

IML 2023: Predicting Saturation Vapour Pressure From Molecular Properties

Helmi Karesti, Janne Penttala, Maija Säisä (‘Group 144’ in Kaggle)

1 Introduction

The objective of this project was to build a machine learning model to predict the saturation vapour pressure (pSat) for atmospherically relevant molecules, based on their interpretable features. The project is based on the GeckoQ dataset that consists of 31,637 molecules with their atomic structures. The molecular data provides relevant information on aerosol particle growth, for which the saturation vapour pressure is a key parameter. Part of this project was to take part in a Kaggle competition, for which we were given a training set of 27,147 molecules and a test set of 4,490 molecules. Our goal was to predict the target variable for the test data, and our prediction accuracy was evaluated by the R2 score.

We started the project by getting familiar with the data through preprocessing and data exploration. After we had verified that our data was clean, we used the data to compare and select our model. After the model selection, we did some parameter tuning and applied a few different feature engineering methods to try and optimize the performance and the accuracy of our model. This report contains a more in-depth explanation of our methods regarding the steps of building the model.

2 Data analysis

2.1 Preprocessing

As a first step to preprocessing we made sure that our data was valid. In other words, we wanted to ensure that there were no missing or mismatched values for any of the variables. This was quite straightforward as Kaggle provided us variable-specific graphs for each variable, meaning we could verify the validity of both the training and test data with the information from Kaggle. We noted that all the data points in both data sets were valid, and thus, the variables did not have missing values.

After validating the data, there were a few minor changes we decided to make to further prepare our data. First, we decided to remove the *Id* column from both data sets, as this would not be needed in predicting our target variable. Second, we applied one-hot encoding for a variable called *parentspecies*, as this was the only categorical variable. This encoding transformed the variable *parentspecies* into eight binary columns (for example, *category_apin*, *category_apin_decane* etc.) for the train set and seven binary columns for the test set. Therefore, after applying one-hot encoding, we needed to create an extra column to the test set for it to match our training set.

Lastly we wanted to ensure the data consistency of both training and test data. To do this, we simply compared the two data sets to each other and made sure that both of them have the same columns in the same order. The number of columns in the preprocessed data sets was 32.

2.2 Data exploration

The data used in the project is based on GeckoQ dataset. Our first step with the project was to get familiar with this dataset and get a preliminary understanding of the features it has. Our plan was to first use material that was already provided by the project descriptions and then continue with some more data exploration by

handling the data ourselves. We first got familiar with the features and their descriptions. Understanding the features perfectly is not necessarily important in this context, however it is useful to have a general understanding of how they might affect our target variable.

2.2.1 Scaling considerations

We noticed that the range of values of the variables varied notably. Because we planned to test several models, some of which would perform better with scaled variables, we decided to apply the `StandardScaler()` on the data sets. We found some information, which claimed that scaling may not be appropriate for one-hot encoded variables. We considered whether we should leave these variables outside of the scaler function. We also found support for applying scaling to one-hot encoded variables to tackle the same challenges that the numerical variables face. Consequently, we tested this matter. Based on the results, not scaling the one-hot encoded variables seemed to have a tiny effect on the R2 score improving the score by around 0.0004-0.0012 for the models we tested, except for Lasso model where the R2 score actually declined by 0.0020. As the effect on the R2 score seemed insignificant based on this analysis, we ended up using the scaled values also for the categorical variables when using scaled values for the numerical variables. We have presented the R2 score for the tested models with and without scaling in Table 1 in the Results section.

2.2.2 Data visualization

To get a better understanding of our training data, we created three plots for each of the 32 explainable variables. For each explainable variable, we plotted

- i) the explainable variable against our target variable as a scatter plot,
- ii) the count of observations per explainable variable value as a histogram, and
- iii) the distribution of explainable variable values as a boxplot.

We did this with and without scaling, and concluded that the shape of the graphs, naturally, remains the same, and the only difference is the value of the explainable variable, that is, the x axis. The three graphs for one of the explainable variables (scaled *MW*) are shown in Figure 1 as an example. We also plotted some graphs for the test data in order to look for any differences between the distributions of the training and test data.

