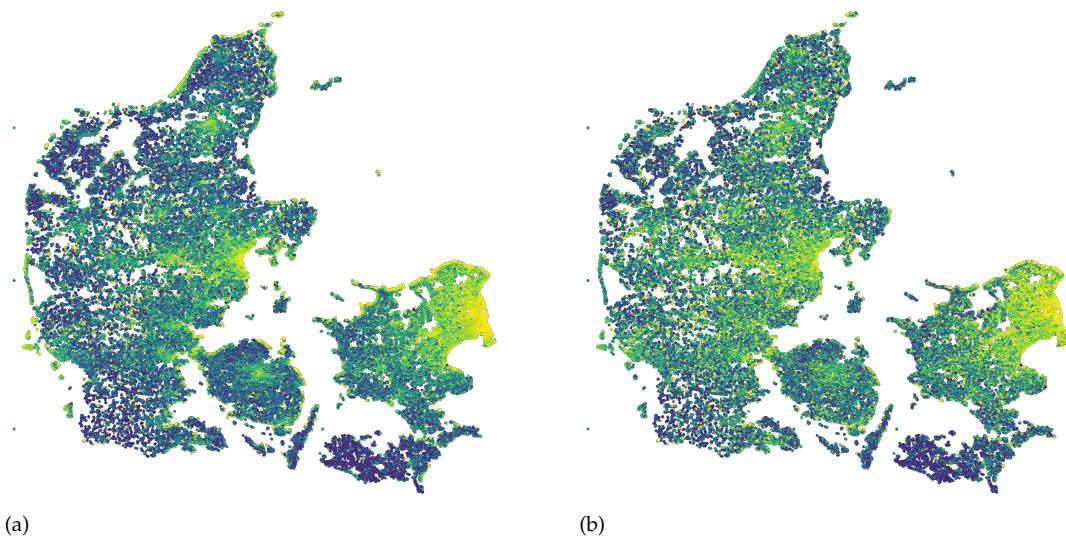
Code	Name
100	Villa
200	Rækkehus
300	Ejerlejlighed
400	Fritidsbolig
401	Kolonihave
500	Andelsbolig
600	Landejendom
700	Helårsgrund
800	Fritidsgrund
900	Villalejlighed
1000	Kvæggård
1100	Svinegård
1200	Planteavlsgård
1300	Skovejendom
1400	Lystejendom
1500	Specialejendom

Table 3.1: Mapping between the code in `SagTypeNr` and the type of residence. The two important types of residences are “villa” (one-family houses) and “ejerlejlighed” (owner-occupied apartments).



mostly interesting due to fact that it clearly shows how the Great Belt and especially the Baltic Sea separates Denmark into three parts; the Western part, the Eastern part, and then Bornholm. Note that more than 5 % of the residences' locations are unknown values, so-called "Not A Number"s (NANs). The distribution of the area, Figure 3.2 (d), shows that most residences are between 50 m<sup>2</sup> and 200 m<sup>2</sup>, as expected in Denmark. However, a relatively large part of the residences, 2.5 %, are listed as having an area of 0 m<sup>2</sup> which are obviously erroneous entries.

Figure 3.2: Distributions of four out of the 168 input variables. Subplot (a) shows the date of the sale, Subplot (b) shows the type of residence, Subplot (c) shows the longitude, Subplot (d) shows the area of the house.

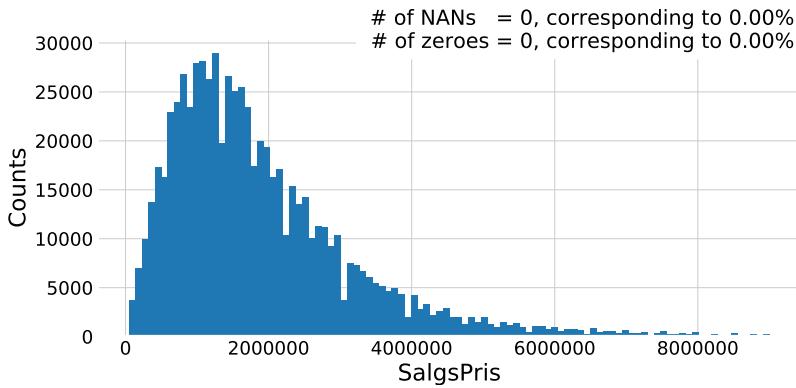


The geographic distribution of sales are shown in Figure 3.3. The residences are coloured according to the square meter price in Figure 3.3 (a) and according to the sales price in Figure 3.3 (b). Notice the strong correlation between the distance to water and the square meter price, a correlation that is less visible when looking at the sales

Figure 3.3: Geographic distribution of the sold residences. In subplot (a) the sales are colored according to their square meter price and in subplot (b) according to the sales price.

price. Since these plots each contain 674 647 points<sup>4</sup>, over-plotting quickly becomes an issue. To circumvent this, the software package called DataShader [8] was used which in a simple, consistent, and not at least computationally efficient manner allows one to plot big data.

The most important of the features is the sales price, called `SalgsPris` in the dataset. Its distribution is shown in Figure 3.4. This is a positively skewed distribution that shares visual similarities with a log-normal distribution. The mode of the all sales prices<sup>5</sup> is 1.1 M.kr. and the median is 1.6 M.kr. The mean is 2.0 M.kr. but this value is heavily influenced by a few very high values.



<sup>4</sup> Only sales with a valid GPS-coordinate and area of residence are shown

<sup>5</sup> Measured in millions DKK, M.kr.

Figure 3.4: Histogram of prices of houses and apartments sold in Denmark.

### 3.1.1 Correlations

Having shown the 1D-distributions of all the different variables in the previous section, the next step would be to look at the correlations between the variables. Since there are 168 input variables, it is almost impossible to understand every inter-variable correlation, however, it is tried in Figure A.15. In this figure, the correlations  $\rho$  between all numerical variables that are not obviously related to other variables (like the GPS-coordinates that are in both latitude-longitude and ETRS89<sup>6</sup> format), are plotted as a  $(86 \times 86)$ -dimensional heatmap<sup>7</sup>.

Even though inter-variable correlations are important in the exploratory data analysis (EDA) phase, what is more important is to get a better understanding of how the input variables correlate with the output variable; the sales price. This is shown in Figure 3.5 for the variables where  $|\rho| > 10\%$ . It is the previous property evaluation, `EjdVurdering_EjendomsVaerdi` that is correlated the most with the sales price, which does not come as any huge surprise. Other positively correlated variables are the cost of ownership<sup>8</sup>, area, number of bathrooms, longitude, and distance to nearest wind mill. In the other end, the local income tax, `Kommune_SkatteProcent`, is the variable that is the most negatively correlated to the sales price, followed by the geographical variables related to province, municipality, and postcode.

The correlation  $\rho$  used above is the linear correlation which only captures linear relationships between variables. All modern machine

<sup>6</sup> European Terrestrial Reference System 1989.

<sup>7</sup> With the condition that each plotted variable has to have at least one inter-variable correlation higher than  $|\rho| > 30\%$

<sup>8</sup> This is a great example of the fact that correlation does not imply causation as simply increasing the cost of ownership does not increase the sales price.

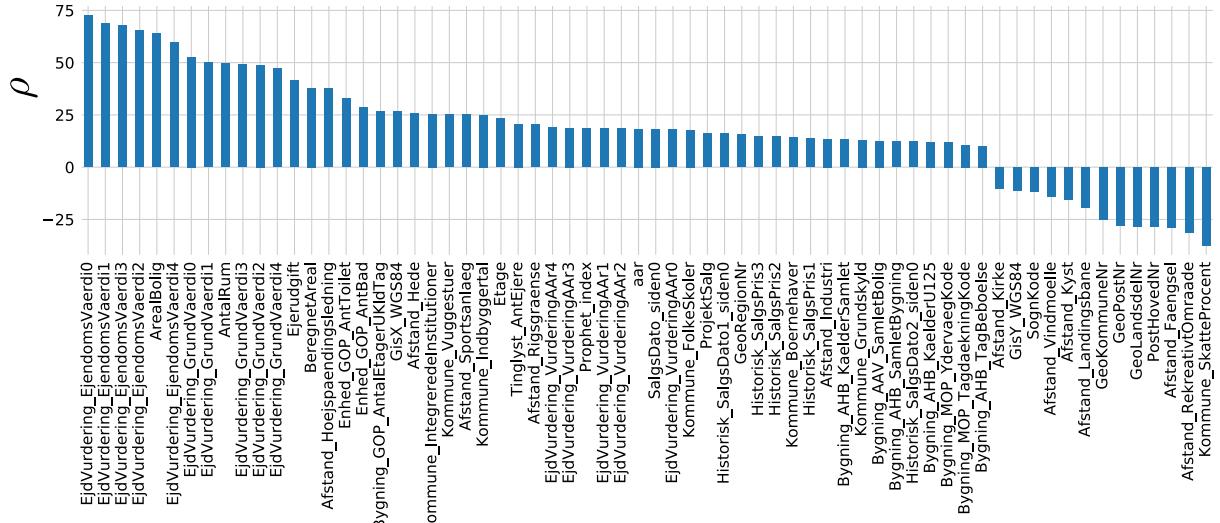


Figure 3.5: Linear correlation  $\rho$  between each variable and the sales price for variables where  $|\rho| > 10\%$ .

learning algorithms, however, are also able to capture higher-order correlations and thus a higher-order correlation measure is needed. The maximal information coefficient (MIC), which is a value between 0 and 1, is such a nonlinear measure<sup>9</sup>. It finds correlation based on the intuition that if two variables are correlated, it should be possible to split the data up into smaller grids where, if they are correlated, the grid that contain points should contain many points and the rest of the grids should be (relatively) empty [80]. This is in comparison to two uncorrelated variables which would simply display noisy behavior and only have few grids with many points in. Albanese et al. [12] extended on this idea and developed the computationally efficient algorithm called MICtools which computes the estimator  $\text{MIC}_e$  for MIC. An example of this nonlinear correlation is seen in Figure 3.6. Here the relationship between the normal linear correlation  $\rho$  and  $\text{MIC}_e$  are shown for four synthetic datasets. Notice that  $\text{MIC}_e$  performs particularly well for the sine wave, and decent for the line and parabola, but only slightly captures the relationship for the exponential growth. The influence of noise on  $\rho$  and  $\text{MIC}_e$  is shown in Figure A.16.

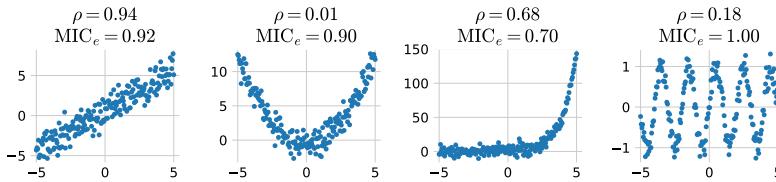


Figure 3.6: Comparison of the linear correlation  $\rho$  and the nonlinear  $\text{MIC}_e$  for a straight line, a parabola, an exponential, and a sine wave, all with noise added. See also Figure A.16.

Using  $\text{MIC}_e$  as the correlation measure between the numerical variables and the sales prices, the variables with a  $\text{MIC}_e$ -score higher than 10 % are shown in Figure 3.7. Again, the previous property evaluations are the most correlated features to the sales price followed by the parish code, `Sognkode`. In general the geographical variables score high here with also the post code and municipality number.

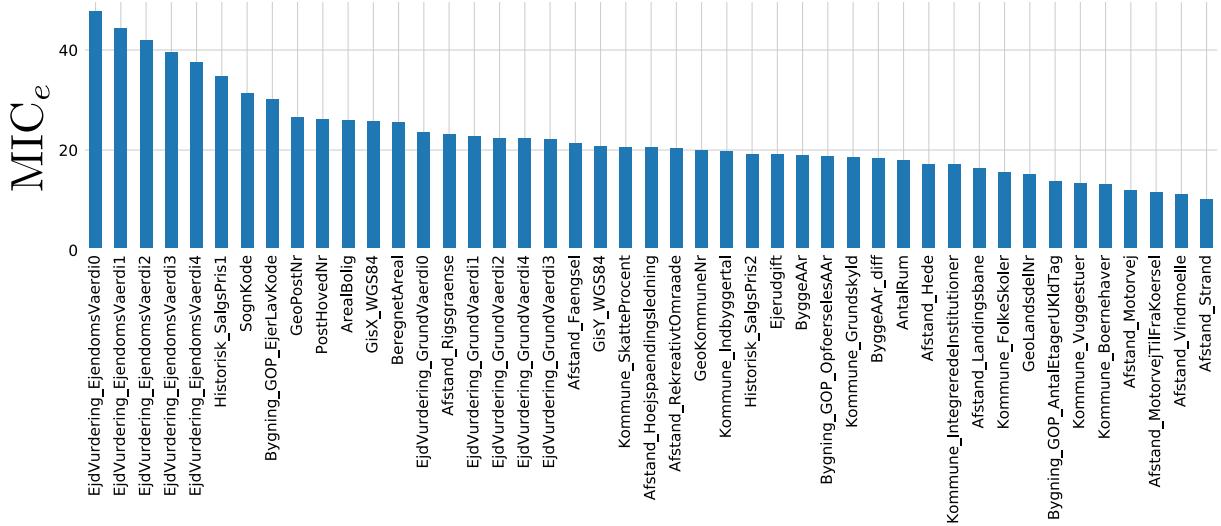
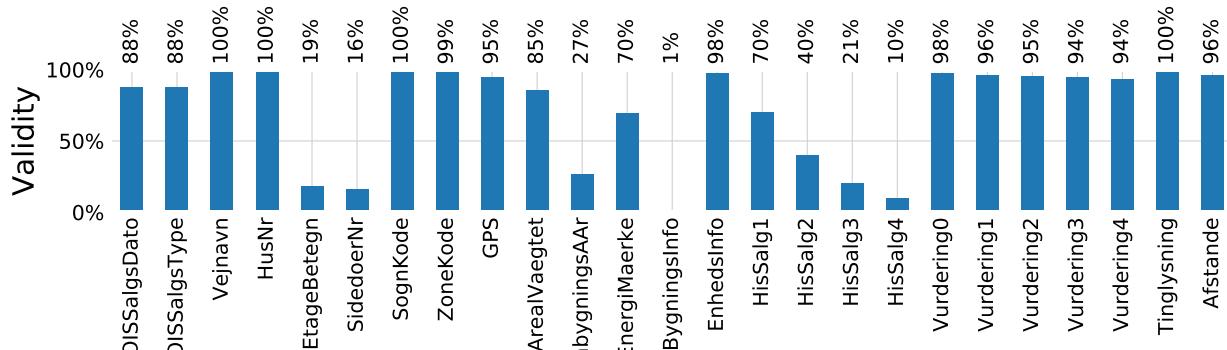


Figure 3.7: Nonlinear correlation between variables and price using Maximal Information Coefficient (MIC) for variables where  $\text{MIC} > 10\%$ .

### 3.1.2 Validity of Input Variables

The fact that some of the variables contain considerable amounts of invalid values, NaNs, requires this to be taken into account before any further analysis. The validity, defined as the percentage of valid observations, of every variable is shown in Figure 3.8. Here the 168 variables are grouped together into 25 groups where each group share the same validity. An example of this is all of the 16 different distance-variables<sup>10</sup>. We see that most of the variables have validities around more than 85 %, however, a few of the variables, especially information about the building, `BygningsInfo`, have validities less than 20 %.



<sup>10</sup> Distance to: prison, heath, high-voltage transmission line, industry, visible railroad, church, churchyard, coast, landing strip, motorway, access to motorway, recreational area, border, sports centre, beach, and windmill.

