

Identifying Optimal Locations for New Home Construction

Matt Cannon, David Kaplan, Dan Quinn, and Roberto Reis

Project 01
July, 22 - 2024

GitHub Repository: https://github.com/dfquinn23/Project_1



Outline

Goals

Data Sources

Census Data, BEA, Zillow, FEMA

Data Integration

Main Outcomes

Business Application

Lessons Learned and Improvements

Project Overview

Home builders face the critical decision of selecting the best locations for new construction projects to maximize profitability and meet market demand.

This project aims to identify and evaluate the most suitable areas for building new homes by analyzing various factors such as market affordability and related fundamentals, demographic trends, and environmental risks.

This analysis will provide actionable insights to guide home builders in making informed decisions about where to build.

Objectives

1. Determine the top locations for new home construction based on comprehensive data analysis.
2. Analyze the impact of key factors such as land cost, market demand, infrastructure, demographics, and environmental risks on the suitability of a location.
3. Provide a comparative analysis of different regions to highlight their strengths and weaknesses for new home construction.
4. Offer actionable recommendations to home builders on the most promising areas for development.

Data Sources

01



American Community Survey (ACS) is the annual survey the Census Bureau conducts in between the decennial Census.

It is a survey that provides information across more than 20 datasets on an ongoing basis.

For our final analysis, we employed county-level data to provide granular data to assess promising areas for development.

In our initial research, we also reviewed metro, state, and national-level data.

ACS datasets we employed in our initial screens:

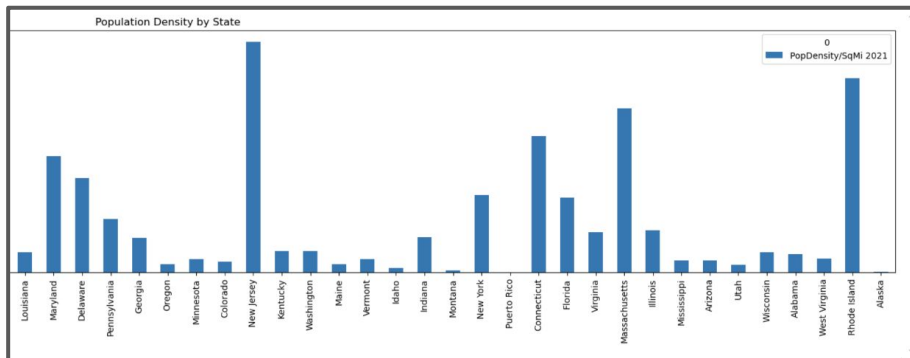
- Total population growth
- Housing Affordability
- Poverty Rates

There were dozens of other datasets that we could have used to further refine our analysis, but these three provided our team with a strong foundation from which to proceed.

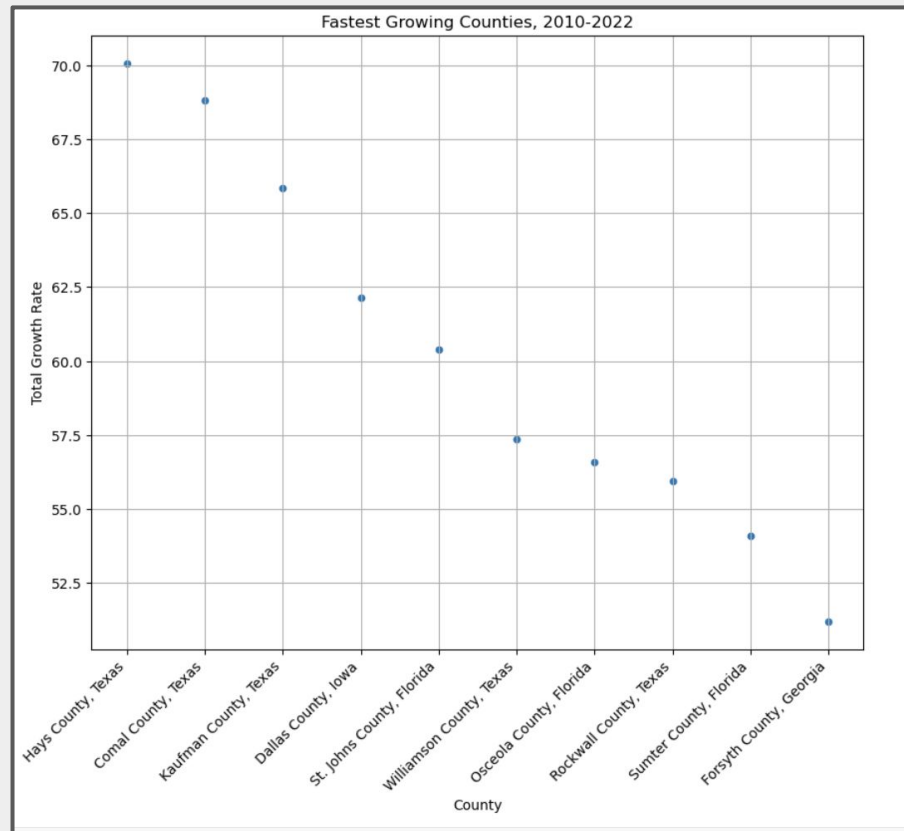
Had to ensure data frame column headings and data types matched across the various tables.

For example, one of our initial screens related to state population density in 2021.

