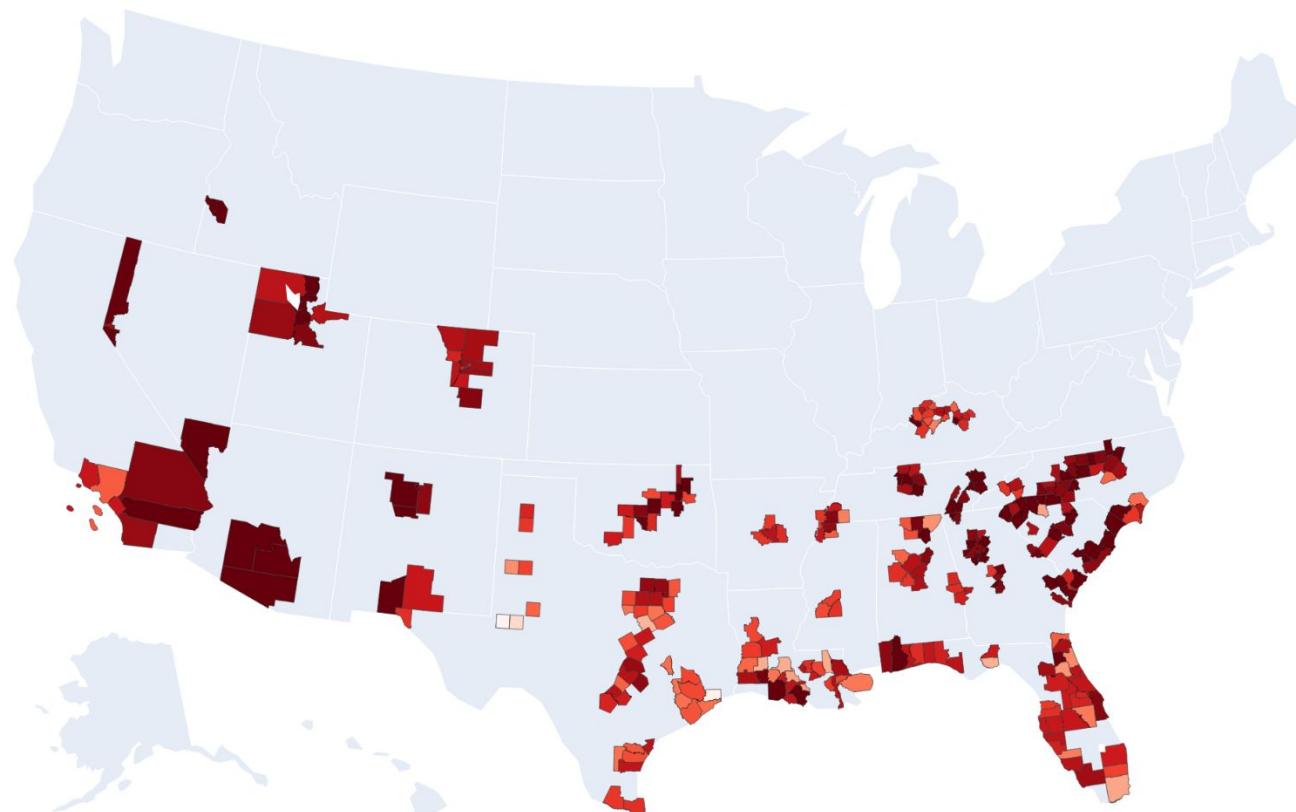


Predicting Real Estate Hot Spots (DSO 597)

Nattawut Kananusorn, Yi-Ching Lin, Chinka Huang, Siheng Rong, Thomas Nguyen, and Keegan O'Neill



Introduction

Goal: Identify the 10 areas (US Census Tract) in the country with the highest expected rent growth in the next three years (2025).

Methodology: Use property level rent data in conjunction with additionally gathered data to train models that predicts the change in rent for the entire country from 2022 to 2025.

Machine Learning Models Used:

- **Traditional ML Models:** Linear Regression, Decision Tree, Random Forest, XGBoost, LightGBM, and Long Short Term Memory.
- **Time-Series Specific Models:** Pooled OLS, Fixed Effects and Random Effects.