To see how closely related the different validity groups are, one can look at the dendrogram in Figure 3.9. The dendrogram is based on a hierarchical clustering algorithm [100] where the different groups are clustered according to the linear correlation of their validity. This diagram is supposed to be read in a top-down approach, where it can be seen that the name of the street, `Vejnavn`, and the number of the residence, `HusNr`, correlate a lot and are thus clustered very early. The year of the last time the residence was rebuilt or greatly

Figure 3.8: Percentage of valid counts for each variable grouped together in categories.

modified, `OmbygningsAAr`, does not correlate strongly with any of the other variables and is thus the last variable to be clustered together. A heatmap of the validity inter-variable correlations is shown in Figure A.17.

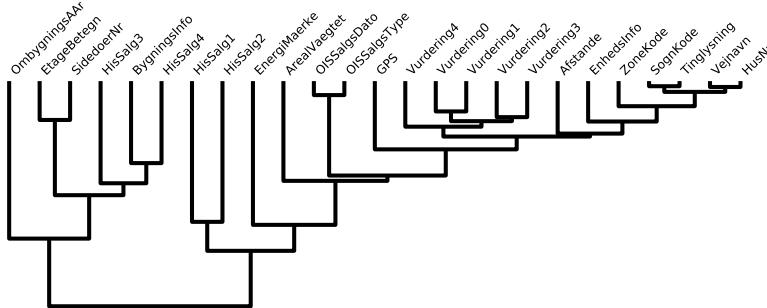


Figure 3.9: Validity Dendrogram based on hierarchical clustering of the linear correlation of validity for the housing price variables clustered together.

### 3.1.3 Selection Criteria

Given the 1D input variable distributions and their validity, we apply some very basic selection criteria before any further analysis. These selection criteria are shown in Table 3.2.

	Description	Remaining	Removed
Area	$20 \text{ m}^2 \leq \text{Area} \leq 500 \text{ m}^2$	689 140	23 666
Price	$0.1 \text{ M.kr.} \leq \text{Price} \leq 100 \text{ M.kr.}$	687 546	1 594
Type	Has sales type 1	605 415	82 131
GPS	Has valid GPS coordinates	578 860	26 555
Private	Only non-business sales	549 140	29 720
Time	Sold in 2009 or later	520 548	28 592

The sales type, `OISSalgsType`, is an OIS<sup>11</sup> code which describes what type of sale it is: when it is 1 it is considered a normal sale, compared to e.g. forced sales. These cuts are seen as the minimum requirements for what constitutes a curated dataset with no obvious outliers. The reason why the time requirement is applied is to reduce the effect of the financial crisis to creep into the model.

## 3.2 Feature Augmentation

The analysis have dealt with different types of residences all together until now. The rest of the analysis will be applied on single-family houses and owner-occupied apartments independently from this point on.

First, invalid counts are dropped such that variables which contain more than 10 % NaNs are dropped, and duplicate rows are also removed. Then some manual features are added based on existing features. The day of the month, the month, and the year are extracted from the sales date. The sales date is converted to the numbers of days since January 1<sup>st</sup>, 2009. From the number of the house, `HusNr`, the value is extracted along with a boolean flag indicating whether

Table 3.2: Overview of the basic selection criteria which define the minimum information needed to predict the price of a sale.

<sup>11</sup> OIS is short for “Den Offentlige Informationsserver”, the Danish public information server, and it collects information about Danish residences [9].

String	Explanation	Code
NAN	No side door	0
TH, TV	Right, left	11
MF	Center	12
-	The rest	15

Table 3.3: Side door mapping. If the side door string contains e.g. “TH” this gets the code 11.

or not it includes a letter (eg. "27B"). The number of the side door, `SidedoerNr`, is formatted according to Table 3.3 and the road name according to Table 3.4. The age of the house is added<sup>12</sup> and the amount of time (in years) since last major modification. The energy rating label, `EnergiMaerke`, is also converted from strings to values according to Table A.1.

Finally, some of the variables in the dataset are not suitable for machine learning (or are simply transformations of other variables) and are thus dropped. These are variables such as the ID, when the house was deleted at Boligsiden, or the cash price.

### 3.2.1 Time-Dependent Price Index

In addition to the manual data augmentation, a time-dependent price index is also added. We make use of the open source package called Prophet made by Taylor and Letham [89] at Facebook. It is based on a decomposable time series model [53] with two<sup>13</sup> components; the trend  $g(t)$  and the seasonality  $s(t)$ :

$$y_{\text{Prophet}}(t) = g(t) + s(t) + \epsilon_t, \quad (3.1)$$

where  $\epsilon_t$  is a normally distributed error term. Taylor and Letham [89] fit this equation with a generalized additive model (GAM) [54] which they argue has several practical advantages compared to ARIMA<sup>14</sup> models which are commonly used in economics [71].

We fit the Prophet model on the weekly median price pr. square meter (PPSM) up until (and including) 2017. The results of the prophet model fitted on (owner-occupied) apartments are seen in Figure 3.10 and Figure 3.11. In Figure 3.10 the weekly median price pr. square meter for apartments are shown as black dots with the fitted Prophet model shown in blue. The  $1\sigma$  uncertainty intervals are shown as the transparent blue band. The Prophet model not only allows predicting previous and future PPSMs, it also return the uncertainty of this prediction. The future predictions for the PPSM are the values after 2018. The trend and seasonality of the model are shown in Figure 3.11 where the top plot is the overall trend  $g(t)$  and the bottom plot is the seasonality  $s(t)$ . Whereas the trend just continues to rise, the seasonality shows that residences are generally sold for a higher price in the Summer months compared to the Winter months. The Prophet model plots for one-family houses are shown in Figure A.18 and A.19.

Using the Prophet model, we define the price index (PI) to be the Prophet-predicted PPSM,  $y_{\text{Prophet}}$ , for each residence normalized by the mean to give values around 1:

$$\text{PI}(t) = \frac{y_{\text{Prophet}}(t)}{\langle y_{\text{Prophet}}(t) \rangle}, \quad (3.2)$$

where  $\langle \cdot \rangle$  refers to the average. The price index thus works as a measure of the national price for houses or apartments at a given time and is added as a variable to the dataset.

<sup>12</sup> In addition to only having the year the house was built.

Contains	Explanation	Code
Vej	Road	0
Gade	Street	1
Alle, Allé	Avenue	2
Boulevard	Boulevard	3
-	The rest	-1

Table 3.4: Street mapping. If the street name contains e.g. "Vej" this gets the code 0.

<sup>13</sup> In their paper, Taylor and Letham [89] include a holiday component in their analysis as well which is not included in this project.

<sup>14</sup> AutoRegressive Integrated Moving Average.

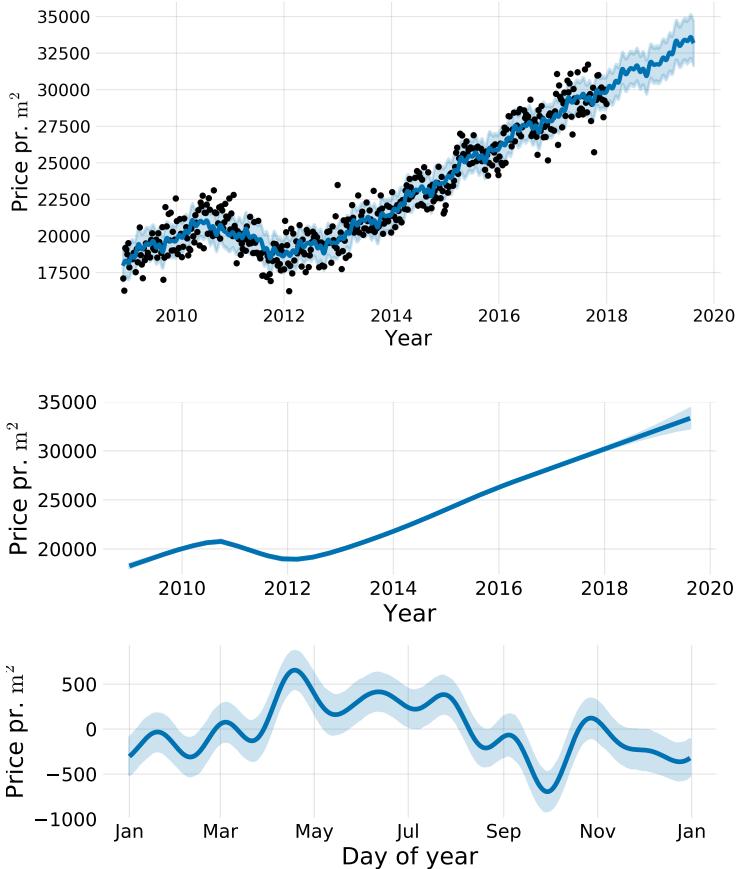


Figure 3.10: The predictions of the Facebook Prophet model trained on square meter prices for owner-occupied apartments sold before January 1st, 2018. The data is down-sampled to weekly bins where the median of each week is used as in input to the Prophet model, shown seen as black dots in the figure. The model's forecasts for 2018 and 2019 are shown in blue with a light blue error band showing the  $1\sigma$  confidence interval.

Figure 3.11: The trends of the Facebook Prophet model trained on square meter prices for owner-occupied apartments sold before January 1st, 2018. In the top plot is the overall trend as a function of year and in the bottom plot is the yearly variation as a function of day of year. It can be seen that the square meter price is higher during the Summer months compared to the Winter months, however, compared to the overall trend this effect is minor (< 10%).

### 3.3 Evaluation Function

The choice of evaluation function  $f_{\text{eval}}$  is an important decision. The evaluation function will be based on the relative prediction  $z_i$ :

$$z_i = \frac{\hat{y}_i - y_i}{y_i}, \quad (3.3)$$

where  $y_i$  is the true price and  $\hat{y}_i$  the predicted one. The relative prediction is defined such that it is positive when  $\hat{y}_i > y_i$ . Due to the outlier selection criteria applied earlier, the denominator is made sure to always be positive (and never 0), and  $\mathbf{z}$  is expected to approximately follow a normal distribution<sup>15</sup>. Initially the mean of  $\mathbf{z}$  was considered as the choice of evaluation function  $f_{\text{eval}} \stackrel{?}{=} \text{mean}(\mathbf{z})$  though this only ensures a minimum of bias, not necessarily a low spread. This lead the discussion on to look at the standard deviation of  $z$  as  $f_{\text{eval}} \stackrel{?}{=} \text{std}(\mathbf{z})$ . The mean and the standard deviation, however, are not very robust estimators since they are heavily influenced by outliers. The mean (and thus also the standard deviation) has an *asymptotic breakdown point* at 0 %, where the breakdown point is defined as the smallest fraction<sup>16</sup> of bad observations that can cause an estimator to become arbitrarily small or large: a single outlier with an arbitrarily large value may cause the mean to diverge to that large value [57]. In comparison, the median has an asymptotic breakdown point of 50 % and is thus a more robust estimator of cen-

<sup>15</sup> Where  $\mathbf{z}$  is the vector of all relative predictions  $\mathbf{z} \in \mathbb{R}^N$ .

<sup>16</sup> Where the “asymptotic” in “asymptotic breakdown point” refers to when the number of samples goes to infinity.

trality. A robust measure of the variability or dispersion of a sample  $\mathbf{x}$  (compared to e.g. the standard deviation  $\sigma$ ) is the median absolute deviation (MAD) written as:

$$\text{MAD}(\mathbf{x}) = c \cdot \text{median}(|\mathbf{x} - \text{median}(\mathbf{x})|), \quad c = \frac{1}{\Phi^{-1}\left(\frac{3}{4}\right)}, \quad (3.4)$$

where  $c$  is a normalization constant to make MAD a consistent estimator of the standard deviation  $\sigma$  assuming normally distributed data and  $\Phi^{-1}$  is the percent point function<sup>17</sup> [65]. The MAD is thus the median of the absolute differences between the data and the median of the data. We are, however, not just interested in having the distribution of  $\mathbf{z}$  as narrow as possible, we also want it centered around 0. We thus continue with the following evaluation function:

$$\begin{aligned} \text{MAD}_0(\mathbf{x}) &\equiv c \cdot \text{median}(|\mathbf{x} - 0|) = c \cdot \text{median}(|\mathbf{x}|) \\ f_{\text{eval}}(\mathbf{z}) &\equiv \text{MAD}_0(\mathbf{z}) = c \cdot \text{median}(|\mathbf{z}|). \end{aligned} \quad (3.5)$$

To get an intuition about the size of a “good” value of  $\text{MAD}_0$ , one could calculate it comparing the asking price with the actual sales price. Doing so, one finds:  $f_{\text{eval}}(\mathbf{z}_{\text{OFH}}) = 11.35\%$  for houses and  $f_{\text{eval}}(\mathbf{z}_{\text{OOA}}) = 5.72\%$  for apartments. In some cases  $f_{\text{eval}}$  will still be referred to as MAD and it will be mentioned explicitly if the form in equation (3.4) is meant.

### 3.4 Initial Hyperparameter Optimization

With the initial cleaning and feature adding done, the shapes of the ML-ready datasets are: (291 317, 144) for houses (OFHs) and (114 166, 144) for apartments (OOA), both sharing the same variables. All of the variables which are used from this point on are shown in Table A.2. The data are split into training and test sets such that training is defined as every sale from before 2018, the test set is every sale from 2018, and, since more data came after the project started, 2019 is a small extra test set.

The number of observations for the different sets are shown in Table 3.5. Since the dataset has been shown to be quite noisy with a lot of invalid counts a *tight* selection of the data is also applied. The tight selection is defined as residences which are within the 1% to 99% quantiles of all<sup>18</sup> numeric variables with more than 3 unique values. The number of observations for the different tight sets are shown in Table 3.6.

A small study into the effect of some various hyperparameters was performed before any further fitting. This study investigated the effect of the old sales by assigning them a lower weight depending on time. It was investigated whether or not the model would perform better if samples got the time-dependent weight  $w(t)$  given by:

$$\begin{aligned} w'(t) &= e^{k \cdot t}, \quad k = \frac{\log 2}{T_{\frac{1}{2}}}, \\ w(t) &= \frac{w'(t)}{\langle w'(t) \rangle}, \end{aligned} \quad (3.6)$$

Here  $|\mathbf{x}|$  is the absolute value of  $x_i$  applied element-wise and not the norm of the vector  $\mathbf{x}$ .

<sup>17</sup> Inverse of the cumulative distribution function.

The robust estimator MAD is derived assuming symmetric distributions. Robust measures of the variability of non-symmetric samples have also been developed, see Rousseeuw and Croux [82] for more details.

	Houses	Apartments
Train	240 070	93 115
Test	34 628	14 183
2019	16 619	6868

Table 3.5: Number of observations for houses and apartments in the training, test, and 2019 set.

Tight	Houses	Apartments
Train	143 179	57 795
Test	20 338	8376
2019	9683	4030

Table 3.6: Number of observations for houses and apartments in the training, test, and 2019 set for the tight selection.

<sup>18</sup> Except the variables that contains the words: “aar” (year), “dato” (date), and “prophet”.

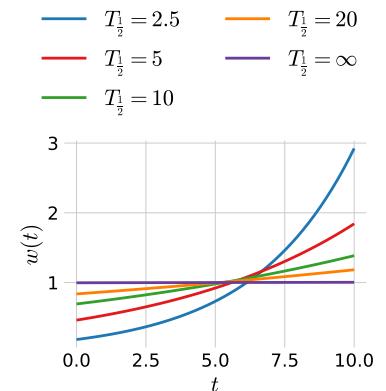


Figure 3.12: The sample weight  $w(t)$  as a function of time  $t$  where the time is in years after January 1<sup>st</sup>, 2009. Here seen plotted for different values of the half-life  $T_{\frac{1}{2}}$ .

