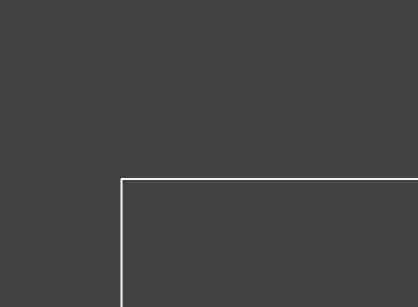


Predicting Housing Prices

By Dustin Stewart



How can we make good
real estate investments if
we can't accurately
predict the selling price of
a property?



The Problem

Problem

How can we make good real estate investments if we can't accurately predict the selling price of a property?

Objective

Create a regression model that can accurately predict the selling price of a home

**Dustin
Stewart**



About Me

Data scientist at General Assembly and have a background in digital marketing as well as finance.

Started doing data analysis back in 2017 while studying at the University of Southern California

Table of Contents

01

The Data

02

Exploring the Data

03

Creating a Model

04

Refining the Model

05

Model Insights

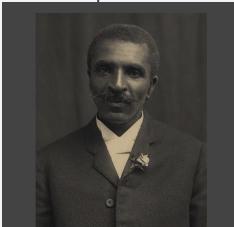
06

Conclusions

01

The Data

Data Background



70
Columns

2006-
2010

2051
Rows



Ames, Iowa

Population: 55,000

Home to Iowa State University, founded in 1858

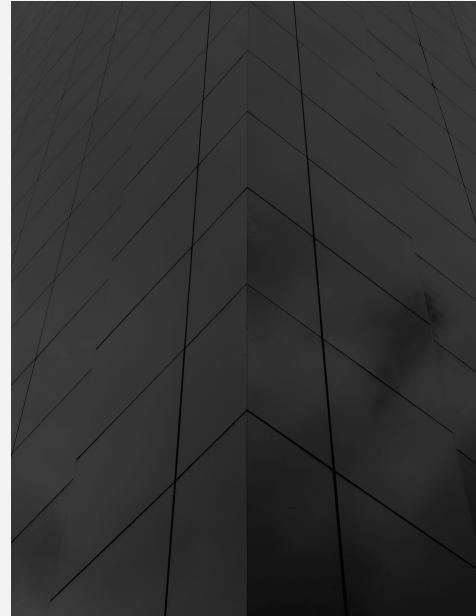
The Data

Details of over 2000 houses sold between 2006 and 2010 in Ames

The data starts during the last years of the housing bubble and continues for two years after it burst

Data Cleaning

1. Replace null values while preserving data integrity
2. Reformat categorical and boolean variables so they can be used in model as effectively as possible
3. Clean further as exploratory analysis demands