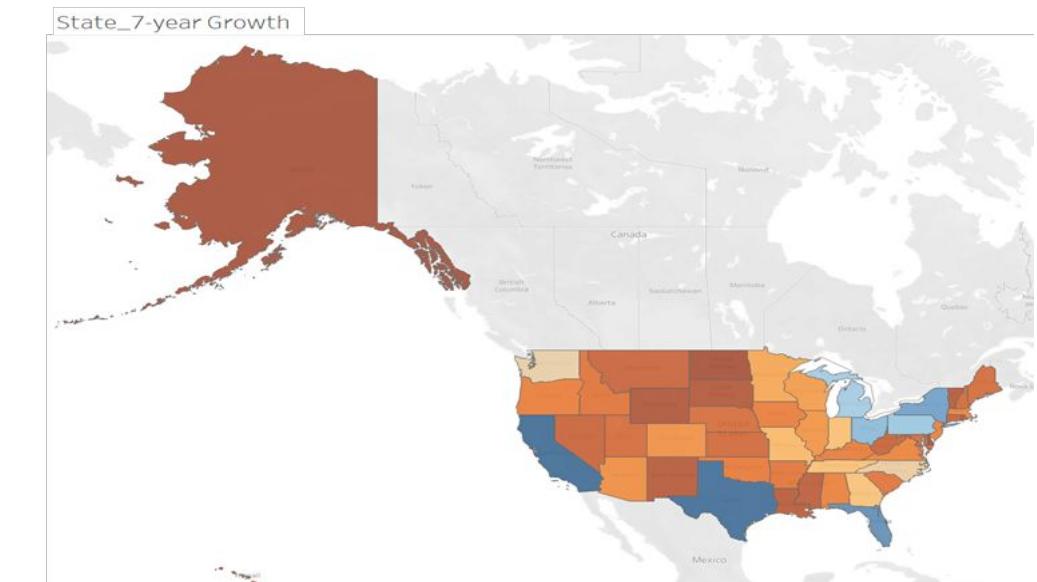
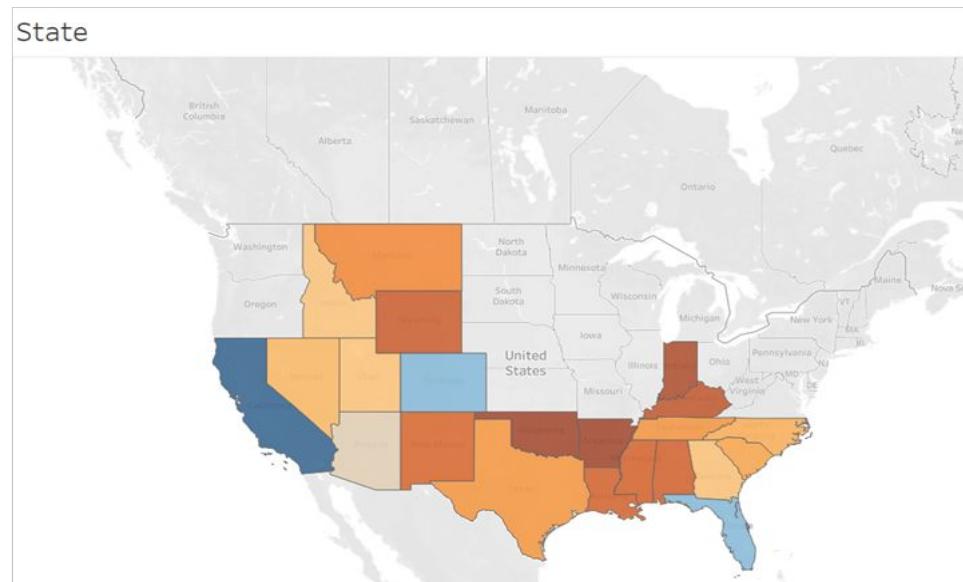
