# Project 8 - Decision trees and neural networks

Nate Taylor & Clark Farnsworth
Slides | Repository | Dataset

## Introduction

Housing prices are on the rise, especially in certain parts of the country, and many first-time buyers and sellers alike are uncertain as to what to expect when approaching a sale. Through decision tree and neural network predictive models, we hoped to answer questions Utah/Colorado buyers and sellers might have, such as: What kind of house can I afford with the money I can realistically put towards a house? What are similar houses going for right now? If I want a house with certain characteristics, what can I expect to pay? Using a dataset of houses on the market in mid-2022, we used each house bedroom, bathroom, size, acreage, and location data to answer these questions and discover trends.

## Dataset

The dataset we used was a USA Real Estate Dataset from Kaggle, which contained data for over 2 million houses on the market in mid-2022 across all US states. Data it provided included broker, status (recently sold or still on the market), price, numbers of bedrooms and bathrooms, acreage, location, size, and last sell date. The numerical data was especially useful in predicting housing prices, so our models used bedrooms, bathrooms, acreage, size, and state. We narrowed our data down to houses in Utah and Colorado only (about 20,000 of the samples), as our stakeholders reside or wish to reside there. In order to increase predictive power, we separated houses into price ranges based on price quantiles.
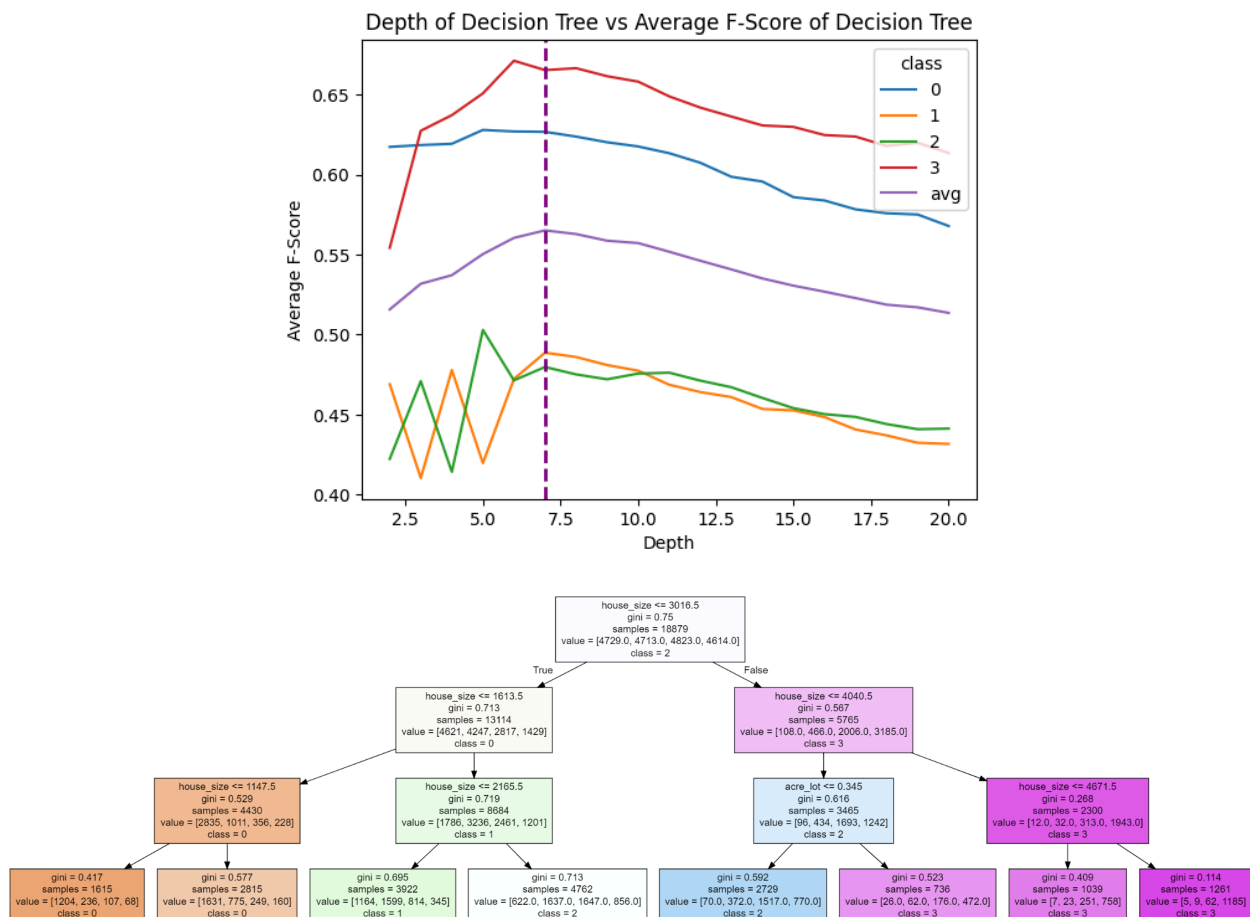
## Analysis technique

