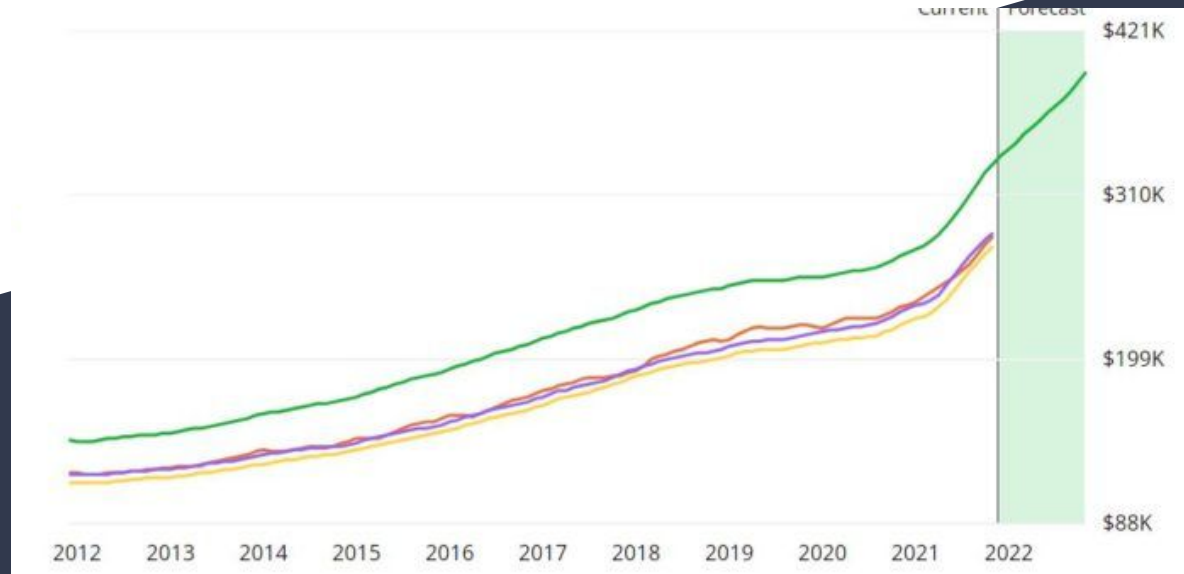# Housing Sales Price Prediction in North Dallas 2021

# Project Background

- The housing market in North Dallas has been growing rapidly over the last few years and even more so in the last 6 months.
  - Several large corporations are relocating their headquarters to the area and moving large numbers of employees with them which is driving the market upward.
- The goal of the project is to solve the problem of optimal housing pricing by using machine learning to estimate a sale price of homes in a small section of North Dallas suburbs.
- This project will use a small subset of the available data as a case study to prove effectiveness of the methodology. If the models show that a home could be sold higher this results in more money for the home sellers, a larger fee for the realty company, and a larger commission for the realtor.
- If the model shows that homes are overpriced then it can be suggested to lower the price if it is needed to move the home quickly.

# Data Wrangling



**List of Features in Raw Dataset**

- Goal of Data Wrangling
  - Remove unnecessary features
  - Investigate and fix null values
  - Each row for each feature has a relevant value
  - Each feature is the correct Dtype for analysis
- The list of features and data types of the raw data can be seen to the left
  - It is clear that there is a lot of work to be done removing features
    - Before even looking at values, features that include personal information were removed
      - Owner names
      - Mortgage information
      - Tax information
      - Contact information

# Data Wrangling

- This was a very well populated data set so all of the numerical values were populated
- The categorical values needed some work to make sure each row had a value for each feature
- Since the values of some of the categorical features were either 'Yes' or 0 it was assumed that 0 was a no, so they were filled as such
  - Has pool
  - Basement
  - Porch Type
- Other features were the opposite so were filled with 'No'
  - Heat Type
  - Air conditioning

```
cleaning['Has Pool'].fillna('None', inplace = True)
cleaning['Basement'].fillna('No Basement', inplace = True)
cleaning['Porch Type'].fillna('None', inplace = True)
cleaning['Heat Type'].fillna('Yes', inplace = True)
cleaning['Roof Material Type'].fillna('Composition Shingle', inplace = True)
cleaning['Air conditioning'].fillna('Yes', inplace = True)
cleaning = cleaning.drop(['Basement SQFT', 'Total Assessed Value', 'Assessed Land', 'Assessed Improvement', 'State', 'Status', 'Z
```

**Code Snippet Showing Filling Null Categorical Data**

| City | object |
|---|---|
| ZIP 5 | float64 |
| ZIP 4 | float64 |
| List Price | int64 |
| Township Name | object |
| Subdivision | object |
| Tax Amount | float64 |
| Total Market Value | int64 |
| Market Value Land | int64 |
| Market Value Improvement | int64 |
| Estimated Value | float64 |
| County Use Code | object |
| Lot SQFT | float64 |
| New Construction | object |
| Beds | float64 |
| Baths | float64 |
| Total Building Area | int64 |
| Living Area SQFT | float64 |
| Garage SQFT | float64 |
| Basement | object |
| Stories | object |
| Year Built | float64 |
| Air conditioning | object |
| Heat Type | object |
| Fireplace indicator | bool |
| Porch Type | object |
| Has Pool | object |
| Parking spaces | int64 |
| Sales Price | float64 |
| dtype: object | |