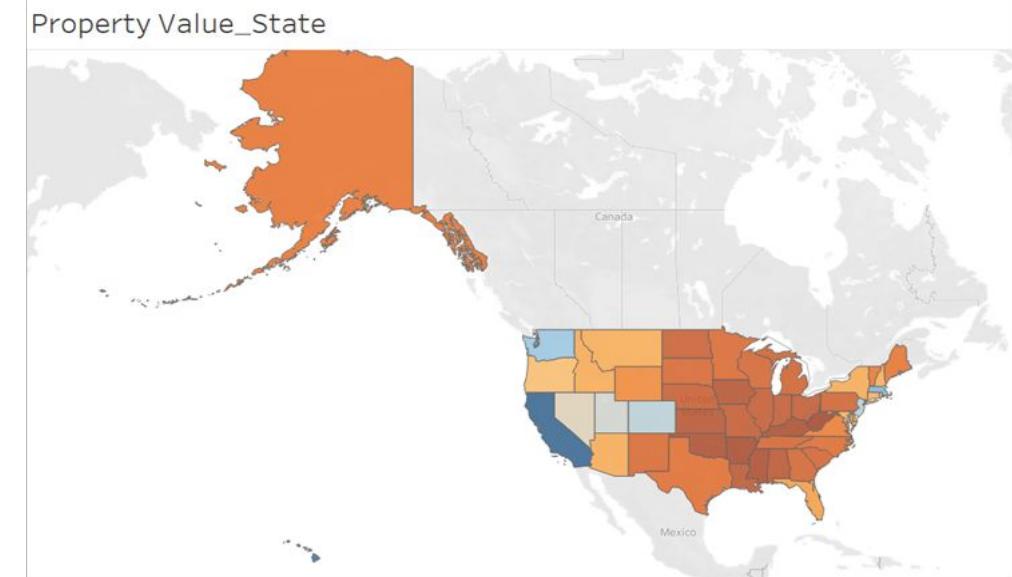
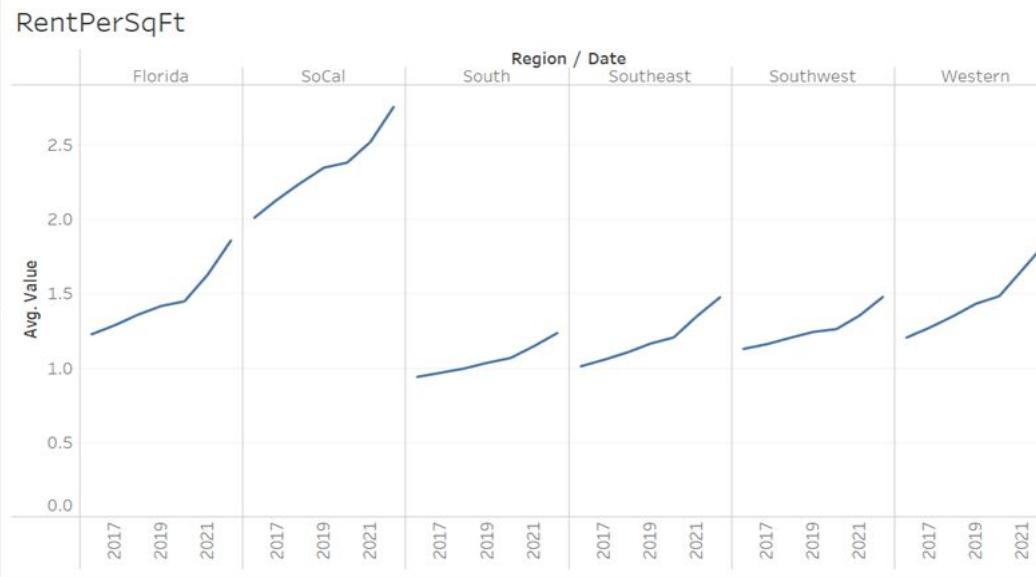
Overview of Methodology