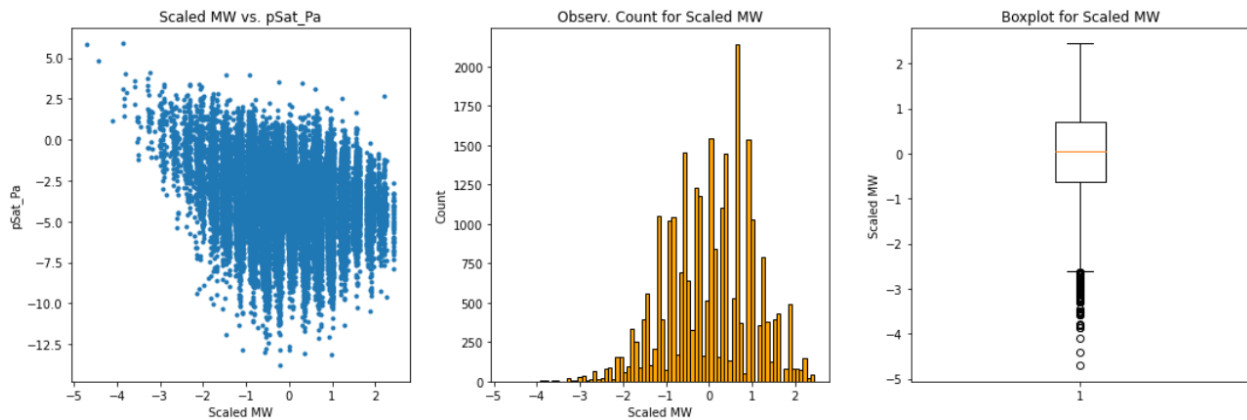


Figure 1: The three graphs for the (scaled) explainable variable *MW*. Equal plots created for each of the explainable variables.

The scatter plot

Analyzing the first graph, that is, the scatter plot, helped us to understand the relationship of each explainable variable and the target variable better. The scatter plots demonstrated well that there was only one continuous

variable while the rest of the variables were discrete (including the one-hot encoded categorical variable). That is, the *MW* or molecule’s molecular weight was the only continuous variable. However, we concluded that this information about discrete vs continuous variables does not guide us to any specific direction or provide any useful insights.

We also noticed that many of the variables had in maximum three different values, so the values were quite concentrated on specific points. Later, this information was used in feature selection, testing the effect of dropping the features with only a few different values.

Furthermore, we could observe a somewhat linear relationship for some of the variables, for example, *MW*, *NumHBondDonors*, *NumOfAtoms*, and *NumOfO* (see Figure 2). Therefore, also including the linear regression model and feature selection based on only these type of features was supported.

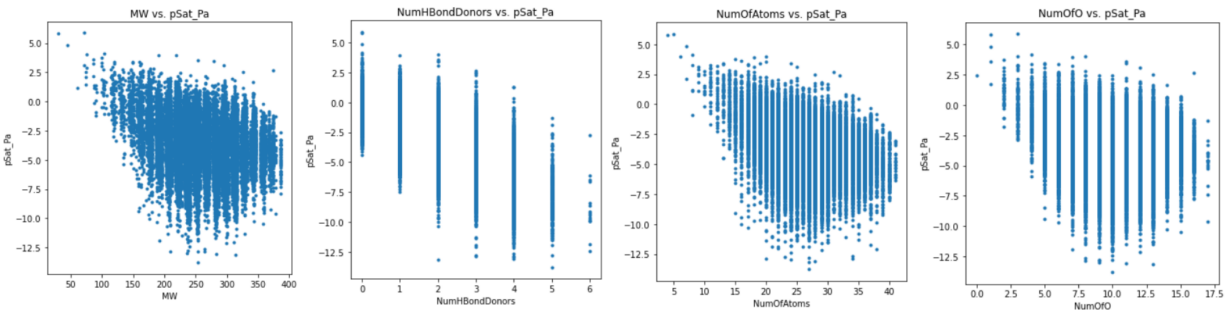


Figure 2: Somewhat linear relationship observed in some of the scatter plots.

