# Project 4

KIMBERLY ROETEN & EVAN WALSH

CIS 445-01: DATA MINING

For this project, the nodes that were used were a file import with the real estate data, a data partition node that only uses training and validation, as the data set is too small to use test data, a statexplore node to produce a statistical summary of the input data, then I used a filter node for the variable selection and transformation nodes leading up to interactive binning, and the 3 models we want in this diagram, neural network, linear regression, and memory base reasoning with 10 neighbors. These three models were also set up without the filtering nodes, to see the difference. So essentially there are two paths in this diagram, one that filters outliers and one that doesn't.

When calling to perform the data cleaning and variable transformation, it appears that the total area variable is likely to be a linear combination of the variables first floor, second floor, and upper area. Assuming the total area variable is a summation of these other variables, using all four in the model is not necessary, unless you want to find out the most significant variable leading up to that point, that makes the total area variable so significant.

The most significant variable that is determined by the linear regression node not connected to the filter node includes garage type. The most significant variable that is determined by the linear regression node that is connected to the filtering nodes includes is total area.

After examining the effects and p values, there is no need to use the impute or replacement node because after variable selection, transformation, and interactive binning are used, their purpose is to replace values, so there should be no missing values in the data set after these stages in the workflow diagram. The most significant f value after analyzing the effects is the age variable, which makes sense considering it was first of all assigned to the interval level, and these types have the most significant distance clearly defined between the values/levels this

variable takes. It is such a broad variable that in terms of results produced, and because age had negative values before it underwent the transformation and filtering node, the variables change based on their performance of the data cleaning.

In terms of errors in the assessment score distributions for all six models, I believe the best model to predict the sale prices more accurately would be the linear regression that does not undergo filtering, because the mean predicted has a much steadier flow than any of the other models, and proves that it will have the most significant results when it comes to sale price, the target variable. Not even removing the outliers after filtering and variable transformation will produce more accurate results with less errors, its only purpose is to improve data preparation and not have as many issues.

In terms of low RMSE, MAE, and high $R^2$, obtaining great results was hard as the older houses in low-end neighborhoods are difficult to appraise accurately due to their lack of homogeneity. In future projects, perhaps using a clustering node to group similar properties together would be the key to creating the models that could possibly yield better prediction results, or smaller errors so to speak.

In terms of how well this specific project got students like me more acquainted with the software, the instructions on how to complete it was a 5 on a scale of 1 to 10, because even though it gave us an idea of what nodes to use, the instructions were pretty vague compared to previous projects and tutorials. And not having used filter and MBR nodes before was challenging, because until we figure out how to set up the diagram and properly connect all the nodes, we were experiencing run errors, but hopefully our results will do. Evan's contribution to this project was the flow diagram, and Kimberly's contribution was the report summary.