We wanted to easily visualize which house features contributed most to predicting house price, so we implemented decision tree models to predict ranges of prices. This allowed us to see basic decision paths (which can be thought of as house feature combinations) that led to certain price ranges. We ran 30 iterations of Monte-Carlo cross-validation with 75%/25% training/testing data splits to ensure that our models didn't overfit the data we were using. We also utilized neural networks to further isolate important house features and to see if we could get more powerful predictions.
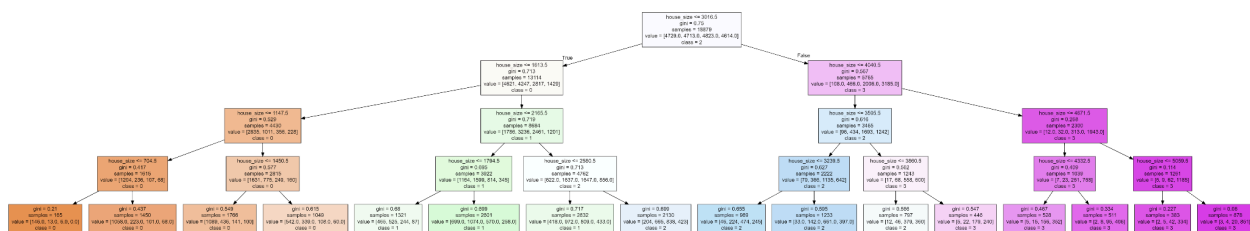
## Results

We found that our predictive models weren't as powerful as we hoped they would be at predicting price ranges of houses. On average, our best predictions could predict houses in the price range of $12,000 to $430,000 with an f-score of about 62.7% and houses in the price range of $775,000 to $3,190,000 with an f-score of 66.5%.
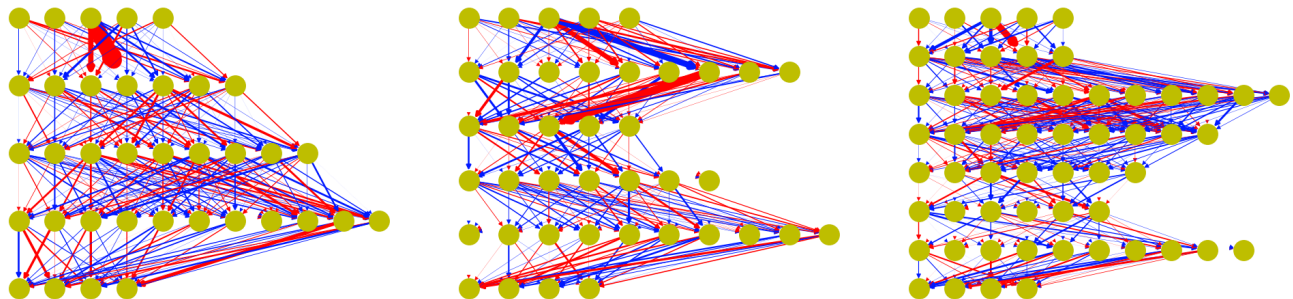
Predictions of houses in the middle two price ranges ($430,000 to $555,000 and $555,000 to $775,000) averaged f-scores of 48.8% and 48.0% respectively at best. We used f-score as an evaluation of our predictions because between precision and recall, neither was more crucial than the other. We concluded that the sheer variety of house feature combinations contributed to this unpredictable quality of our data. Our models overall were better at predicting houses in the lowest and highest price range, but surprisingly, predictive power did not vary significantly between decision tree models and neural networks, and even different tree depths or network layer counts did little to improve it.