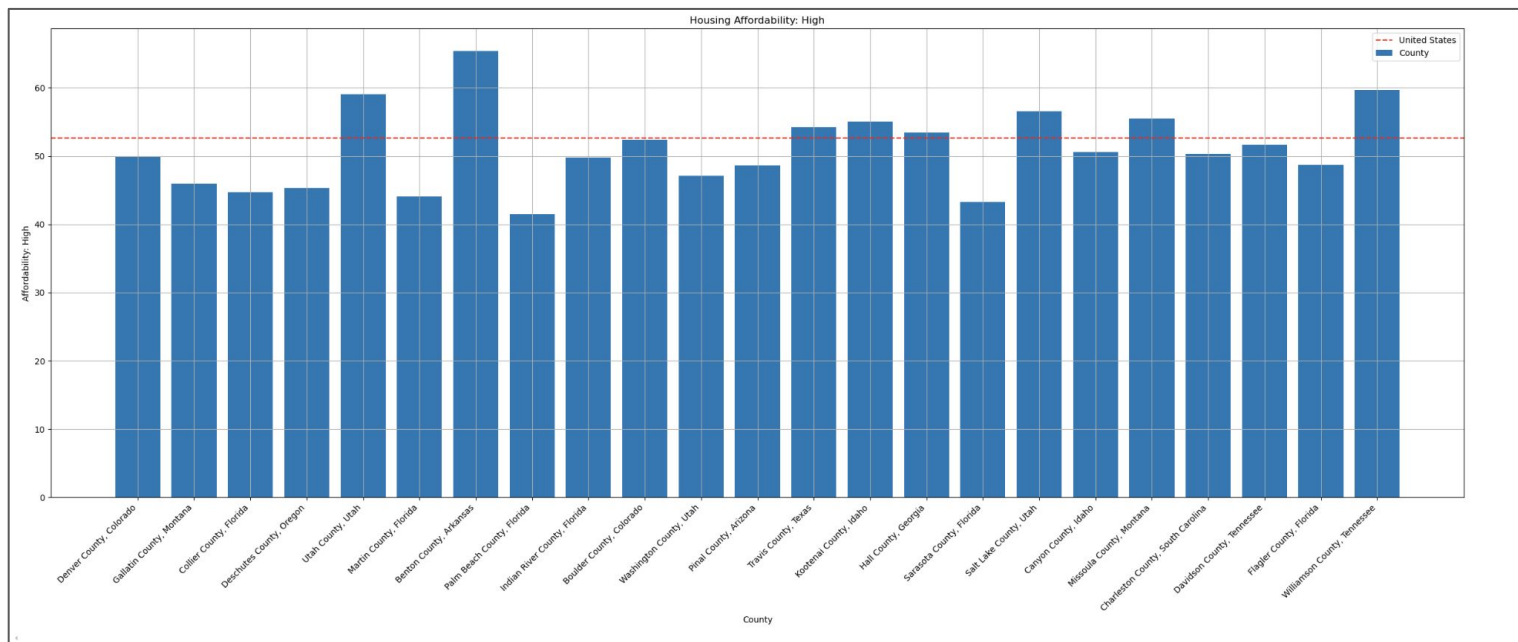
As you can see in the cropped image below, we were able to easily ascertain which states had the highest population density per square mile:



In another, we identified the fastest growing counties in the US, from 2010-2022:



Finally, these are a bit more interesting - housing affordability graphs (*High Affordability shown here - also have Mid- and Low*). These were based off of the percentage of household income spent on owner-occupied housing in our selected counties (2010-2022):