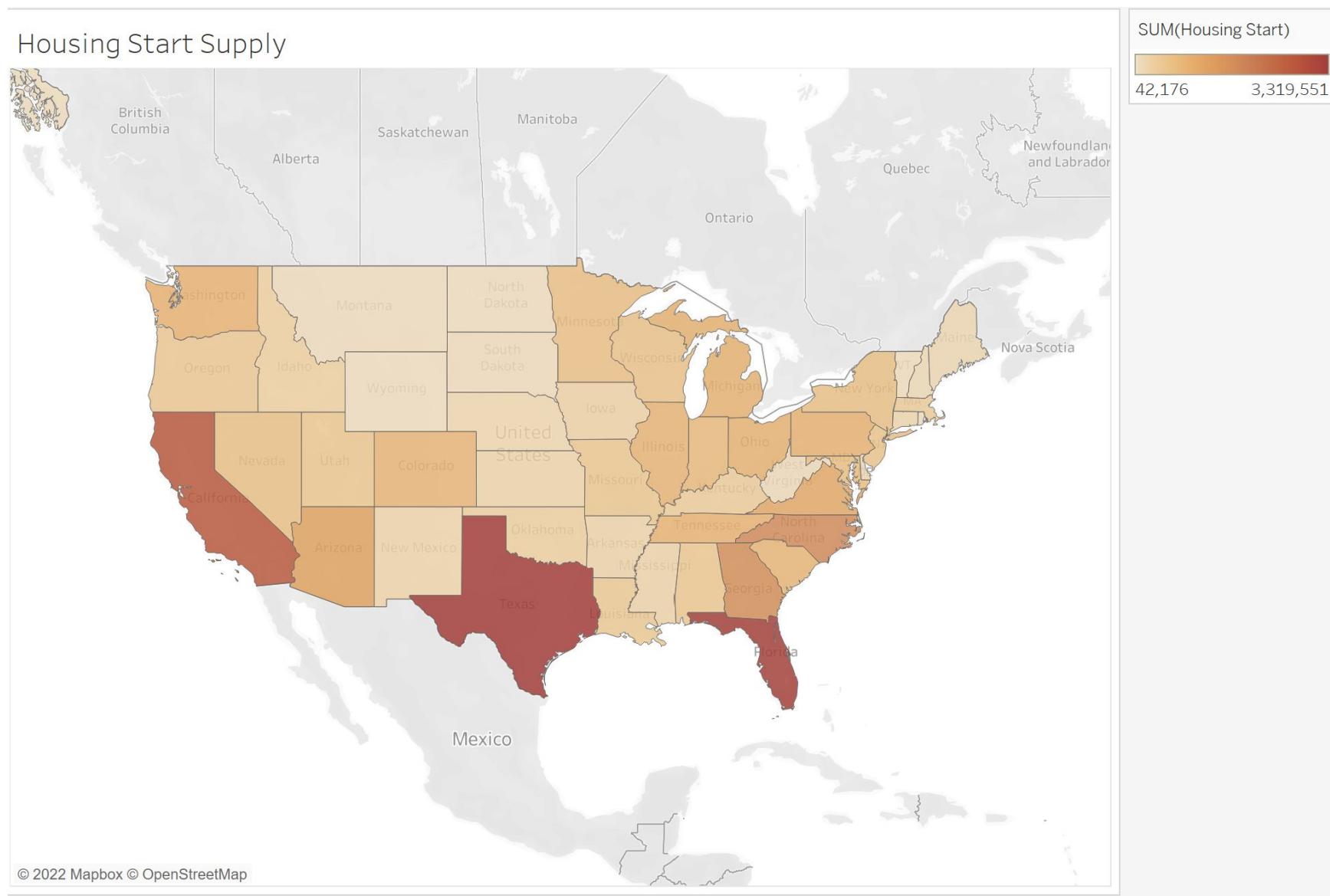
Factors Affecting Rent Growth

- **Housing data**
 - Occupancy Rate
 - Housing Supply
 - Property Values
- **Economics data**
 - Tax Rates (Sales, Income, Corporate, Combined)
 - Debt
 - National Indicators (Bond, SP 500, Inflation Rate)
- **Location data**
 - Starbucks, Banks, Public School, Private Schools, Hospitals (Normalized by Location Density)
 - Zip Code Area
- **Demographic data**
 - Population (gender, age, race)
 - Median Household Income
 - Employment Rate (16+ employed population / 16+ labor force)
- **Crime data**
 - Crime Rate