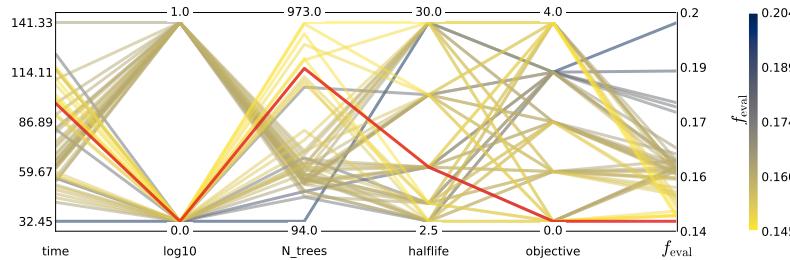
where  $T_{\frac{1}{2}}$  is the half-life. This is illustrated in Figure 3.12 for the different values of  $T_{\frac{1}{2}}$  used in the study. In addition to the weight, it was also investigated whether or not a  $\log_{10}$ -transformation of the price would increase performance. Finally the choice of loss function was also added to the study for the five different loss functions defined in section 2.5. A grid search<sup>19</sup> was performed for:

$$\begin{aligned} T_{\frac{1}{2}} &\in \{2.5, 5, 10, 20, \infty\} \text{ years} \\ \log_{10} &\in \{\text{True}, \text{False}\} \\ \ell &\in \{\ell_{\text{Cauchy}}, \ell_{\text{Fair}}, \ell_{\text{LogCosh}}, \ell_{\text{SE}}, \ell_{\text{Welsch}}\}. \end{aligned} \quad (3.7)$$

The grid search is performed for two boosted decision tree models based on XGBoost [36] (XGB), one for apartments and one for houses. They were each training with 5-fold cross validation and early stopping was applied with a patience of 100.

For apartments, the loss function with the lowest  $f_{\text{eval}}$  was the Cauchy loss with  $T_{\frac{1}{2}} = 10$  years and no  $\log_{10}$ -transformation. This model terminated by early stopping after 770 trees, see Table 3.7. For houses the loss function with the lowest  $f_{\text{eval}}$  was (also) the Cauchy loss with  $T_{\frac{1}{2}} = 20$  years and (also) no  $\log_{10}$ -transformation. This model terminated by early stopping after 1514 trees, see Table 3.8. It is interesting to note, that the best results for both houses and apartments share the same loss function and no  $\log_{10}$ -transformation.

The  $\log_{10}$ -transformation was included as some machine learning methods assume that the dependent variable,  $y$ , is normally distributed. Initially, it also showed better results than no transformation, however, this turned out to be a numerical consequence which was alleviated by dividing all prices with a million, such that  $y$  had units of M.kr instead of kr. All of the results for the apartments are shown in Tables A.3 to A.7 along with all of the results for the houses in Tables A.8 to A.12.



The results of the hyperparameter optimization are shown in the parallel coordinate plot in Figure 3.13. The hyperparameters (along the time taken and the number of trees) are plotted along the  $x$ -axis and the value of the hyperparameter on the  $y$ -axis. Every iteration of the grid search is thus a line on the plot. The lines are colored according to their evaluation score; the best hyperparameter is shown in red. For the hyperparameter `log10`, 0 means False and 1 means True, for `Halftime`  $\infty$  is mapped to 30, and for `objective` the functions Cauchy (0), Fair (1), LogCosh (2) SquaredError (3), and Welsch

<sup>19</sup> Grid search was acceptable since the parameter space is small and two of the three dimensions are non-numerical.

Half-life	$\log_{10}$	$N_{\text{trees}}$	$f_{\text{eval}}$
2.5	True	293	0.1598
2.5	False	814	0.1466
5	True	304	0.1610
5	False	923	0.1468
10	True	266	0.1610
<b>10</b>	<b>False</b>	<b>770</b>	<b>0.1450</b>
20	True	288	0.1613
20	False	967	0.1467
$\infty$	True	340	0.1601
$\infty$	False	807	0.1480

Table 3.7: Results of the initial hyperparameter optimization for apartments for the best loss function  $\ell_{\text{Cauchy}}$ . The best hyperparameter is shown in bold.

Half-life	$\log_{10}$	$N_{\text{trees}}$	$f_{\text{eval}}$
2.5	True	434	0.1991
2.5	False	1007	0.1872
5	True	350	0.1999
5	False	1130	0.1858
10	True	436	0.1992
10	False	1183	0.1850
20	True	397	0.2003
<b>20</b>	<b>False</b>	<b>1514</b>	<b>0.1833</b>
$\infty$	True	449	0.1992
$\infty$	False	1351	0.1844

Table 3.8: Results of the initial hyperparameter optimization for houses for the best loss function  $\ell_{\text{Cauchy}}$ . The best hyperparameter is shown in bold.

Figure 3.13: Hyperparameter optimization results for apartments. The results are shown as parallel coordinates with each hyperparameter along the  $x$ -axis and the value of that parameter on the  $y$ -axis. Each line is an iteration of the grid search, colored according to the performance of that hyperparameter as measured by the MAD. The **single best hyperparameter** is shown in red. For the hyperparameter `log10` 0 means False and 1 means True, for `Halftime`  $\infty$  is mapped to 30, and for `objective` the functions Cauchy (0), Fair (1), LogCosh (2) SquaredError (3), and Welsch

(4) are mapped to the integers in the parentheses.

(4) are mapped to the integers in the parentheses. The results of the grid search for houses are shown in Figure A.21.

What can be concluded from Figure 3.13 is that A) there is a clear preference to not log-transform the data, B) the models with many trees (selected by early stopping) generally performed better than the ones with fewer trees, C) it is difficult to see a clear pattern for the halflife, and that D) there seems to be a tendency for the Cauchy loss to be the best. What is not seen in the figure, however, are how the uncertainties of the different iterations also matter, where the uncertainties are the standard deviations (not of the mean) of the 5 folds in the cross validation.

For the half-life and the choice of objective function these are shown in Figure 3.14 and 3.15. In the first of the two figures it is easily seen that even though  $T_{\frac{1}{2}} = 10 \text{ yr}$  is the minimum, the uncertainties are so large that nothing can be concluded definitely with regards to the half-life parameter related to the weights  $w(t)$ . On the contrary, in the second figure there is a clear performance difference between the different loss functions where the Cauchy loss achieves the best (lowest) value of  $f_{\text{eval}}$  and especially the Squared Error is disregarded.

The rest of the machine learning models thus continue with the following hyperparameters:

	$\log_{10}$	Half-life	Loss function
Apartments	False	10 yr	Cauchy
Houses	False	20 yr	Cauchy

### 3.5 Hyperparameter Optimization

With the initial hyperparameter optimization finished, the actual training of the models was started. Again, two XGBoost models were fitted (one for apartments and one for houses) to the training set the hyperparameters were optimized using both random search and Bayesian optimization. These were each run for 100 iterations with 5-fold cross validation and early stopping<sup>20</sup>. The hyperparameters to be optimized were the following:

The `subsample` controls the fraction of row subsampling and is a number between 0 and 1.

The `colsample_bytree` controls the fraction of column down-sampling for each tree, so how many columns (or variables) each tree are allowed to fit to. Is a number between 0 and 1.

The `max_depth` controls the maximum depth of every tree. Is a positive integer (negative values means no maximum depth).

The `min_child_weight` variable controls when the decision tree algorithm should stop splitting a node into further nodes (and will thus turn it into a leaf). Is a positive integer.

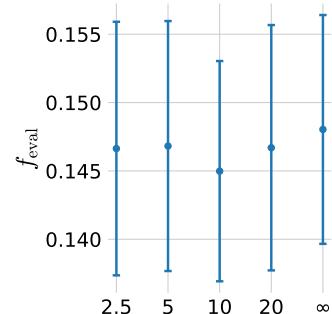


Figure 3.14: Evaluation score as a function of the weight half-life  $T_{\frac{1}{2}}$  with the standard deviation over the 5 folds as errorbars for apartments.

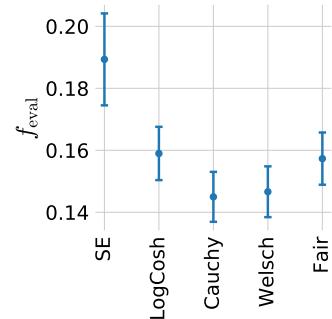


Figure 3.15: Evaluation score as a function of the loss function with the standard deviation over the 5 folds as errorbars for apartments.

<sup>20</sup> This takes a good 4 hours for the random search, 7 hours for the Bayesian optimization, and 90 minutes to optimize the learning rate with early stopping for apartments when run on the local computing cluster, HEP. For the houses this takes a good 24 hours for the random search, 34 hours for the Bayesian optimization, and almost 5 hours to optimize the number of estimators with early stopping.

`reg_lambda` controls the L<sub>2</sub> regularization term. Is a positive number.

`reg_alpha` controls the L<sub>1</sub> regularization term. Is a positive number.

The ranges of the hyperparameters were chosen by a manual, iterative process of fitting a subset of the data (1%–10%) and making sure that the best hyperparameter is not sufficiently close to the boundary of the range; if it is, then the range is extended. The final ranges chosen are shown in Table 3.9. Here  $\mathcal{U}(a, b)$  refers to a uniform distribution from  $a$  to  $b$ , and  $\mathcal{U}_{\text{int}}(a, b)$  is the same only an integer distribution.

The fitting pipeline for this subproject is to first run both random search and Bayesian optimization to compare their results. The best of the two models is chosen and finally learning rate optimized by early stopping where the learning rate  $\eta$  is reduced from  $\eta = 0.1$  to  $\eta = 0.01$ . In the end one ends up with a model that has been hyperparameter optimized for preprocessing optimizations, loss functions, sample weights, normal BDT hyperparameters and finally the learning rate. I have manually implemented this pipeline in Python since no other packages provide the same flexibility as a custom implementation that works fully automated.

The results of the random search and the Bayesian optimization are shown in Figure 3.16 and A.23. The corresponding plots for houses are shown in Figure A.24 and A.25.

Hyperparameter	Range
<code>subsample</code>	$\mathcal{U}(0.5, 0.9)$
<code>colsample_bytree</code>	$\mathcal{U}(0.1, 0.99)$
<code>max_depth</code>	$\mathcal{U}_{\text{int}}(1, 20)$
<code>min_child_weight</code>	$\mathcal{U}_{\text{int}}(1, 30)$
<code>reg_lambda</code>	$\mathcal{U}(0.1, 4)$
<code>reg_alpha</code>	$\mathcal{U}(0.1, 4)$

Table 3.9: Probability Density Functions used in the random search to draw new sets of values for the hyperparameters. Each hyperparameter is drawn from the distribution seen in the table.

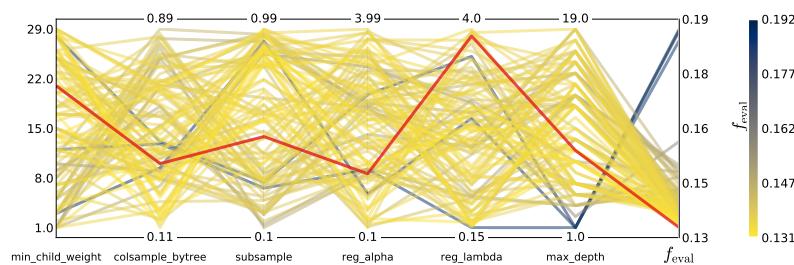
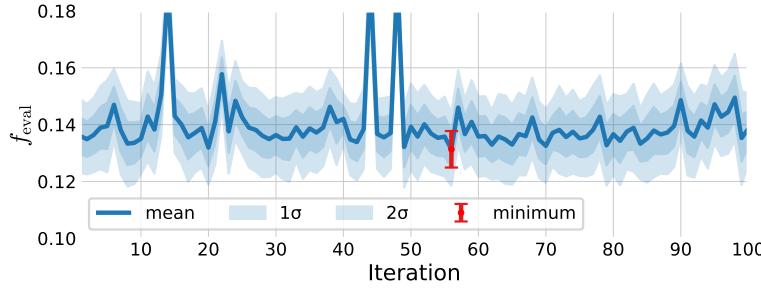


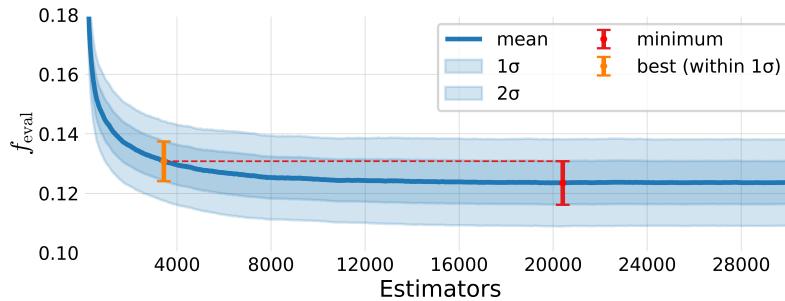
Figure 3.16: Hyperparameter optimization results of XGBoost parameters of the housing model for apartments shown as parallel coordinates. Here shown for random search as hyperparameter optimization.

As with the initial grid search, it is important to compare the results in Figure 3.16 with their uncertainties. This is shown in Figure 3.17. Here the value of the evaluation function along with its  $1\sigma$  and  $2\sigma$  uncertainties are plotted as a function of the iteration along the abscissa. The minimum value of  $f_{\text{eval}}$  is shown in red. Even though this is the minimum value, notice how most of the other iterations are within  $1\sigma$ . The plot for the Bayesian optimization in Figure A.27 shows the exploration phase in the first half of the iterations where it afterwards converge to a more flat minimum, however, this minimum was still worse than the one found with random search.

The best of the random search and the Bayesian optimization models are chosen for subsequent analysis by first reducing the learning rate to  $\eta = 0.01$  and then find the best number of estimators by early stopping. The evaluation function as a function of number of estimators, also known as the *learning curve*, is seen in Figure 3.18.



This curve is the realization of Figure 2.2 in real data, where it first improves a lot and then finds a stable plateau. The minimum is shown in red with its uncertainty. To reduce the risk of overfitting and model complexity – with the further advantage of resulting in a faster model at prediction time – we keep the model with the lowest number of estimators that are still within  $1\sigma$  of the minimum: see the orange point in the figure. This results in a model that contains less than a fifth of the number of trees and is thus also significantly simpler and faster at inference time. This is the final hyperparameter optimized model that will be used for the further analysis.



### 3.6 Results

The performances of the different models (the random search optimized, the bayesian optimization optimized and the best of the two  $\eta$  optimized with early stopping) are shown in Figure 3.19. Here the distribution of the relative price prediction  $z$  of the model evaluated on the test set, apartments sold in 2018, is shown for the three models. In addition to the distributions, also the performance metrics are shown: the value of the evaluation function<sup>21</sup> along with the fraction of relative price predictions that are within the specified percentage. In this particular instance it is seen that 41.3 % of the predictions by the final early stopping model are less than  $\pm 5\%$  wrong, 69.8 % within  $\pm 10\%$ , and 91.9 % within  $\pm 20\%$ . Note that the differences between the three models are minor and that the distributions almost cannot be distinguished from each other.

An interesting observation from the performance metrics of the test set is the low value of  $f_{\text{eval}} = 9.289\%$  (`MAD`). By looking at Figure 3.18 one would expect a test loss of around 12 % assuming iid. samples. However, this assumption does not seem to be fulfilled for the test set. The performances of the realtors are also better for

Figure 3.17: The results of running random search on apartments using the XGB-model. The **minimum (mean)** loss along with its uncertainty is shown in red, the **means** for the different iterations of random search in blue, and as light blue bands are the **one (and two) standard deviation(s)**, all as a function of iteration number.