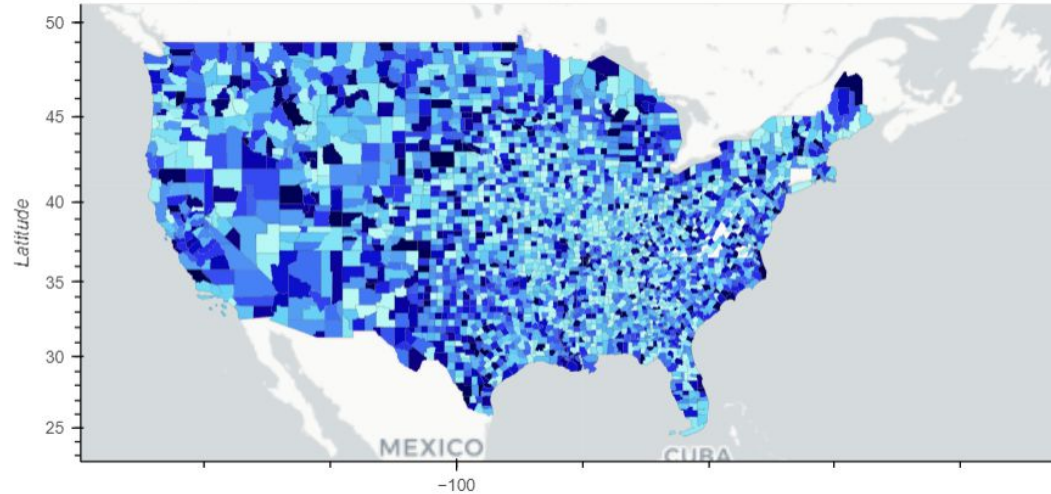


Data Sources 02

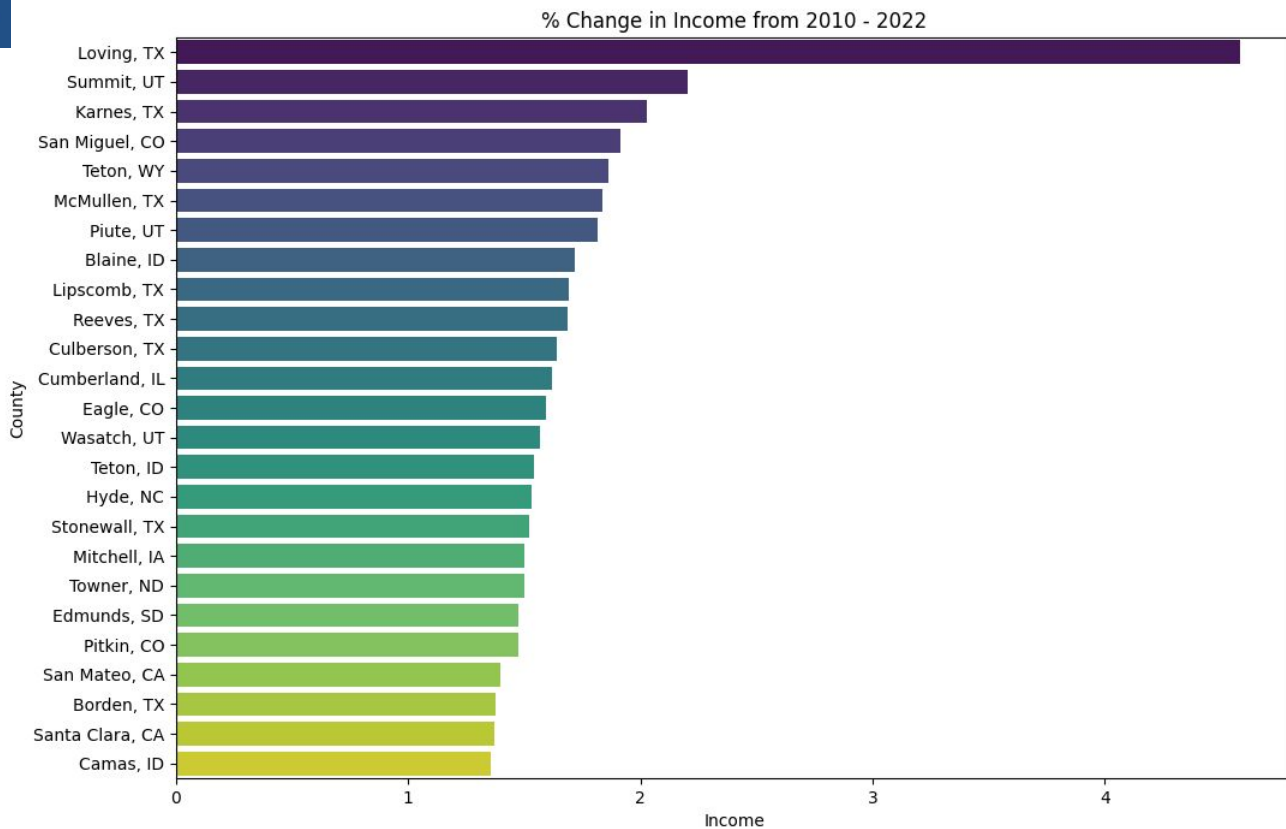
Bureau of Economic Analysis (BEA)



Percentage Increase in Income by County 2010 - 2022



- **assess** the nations income growth by county
- **identify** counties with positive income growth
- **help** showcase counties with higher income growth