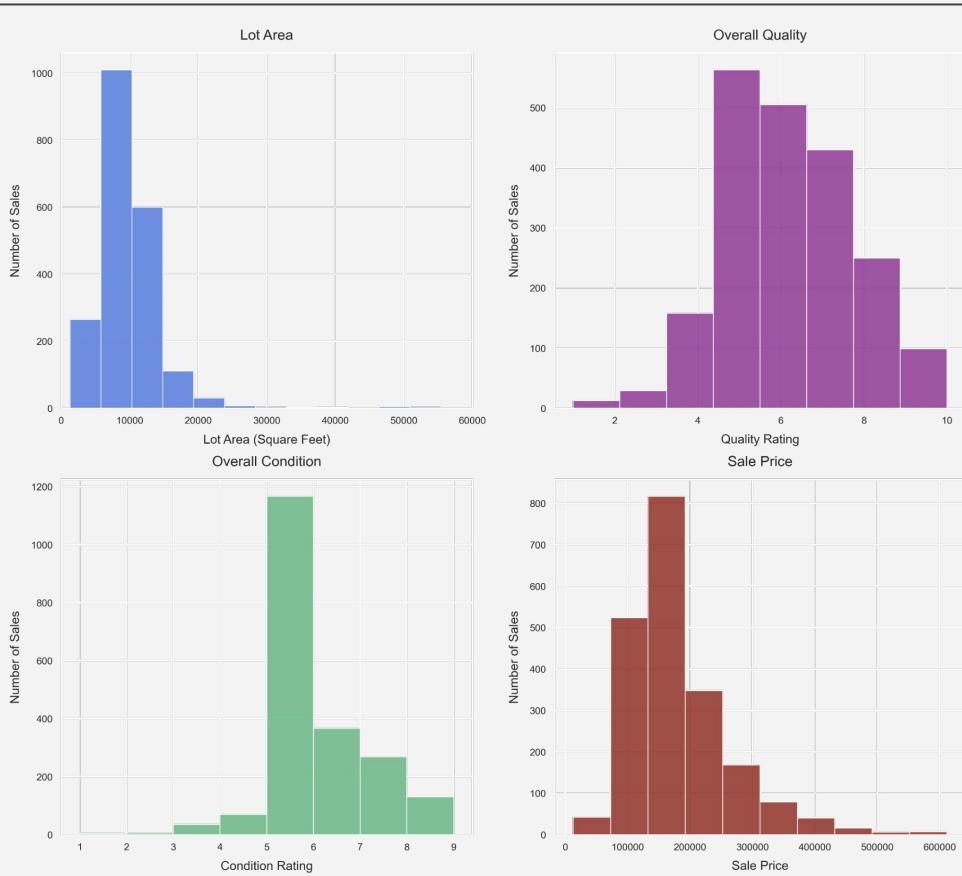




Exploring the Data

02

Methodology

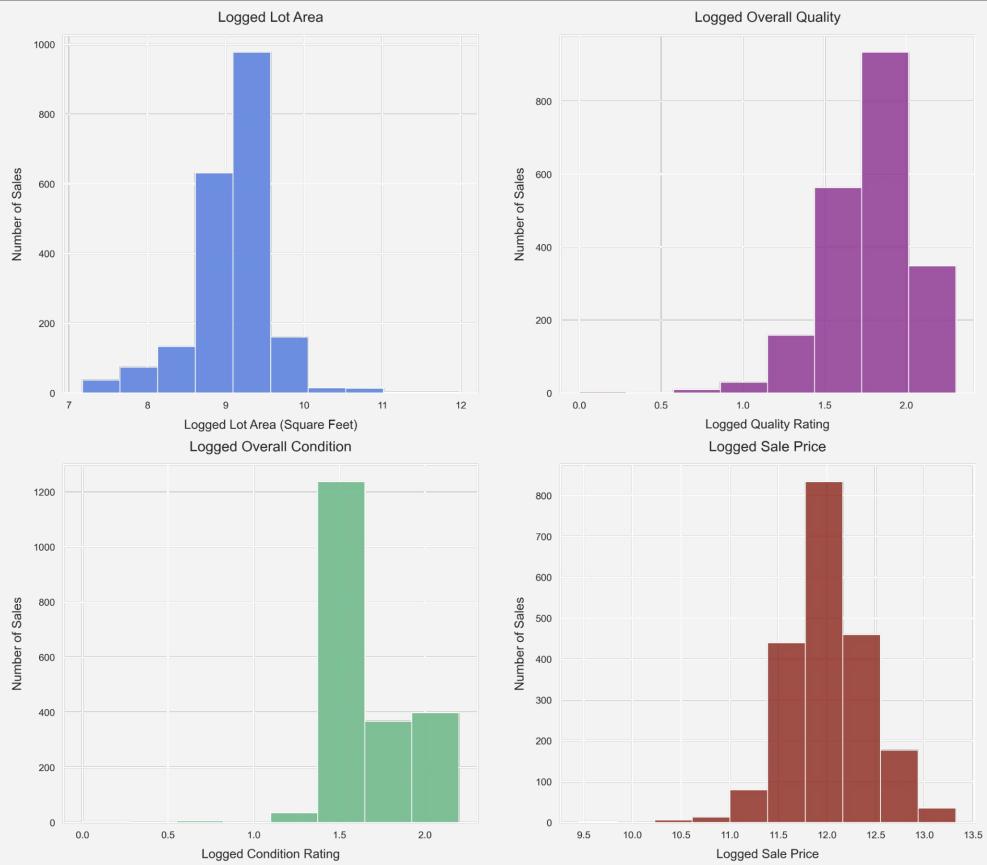


Analyzing Distributions

We see that certain distributions appear to be skewed right

A common way to make the distribution more normal is to use a logarithmic transformation