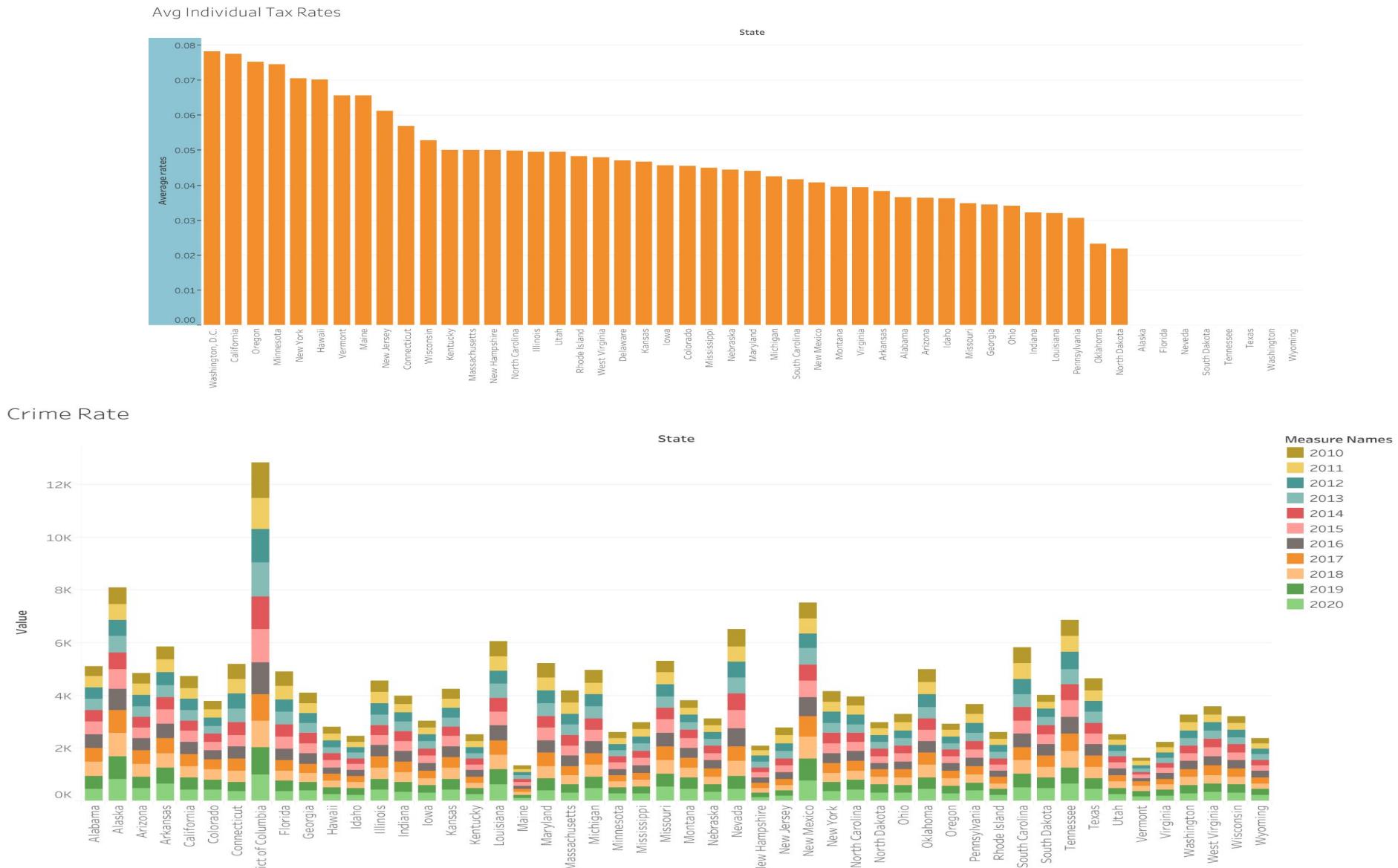
EDA: Rent and Property Value



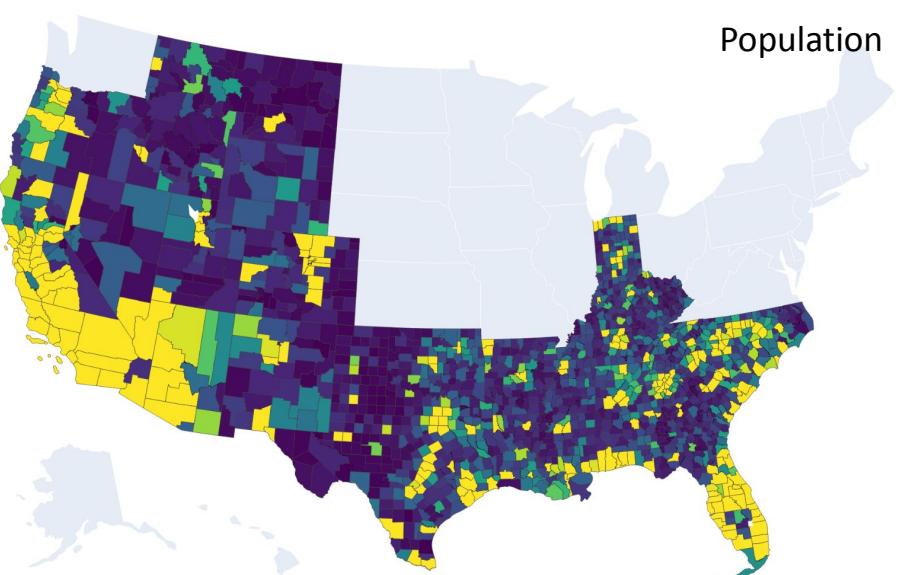
EDA: Housing Supply



EDA: Economics Factors



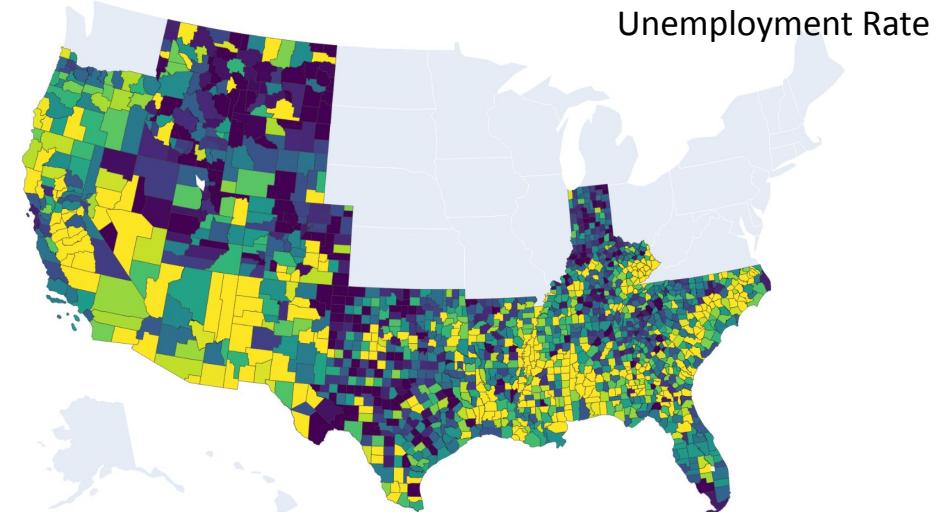