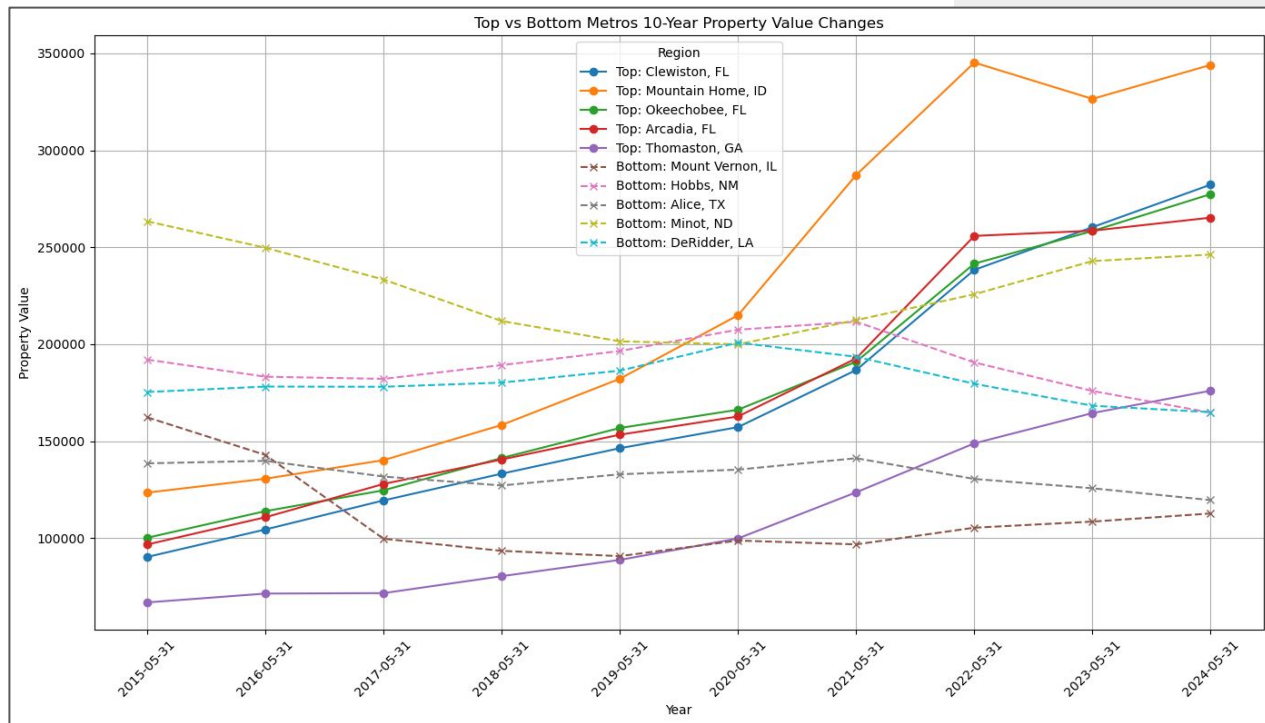




Zillow Home Values

Zillow Home Value Index (ZHVI):

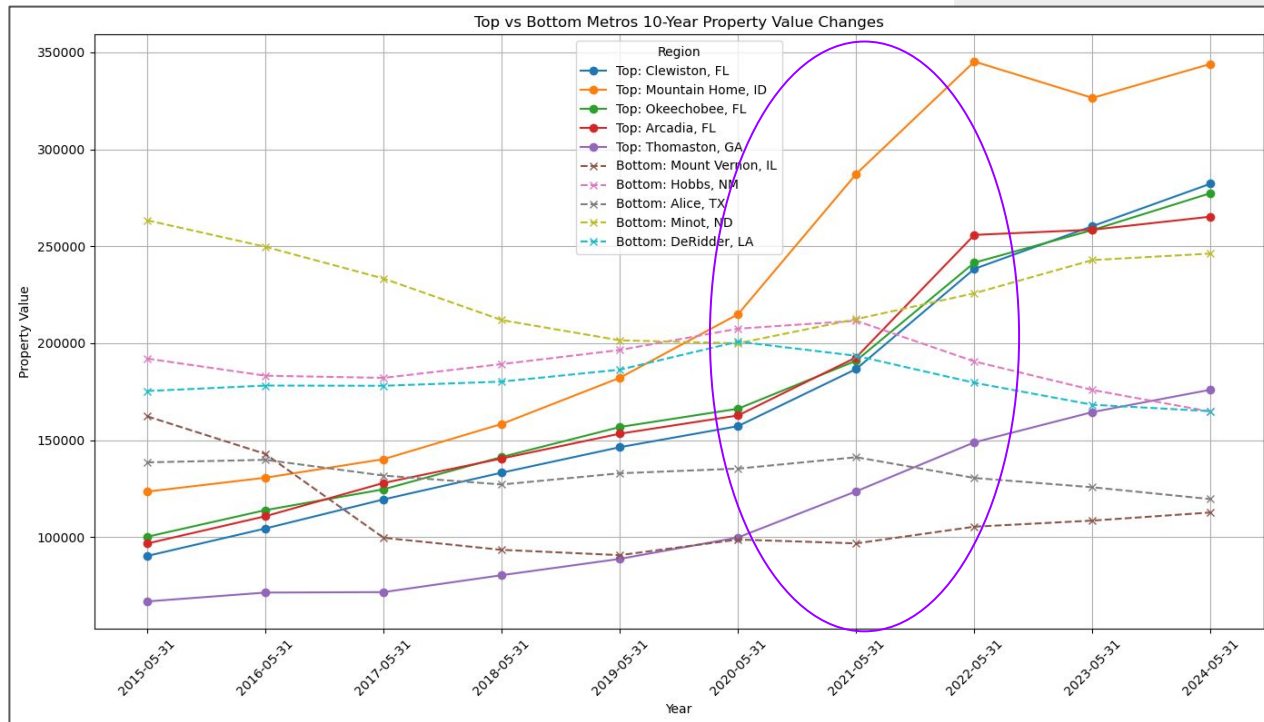
- A measure of typical home value and market changes for a region
- Reflects the typical value for homes in the 35th to 65th percentile range
- The data is made available as a smoothed, seasonally adjusted measure





Zillow Home Values

The Covid Bump
03/2020 - 06/2021



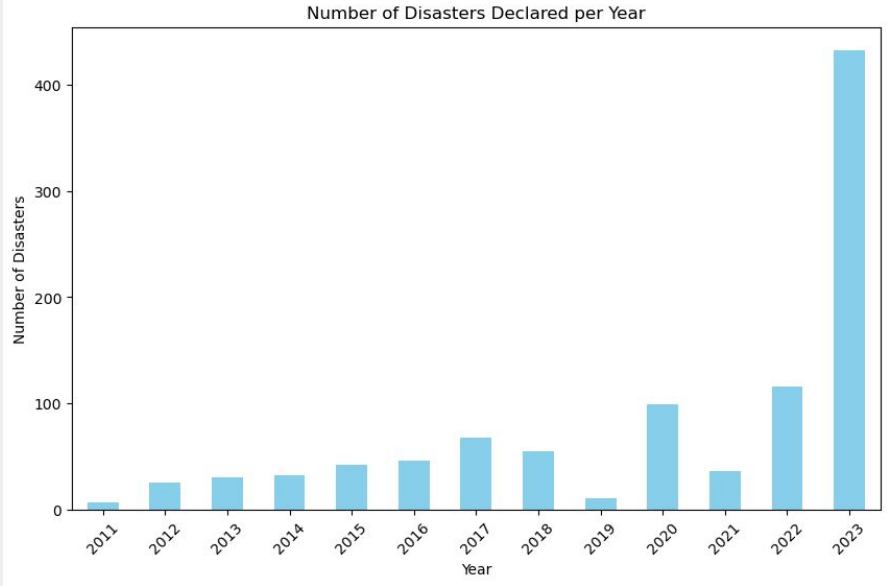
Data Sources

04

Federal Emergency Management Agency (FEMA)