Depth of Decision Tree vs Average F-Score of Decision Tree



Our decision trees seemed to indicate that a house's price was almost entirely accounted for by the house's size. Tree visualizations revealed that relatively clear size boundaries could be drawn between different price ranges. $12,000 to $430,000 homes had square footages up to around 1,600 square feet, while $430,000 to $555,000 homes ranged up to 2,580 square feet and $555,000 to $775,000 homes ranged up to 3,680 square feet. Acreage played an extremely minor role in determining price, and bathrooms an even smaller role. This indicates to a seller that house size is a much more important factor in price-setting than number of bedrooms or other features.

Conversely, neural networks relied heavily on acreage in predictive models, with house size and number of bathrooms playing a limited role in strengthening predictions. Summarily, between the two model types, it unsurprisingly seems like acreage and house size consistently predicted the price of houses the best. This means that buyers more concerned about other house features such as number of bedrooms or bathrooms will have an easier time finding affordable houses by searching for these features amongst houses of smaller size and acreage.



## Technical

We wanted to focus our analyses on houses in or near Utah, so we first narrowed our data down to several states in the general area. Distribution comparisons revealed that Colorado's housing data most closely matched Utah's while other states varied, so we further limited the data down to Utah and Colorado only. We dropped any houses that were missing important feature values (the feature combinations varied so widely that we felt like trying to estimate values would lead to problems), and we also removed price outliers (more than 2 standard deviations from the mean price), as we wanted to focus more on houses affordable to the average American than the grossly expensive houses that were severely right-skewing the data. Lastly, we separated houses into price ranges based on price quartiles to limit the number of classes to predict. Despite the difference in ranges of the quartiles, we felt like quartile ranges were a suitable choice because it provided even classes that did not need to be weighted in predictive models to account for class imbalance. As we should have predicted, our choice of price ranges led to the middle price ranges with narrower ranges being harder to predict, so an alternative approach we might've tried would be to separate the data into

uniform price ranges and either random sample each price range to produce a balanced number for each class or weight classes in predictive models.

We felt this dataset would do well with decision trees because we thought house feature combinations would be easy to visualize along tree paths. Trees did indeed show some interesting trends and provided some easily-visualizable separations of the data. One such observable pattern was the effect of tree depth. Our tree of depth 3 relied almost exclusively on house size (seemingly high bias), but the tree of depth 5 only relied on two additional features to only a very minor extent. Trees using only housing price performed about the same as trees with additional features, showing that the simpler trees were not overly simple or biased. Testing depths ranging from 2 to 20 showed that a depth of only 7 maximized predictive power, while depths greater than 7 showed the effects of high variance and overfitting with decreasing f-scores. Deeper trees showed splits based on house size at the highest levels, followed by acreage and then number of bathrooms. Interestingly, simpler trees showed that acreage was key in distinguishing between the highest two price ranges, but not as key for the other price ranges.

We hoped neural networks would perform better than our trees, but they performed about the same. We did not maximize the usefulness of neural networks because we had few features to feed it as inputs, so finding a housing dataset with an increased variety of features would've been more effective. Interestingly, neural networks relied much heavier on house acreage than other features. We were expecting to see a reliance on house size like the decision tree models, but varying hidden layer sizes and counts showed a consistent pattern of preference for acreage in the models. Between the two types of models, it makes sense that house size and acreage are the two biggest factors in house price, so we might assume that our models were functioning correctly.