EDA: Demographic Factors



Population

Population

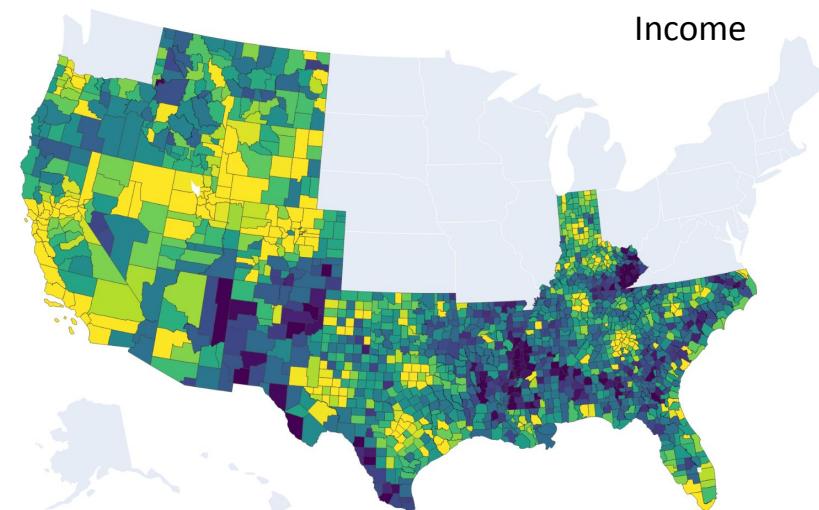
140k
120k
100k
80k
60k
40k
20k
0k



Unemployment Rate

Unemployment Rate

0.08
0.07
0.06
0.05
0.04