FEMA



- To **assess** an area's risk profile
- To **identify** locations where lower disaster risk correlates with stable or increasing property values
- To **help** balance the trade-off between safety and profitability

Data Integration

Top 23 counties

- Combined BEA, Census, and Zillow data
- Sorted data sets to get top counties
- sliced data sets for 10 percent of counties
- Merged data sets to find county overlap

Data Cleaning

BEA

	GeoFips	GeoName	Income Change
0	48301	Loving, TX	4.581468
1	49043	Summit, UT	2.204708

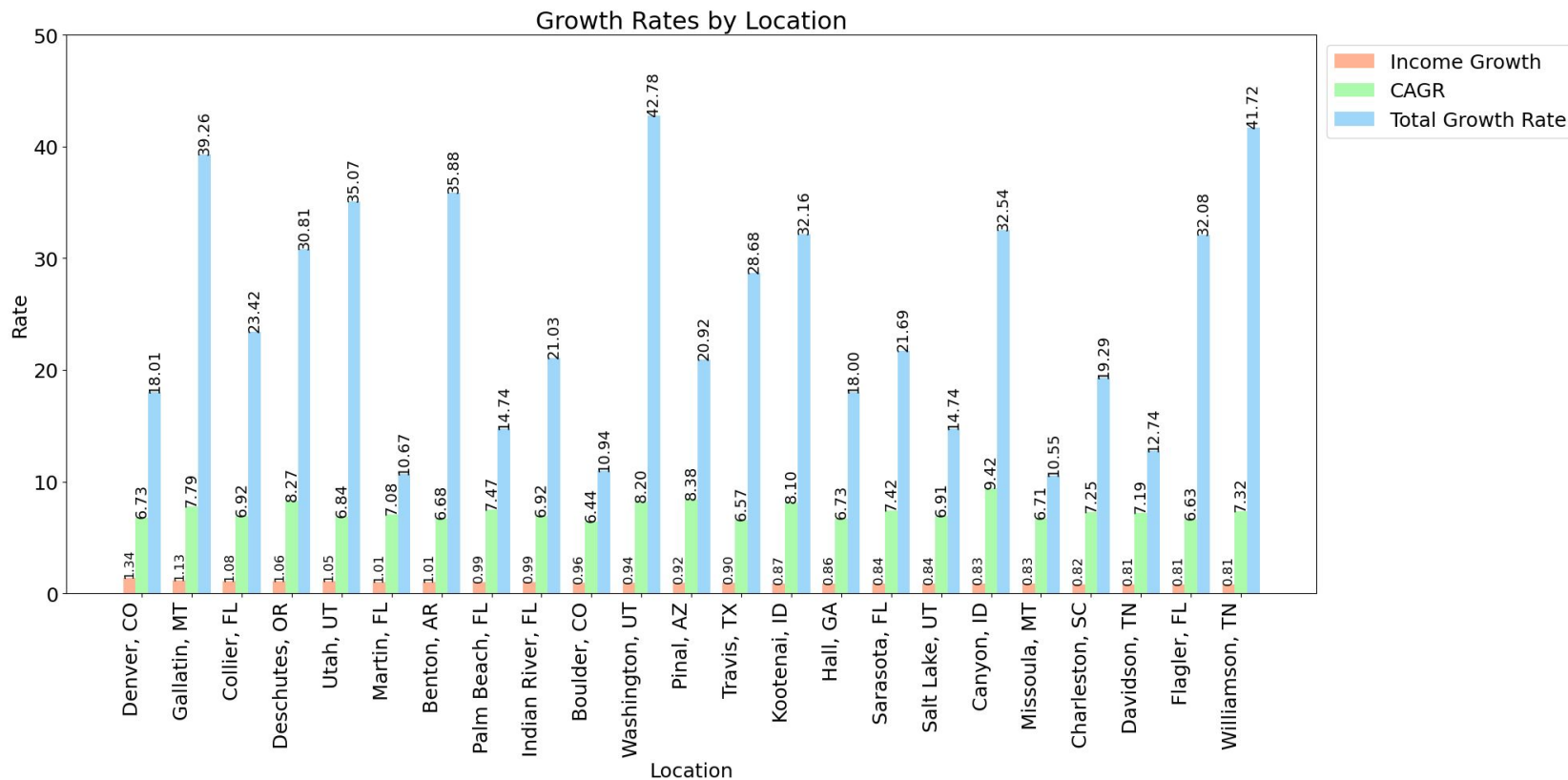
Zillow

	County_State	CAGR
0	Petersburg City, VA	11.300335
1	Clayton County, GA	10.520823

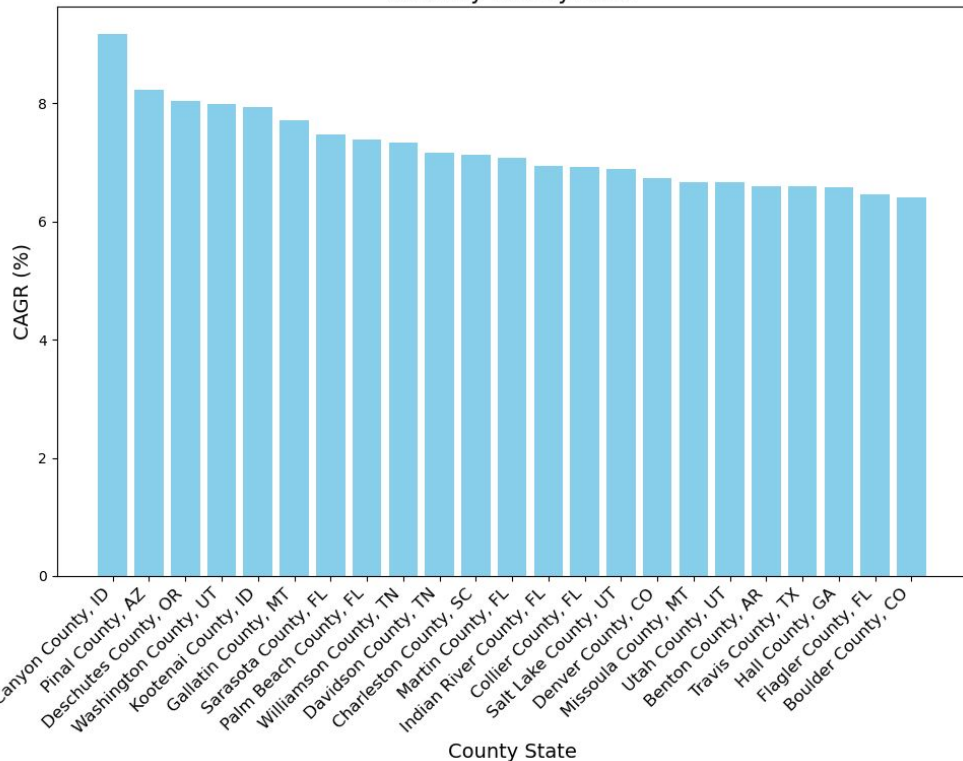
Census

	County	State Code	County Code	Total Growth Rate
0	Tuscaloosa County, Alabama	1	125	0.214032
1	St. Clair County, Alabama	1	115	0.121429

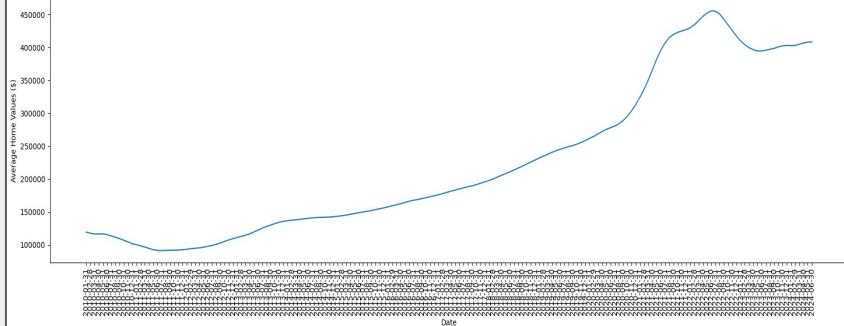
Output



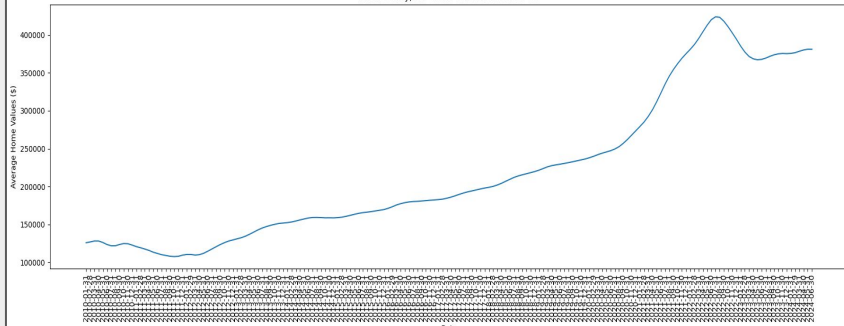
CAGR by County State



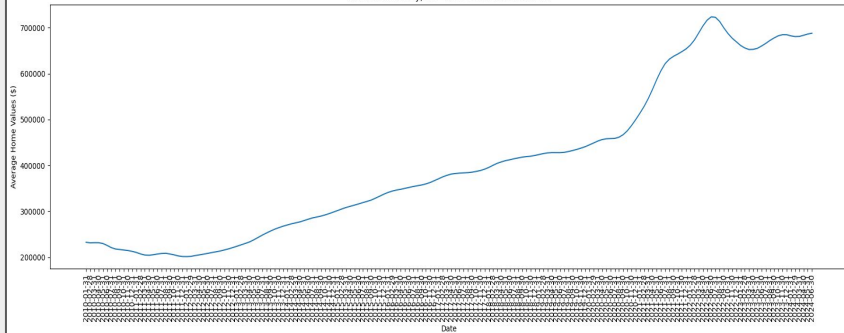
Canyon County, ID - 2010-01-31 to 2024-06-30