Figure 3.18: The results of early stopping on apartments using the XGB-model. The **minimum (mean)** loss along with its uncertainty is shown in red, the **means** for the different iterations of random search in blue, and as light blue bands are the **one (and two) standard deviation(s)**, all as a function of number of estimators (trees). In orange the “**best**” number of estimators is shown, defined as the lowest number of estimators which are still within  $1\sigma$  of the minimum value.

<sup>21</sup> Written as `MAD`.

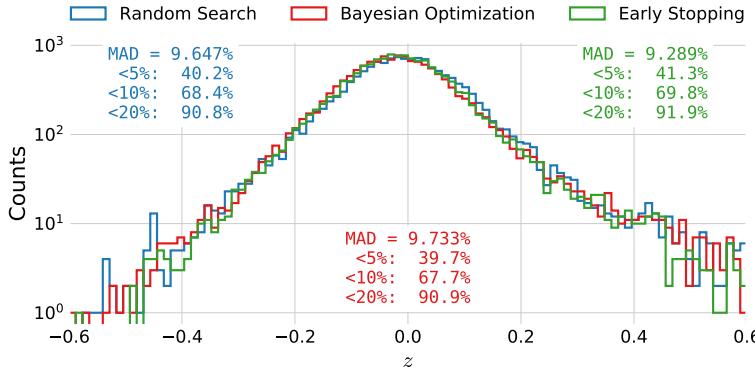
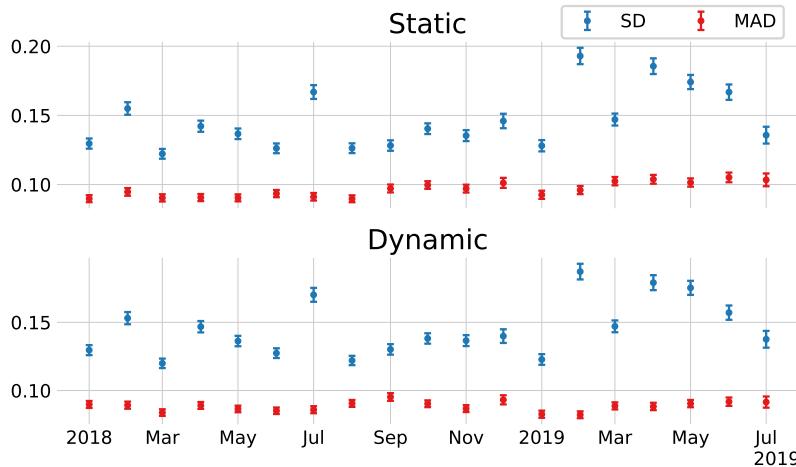


Figure 3.19: Histogram of  $z$ -values of the XGB-model trained on apartments. The performance after hyperparameter optimization using `random search` is shown in blue, for `Bayesian optimization` in red, and the learning rate optimized with `early stopping` in green.

the test set than the training set, and by comparing the test set with the extra 2019 set it seems to be 2018 that was an extra homogenous year, see Table 3.10. The “Tight” in the the table corresponds to the realtors performance on the tight version of the different datasets. The performance of the XGB models on the tight test set is shown in Figure A.30, where it can be seen that the final model has  $f_{\text{eval}} = 8.383\%$ .

To gauge the predictive power of the model over time, we applied the model to the next month’s data and evaluated the results for that month, for all the months in the test set (2018) and 2019. We applied two different methods of forecasting: *static* forecasting where the model is only trained once, and *dynamic* forecasting where the model is retrained after each month on all of the previous sales. These predictions allow the relative predictions  $z$  to be calculated and the MAD and the standard deviation (SD) of  $z$  are shown in Figure 3.20.



In the top subplot the results for the static forecast are shown, whereas the dynamic results are shown in the bottom subplot. The errorbars are calculated using the usual variance of the standard deviation:

$$\sigma_\mu = \frac{\sigma}{\sqrt{N}}, \quad \sigma_\sigma = \frac{\sigma}{\sqrt{2N}}, \quad (3.8)$$

where  $\sigma_\mu$  is the standard deviation of the mean and  $\sigma_\sigma$  is the standard deviation of the standard deviation<sup>22</sup> [18]. Notice the large fluctua-

	Train	Test	2019
Normal	5.80 %	4.97 %	6.19 %
Tight	5.69 %	4.94 %	6.19 %

Table 3.10: The MAD of the realtors’ predictions for the normal and tight selections in the training, test, and 2019 datasets.

Figure 3.20: Performance of 1-month forecasts for apartments sold in 2018 and 2019. For both plots the XGB model is trained on data up to (but excluding) 2018. Top) The performance of the static model’s prediction for both the standard deviation (SD) and MAD of the  $z$ -scores. Bottom) Same as above, however, based on a dynamic model, i.e. a model which is retrained after every month to include the previous month’s sales.

<sup>22</sup> That this estimator is biased does not matter since we are in the large  $N$  limit,  $N \sim 1000$ .

tions in the standard deviations over time compared to MADs which is an effect of MAD being a robust estimator. What is also interesting to note is that the MAD seems to increase over time for the static model, albeit slowly. In comparison, for the dynamic model this seem less pronounced. This figure not only shows the time dependence of the performance of the model, but also that the model is quite stable over time, at least for the dynamic model.

Using the relative price predictions  $z$ , we construct the Market Index, MI. This is an index which measures the overall level of the Danish housing market based on the assumption that if the houses sold in a month are generally sold at a higher price than was predicted by the model, there can be two reasons for it: either the model was wrong or the market simply changed in the time span. With the latter assumption, the ratio between the prediction and the actual price of a residence is thus a measure of the market index:

$$\begin{aligned} z_{mi} &\equiv \frac{\hat{y}}{y} = 1 - \frac{\hat{y} - y}{y} = 1 - z \\ MI_{\text{mean}} &= \text{mean}(z_{mi}) \\ MI_{\text{median}} &= \text{median}(z_{mi}). \end{aligned} \quad (3.9)$$

Here  $z_{mi}$  is the vector of ratios between prediction and actual price, and the market index can then be estimated using either the mean  $MI_{\text{mean}}$  or the median  $MI_{\text{median}}$ . The market indices for every month of the forecast described in the previous figure are shown in Figure 3.21. In the top panel the market index is shown for the static model. Here the median  $MI_{\text{median}}$  is consistently higher than 1, around 1%. Compare this to the dynamic  $MI_{\text{median}}$  in the lower plot which fluctuates around 1. The mean  $MI_{\text{mean}}$  is added to show that its fluctuations are higher than the medians.

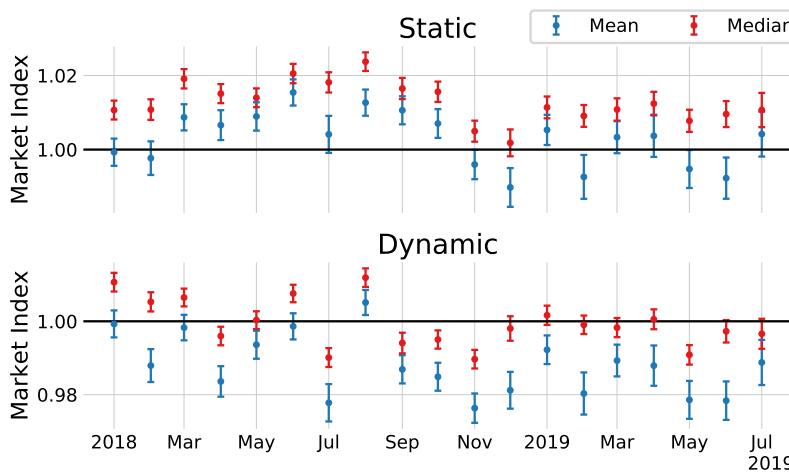


Figure 3.21: The Market Index as defined in equation (3.9) for the static and dynamic 1-month forecasts for 2018 and 2019. The `mean` is plotted in red and the `median` in blue.

The final results for the model is seen in Table 3.11 for owner-occupied apartments and in Table 3.12 for one family houses. The MAD is around 9 % for apartments and 16 % for houses, which is still worse than the realtors' prediction, yet still acceptable for a model that does not have any variables describing the indoor conditions.

For apartments around 40 % of all the predictions are within  $\pm 5\%$  and more than 90 % are within  $\pm 20\%$  which is similar to the performance of the professional automated property evaluations from e.g. Bolighed [25]. The performance on the tight cuts is shown in Table A.13 and A.14.

	MAD (%)	$\leq 5\%(\%)$	$\leq 10\%(\%)$	$\leq 20\%(\%)$
Train	7.83	47.90	75.74	93.97
Test	9.29	41.33	69.77	91.91
2019	9.89	38.85	66.76	90.04

	MAD (%)	$\leq 5\%(\%)$	$\leq 10\%(\%)$	$\leq 20\%(\%)$
Train	14.12	28.95	51.89	78.04
Test	15.79	25.61	47.52	75.14
2019	16.50	24.01	45.89	74.22

Table 3.11: Performance metrics for the final housing model trained and evaluated on apartments.

Table 3.12: Performance metrics for the final housing model trained and evaluated on houses.

### 3.7 Model Inspection

One of the most important aspects of applying advanced machine learning methods, in addition to getting accurate predictions is to understand the model. As mentioned in section 2.8, it is possible to use SHAP values to inspect the trained model for both local predictions and for global feature importances  $\phi_i^{\text{tot}}$ . An example of how to use SHAP to better understand a local prediction is seen in Figure 3.22. Here the SHAP values for a particular sale are visualized as a bar chart where the green colors are positive values, red values negative values and the blue is the final prediction  $\hat{y}$ . Remember that for SHAP values the prediction is the sum of all of the SHAP values, see equation (2.40). Here  $\phi_0$  is the expectation value denoted as Bias in the plot. To show all 143 variables would make the figure excessively large, so two extra bins have been added: the Overflow bar which is the sum of the remaining positive SHAP values and likewise with Underflow for negative values. The sum of all the green and red bars adds up to  $\hat{y} = 6.86 \text{ M.kr.}$  in this particular instance and the actual sold value was  $y = 6.35 \text{ M.kr.}$  Thus, in cases where there is a large discrepancy between the predicted and actual sales prices, one can use this tool to better understand why the prediction was estimated as it was.

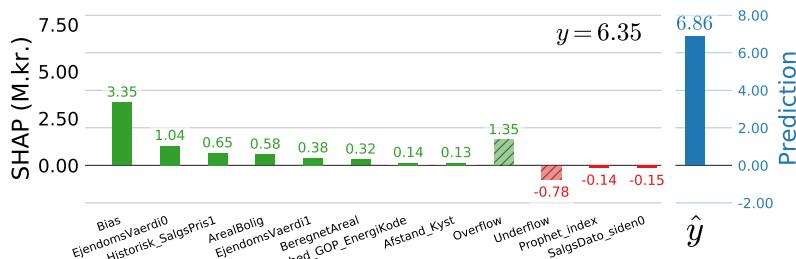


Figure 3.22: Model explanation for XGB model for a specific apartment. The bars are the variables in the dataset that the model found most important sorted after their importance for this particular apartment. The bias bar refers to the expected value of the model, which is simply the mean of the training set which acts as the naive prediction baseline. The “cutoff positive (negative)” bars are the sum of the remaining positive (negative) values that are not shown. On the right hand side of the plot is the model prediction shown. The model prediction is the sum of all of the bars in the left part (6.86 M.kr. in this example). The negative values are shown in red, positive ones in green, and the bias in blue.

Figure 3.23 shows which variables are most important on a global<sup>23</sup> scale as measured by  $\phi_i^{\text{tot}}$ . Here the variables are sorted according to the normalized  $\phi_i^{\text{tot}}$  which is shown in parentheses after each variable name. In the center of the plot is shown a dot-plot of the dataset plotted with the SHAP value on the abscissa and colored according to the feature value. The way to interpret this plot is as follows. Take a variable of interest, e.g. the area of the apartment `ArealBolig` with  $\phi_{\text{ArealBolig}}^{\text{tot}} = 5.35\%$ . Every dot is a sale plotted as a function of their SHAP value with a spread such that the height corresponds to the SHAP distribution of that specific feature. For the area it can be seen that there is a long tail towards high SHAP values, however, most of the samples have slightly negative SHAP values. The dots are colored according to their feature value and it can thus be seen that large apartments (red) are given a higher SHAP value than small apartments (blue); precisely as expected from the model. In contrary, when the total days on market (`LiggetidSamlet`) is large it pushes the prediction in the negative direction.

<sup>23</sup> Global here meaning for the entire dataset.

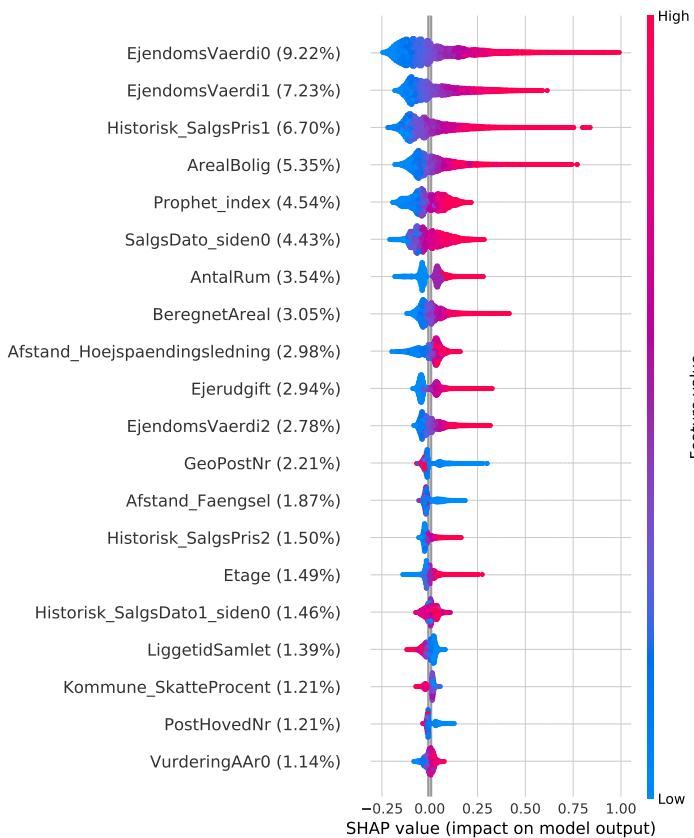


Figure 3.23: Feature importance of apartment prices using the XGB-model. The feature importance is measured using SHAP values. The variables are sorted top to bottom according to their overall feature importance, i.e. the previous public property valuation `EjendomsVaerdi0` is the most important single feature. Along the  $x$ -axis is the impact on model output, in this example the price in M.kr. This axes is colored by the value of the feature, from **low** (blue) to **high** (red). In this particular example we see that high values of the previous public property valuation has high, positive impact on the model prediction – exactly as expected. This is exactly opposite the total days on market described by the variable `LiggetidSamlet` where a high value has a negative impact.

The SHAP software [67] not only allows for 1D dependencies to be gauged, it also allows for so-called *interaction plots*. These plots shows the 2D-dependence between the variable and the SHAP value colored according to a second variable. Since the previous public property valuation (PPPV) (`EjendomsVaerdi0`) is the most important of the features, the interaction plot of this variable is seen in Figure 3.24.