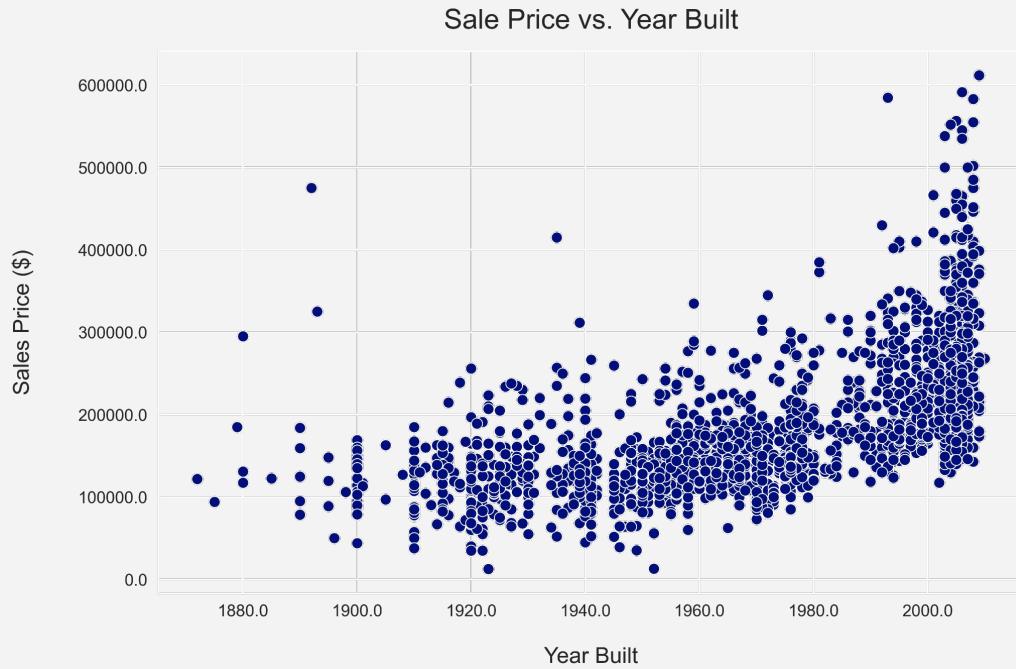
Methodology



Sales Price and Lot Area are definitely more normally distributed when logged

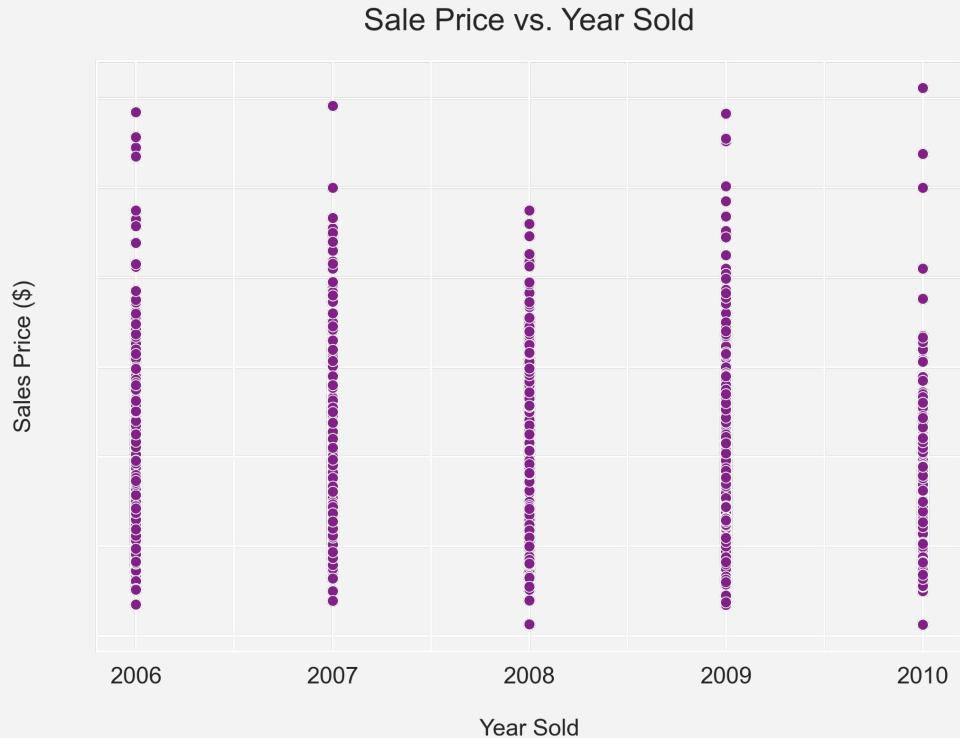
Scatter Plot Analysis

Why are there so many outliers towards the bottom of this distribution?