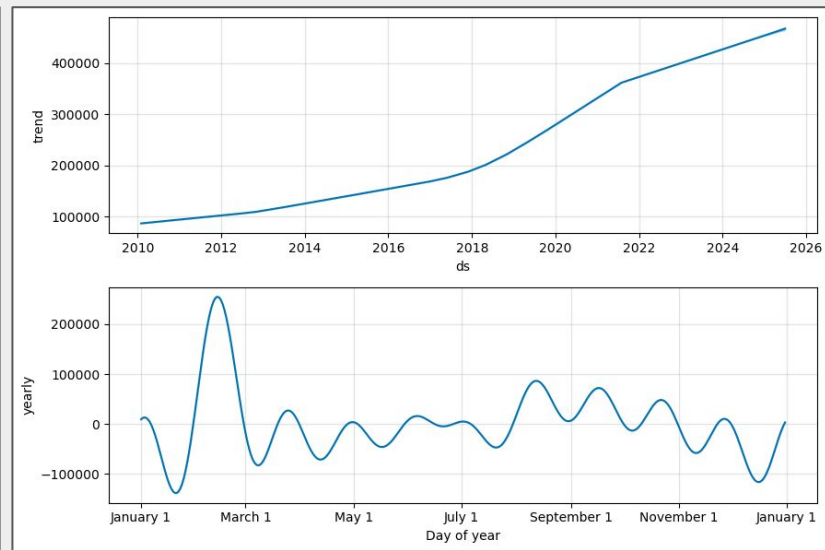
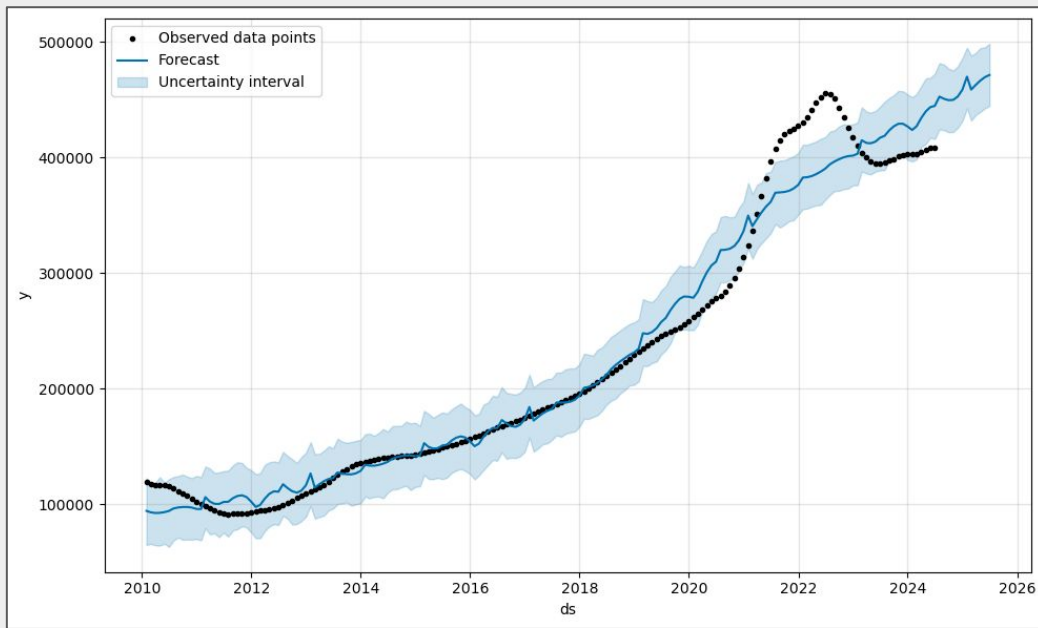


Pinal County, AZ - 2010-01-31 to 2024-06-30



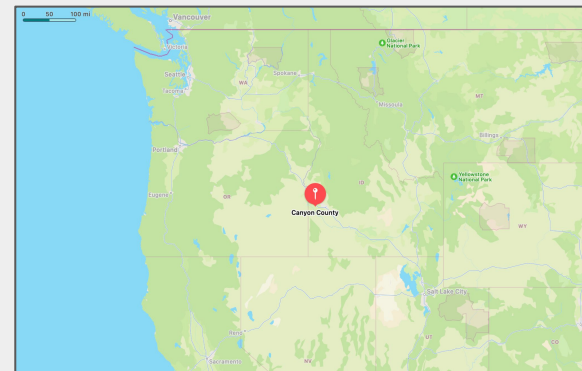
Deschutes County, OR - 2010-01-31 to 2024-06-30