The histogram

Analyzing the second graph, that is, the histogram, helped us to observe the occurrence of the various variable values. A selection of these different histograms is shown in Figure 3. We noticed that the histograms of the various variables remarkably differed from each other. Some of the count distributions reminded us of a half or full normal distribution while in the other extreme the values were concentrated around only one value.

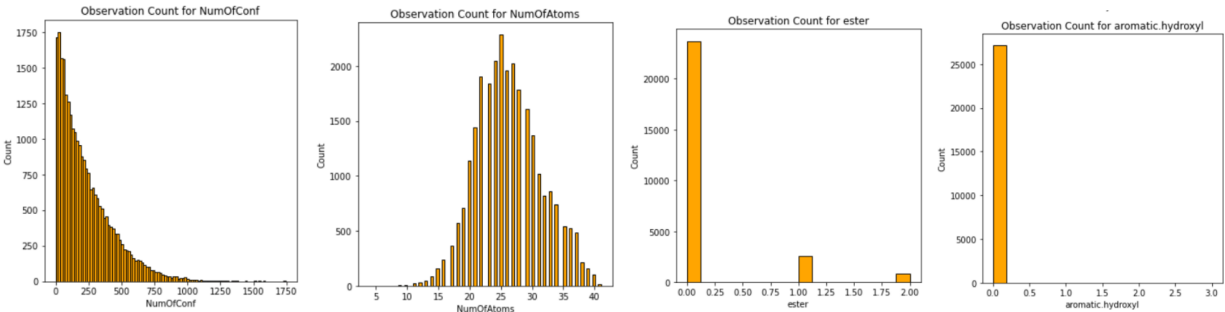


Figure 3: Various distributions for the count of observation values.

The boxplot

The third graph, that is the boxplot, complemented well the first two graphs. As expected, the distributions of the (scaled) values notably differed from each other between the different variables. We noticed that some distributions did not have any outliers while others had plenty of them (two graphs on the left in Figure 4). The variables, which were concentrated around a single value, demonstrated this behavior also in the boxplot; outside of the single value there were only outliers (third graph from the left in Figure 4). Furthermore, we noticed that from the one-hot encoded categorical variable, the *category_toluene* seemed to be the only one

with a remarkable amount of two values (the graph on the right in Figure 4). Due to this observation, we tested to drop the other categorical variables (except for *category_toluene*) to see the effect on the results.

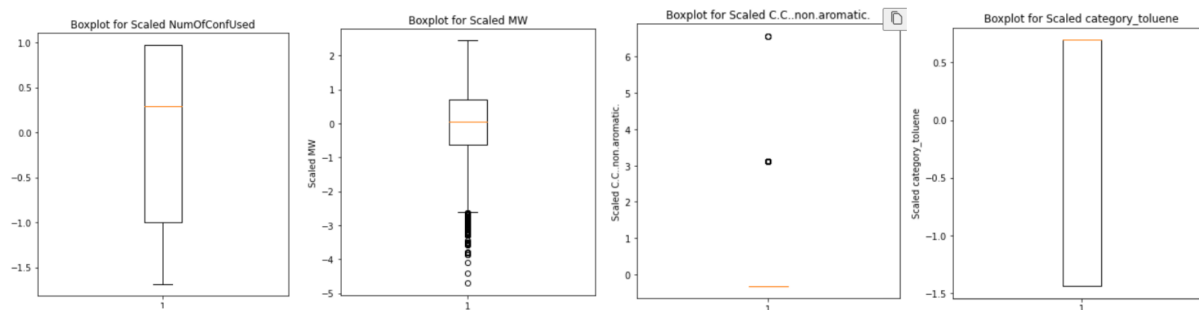


Figure 4: Distributions of observation values as a boxplot.

Comparison of training and test data

We plotted a histogram and a boxplot also for the variables in the test data in order to look for any differences between the distributions of the training and test data. The graphs of the training and test data for the variable *MW* are presented in Figure 5. When comparing the data in the graphs for the variable *MW*, for example, we noticed that the values in the test data started from 100 while the training data also included values below 100. Consequently, we used this information in feature engineering.