0.03



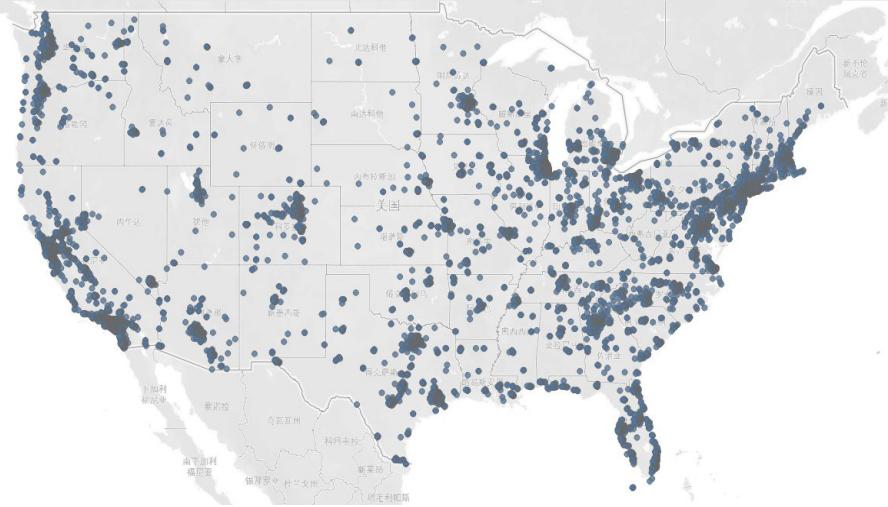
Income

Median Household Income

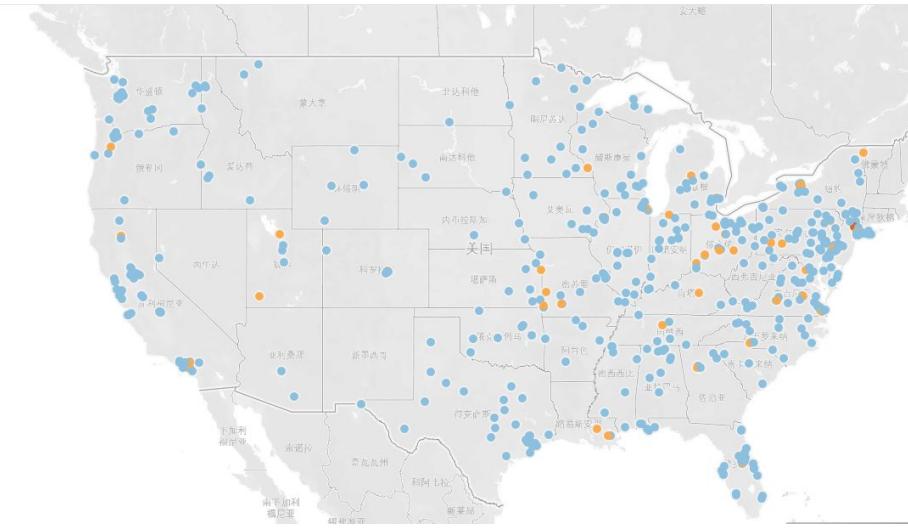
\$60k
\$55k
\$50k
\$45k
\$40k
\$35k
\$30k
\$25k
\$20k
\$15k
\$10k
\$5k
\$0k