Scatter Plot Analysis

How did the bursting of the housing bubble affect Ames?



Methodology

High Positive Correlations

Correlations with Sale Price

	Sale Price
Overall Qual	0.8
Gr Liv Area	0.71
Kitchen Qual	0.69
Garage Area	0.65
Total Bsmt SF	0.64
nice_exterior	0.63
1st Flr SF	0.63
Year Built	0.57
Year Remod/Add	0.55
Full Bath	0.54
has_poured_concrete.foundation	0.53
TotRms AbvGrd	0.51
has_fireplace	0.49
has_attached_garage	0.47
Fireplaces	0.47
BsmtFin SF 1	0.43
newer	0.36
Open Porch SF	0.34
Wood Deck SF	0.33
good_hood	0.3
Lot Area	0.3
has_paved_drive	0.29
Bsmt Full Bath	0.28
Central Air	0.28
has_decent_garage	0.27
has_hip_roof	0.27
	Sale Price
has_brick_face	0.26
2nd Flr SF	0.25
low_density_residential	0.23
two_story	0.23
on_hill	0.21
in_culdesac	0.16
has_basement	0.15
Bedroom AbvGr	0.14
Screen Porch	0.13
floating_village	0.11
3Ssn Porch	0.049
planned_development	-0.012
post_recession_sale	-0.017
Low Qual Fin SF	-0.044
remodeled	-0.047
bad_hood	-0.056
split	-0.066
gravel_street	-0.07
has_nice_fence	-0.081
MS SubClass	-0.089
poor_functionality	-0.095
Overall Cond	-0.097
Enclosed Porch	-0.14
HasAlley	-0.14
near_artery_or_feeder	-0.18
pre_war	-0.24
	Sale Price

Feature Selection

Looking at correlations between different variables and sales to help choose which to use for the baseline model.

Continually return to these correlations to figure out why a variable's weight in the model doesn't reflect its level of correlation with sales price.

High Negative Correlations



03 Creating a Model

Model Selection

Linear Regression



R2 Score

=

0.904

This is the basic regression model

Lasso



0.905

Lasso adds a penalty that increases as the weight of the variable increases

Ridge



0.904

Ridge also adds a penalty but will never make a coefficient absolute zero like Lasso does

Baseline Score

We calculate a baseline error to compare our model to by scoring the mean sale price from our data

Error = **\$79,282**

First Model Score

We run all of our numeric variables through the model:

Training Error = **\$20,801**
R2 Score = **0.91**

Wow what an improvement!

Testing Error = **\$95,373**
R2 Score = **0.83**



Refining the Model

04

Methodology

	Features	Coefficient
0	Overall Qual	0.117359
1	Year Built	0.079079
2	BsmtFin SF 1	0.035522
3	1st Flr SF	0.105314
4	2nd Flr SF	0.085343
5	Kitchen Qual	0.031546
6	remodeled	0.002605
7	post_recession_sale	-0.002916
8	HasAlley	-0.003548
9	has_brick_face	0.003637
10	has_decent_garage	0.021931
11	floating_village	0.011741
12	low_density_residential	0.008183
13	gravel_street	-0.008023
14	on_hill	0.004157
15	in_culdesac	0.001034
16	near_artery_or_feeder	-0.016018
17	has_hip_roof	0.006775
18	nice_exterior	0.000684
19	has_poured_concrete.foundation	0.010510
20	poor_functionality	-0.020989
21	has_fireplace	0.023629
22	has_porch_or_deck	0.006814
23	log_overall_cond	0.055932
24	log_lot_area	0.053222
25	Neighborhood_Blueste	-0.000000
26	Neighborhood_BrDale	-0.003984

Recursive Feature Selection & Creation

By:

- A. Calculating the significance of each variable's effect on the model
- B. Comparing their effect on the model to correlations
- C. Using a bit of trial and error