Here the SHAP value of the `EjendomsVaerdi0` is plotted as a function of `EjendomsVaerdi0`. This plot shows that the higher the PPPV, the higher the model output. The colors show how this trend depends on time by the variable `SalgsDato_siden0` which is the number of days since January 1<sup>st</sup> 2009 that the apartment was sold. The plots shows that newer sales with a high PPPV has even higher SHAP values than older sales with the same high PPPV. This agrees with the fact that the market has gone up since 2009. On the other hand, for low PPPV apartments the relationship with time is inverse, however, the effect is much smaller here. The SHAP software chose to color by `SalgsDato_siden0` since this is the variable which explains most of the variation for a given PPPV.

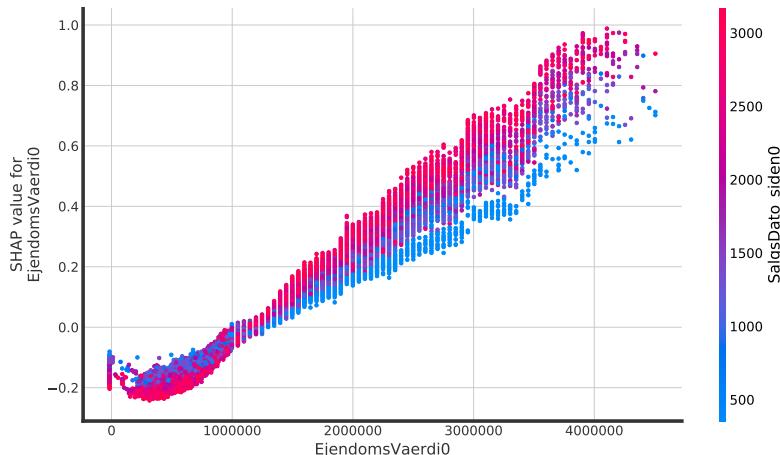


Figure 3.24: Interaction plot of the feature importances for the XGB model on apartments. Here shown for the previous public property valuation (`EjendomsVaerdi0`) colored according to the sales date (`SalgsDato_siden0`).

What can be concluded from this plot is that for apartments with a high PPPV and that were sold recently were sold for more than apartments with the same PPPV that were sold longer time ago; effectively the time dependence of sales prices. These interaction plots serves as helpful sanity checks to see that the model is actually learning reasonable relationships between the variables.

### 3.8 Multiple Models

In addition to the XGBoost [36] (XGB) BDT model, several other models were also tested. During the sub-project the LightGBM package [61] (LGB), also a BDT model, was released and started gaining traction in the ML community, especially for large-scale data analysis. In comparison to XGBoost, LightGBM implements some extra binning and categorical assumptions that greatly speeds up the fitting process. It was trained in the same way as the XGBoost model with the same hyperparameter optimization process and range for the hyperparameters.

To compliment the boosted decision trees, a simple linear model (LIN) with  $L_2$  loss (ridge regression) was fitted. Before the fit all of the NaNs were median-imputed which means that all invalid values were replaced with the median along each column. All the features were scaled<sup>24</sup> with a robust scaler from Scipy [100] in the (25, 75)

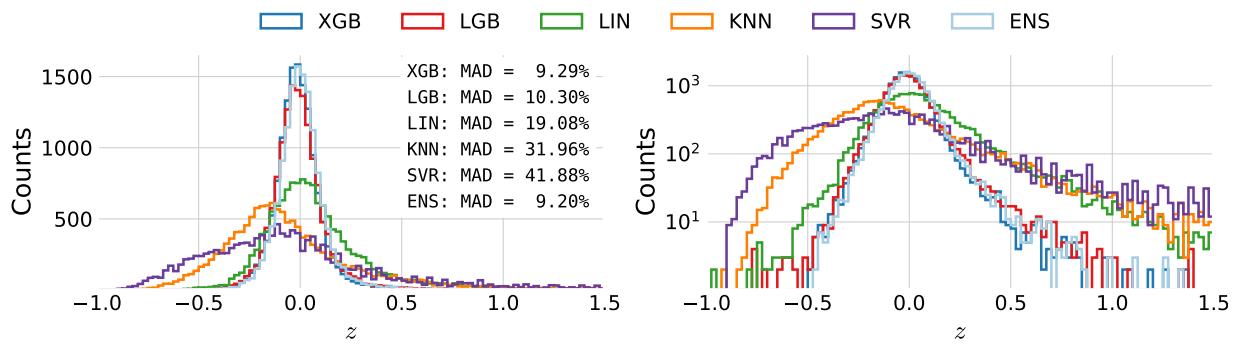
<sup>24</sup> Scaling of input features is generally an important preprocessing step for non-tree based ML models.

quantiles range and the regularization parameter was hyperparameter optimized. This linear model is quick and easy to both implement and fit, and can be seen as the simplest baseline model.

The K-nearest neighbors (KNN) algorithm was fitted to the data with a data preprocessing pipeline similar to the linear model with median imputation and robust scaling of the input features. The number of neighbors,  $K$ , was hyperparameter optimized together with the  $p$ -norm of the metric,  $p \in \{1, 2\}$  (Manhattan, Euclidean, see also subsection 2.4.1).

Finally, support vector machines<sup>25</sup> (SVR) were used with the same preprocessing pipeline as the previous two methods and hyperparameter optimized for the  $L_2$  regularization parameter  $C$  and the kernel coefficient  $\gamma$  for the radial basis function (RBF) kernel.

<sup>25</sup> Known as support vector regression when applied to regression problems [17].



The results for the five different models are shown in Figure 3.25, where the two subplots are the same up to a logarithmic scaling of the ordinate axis in the right plot. It is easily seen how XGBoost and LightGBM are the best-performing models together with the ensemble (ENS) of the different models, see paragraph below for further explanation. In the figure the MAD is also shown on the test set, which confirms the visual clue: that XGBoost performs best, followed by the ensemble model.

The five different models – XGB, LGB, LIN, KNN, SVR – should be able to each capture different parts of the hyperdimensional phase space and an ensemble of these models would thus be expected to be as good as or better than the best of the individual models. This kind of ensemble model is sometimes called a *super learner* in the statistics community [76, 95]. To make sure that the ensemble model is not just retraining on the training set, and thus end up overfitting, we follow the process from Polley and van der Laan [76]. Using cross-validation for time-series<sup>26</sup> data with 10 folds, the training data is split up into folds sorted by time. Each fold is fitted with all five models, and the prediction of the next fold is made for all five models. This is repeated for the remaining folds until one ends up with a matrix of predictions  $Z \in \mathbb{R}^{(N \times 5)}$  for  $N$  training samples. Since all the folds in  $Z$  consists of predictions on unseen data<sup>27</sup> this prevents overfitting. The meta learner then fits  $Z$  to the actual predictions of the training data  $y$  in the usual way. The combination of a meta

Figure 3.25: Distribution of the relative predictions  $z$  for apartments for the five different models and the ensemble model. The left plot shows the normal histogram (with a linear scale), whereas the right one has a log-scaled  $y$ -axis to better see the tails of the distributions. The MAD scores for the different models is further shown in the left plot.

<sup>26</sup> See subsection 2.4.2.

<sup>27</sup> The predictions for the very first fold is based on training data.

learner fitted to the predictions of individual models is called an ensemble model.

At first an XGBoost model was used as the meta learner yielding decent results, however, still performing worse than the single XGB model ( $MAD = 9.57\%$ ). To better understand the issue, the meta learner was changed to a linear model which would basically just compute a weighted average of the different models:

$$\Psi_{\text{meta}}(\mathbf{x}) = \sum_{i=1}^5 \alpha_i \Psi_i(\mathbf{x}), \quad (3.10)$$

where  $\alpha$  is a vector of the weights for the meta learner and  $\Psi_i$  is each of the individual ML models. The linear model performed even worse than using XGB as the meta learner ( $MAD = 10.48\%$ ), yet it was more transparent. During the debugging process it was realized that none of these models actually optimize the evaluation function, MAD, directly. The XGBoost model was trained using the Cauchy loss found earlier and the linear regression model using a simple squared error loss. Since a simple weighted average should work as the meta learner [76], a custom algorithm for finding  $\alpha$  according to MAD was implemented.

Given the training data, the evaluation function as a function of  $\alpha$  was minimized using the MINUIT algorithm [59] via the iminuit [90] Python interface. It yielded decent results, yet they were all very dependent on the initial parameter of the fit indicating many local minima. A scan over the 5-dimensional hyperspace in steps of 0.01 was thus conducted and the result of this scan was used as the new initial parameter in the minimization routine. This yielded the following result:

$$\alpha = \begin{bmatrix} \alpha_{\text{LIN}} \\ \alpha_{\text{KNN}} \\ \alpha_{\text{SVR}} \\ \alpha_{\text{XGB}} \\ \alpha_{\text{LGB}} \end{bmatrix} = \begin{bmatrix} 0.202 \\ 0.002 \\ 0.001 \\ 81.302 \\ 20.002 \end{bmatrix} \times 10^{-2}. \quad (3.11)$$

The fact that it sums to more than 1 just corresponds to an overall scaling. When using the found  $\alpha$  in equation (3.10), one gets the ensemble model (ENS) shown in Figure 3.25 with a  $MAD = 9.20\%$ . This value is the evaluation loss on the test set based on only the training data and outperforms all of the individual models.

The paragraphs above refer to apartments, however, the intermediate results for houses showed the same pattern. The combined model along with its performance is shown in Figure A.31

### 3.9 Discussion

The subproject of estimating housing prices has focussed a lot on experimenting with different machine learning models and how to optimize them. As shown in the previous sections, the choice of

ML model is by far the most important. Actually, the gain from hyperparameter optimization is quite small, especially considering the amount of time spent on it<sup>28</sup>. With the dataset at hand, decent results were achieved, however, they were nowhere near the performance of the realtors' predictions.

There are two main reasons for this, the first being that realtors are educated within this field and thus have developed the skills required for estimating the price of a house over many years of hard work. The second reason is the fact that the realtor has access to a lot more information than the ML models have. We are not in possession of any *indoor variables*. The area of the house, the number of rooms, the name of the street, and the distance to a highway are all variables that are in the data set but none of them describe the overall quality of the house, the maintenance level, the age of the kitchen or bathroom. These indoor features are invisible to the ML model.

During the project it was investigated how to get access to these variables. At first the online images from each sale was suggested, however, it turned out that Boligsiden only have the right to use them while a residence is for sale; when it is sold all rights return to the photographer. The images are not the only thing that provide more information about the condition of the residence, also the descriptions does that.

The descriptions turned out to be available for most of the sales and was investigated for a short period. At that time of the project, the MAD for (a subset of the) apartments was around  $\pm 14\%$  and  $\pm 20\%$  for houses. By using methods from the big natural language processing (NLP) community within the field of machine learning, it was possible to reduce the MAD to around  $\pm 12\%$  and  $\pm 15\%$  for apartments and houses respectively. From the improvement in performance it is noticeable how apartments in general are much more uniform compared to houses where the “inside” is more decisive regarding the price.

The methods for translating the text to numerical variables decipherable by classical ML models were for instance simple *bag of words* (BOW) models and *Term Frequency, Inverse Document Frequency* (TF-IDF) but also slightly more advanced statistical tools such as *Latent Dirichlet Allocation* (LDA).

An old example of this were a housing-based model that was trained with the five numerical variables<sup>29</sup>: the PPPV, postal code, year of construction, distance to shore, and the weighted area. and the text descriptions (encoded with TF-IDF). Figure 3.26 shows the SHAP summary plot of the trained model. As expected, the most important features are the numerical ones, however, the word `flot` (“pretty”) was in top five. The model also learnt that `flot` has a positive impact on the price compared to `trænger` (“requires”) which has a negative impact.

The descriptions turned out to be more time-consuming to extract for Boligsiden and along with the fact that overall deadline was

<sup>28</sup> Not only user-time programming it, but also the computational resources spent.

<sup>29</sup> With the Danish names: `Ejendomsvaerdio`, `GeoPostNr`, `ByggeAAr`, `Afstand_Kyst`, and `BeregnetAreal`.

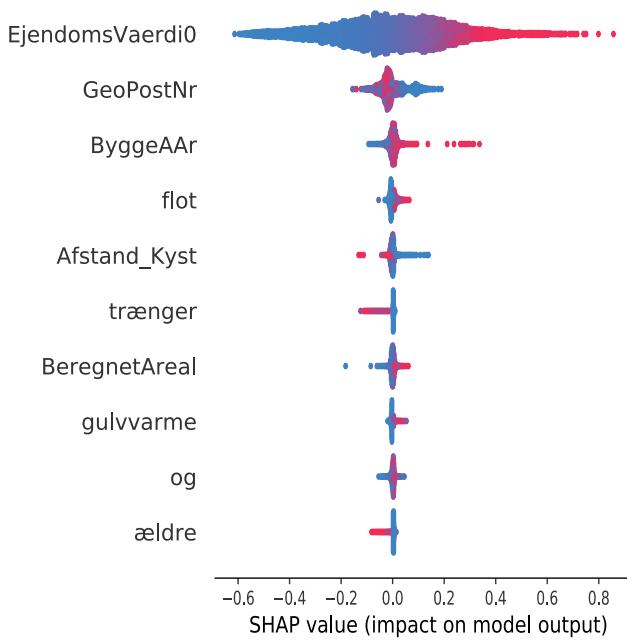


Figure 3.26: Feature importance of villas when the descriptions are also used. The variables are sorted top to bottom according to their overall feature importance, i.e. the previous public property valuation `EjendomsVaerdi0` is the most important single feature. Along the  $x$ -axis is the impact on model output, in this example the price in M.kr. This axis is colored by the value of the feature, from **low** (blue) to **high** (red). In this particular example we see that high values of the previous public property valuation has high, positive impact on the model prediction – exactly as expected. This is in contrast to the word “requires” (`trænger`) where a high value has a negative impact.

quickly approaching, the remaining time was focussed on the main part of the project, the quark gluon discrimination. Given more time, the text analysis would definitely be the first step for further improving the accuracy and precision of the price predictions.

Another step would be to apply more modern deep learning<sup>30</sup> methods. These methods were briefly experimented with in the initial stages of this subproject but showed inferior performance compared to BDTs. It is generally accepted in the ML community (with some modifications) that neural networks underperform, or at least not outperform, classic ML methods on structured<sup>31</sup> data [63]. Most often they have the inherent complexity needed to perform as well as other ML methods, however, this requires extensive architecture optimization, or, in short; the hypothesis space for neural networks is much larger than for classical ML methods and thus requires more care to avoid overfitting.

### 3.10 Conclusion

**XXX TODO!**

<sup>30</sup> Basically advanced neural networks with many layers.

<sup>31</sup> In general data that can be described by a spread sheet, i.e. has a well-defined number of variables and observations which is why it is also known as *tabular data*.