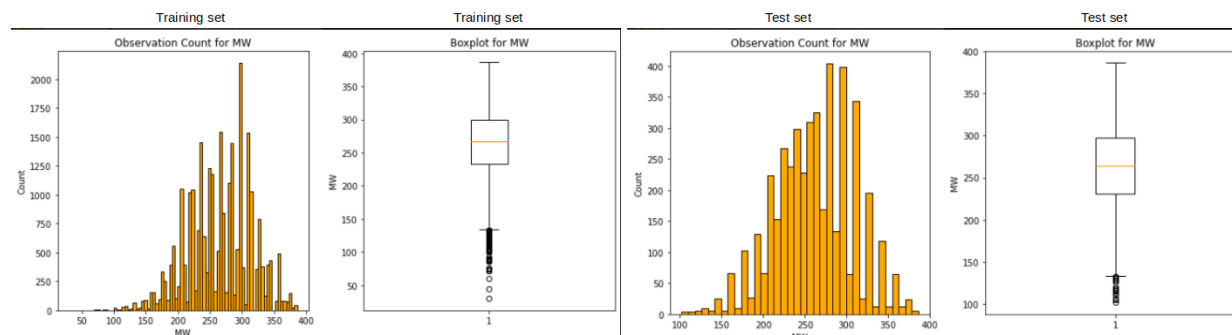


Figure 5: Histogram and boxplot for the MW variable in the training and test data.

2.2.3 Outlier considerations

The observations from the boxplots encouraged us to test dropping outliers from the data set.

First, we did this by dropping all the rows with at least one outlier based on the distance from the first quartile (Q1) or the third quartile (Q3); we defined the distance as a multiplier (for example, 1.5) times the interquartile range (IQR). However, this analysis did not improve the results because the number of the resulting rows was in the magnitude of 10% of the original number of rows. This analysis helped us to understand that there seemed to be a trade-off between the number of rows in the data set and excluding the outliers from the data set. Because this approach did not work well for our purposes, we tested another approach to drop outliers based on a threshold for z score.

We tested dropping out rows which included at least one variable that could be regarded as an outlier based on a threshold level (tested a threshold range from 2 to 4). That is, a variable that had a z score higher than the threshold level was considered as an outlier. The results made more sense because the number of rows in a data set after removing the outliers varied from 10,561 ($z < 2$) to 25,166 ($z < 4$). However, the conclusions were the same; a larger number of rows in the data set outweighed any benefit received from removing

outliers with this method for most of the models. Worth noting is that for the OLS Linear Regression model, the R2 score slightly increased when moving from a threshold level 4 to 3; regardless, the R2 score of this model was below SVR, so it would not be worth removing outliers from the data set based on this method.

3 Methods

The methods discussed in this section aim to improve the model R2 score in the Kaggle competition. However, the evaluation of the model optimization process for the Kaggle competition turned out to be challenging, as the R2 score in Kaggle did not correlate with local results. The cross-validated R2 score was locally at best varying between 0.71 and 0.76, while the results in Kaggle were between 0.62 and 0.66. The changes in the R2 score were often somewhat marginal both locally and in Kaggle, which made it difficult to evaluate the true effect of the selected model and the feature engineering steps. We were later informed that the test data we had for the Kaggle competition was actually somewhat skewed when compared to the training data. This could explain at least part of the differences we had between Kaggle & local scores. In order to stay consistent, the R2 scores discussed in this and the latter sections were obtained locally with the cross-validation method. For reader reference, the models we used were imported from the Python's `sklearn` and `xgboost` libraries. The tools needed in the model selection & parameter tuning, as well as in the feature engineering steps, were from the `sklearn`, `pandas` and `numpy` libraries.

3.1 Model selection

We started out by comparing different regression models on the preprocessed training data. Recall that in the preprocessing step we applied one-hot encoding on the *parentspecies* and did not yet drop any features apart from the *Id* column, meaning that we had 32 features in total. Initially, we tested with OLS Linear Regression, Random Forest, Support Vector Regressor, Lasso Regression and eXtreme Gradient Boosting. Later, we also added MLP Regression and Elastic Net to the comparison.