**Final List of Features and Type**

# Exploratory Data Analysis



**Pearson Correlation Heatmap**

- Exploratory data analysis is an important part of getting to know your data before you move on to modeling
  - Summarize and visualize important features
  - Identify patterns
  - Identify correlations
- To the left you can see the Pearson correlation heatmap
  - This is a quick way to view correlations of all of your features
- Another useful library is Pandas Profiling
  - Histograms of each feature
  - Number of null values remaining
  - Identify outliers
- After viewing histograms and outliers you can circle back to your data wrangling to fix any features that are still problematic

# Exploratory Data Analysis

- Another way to visualize the correlations is to scatter plot each of the features vs the target variable
- After looking through all of the variables it seems a little unfair to use some of them due to collinearity. The features that may be removed are Tax Amount, Total Market Value, Market Value Land, Market Value Improvement, and Estimated Value
  - In a future update to the project I would use t-SNE with different thresholds to see how many of the correlated features would be removed

# Preprocessing

- One Hot encoding variables
  - Massively expanded the number of features, 29-5893
  - Looked into some of the encoded features more and found the 4 digit ZIP codes had very few duplicates
  - These were removed to prevent over complication
- Different X/Y datasets for modeling
  - Tried different combinations of features for modeling
  - Just house features vs including location or price estimates
- Create new features for modeling
  - It was determined that the predicted sales price had the largest effect on the actual sales price
  - I tested creating new features by multiplying the predicted price by 90%-110%  to see how it would change the actual sales price
- Split into train/test subsets

# Modeling

```
start = time.time()

steps = [('scaler', StandardScaler()), ('xgbr', XGBRegressor())]
pipe = Pipeline(steps)

n_estimators = [int(x) for x in np.linspace(start = 200, stop = 2000, num = 10)]
max_depth = [1,2,3,4,5,6,7,8,9,10]
eta = [.001,.005,.01,.025,.05,.1,.2,.3]
subsample = [.25,.5,.75,1]
colsample_bytree = [.25,.5,.75,1]

random_grid = {'xgbr__n_estimators': n_estimators,
               'xgbr__max_depth': max_depth,
               'xgbr__eta': eta,
               'xgbr__subsample': subsample,
               'xgbr__colsample_bytree': colsample_bytree}

xgboost = RandomizedSearchCV(
          estimator=pipe,
          param_distributions = random_grid, n_iter = 100, cv = 5, verbose=2, random_state=42, n_jobs = -1)

grid_result = xgboost.fit(X_big_train, y_big_train)
best_params = xgboost.best_params_
print(best_params)

y_pred_xgboost = xgboost.predict(X_big_test)
y_tr_pred_xgboost = xgboost.predict(X_big_train)


print('It takes %s minutes' % ((time.time() - start)/60))
median_mae_xgboost = mean_absolute_error(y_big_train, y_tr_pred_xgboost), mean_absolute_error(y_big_test, y_pred_xgboost)
median_mae_xgboost

print(median_mae_xgboost)


r2_score(y_big_train, y_tr_pred_xgboost), r2_score(y_big_test, y_pred_xgboost)

Fitting 5 folds for each of 100 candidates, totalling 500 fits
{'xgbr__subsample': 0.5, 'xgbr__n_estimators': 1800, 'xgbr__max_depth': 2, 'xgbr__eta': 0.005, 'xgbr__colsample_bytree': 0.75}
It takes 75.25825051069259 minutes
(60930.19888094893, 64929.436901653615)
```

**Code Snippet of Best Performing Model Tuning - XGBoost**

- When I got to the modeling phase I planned to test a wide variety of model types in order to compare results
  - Support Vector Regression
  - Random Forest Regression
  - Lasso Regression
  - Ridge Regression
  - XGBoost Regression
- For each of the models I used GridSearchCV for hyperparameter optimization and cross validation in order to find the best combination of parameters
- Metrics Used
  - R2
  - Mean Average Error
  - Both metrics were used to assess the training and test datasets to see how well the model learned the training data and how well it generalized to the testing data

# Modeling Results

| Model | Train MAE ($) | Test MAE ($) |
|---|---|---|
| Support Vector Regression | 116,322 | 109,472 |
| Random Forest Regression | 67,369 | 74,960 |
| Lasso Regression | 64,532 | 66,370 |
| Ridge Regression | 64,574 | 66,138 |
| XGBoost Regression | 60,930 | 64,929 |

- It is plain to see the XGBoost model had the lowest error and would be the recommended model to use for production.
- All of the models seem to generalize well to the testing data and did not overfit on the training data.
- The five most important features used by the model can be seen here

| Feature | Importance |
|---|---|
| Estimated Value | 0.181957 |
| Total Market Value | 0.107437 |
| List Price | 0.079475 |
| Tax Amount | 0.029974 |
| Market Value Improvement | 0.023867 |

# Conclusions

- This project was a good overview of the entire data science sequence with practise in each phase
    - Data manipulation/wrangling
    - Exploratory Data Analysis
    - Preprocessing
    - Modeling
- Future work or changes to process
    - As a continuation of the process I would include more in the feature selection/extraction process
        - t-SNE for feature selection
        - Principal component analysis for feature extraction
    - For a different direction for modeling I would try a neural network to see if better results can be achieved