*(This page is intentionally left blank).*

## Part II



*(This page is intentionally left blank).*

# *List of Figures*

2.1	The Learning Problem . . . . .	6
2.2	Approximation-Estimation Tradeoff . . . . .	10
2.3	Regularization Strength . . . . .	11
2.4	Regularization Effect of $L_2$ . . . . .	12
2.5	Regularization Effect of $L_1$ . . . . .	12
2.6	$k$ -Fold Cross Validation . . . . .	13
2.7	$k$ -Fold Cross Validation for Time Series Data . . . . .	13
2.8	Objective Functions . . . . .	15
2.9	Decision Tree Cuts In Feature Space . . . . .	16
2.10	Decision Tree . . . . .	16
2.11	Grid Search . . . . .	20
2.12	Random Search . . . . .	21
2.13	Bayesian Optimization . . . . .	22
3.1	Price Index of the Danish Housing Market . . . . .	27
3.2	Distributions of the Variables in the Housing Prices . . . . .	29
3.3	Geographic Distribution of the Sold Residences . . . . .	29
3.4	Histogram of Prices of Houses and Apartments Sold in Denmark . . . . .	30
3.5	Linear Correlation Between Variables and Price . . . . .	31
3.6	Comparison of the Linear Correlation $\rho$ and the Nonlinear MIC . . . . .	31
3.7	Nonlinear Correlation Between Variables and Price . . . . .	32
3.8	Validity of Input Features . . . . .	32
3.9	Validity Dendrogram . . . . .	33
3.10	Prophet Forecast for Apartments . . . . .	35
3.11	Prophet Trends . . . . .	35
3.12	Sample Weight as a Function of Time for Different Half-Lives. . . . .	36
3.13	Parallel Coordinate Plot of the Initial Hyperparameter Optimization for Apartments . . . . .	37
3.14	Initial HPO Results for the Weight Half-life $T_{\frac{1}{2}}$ . . . . .	38
3.15	Initial HPO Results for the Loss Function . . . . .	38
3.16	Parallel Coordinate Plot of the Random Search Hyperparameter Optimization Results of XGBoost for Apartments . . . . .	39
3.17	Hyperparameter Optimization: Random Search Results . . . . .	40
3.18	Early Stopping Results . . . . .	40
3.19	Performance of XGB-model for Apartments . . . . .	41
3.20	Standard Deviation and MAD of the Static and Dynamic XGB Forecasts . . . . .	41
3.21	Market Index based on the Static and Dynamic XGB Forecasts . . . . .	42
3.22	SHAP Prediction Explanation for Apartments . . . . .	43
3.23	Feature importance of Apartments Prices . . . . .	44
3.24	Feature Importance Interaction Plot for Apartments . . . . .	45

3.25 Performance Comparison of Multiple Models for Apartments . . . . .	46
3.26 Feature Importance of Villas With Descriptions . . . . .	49
4.1 The Standard Model . . . . .	54
4.2 Feynman Diagram for the Jet Production at LEP . . . . .	55
4.3 Quark Splitting . . . . .	55
4.4 Hadronization Process . . . . .	56
4.5 The ALEPH Detector . . . . .	57
4.6 Polar Angle . . . . .	57
4.7 Azimuthal Angle . . . . .	57
4.8 Track Significance . . . . .	59
5.1 Histograms of the Vertex Variables . . . . .	66
5.2 UMAP Visualization of the Vertex Variables for 4-Jet Events . . . . .	67
5.3 UMAP Visualization of the Vertex Variables for 3-Jet Events . . . . .	67
5.4 Correlations Between of the Vertex Variables . . . . .	67
5.5 Plot of the Log-Loss $\ell_{\log}$ . . . . .	68
5.6 Hyperparameter Optimization of $b$ -Tagging . . . . .	69
5.7 Parallel Plot of HPO Results for 4-Jet $b$ -Tagging . . . . .	70
5.8 $b$ -Tag Scores in 4-Jet Events . . . . .	70
5.9 ROC Curve for 4-Jet $b$ -Tagging . . . . .	71
5.10 Distribution of $b$ -Tags in 4-Jet Events . . . . .	71
5.11 Global Feature Importances for the LGB $b$ -Tagging Algorithm on 4-Jet Events . . . . .	72
5.12 The Expit Function . . . . .	72
5.13 The Logit Function . . . . .	72
5.14 SHAP 3-Jet Model Explanation for $b$ -Like Jet . . . . .	72
5.15 $b$ -Tagging Efficiency $\varepsilon_b^{b\text{-sig}}$ as a Function of Jet Energy . . . . .	74
5.16 $b$ -Tagging Efficiency $\varepsilon_g^{g\text{-sig}}$ as a Function of Jet Energy . . . . .	75
5.17 Hyperparameter Optimization of $g$ -Tagging . . . . .	77
5.18 1D Sum Models Predictions and Signal Fraction for 4-jets events . . . . .	79
5.19 $g$ -Tag Scores in 4-Jet Events . . . . .	79
5.20 ROC Curve for $g$ -Tag in 4-Jet Events . . . . .	80
5.21 Distribution of $g$ -Tag Scores in 4-Jet Events for Signal and Background . . . . .	80
5.22 Distribution of $b$ -Tag Scores in 3-Jet $l$ -Quark Events for Low and High $g$ -Tag Values . . . . .	81
5.23 3D Scatter Plot of $\beta_{\text{tag}}$ -Values for High and Low $\gamma_{\text{tag}}$ $l$ -Quark Events . . . . .	81
5.24 $g$ -Tagging Pseudo Efficiency for $b\bar{b}g$ -Events as a Function of $g$ -Tag . . . . .	82
5.25 $g$ -Tagging Pseudo Efficiency for $b\bar{b}g$ -Events as a Function of The Mean Invariant Mass . . . . .	83
5.26 Generalized Angularities . . . . .	83
5.27 Generalized Angularities for Charged Gluons Jets in 3-Jet Events: $\lambda_0^2$ . . . . .	84
5.28 Soft Wide Angle Gluons in 4-Jet Events . . . . .	87
5.29 Soft Collinear Gluons in 4-Jet Events . . . . .	87
5.30 Hard Non $g \rightarrow gg$ Gluons in 4-Jet Events . . . . .	87
5.31 $g$ -Tagging Efficiency for 4-Jet Events in MC as a Function of the Normalized Gluon-Gluon Jet Energy Difference Asymmetry $E_{\text{diff}}$ . . . . .	88
5.32 Closure Plot Comparing MC Truth and the Efficiency Corrected $g$ -Tagging Model in 4-Jet Events for the Normalized Gluon Gluon Jet Energy Asymmetry . . . . .	89
5.33 Overview of the Four Regions in the $R_{gg}^{k_t}$ - $R_{gg}^{\text{CA}}$ Phase Space . . . . .	91
5.34 Gluon Splitting Distribution Comparison in MC and Data for $R_{gg}^{k_t}$ - $R_{gg}^{\text{CA}}$ Phase Space Region A . . . . .	92
5.35 Ratio Plot of $E_{\text{diff}}$ in MC and Data for $R_{gg}^{k_t}$ - $R_{gg}^{\text{CA}}$ for Region A . . . . .	93

5.36 Ratio Plot of $E_{\text{diff}}$ in MC and Data for $R_{gg}^{k_t}$ - $R_{gg}^{\text{CA}}$ for Region D . . . . .	93
A.1 Distributions of the Input Variables – Part 1 . . . . .	96
A.2 Distributions of the Input Variables – Part 2 . . . . .	97
A.3 Distributions of the Input Variables – Part 3 . . . . .	98
A.4 Distributions of the Input Variables – Part 4 . . . . .	99
A.5 Distributions of the Input Variables – Part 5 . . . . .	100
A.6 Distributions of the Input Variables – Part 6 . . . . .	101
A.7 Distributions of the Input Variables – Part 7 . . . . .	102
A.8 Distributions of the Input Variables – Part 8 . . . . .	103
A.9 Distributions of the Input Variables – Part 9 . . . . .	104
A.10 Distributions of the Input Variables – Part 10 . . . . .	105
A.11 Distributions of the Input Variables – Part 11 . . . . .	106
A.12 Distributions of the Input Variables – Part 12 . . . . .	107
A.13 Distributions of the Input Variables – Part 13 . . . . .	108
A.14 Distributions of the Input Variables – Part 14 . . . . .	109
A.15 Linear Correlations . . . . .	110
A.16 Comparison of the Linear Correlation $\rho$ and the Nonlinear MIC for Different Levels of Noise . . . . .	111
A.17 Validity Heatmap . . . . .	111
A.18 Prophet Forecast for apartments . . . . .	113
A.19 Prophet Trends . . . . .	113
A.20 Parallel Coordinate Plot of the Initial Hyperparameter Optimization for Apartments . . . . .	120
A.21 Parallel Coordinate Plot of the Initial Hyperparameter Optimization for Houses . . . . .	120
A.22 Parallel Coordinate Plot of the Hyperparameter Optimization for Apartments Using Random Search . . . . .	121
A.23 Parallel Coordinate Plot of the Hyperparameter Optimization for Apartments Using Bayesian Optimization . . . . .	121
A.24 Parallel Coordinate Plot of the Hyperparameter Optimization for Houses Using Random Search . . . . .	122
A.25 Parallel Coordinate Plot of the Hyperparameter Optimization for Houses Using Bayesian Optimization . . . . .	122
A.26 Random Search Results as a Function of Iteration for Apartments . . . . .	123
A.27 Bayesian Optimization Results as a Function of Iteration for Apartments . . . . .	123
A.28 Random Search Results as a Function of Iteration for Houses . . . . .	124
A.29 Bayesian Optimization Results as a Function of Iteration for Houses . . . . .	124
A.30 Performance of XGB-model on apartment prices . . . . .	125
A.31 Performance Comparison of Multiple Models for Houses . . . . .	125
B.1 UMAP Parameter Grid Search . . . . .	128
B.2 Visualization of the t-SNE algorithm . . . . .	128
B.3 Parallel Coordinate Plot of HPO results for 3-Jet $b$ -Tagging . . . . .	129
B.4 Parallel Coordinate Plot of HPO results for 4-Jet $b$ -Tagging . . . . .	129
B.5 $b$ -Tag Scores in 3-Jet Events . . . . .	130
B.6 $b$ -Tag Scores in 4-Jet Events . . . . .	130
B.7 ROC curve for 3-jet $b$ -tagging . . . . .	131
B.8 ROC curve for 4-jet $b$ -tagging . . . . .	131
B.9 Distribution of $b$ -Tags in 3-Jet Events . . . . .	132
B.10 Distribution of $b$ -Tags in 4-Jet Events . . . . .	132
B.11 Global Feature Importances for the LGB $b$ -Tagging Algorithm on 3-Jet Events . . . . .	133

B.12	Global Feature Importances for the LGB $b$ -Tagging Algorithm on 4-Jet Events	133
B.13	$b$ -Tagging Efficiency $\varepsilon_b^{b\text{-sig}}$ as a Function of Jet Energy	134
B.14	$b$ -Tagging Efficiency $\varepsilon_g^{g\text{-sig}}$ as a Function of Jet Energy	134
B.15	$b$ -Tagging Efficiency $\varepsilon_b^{g\text{-sig}}$ as a Function of Jet Energy	135
B.16	$b$ -Tagging Efficiency $\varepsilon_g^{b\text{-sig}}$ as a Function of Jet Energy	135
B.17	Parallel Coordinate Plot of HPO Results for 3-Jet $g$ -Tagging for Energy Ordered Jets	136
B.18	Parallel Coordinate Plot of HPO Results for 3-Jet $g$ -Tagging for Shuffled Jets	136
B.19	Parallel Coordinate Plot of HPO Results for 4-Jet $g$ -Tagging for Energy Ordered Jets	137
B.20	Parallel Coordinate Plot of HPO Results for 4-Jet $g$ -Tagging for Shuffled Jets	137
B.21	PermNet Architecture	138
B.22	1D Sum Models Predictions and Signal Fraction for 4-jets events	139
B.23	1D LGB Model Cuts for 4-jets events	139
B.24	1D Sum Models Predictions and Signal Fraction for 3-jets events	140
B.25	1D LGB Model Cuts for 3-jets events	140
B.26	$g$ -Tag Scores in 4-Jet Events	141
B.27	$g$ -Tag Scores in 3-Jet Events	141
B.28	ROC Curve for $g$ -Tag in 4-Jet Events	142
B.29	ROC Curve for $g$ -Tag in 3-Jet Events	142
B.30	Distribution of $g$ -Tag Scores in 4-Jet Events for Signal and Background	143
B.31	Distribution of $g$ -Tag Scores in 3-Jet Events for Signal and Background	143
B.32	Generalized Angularities for Charged Gluons Jets: $\lambda_0^2$	144
B.33	Generalized Angularities for Charged Gluons Jets: $\lambda_{1/2}^1$	144
B.34	Generalized Angularities for Charged Gluons Jets: $\lambda_1^1$	144
B.35	Generalized Angularities for Charged Gluons Jets: $\lambda_1^2$	145
B.36	Generalized Angularities for Charged Gluons Jets: $\lambda_0^0$	145
B.37	Generalized Angularities for Neutral Gluons Jets: $\lambda_0^2$	146
B.38	Generalized Angularities for Neutral Gluons Jets: $\lambda_{1/2}^1$	146
B.39	Generalized Angularities for Neutral Gluons Jets: $\lambda_1^1$	146
B.40	Generalized Angularities for Neutral Gluons Jets: $\lambda_1^2$	147
B.41	Generalized Angularities for Neutral Gluons Jets: $\lambda_0^0$	147
B.42	$g$ -Tagging Efficiency for 4-Jet Events in MC as a Function of $E_{\text{diff}}$	148
B.43	$g$ -Tagging Efficiency for 4-Jet Events in MC as a Function of $E_{\text{rel,min}}$	148
B.44	$g$ -Tagging Efficiency for 4-Jet Events in MC as a Function of $E_{\text{rel}}$	148
B.45	$g$ -Tagging Efficiency for 4-Jet Events in MC as a Function of $\Delta_\theta$	149
B.46	$g$ -Tagging Efficiency for 4-Jet Events in MC as a Function of $m_{gg}$	149
B.47	$g$ -Tagging Efficiency for 4-Jet Events in MC as a Function of $\phi_{  }$	149
B.48	$g$ -Tagging Efficiency for 4-Jet Events in MC as a Function of $\ln(k_t^2/m_{\text{vis}}^2)$	149
B.49	$g$ -Tagging Efficiency for 4-Jet Events in MC as a Function of $p_{\perp,A}^2$	149
B.50	Closure Plot Comparing MC Truth and the Efficiency Corrected $g$ -Tagging Model in 4-Jet Events for $E_{\text{diff}}$	150
B.51	Closure Plot Comparing MC Truth and the Efficiency Corrected $g$ -Tagging Model in 4-Jet Events for $E_{\text{rel,min}}$	150
B.52	Closure Plot Comparing MC Truth and the Efficiency Corrected $g$ -Tagging Model in 4-Jet Events for $E_{\text{rel}}$	151
B.53	Closure Plot Comparing MC Truth and the Efficiency Corrected $g$ -Tagging Model in 4-Jet Events for $\Delta_\theta$	151
B.54	Closure Plot Comparing MC Truth and the Efficiency Corrected $g$ -Tagging Model in 4-Jet Events for $m_{gg}$	151

B.55 Closure Plot Comparing MC Truth and the Efficiency Corrected $g$ -Tagging Model in 4-Jet Events for $\phi_{\parallel}$ . . . . .	152
B.56 Closure Plot Comparing MC Truth and the Efficiency Corrected $g$ -Tagging Model in 4-Jet Events for $\ln(k_t^2/m_{\text{vis}}^2)$ . . . . .	152
B.57 Closure Plot Comparing MC Truth and the Efficiency Corrected $g$ -Tagging Model in 4-Jet Events for $p_{\perp,A}^2$ . . . . .	152
B.58 Gluon Splitting Distribution Comparison in MC and Data for Phase Space Region A . . . . .	153
B.59 Gluon Splitting Distribution Comparison in MC and Data for Phase Space Region B . . . . .	154
B.60 Gluon Splitting Distribution Comparison in MC and Data for Phase Space Region C . . . . .	155
B.61 Gluon Splitting Distribution Comparison in MC and Data for Phase Space Region D . . . . .	156