At this point, we did not set any model parameters apart from the random seed. We used cross-validation ($k=5$) to evaluate the R2 scores and fitted the models with and without scaling. The models that clearly benefited from the scaling were SVR and MLP. For Lasso and Elastic Net, the results actually worsened with scaling and for OLS Linear Regression, Random Forest and XGBoost the results remained more or less the same. The full results are presented later in section 4. At this stage, we obtained the best overall R2 score with MLP (0.7421) and SVR (0.7383) using scaled data. XGBR also performed well (0.7299), and due to its significantly faster performance, we used it in the latter feature engineering phases among the two other best performing models.

3.2 Parameter tuning

After model selection, we tried to improve the performance of the best models by tuning the hyperparameters. We tried using both `GridSearchCV()` and `RandomizedSearchCV()` to find the best hyperparameter combinations. For SVR, we tried different kernels (linear, rbf, poly), C values (0.1, 1, 10, 100), epsilon values (0.1, 0.2, 0.5) and gamma values (scale, auto). Unfortunately, we quickly realized that in this case the GridSearch with $k=5$ cross-validation took too long to run and we had to settle for a small RandomizedSearch and manual selection. In the end, we were able to obtain at least one hyperparameter combination that slightly improved the result for SVR. R2 score of 0.7444 was achieved with $C=10$ (default 1), $\text{kernel}=\text{'rbf'}$, $\text{gamma}=\text{'scale'}$ & $\text{epsilon}=0.5$ (default 0.1).

With XGBR, we managed to improve the R2 score from 0.7299 to 0.7388 using RandomizedSearch with 20 iterations and $k=5$ cross-validation. The best hyperparameters were: $\text{colsample_bytree}=0.983$, $\text{gamma}=0.162$, $\text{learning_rate}=0.166$, $\text{max_depth}=3$, $\text{min_child_weight}=17$, $\text{n_estimators}=596$, $\text{reg_alpha}=0.127$, $\text{reg_lambda}=0.123$ & $\text{subsample}=0.848$. We tried to find the best hyperparameters again after feature engineering phase, but the hyperparameters remained the same. As the MLP model was only added later on, we did not perform any parameter search for it.

3.3 Feature engineering

3.3.1 Scaling features

We had already scaled the train and test data in the preprocessing phase, but soon realized that when doing feature engineering and calculating R2 scores, we naturally had to use unscaled data and only scale in the cross-validation phase. In other words, the features were scaled separately in each engineering step.

3.3.2 Feature selection

As mentioned in the data exploration section, some of the observations encouraged us to test dropping certain features. We tested keeping only the features that showed a somewhat linear relationship with the target variable; this test did not improve the results for the linear regression model nor other models. We also tested, for example, dropping all the other categorical variables except for *category_toluene* (the only categorical variable notably having the two type of observations); however, the results slightly deteriorated with this change, so dropping only the seven one-hot encoded variables was not supported.

Furthermore, we tested to drop features with value distributions that were extremely concentrated on a certain value. These features included most of the one-hot encoded *parentspecies* variables: *category_None*, *category_apin_decane*, *category_apin_decane_toluene*, *category_apin_toluene* and *category_decane_toluene*. In addition, other such features were *C.C.C.O.in.non.aromatic.ring*, *aromatic.hydroxyl* and *nitroester*.

We also tried a few different combinations of different features, but the best result was obtained by dropping the features mentioned above. At this point, we were still trying out with both SVR & XGBR and they both had slight increases in their R2 scores. However, it is questionable if the improvements were actually significant since the changes in R2 were so marginal. With the features that were left, it was more difficult to evaluate their effect on the target variable. In the hindsight, we probably could have used some methods to try and evaluate the importance of the other features as well.

3.3.3 Creating new features

When reading the related article (Besel et al.), we found some hints about which features could possibly be more related to the target variable. There was a mention in the article that *‘Generally, a large number of functional groups correlates with a low pSat.’* From this we concluded that we could create a new feature to describe the amount of different functional groups in a molecule. After this, we dropped the original functional group features and calculated the R2 score. Unfortunately this did not produce better results, and the R2 score decreased significantly. When keeping the original features and the new feature, the R2 score increased ever so slightly. Again, it was questionable whether the improvement was actually significant.