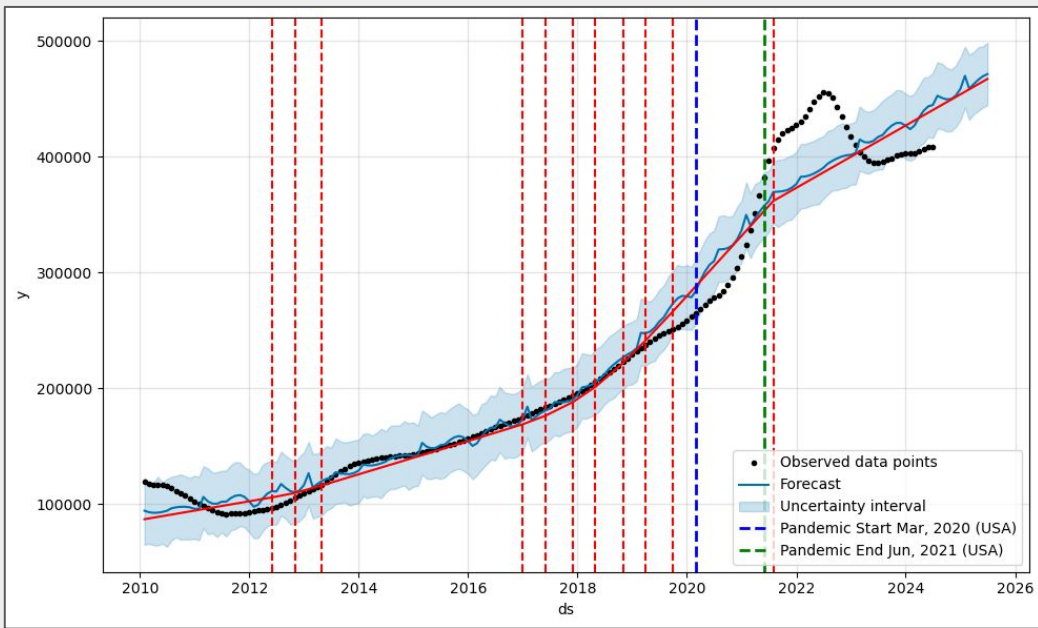




Characteristics for Canyon County, ID:

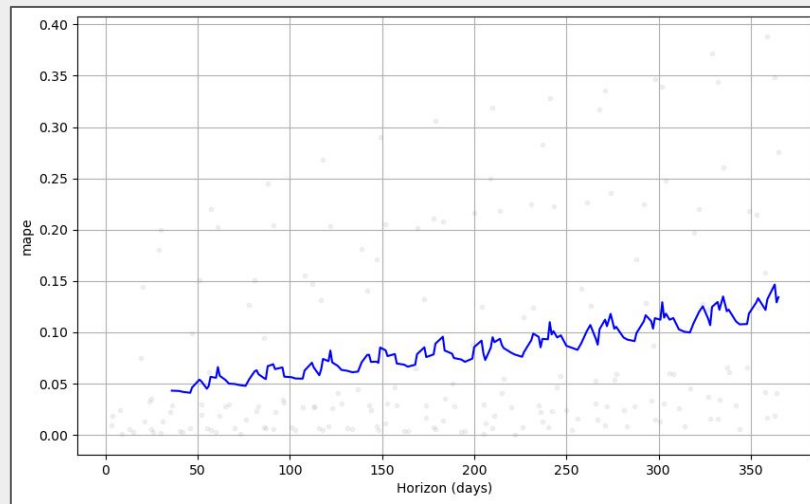
- Experienced a Covid bump
- Forecast is up and to the right
- Yearly market changes are more volatile in Q1, and then trend down from Aug/Sept through the end of the year





Trend Changepoints:

- Identify changes in the trend's trajectory
- Picks up the end of Covid
- Only inferred for first 80% of the time series



Prediction Accuracy

Mean Absolute Percentage Error (MAPE)

- The average deviation between the forecasted value and the actual value.
- This forecast typically errors around 5% for the first 100 days and then errors increase up to around 15% for predictions that are a year out.

Summary

- There are >3,200 counties in the US
- Based on the parameters used to merge our Census, Zillow, and BEA data, we were able to narrow a select client list of 23 counties
- Within the selected list, the environmental risk is negligible

Business Application

The project has to be saleable to a realistic client

Base assumption - 2-week project

Let's assume our "company" does 20 2-week projects per year:

- $[(\text{Salaries } (\$150\text{K} \times 5) / 20) + \$1,000]$ (\$20K in company overhead annually, divided by 20 – we run a tight ship.)
- Add in a profit margin of 25%
- We need to do a cost-benefit analysis for the client:
 - What problem are we solving?
 - How does our solution enable more precise decision-making for our client?
 - Can we justify our total fee against what the ultimate value is for the client?

Total Project Fee: \$50,000

Job Market Analysis

linkedin.com
Glassdoor.com
ChatGPT.com

7,999 Data Scientist jobs in the US

Requirements:

- Educational Background: Typically requires a Bachelor's degree in data science, statistics, computer science.
- Technical Skills: Proficiency in programming languages such as Python or R. Familiarity with data manipulation tools and libraries like Pandas and NumPy, and experience with data visualization tools.
- Specialized Knowledge: Understanding of time series analysis, experience with machine learning models for forecasting.

Experience:

- Experience with large datasets and real-time data processing.
- Background in statistical modeling and regression analysis.
- Prior work with cross-validation and predictive modeling techniques.

Expectations:

- Ability to translate complex datasets into actionable insights.
- Strong analytical skills to monitor and forecast market trends.
- Capable of developing models that predict real estate values and market dynamics.

Salary:

According to Glassdoor, the average base pay for data scientists in the U.S. is around \$113,000 to \$160,000 a year. Senior data scientists might earn between \$160,000 and \$200,000 annually.

Job Outlook:

Employment of data scientists is expected to grow significantly due to the high demand for data-driven decision-making and the advancement in AI and machine learning technologies.

Lessons Learned and Improvements

Lessons Learned:

1. Getting sync'd up earlier on data integration
2. Picking the right data from a mountain of data

Improvements:

1. Fail faster and better expectations of data sifting
2. Clear integration plans for the data

Questions?