*(This page is intentionally left blank).*

# *List of Tables*

3.1	Mapping between the Code in <code>SagTypeNr</code> and the Type of Residence . . . . .	28
3.2	Basic Selection Criteria . . . . .	33
3.3	Side Door Mapping. . . . .	33
3.4	Street Mapping . . . . .	34
3.5	Number of Observations in the Housing Dataset . . . . .	36
3.6	Number of Observations in the Housing Dataset for the Tight Selection . . . . .	36
3.7	Results from the Initial Hyperparameter Optimization for Apartments . . . . .	37
3.8	Results from the Initial Hyperparameter Optimization for Houses . . . . .	37
3.9	PDFs Used in the Random Search . . . . .	39
3.10	Realtors' MAD . . . . .	41
3.11	Performance Metrics for Apartments . . . . .	43
3.12	Performance Metrics for Houses . . . . .	43
5.1	The number of $n$ -jet events for Data and MC. . . . .	64
5.2	Mapping Between the Flavor and the Variable Flevt . . . . .	64
5.3	Number of Different Types of Jets for MC and MCb for $n$ -Jet Events . . . . .	65
5.4	Random Search PDFs for LGB . . . . .	69
5.5	Global SHAP Feature Importances for the $g$ -Tagging Models in 4-Jet Events . . . . .	77
5.6	Generalized Angularities Systematic Errors, Charged Tracks . . . . .	85
5.7	Generalized Angularities Systematic Errors, Neutral Clusters . . . . .	85
5.8	Gluon Splitting Systematic Errors . . . . .	90
5.9	Rgion Definition in the $R_{gg}^{k_t}$ - $R_{gg}^{\text{CA}}$ Phase Space . . . . .	90
A.1	Energy Rating Mapping . . . . .	112
A.2	Final Variable List . . . . .	115
A.3	Initial Hyperparameter Optimization Results for Apartments – SE Loss Function . . . . .	116
A.4	Initial Hyperparameter Optimization Results for Apartments – Log-Cosh Loss Function . . . . .	116
A.5	Initial Hyperparameter Optimization Results for Apartments – Cauchy Loss Function . . . . .	116
A.6	Initial Hyperparameter Optimization Results for Apartments – Welsch Loss Function . . . . .	117
A.7	Initial Hyperparameter Optimization Results for Apartments – Fair Loss Function . . . . .	117
A.8	Initial Hyperparameter Optimization Results for Houses – SE Loss Function . . . . .	118
A.9	Initial Hyperparameter Optimization Results for Houses – Log-Cosh Loss Function . . . . .	118
A.10	Initial Hyperparameter Optimization Results for Houses – Cauchy Loss Function . . . . .	118
A.11	Initial Hyperparameter Optimization Results for Houses – Welsch Loss Function . . . . .	119
A.12	Initial Hyperparameter Optimization Results for Houses – Fair Loss Function . . . . .	119
A.13	Performance Metrics for Apartments with the Tight Selection . . . . .	126
A.14	Performance Metrics for Houses with the Tight Selection . . . . .	126

B.1	Number of different types of jets for MC and MC <sub>b</sub> written in relative numbers such that each row sum to 100 %. See also Table 5.3.	127
B.2	Random Search PDFs for XGB	127
B.3	Global SHAP Feature Importances for the <i>g</i> -Tagging Models in 3-Jet Events	127

# Bibliography

- [1] Advanced Topics in Machine Learning (ATML). URL <https://kurser.ku.dk/course/ndak15014u>.
- [2] Allstate Claims Severity - Fair Loss. URL <https://kaggle.com/c/allstate-claims-severity>.
- [3] Dmlc/xgboost. URL <https://github.com/dmlc/xgboost>.
- [4] HEP meets ML award | The Higgs Machine Learning Challenge. URL <https://higgsml.lal.in2p3.fr/prizes-and-award/award/>.
- [5] The Large Electron-Positron Collider | CERN. URL <https://home.cern/science/accelerators/large-electron-positron-collider>.
- [6] Microsoft/LightGBM. URL [https://github.com/microsoft/LightGBM/blob/b397555d7023fd05f8e56326905fe7b185109de/src/treelearner/serial\\_tree\\_learner.cpp#L282](https://github.com/microsoft/LightGBM/blob/b397555d7023fd05f8e56326905fe7b185109de/src/treelearner/serial_tree_learner.cpp#L282).
- [7] Scikit-hep/uproot. URL <https://github.com/scikit-hep/uproot>.
- [8] Datashader: Revealing the Structure of Genuinely Big Data. URL <https://github.com/holoviz/datashader>.
- [9] O . Wwww.OIS.dk - Din genvej til ejendomsdata. URL <https://www.ois.dk/>.
- [10] M. Abadi et al. TensorFlow: Large-scale machine learning on heterogeneous systems. URL <http://tensorflow.org/>.
- [11] Y. S. Abu-Mostafa, M. Magdon-Ismail, and H.-T. Lin. *Learning From Data*. AMLBook. ISBN 978-1-60049-006-4.
- [12] D. Albanese, S. Riccadonna, C. Donati, and P. Franceschi. A practical tool for maximal information coefficient analysis. 7. ISSN 2047-217X. doi: 10.1093/gigascience/giy032. URL <https://doi.org/10.1093/gigascience/giy032>.
- [13] H. Albrecht, H. Ehrlichmann, T. Hamacher, R. P. Hofmann, T. Kirchhoff, A. Nau, S. Nowak, H. Schröder, H. D. Schulz, M. Walter, R. Wurth, C. Hast, H. Kolanoski, A. Kosche,

- A. Lange, A. Lindner, R. Mankel, M. Schieber, T. Siegmund, B. Spaan, H. Thurn, D. Töpfer, D. Wegener, M. Bittrner, P. Eckstein, M. Paulini, K. Reim, H. Wegener, R. Eckmann, R. Mundt, T. Oest, R. Reiner, W. Schmidt-Parzefall, W. Funk, J. Stiewe, S. Werner, K. Ehret, W. Hofmann, A. Hüpper, S. Khan, K. T. Knöpfle, M. Seeger, J. Spengler, D. I. Britton, C. E. K. Charlesworth, K. W. Edwards, E. R. F. Hyatt, H. Kapitza, P. Krieger, D. B. MacFarlane, P. M. Patel, J. D. Prentice, P. R. B. Saull, K. Tzamariudaki, R. G. Van de Water, T. S. Yoon, D. Reßing, M. Schmidtler, M. Schneider, K. R. Schubert, K. Strahl, R. Waldi, S. Weseler, G. Kernel, P. Križnič, T. Podobnik, T. Živko, V. Balagura, I. Belyaev, S. Chechelnitsky, M. Danilov, A. Droutskoy, Y. Gershtein, A. Golutvin, G. Kostina, D. Litvintsev, V. Lubimov, P. Pakhlov, F. Ratnikov, S. Semenov, A. Snizhko, V. Soloshenko, I. Tichomirov, and Y. Zaitsev. A model-independent determination of the inclusive semileptonic decay fraction of B mesons. 318(2): 397–404. ISSN 0370-2693. doi: 10.1016/0370-2693(93)90146-9. URL <http://www.sciencedirect.com/science/article/pii/0370269393901469>.
- [14] E. Anderson. The Species Problem in Iris. 23(3):457–509. ISSN 00266493. doi: 10.2307/2394164. URL [www.jstor.org/stable/2394164](http://www.jstor.org/stable/2394164).
  - [15] B. Andersson, G. Gustafson, G. Ingelman, and T. Sjöstrand. Parton fragmentation and string dynamics. 97(2): 31–145. ISSN 0370-1573. doi: 10.1016/0370-1573(83)90080-7. URL <http://www.sciencedirect.com/science/article/pii/0370157383900807>.
  - [16] S. R. Armstrong. A Search for the standard model Higgs boson in four jet final states at center-of-mass energies near 183-GeV with the ALEPH detector at LEP. URL <http://wwwlib.umi.com/dissertations/fullcit?p9910371>.
  - [17] M. Awad and R. Khanna. Support Vector Regression. In M. Awad and R. Khanna, editors, *Efficient Learning Machines: Theories, Concepts, and Applications for Engineers and System Designers*, pages 67–80. Apress. ISBN 978-1-4302-5990-9. doi: 10.1007/978-1-4302-5990-9\_4. URL [https://doi.org/10.1007/978-1-4302-5990-9\\_4](https://doi.org/10.1007/978-1-4302-5990-9_4).
  - [18] R. J. Barlow. *Statistics: A Guide to the Use of Statistical Methods in the Physical Sciences (Manchester Physics Series)*. WileyBlackwell, reprint edition. ISBN 0-471-92295-1. URL <http://www.amazon.co.uk/Statistics-Statistical-Physical-Sciences-Manchester/dp/0471922951%3FSubscriptionId%3D13CT5CVB80YFWJEPWS02%26tag%3Dws%26linkCode%3Dxm%26camp%3D2025%26creative%3D165953%26creativeASIN%3D0471922951>.

- [19] J. T. Barron. A General and Adaptive Robust Loss Function. URL <http://arxiv.org/abs/1701.03077>.
- [20] W. Bartel et al. Experimental study of jets in electron-positron annihilation. 101(1):129–134. ISSN 0370-2693. doi: 10.1016/0370-2693(81)90505-0. URL <http://www.sciencedirect.com/science/article/pii/0370269381905050>.
- [21] E. Becht, C.-A. Dutertre, I. W. H. Kwok, L. G. Ng, F. Ginhoux, and E. W. Newell. Evaluation of UMAP as an alternative to t-SNE for single-cell data. page 298430, . doi: 10.1101/298430. URL <https://www.biorxiv.org/content/10.1101/298430v1>.
- [22] E. Becht, L. McInnes, J. Healy, C.-A. Dutertre, I. W. H. Kwok, L. G. Ng, F. Ginhoux, and E. W. Newell. Dimensionality reduction for visualizing single-cell data using UMAP. 37(1):38–44, . ISSN 1546-1696. doi: 10.1038/nbt.4314. URL <https://www.nature.com/articles/nbt.4314>.
- [23] J. Bergstra and Y. Bengio. Random Search for Hyper-parameter Optimization. 13:281–305. ISSN 1532-4435. URL <http://dl.acm.org/citation.cfm?id=2188385.2188395>.
- [24] C. Bierlich. Rope hadronization, geometry and particle production in pp and pA Collisions. URL <https://lup.lub.lu.se/search/ws/files/18474576/thesis.pdf>.
- [25] Bolighed. Bolighed - usikkerhed i data-vurderingen. URL <https://bolighed.dk/om-bolighed/spoergsmaal-og-svar/#boligvaerdi>.
- [26] L. Breiman. Random Forests. 45(1):5–32. ISSN 1573-0565. doi: 10.1023/A:1010933404324. URL <https://doi.org/10.1023/A:1010933404324>.
- [27] R. P. Brent. *Algorithms for Minimization without Derivatives*. Prentice-Hall Series in Automatic Computation. Prentice-Hall. URL <https://cds.cern.ch/record/113464>.
- [28] E. Brochu, V. M. Cora, and N. de Freitas. A Tutorial on Bayesian Optimization of Expensive Cost Functions, with Application to Active User Modeling and Hierarchical Reinforcement Learning. URL <http://arxiv.org/abs/1012.2599>.
- [29] R. Brun and F. Rademakers. ROOT — An object oriented data analysis framework. 389(1):81–86. ISSN 0168-9002. doi: 10.1016/S0168-9002(97)00048-X. URL <http://www.sciencedirect.com/science/article/pii/S016890029700048X>.
- [30] A. Buckley, J. Butterworth, S. Gieseke, D. Grellscheid, S. Hoche, H. Hoeth, F. Krauss, L. Lonnblad, E. Nurse, P. Richardson, S. Schumann, M. H. Seymour, T. Sjostrand,

- P. Skands, and B. Webber. General-purpose event generators for LHC physics. 504(5):145–233. ISSN 03701573. doi: 10.1016/j.physrep.2011.03.005. URL <http://arxiv.org/abs/1101.2599>.
- [31] C. Burgard. Standard model of physics | TikZ example. URL <http://www.texample.net/tikz/examples/model-physics/>.
  - [32] D. Buskulic et al. An investigation of B<sub>d</sub> and B<sub>s</sub> oscillation. 322(4):441–458. ISSN 0370-2693. doi: 10.1016/0370-2693(94)91177-0. URL <http://www.sciencedirect.com/science/article/pii/0370269394911770>.
  - [33] M. Cacciari, G. P. Salam, and G. Soyez. FastJet user manual. 72(3):1896. ISSN 1434-6044, 1434-6052. doi: 10.1140/epjc/s10052-012-1896-2. URL <http://arxiv.org/abs/1111.6097>.
  - [34] S. Catani, Y. L. Dokshitzer, M. H. Seymour, and B. R. Webber. Longitudinally-invariant kt-clustering algorithms for hadron-hadron collisions. 406(1):187–224, . ISSN 0550-3213. doi: 10.1016/0550-3213(93)90166-M. URL <http://www.sciencedirect.com/science/article/pii/055032139390166M>.
  - [35] S. Catani, G. Turnock, and B. R. Webber. Jet broadening measures in e+ e- annihilation. B295:269–276, . doi: 10.1016/0370-2693(92)91565-Q.
  - [36] T. Chen and C. Guestrin. XGBoost: A Scalable Tree Boosting System. pages 785–794. doi: 10.1145/2939672.2939785. URL <http://arxiv.org/abs/1603.02754>.
  - [37] A. Collaboration. Electron efficiency measurements with the ATLAS detector using 2012 LHC proton-proton collision data. 77(3):195, . ISSN 1434-6044, 1434-6052. doi: 10.1140/epjc/s10052-017-4756-2. URL <http://arxiv.org/abs/1612.01456>.
  - [38] C. Collaboration. Search for a Higgs boson in the decay channel H to ZZ(\*) to q qbar l-l+ in pp collisions at sqrt(s) = 7 TeV. 2012(4):36, . ISSN 1029-8479. doi: 10.1007/JHEP04(2012)036. URL <http://arxiv.org/abs/1202.1416>.
  - [39] T. A. Collaboration. Observation of a new particle in the search for the Standard Model Higgs boson with the ATLAS detector at the LHC. 716(1):1–29, . ISSN 03702693. doi: 10.1016/j.physletb.2012.08.020. URL <http://arxiv.org/abs/1207.7214>.
  - [40] T. C. Collaboration. Observation of a new boson at a mass of 125 GeV with the CMS experiment at the LHC. 716(1):30–61, . ISSN 03702693. doi: 10.1016/j.physletb.2012.08.021. URL <http://arxiv.org/abs/1207.7235>.
  - [41] M. Dam. An upper limit for Br(Z->ggg) from symmetric 3-jet zo hadronic decays. 389(2):405–415. ISSN 0370-2693. doi: 10.1016/S0370-2693(96)01450-5.