We can find the variables that best predict sale price

Final Model Criteria

Independent Variables

- Overall Quality
- Overall Condition (log)
- Lot Area (log)
- Year Built
- Finished Basement SF
- 1st Floor SF
- 2nd Floor SF
- Kitchen Quality
- Remodeled
- Sold Post Recession
- Has an Alley
- Has Brick Face
- Has a Decent Garage
- On a Gravel Street
- On a Hill
- In a Cul-de-Sac
- Near an Artery or Feeder Street
- Has a 'Hip' Roof
- Has a Nice Exterior
- Has a Poured Concrete Foundation
- Has Poor Functionality
- Has a Fireplace
- Has a Porch or Deck
- In the Floating Village
- In a Low Density Residential Area
- Neighborhoods





05 Model Insights

R2 Scores

Training = **0.89**

Test = **0.90**

RMSE Scores

Training Error = **\$36,301**

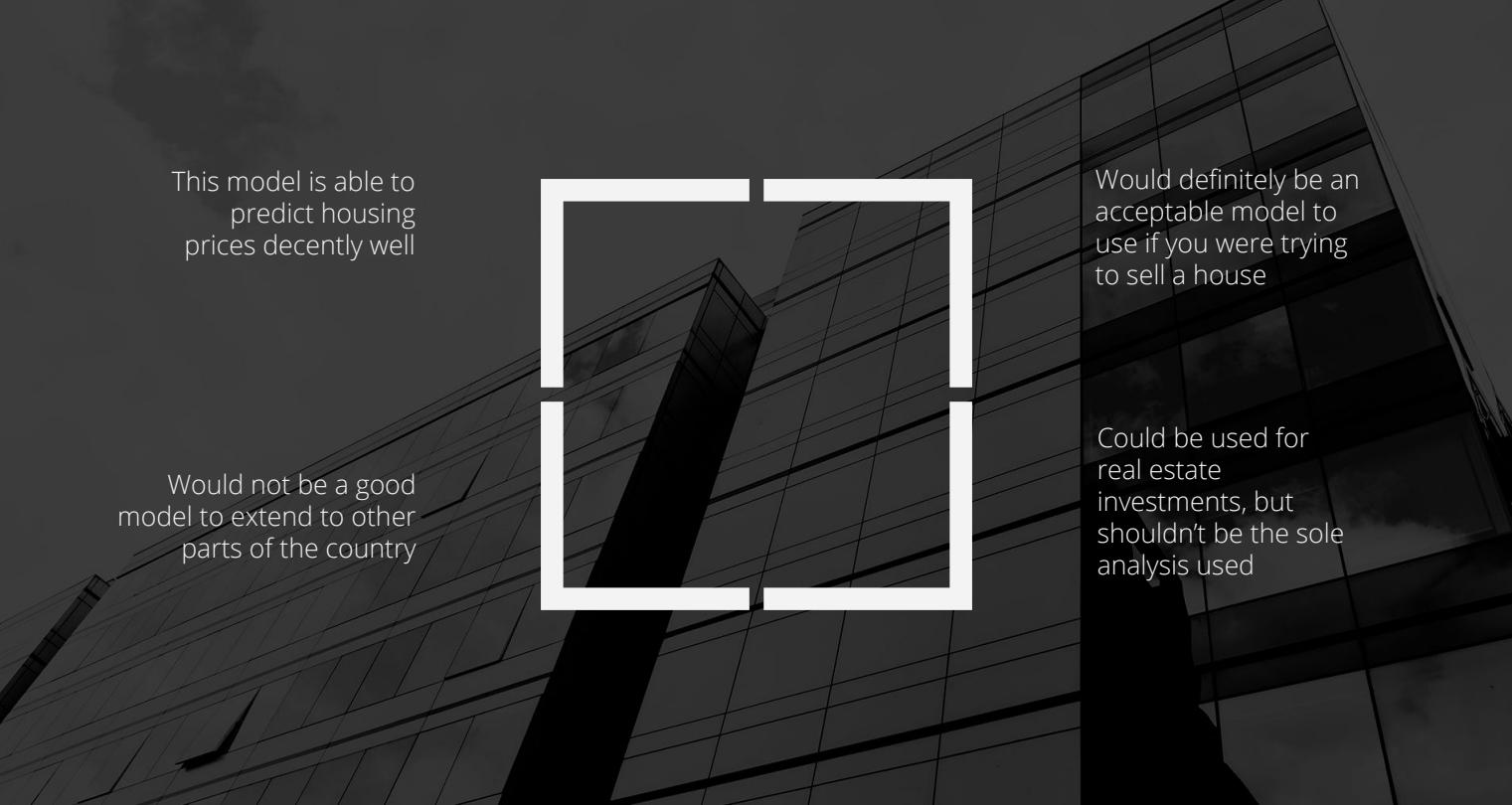
Testing Error = **\$20,691**



06

Conclusions

Conclusions

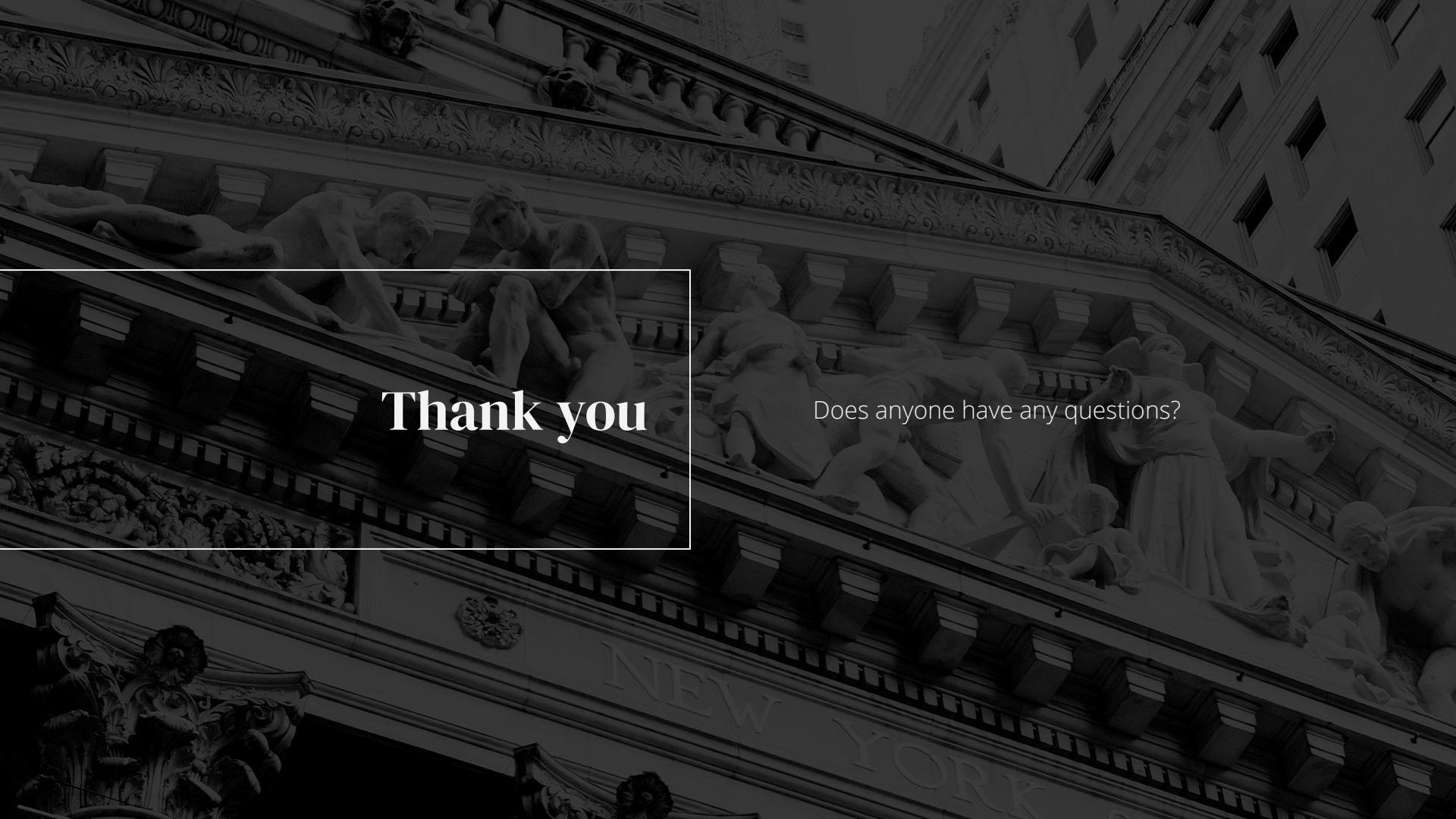


This model is able to predict housing prices decently well

Would not be a good model to extend to other parts of the country

Would definitely be an acceptable model to use if you were trying to sell a house

Could be used for real estate investments, but shouldn't be the sole analysis used



Thank you

Does anyone have any questions?