3.3.4 Transforming features

Another hint we found in the article was related to nitro and nitrate groups and their molecular weights. The article states: *‘Beyond 220g/mol the abundance of nitrate (62 g/mol) and nitro (46 g/mol) groups dominates. These groups have a comparatively large mass, but their contribution to a lower pSat is small, which explains the saturation of pSat to a low value.’* Indeed, when we look at the scatter plot of MW vs. pSatPa (Figure 2), we can see that around MW value of 220 g/mol, the pSat values level out. Since the amount of nitro and nitrate groups was known, we tried to diminish their effect on the regression model by reducing their molecular weights from the total MW values. Unfortunately, we were once again disappointed as our model still didn’t seem to improve and the R2 score remained pretty much the same.

4 Results

Our results consist of initial results using seven different models and our final R2 scores using the best performing models with machine learning methods.

4.1 Initial results with different models

Our initial results include the R2 scores for seven different regression models with no hyperparameter tuning apart from the random seed. In Table 1 we can see how the models compare to each other with and without scaling.

Table 1: Results of the basic models with and without scaling

Model	R2 without scaling	R2 with scaling
MLP Regressor	0.5736	0.7421
SVR	0.4713	0.7383
XGBRegressor	0.7299	0.7299
Random Forest	0.7105	0.7106
OLS Linear Regressor	0.6957	0.6944
Elastic Net	0.4446	0.3990
Lasso	0.3295	0.2664

As we can see from Table 1, Support Vector Regressor and MLP Regressor performed the best out of our basic models. Random Forest, OLS Linear Regressor and XGBRegressor also gave quite good results. Lasso and Elastic Net performed significantly worse than the other models.

4.2 Final R2 scores

After seeing our initial results with different models, we applied some machine learning methods to the best performing models to see if we could improve the R2 score. The machine learning methods applied to achieve these results are discussed more in-depth in the previous section. The final results are presented in figure 6 below.

Model	R2score
SVR	0.7538
XGBRegressor	0.7505
MLP Regressor	0.7493

Figure 6: Results of the best models

The impact of the machine learning methods was not as big as expected. Support Vector Regressor and XGBRegressor had noticeable improvement over their basic models.

5 Conclusions

In conclusion we build a machine learning model and evaluated the performance of our model by R2 score. To build our model we experimented on several basic models and chose the best performing one. By applying machine learning methods we were able to improve our model further. In the model building process we tried out many different approaches to improve our model, however not all of them had noticeable impact.

6 Self-grading

6.1 Grade for the deliverables

Grade: 4

Our workflow for building a machine learning model and the methods we used were in our opinion appropriate for the task. We used multiple machine learning methods in the process of building our model, such as removing outliers, tuning hyperparameters and selecting features. Even though some of the methods did not end up being used to achieve the best results, exploring them gave us information about how the model behaves. A criticism of our approach to building the model would be the lack of domain knowledge we used. A better understanding of the subject could have helped us determine how the features affect our target variable, and our model building wouldn't have relied so much on experimentation.

Regarding the challenge submission our results were reasonable. Our initial results with simply selecting the right model were very promising. By applying machine learning methods to our model, we were able to improve the results; however, we were slightly underwhelmed by the impact our methods had.

We wrote the final report with the purpose of presenting our approach and solutions to the term project in a way that is easy to understand but still explains everything thoroughly. We decided to not show any code in our report. Instead, we present our methods by explaining the logic behind them and use figures and tables to help visualize our results. The main purpose of the report is to communicate what was done, how it was done and why it was done, and we think that our report does a good job at that.

6.2 Grade for the group as a whole

Grade: 4

The atmosphere within the group was positive. What we did well was discuss the subject matter and methods that we used in the project. Any concerns were discussed as a group and having multiple perspectives on a subject was helpful. What could have been improved upon is overall teamwork. Practical matters such as busy schedules and working remotely led to the group functioning more as individuals than as a collective team.