- [42] A. Diaz-Papkovich, L. Anderson-Trocmé, C. Ben-Eghan, and S. Gravel. UMAP reveals cryptic population structure and phenotype heterogeneity in large genomic cohorts. 15(11). ISSN 1553-7390. doi: 10.1371/journal.pgen.1008432. URL <https://www.ncbi.nlm.nih.gov/pmc/articles/PMC6853336/>.
- [43] Y. L. Dokshitzer, G. D. Leder, S. Moretti, and B. R. Webber. Better Jet Clustering Algorithms. 1997(08):001–001. ISSN 1029-8479. doi: 10.1088/1126-6708/1997/08/001. URL <http://arxiv.org/abs/hep-ph/9707323>.
- [44] S. D. DST. Price Index (EJ14) - Statistics Denmark. URL <https://www.dst.dk/en/Statistik/emner/priser-og-forbrug/ejendomme>.
- [45] S. D. Ellis and D. E. Soper. Successive combination jet algorithm for hadron collisions. 48(7):3160–3166. doi: 10.1103/PhysRevD.48.3160. URL <https://link.aps.org/doi/10.1103/PhysRevD.48.3160>.
- [46] D. et al. Buskulic. A precise measurement of hadrons. 313(3):535–548. ISSN 0370-2693. doi: 10.1016/0370-2693(93)90028-G. URL <http://www.sciencedirect.com/science/article/pii/037026939390028G>.
- [47] F. Faye. Frederik Faye / deepcalo. URL <https://gitlab.com/ffaye/deepcalo>.
- [48] R. A. Fisher. The Use of Multiple Measurements in Taxonomic Problems. 7(2):179–188. ISSN 2050-1439. doi: 10.1111/j.1469-1809.1936.tb02137.x. URL <https://onlinelibrary.wiley.com/doi/abs/10.1111/j.1469-1809.1936.tb02137.x>.
- [49] Y. Freund and R. E. Schapire. A desicion-theoretic generalization of on-line learning and an application to boosting. In P. Vitányi, editor, *Computational Learning Theory*, pages 23–37. Springer Berlin Heidelberg. ISBN 978-3-540-49195-8. Adaboost.
- [50] S. L. Glashow. Partial-symmetries of weak interactions. 22(4):579–588. ISSN 0029-5582. doi: 10.1016/0029-5582(61)90469-2. URL <http://www.sciencedirect.com/science/article/pii/0029558261904692>.
- [51] P. Gras, S. Höche, D. Kar, A. Larkoski, L. Lönnblad, S. Plätzer, A. Siódmok, P. Skands, G. Soyez, and J. Thaler. Systematics of quark/gluon tagging. 2017(7):91. ISSN 1029-8479. doi: 10.1007/JHEP07(2017)091. URL <http://arxiv.org/abs/1704.03878>.
- [52] N. Guttenberg, N. Virgo, O. Witkowski, H. Aoki, and R. Kanai. Permutation-equivariant neural networks applied to dynamics prediction. URL <http://arxiv.org/abs/1612.04530>.

- [53] A. E. Harvey and S. Peters. Estimation Procedures for Structural Time Series Models. doi: [10.1002/for.3980090203](https://doi.org/10.1002/for.3980090203).
- [54] T. Hastie and R. Tibshirani. Generalized Additive Models: Some Applications. 82(398):371–386. ISSN 01621459. doi: [10.2307/2289439](https://doi.org/10.2307/2289439). URL [www.jstor.org/stable/2289439](http://www.jstor.org/stable/2289439).
- [55] T. Hastie, R. Tibshirani, and J. Friedman. *The Elements of Statistical Learning: Data Mining, Inference, and Prediction, Second Edition*. Springer Series in Statistics. Springer-Verlag, 2 edition. ISBN 978-0-387-84857-0. URL [//www.springer.com/la/book/9780387848570](http://www.springer.com/la/book/9780387848570).
- [56] K. Hornik. Approximation capabilities of multilayer feedforward networks. 4(2):251–257. ISSN 0893-6080. doi: [10.1016/0893-6080\(91\)90009-T](https://doi.org/10.1016/0893-6080(91)90009-T). URL <http://www.sciencedirect.com/science/article/pii/089360809190009T>.
- [57] P. Huber and E. Ronchetti. *Robust Statistics*. Wiley Series in Probability and Statistics. Wiley. ISBN 978-1-118-21033-8. URL [https://books.google.dk/books?id=j10hquR\\_j88C](https://books.google.dk/books?id=j10hquR_j88C).
- [58] S. Hviid, Juul. Working Paper: A regional model of the Danish housing market. URL <http://www.nationalbanken.dk/en/publications/Pages/2017/11/Working-Paper-A-regional-model-of-the-Danish-housing-market.aspx>.
- [59] F. James and M. Roos. Minuit – a system for function minimization and analysis of the parameter errors and correlations. 10:343–367. doi: [10.1016/0010-4655\(75\)90039-9](https://doi.org/10.1016/0010-4655(75)90039-9).
- [60] R. E. Kalman. A new approach to linear filtering and prediction problems. 82:35–45.
- [61] G. Ke, Q. Meng, T. Finley, T. Wang, W. Chen, W. Ma, Q. Ye, and T.-Y. Liu. LightGBM: A Highly Efficient Gradient Boosting Decision Tree. In I. Guyon, U. V. Luxburg, S. Bengio, H. Wallach, R. Fergus, S. Vishwanathan, and R. Garnett, editors, *Advances in Neural Information Processing Systems 30*, pages 3146–3154. Curran Associates, Inc. URL <http://papers.nips.cc/paper/6907-lightgbm-a-highly-efficient-gradient-boosting-decision-tree.pdf>.
- [62] D. P. Kingma and J. Ba. Adam: A Method for Stochastic Optimization. URL <http://arxiv.org/abs/1412.6980>.
- [63] G. Klambauer, T. Unterthiner, A. Mayr, and S. Hochreiter. Self-Normalizing Neural Networks. URL <http://arxiv.org/abs/1706.02515>.
- [64] A. J. Larkoski, J. Thaler, and W. J. Waalewijn. Gaining (Mutual) Information about Quark/Gluon Discrimination. 2014 (11). ISSN 1029-8479. doi: [10.1007/JHEP11\(2014\)129](https://doi.org/10.1007/JHEP11(2014)129). URL <http://arxiv.org/abs/1408.3122>.

- [65] C. Leys, C. Ley, O. Klein, P. Bernard, and L. Licata. Detecting outliers: Do not use standard deviation around the mean, use absolute deviation around the median. 49(4):764–766. ISSN 0022-1031. doi: 10.1016/j.jesp.2013.03.013. URL <http://www.sciencedirect.com/science/article/pii/S0022103113000668>.
- [66] S. M. Lundberg and S.-I. Lee. A unified approach to interpreting model predictions. In *Proceedings of the 31st International Conference on Neural Information Processing Systems*, NIPS’17, pages 4768–4777. Curran Associates Inc. ISBN 978-1-5108-6096-4. URL <http://dl.acm.org/citation.cfm?id=3295222.3295230>.
- [67] S. M. Lundberg, G. G. Erion, and S.-I. Lee. Consistent Individualized Feature Attribution for Tree Ensembles - SHAP. . URL <http://arxiv.org/abs/1802.03888>.
- [68] S. M. Lundberg, G. G. Erion, and S.-I. Lee. Consistent Individualized Feature Attribution for Tree Ensembles. . URL <http://arxiv.org/abs/1802.03888>.
- [69] A. L. Maas. Rectifier nonlinearities improve neural network acoustic models.
- [70] L. McInnes and J. Healy. UMAP: Uniform Manifold Approximation and Projection for Dimension Reduction. URL <http://arxiv.org/abs/1802.03426>.
- [71] T. C. Mills. *Time Series Techniques for Economists / Terence c. Mills*. Cambridge University Press Cambridge [England] ; New York. ISBN 0-521-34339-9 0-521-40574-2. URL <http://www.loc.gov/catdir/toc/cam031/89007187.html>.
- [72] N. Mohd Razali and B. Yap. Power Comparisons of Shapiro-Wilk, Kolmogorov-Smirnov, Lilliefors and Anderson-Darling Tests. 2.
- [73] I. Mulalic, H. Rasmussen, J. Rouwendal, and H. H. Woltmann. The Financial Crisis and Diverging House Prices: Evidence from the Copenhagen Metropolitan Area. ISSN 1556-5068. doi: 10.2139/ssrn.3041272. URL <https://www.ssrn.com/abstract=3041272>.
- [74] Particle Data Group et al. Review of Particle Physics. 98(3):030001. doi: 10.1103/PhysRevD.98.030001. URL <https://link.aps.org/doi/10.1103/PhysRevD.98.030001>.
- [75] F. Pedregosa, G. Varoquaux, A. Gramfort, V. Michel, B. Thirion, O. Grisel, M. Blondel, P. Prettenhofer, R. Weiss, V. Dubourg, J. Vanderplas, A. Passos, D. Cournapeau, M. Brucher, M. Perrot, and E. Duchesnay. Scikit-learn: Machine learning in Python. 12:2825–2830.

- [76] E. Polley and M. van der Laan. Super Learner In Prediction. URL <https://biostats.bepress.com/ucbbiostat/paper266>.
- [77] J. Proriol, J. Jousset, C. Guicheney, A. Falvard, P. Henrard, D. Pallin, P. Perret, and B. Brandl. TAGGING B QUARK EVENTS IN ALEPH WITH NEURAL NETWORKS (comparison of different methods : Neural Networks and Discriminant Analysis). page 27.
- [78] A. Purcell. Go on a particle quest at the first CERN webfest. URL <https://cds.cern.ch/record/1473657>.
- [79] S. Ravanbakhsh, J. Schneider, and B. Poczos. Deep Learning with Sets and Point Clouds. URL <http://arxiv.org/abs/1611.04500>.
- [80] D. N. Reshef, Y. A. Reshef, H. K. Finucane, S. R. Grossman, G. McVean, P. J. Turnbaugh, E. S. Lander, M. Mitzenmacher, and P. C. Sabeti. Detecting Novel Associations in Large Data Sets. 334(6062):1518–1524. ISSN 0036-8075, 1095-9203. doi: 10.1126/science.1205438. URL <http://science.sciencemag.org/content/334/6062/1518>.
- [81] J. W. Rohlfs. *Modern Physics from A to Z*. John Wiley and Sons. ISBN 978-0-471-57270-1.
- [82] P. J. Rousseeuw and C. Croux. Alternatives to the Median Absolute Deviation. 88(424):1273–1283. ISSN 0162-1459. doi: 10.1080/01621459.1993.10476408. URL <https://www.tandfonline.com/doi/abs/10.1080/01621459.1993.10476408>.
- [83] A. Salam. Weak and electromagnetic interactions. In *Selected Papers of Abdus Salam*, volume Volume 5 of *World Scientific Series in 20th Century Physics*, pages 244–254. WORLD SCIENTIFIC. ISBN 978-981-02-1662-7. doi: 10.1142/9789812795915\_0034. URL [https://www.worldscientific.com/doi/abs/10.1142/9789812795915\\_0034](https://www.worldscientific.com/doi/abs/10.1142/9789812795915_0034).
- [84] S. Schmitt. Data Unfolding Methods in High Energy Physics. 137:11008. ISSN 2100-014X. doi: 10.1051/epjconf/201713711008. URL <http://arxiv.org/abs/1611.01927>.
- [85] L. Scodellaro. B tagging in ATLAS and CMS. URL <http://arxiv.org/abs/1709.01290>.
- [86] L. Shapley. A value for n-person games. In *The Shapley Value*, volume 28 of *Annals of Math Studies*, pages 307–317. doi: 10.1017/CBO9780511528446.003.
- [87] T. Sjöstrand, S. Ask, J. R. Christiansen, R. Corke, N. De-sai, P. Ilten, S. Mrenna, S. Prestel, C. O. Rasmussen, and P. Z. Skands. An Introduction to PYTHIA 8.2. 191:159–177. ISSN 00104655. doi: 10.1016/j.cpc.2015.01.024. URL <http://arxiv.org/abs/1410.3012>.

- [88] P. Skands. Peter Skands.
- [89] S. J. Taylor and B. Letham. Forecasting at scale. doi: 10.7287/peerj.preprints.3190v2. URL <https://peerj.com/preprints/3190>.
- [90] i. team. Iminuit – A python interface to minuit. URL <https://github.com/scikit-hep/iminuit>.
- [91] J. Thaler. Report of the Les Houches Quark/Gluon Subgroup. (1):28.
- [92] R. Tibshirani. Regression Shrinkage and Selection via the Lasso. 58(1):267–288. ISSN 0035-9246. URL [www.jstor.org/stable/2346178](http://www.jstor.org/stable/2346178).
- [93] A. Tikhonov. *On the Stability of Inverse Problems*, volume vol. 39 of *Doklady Akademii Nauk SSSR*.
- [94] S. Toghi Eshghi, A. Au-Yeung, C. Takahashi, C. R. Bolen, M. N. Nyachienga, S. P. Lear, C. Green, W. R. Mathews, and W. E. O’Gorman. Quantitative Comparison of Conventional and t-SNE-guided Gating Analyses. 10. ISSN 1664-3224. doi: 10.3389/fimmu.2019.01194. URL <https://www.frontiersin.org/articles/10.3389/fimmu.2019.01194/full>.
- [95] d. L. M. J. van, E. C. Polley, and A. E. Hubbard. Super Learner. 6(1). ISSN 1544-6115. doi: 10.2202/1544-6115.1309. URL <https://www.degruyter.com/view/j/sagmb.2007.6.issue-1/sagmb.2007.6.1.1309/sagmb.2007.6.1.1309.xml>.
- [96] J. J. van der Bij and E. W. N. Glover. Z boson production and decay via gluons. 313(2):237–257. ISSN 0550-3213. doi: 10.1016/0550-3213(89)90317-9. URL <http://www.sciencedirect.com/science/article/pii/0550321389903179>.
- [97] L. van der Maaten and G. Hinton. Visualizing Data using t-SNE. 9:2579–2605. ISSN ISSN 1533-7928. URL <http://www.jmlr.org/papers/v9/vandermaaten08a.html>.
- [98] F. van Veen. The Neural Network Zoo. URL <http://www.asimovinstitute.org/neural-network-zoo/>.
- [99] V. Vapnik. Principles of Risk Minimization for Learning Theory. In *Proceedings of the 4th International Conference on Neural Information Processing Systems*, NIPS’91, pages 831–838. Morgan Kaufmann Publishers Inc. ISBN 978-1-55860-222-9. URL <http://dl.acm.org/citation.cfm?id=2986916.2987018>.
- [100] P. Virtanen, R. Gommers, T. E. Oliphant, M. Haberland, T. Reddy, D. Cournapeau, E. Burovski, P. Peterson, W. Weckesser, J. Bright, S. J. van der Walt, M. Brett, J. Wilson, K. Jarrod Millman, N. Mayorov, A. R. J. Nelson, E. Jones, R. Kern, E. Larson, C. J. Carey, I. Polat, Y. Feng, E. W. Moore,

- J. VanderPlas, D. Laxalde, J. Perktold, R. Cimrman, I. Henriksen, E. A. Quintero, C. R. Harris, A. M. Archibald, A. H. Ribeiro, F. Pedregosa, P. van Mulbregt, and S. . . Contributors. SciPy 1.0—Fundamental Algorithms for Scientific Computing in Python. page arXiv:1907.10121. URL <https://ui.adsabs.harvard.edu/abs/2019arXiv190710121V/abstract>.
- [101] I. Wallach and R. Lilien. The protein–small-molecule database, a non-redundant structural resource for the analysis of protein-ligand binding. 25(5):615–620. ISSN 1367-4803. doi: 10.1093/bioinformatics/btp035. URL <https://academic.oup.com/bioinformatics/article/25/5/615/183421>.
  - [102] S. Weinberg. A Model of Leptons. 19(21):1264–1266. doi: 10.1103/PhysRevLett.19.1264. URL <https://link.aps.org/doi/10.1103/PhysRevLett.19.1264>.
  - [103] H. Wickham. Tidy data. 59(10):1–23. ISSN 1548-7660. doi: 10.18637/jss.v059.i10. URL <https://www.jstatsoft.org/v059/i10>.
  - [104] M. Wobisch and T. Wengler. Hadronization Corrections to Jet Cross Sections in Deep-Inelastic Scattering. URL <http://arxiv.org/abs/hep-ph/9907280>.
  - [105] M. Zaheer, S. Kottur, S. Ravanbakhsh, B. Poczos, R. Salakhutdinov, and A. Smola. Deep Sets. URL <http://arxiv.org/abs/1703.06114>.