EDA: Location Factors

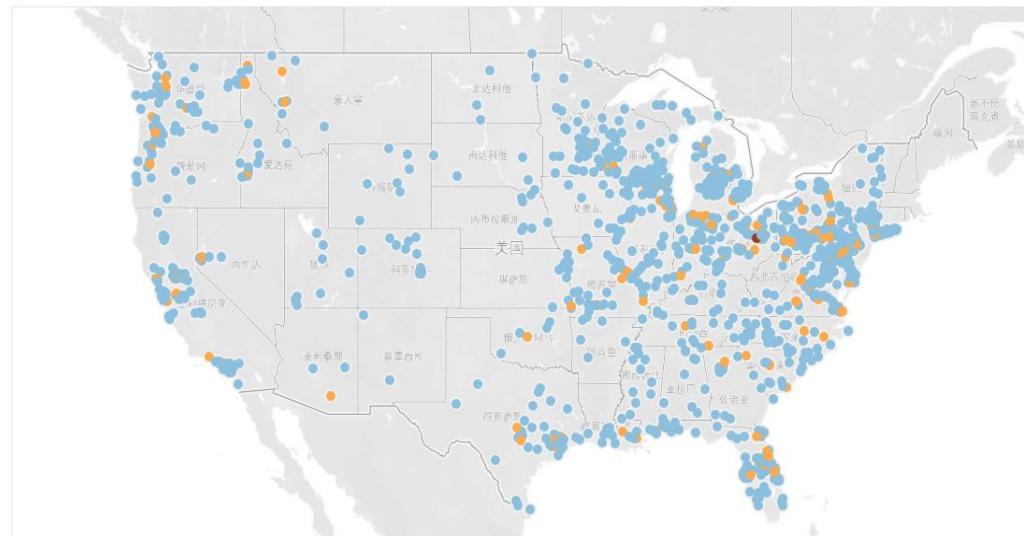
starbucks



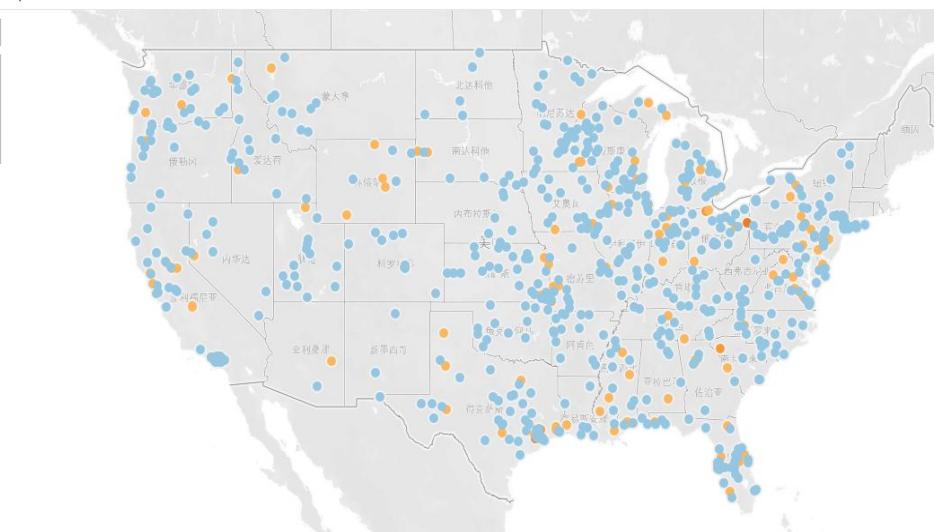
college



private school



hospital



- Starbucks stores
- Public schools
- Private schools
- Colleges
- Hospitals
- Banks
- Grocery visits
- Manufacturing facilities

Data Cleaning

- Map Property Lat/Lon to Census tracts
- Add Third Party Data
 - Location data, Economics data, etc.
- Group Dataset by Year and Tract
- **Feature Enhancement**
 - Ex: Employment Rate, Male to Female Ratio
- **Missing Values**
 - Drop data
 - Missing time series rent data (15% of tracts)
 - Linear Interpolate
 - Corporate Finance Data (SP500, Corporate Bonds, T Bills), Crime Rate, Demographic Data
 - KNN
 - Employment Rate (2 tracts), Property Value (100 Properties)
 - County Median
 - Population, Income, Employment Data (33%)

Final Panel Dataset

	A	B	C	D	E	F	G	H	I	J	K	L	M	N
1	tract	year	id	growth	growth_per_sqft	mean_rent	mean_rent_per	total_units	region	state_code	county	property_sqft	current_rent	state_z
2	1003010500	2016	1.00301E+13	3.47909791	3.968253968	627.0833333	0.655	68 s	1	3	65200	633 AL		
3	1003010500	2017	1.00301E+13	4.598006645	4.071246819	655.9166667	0.681666667	68 s	1	3	65200	633 AL		
4	1003010500	2018	1.00301E+13	-5.107356117	-4.88997555	622.4166667	0.648333333	68 s	1	3	65200	633 AL		
5	1003010500	2019	1.00301E+13	1.017539162	1.285347044	628.75	0.656666667	68 s	1	3	65200	633 AL		
6	1003010500	2020	1.00301E+13	5.990722333	5.837563452	666.4166667	0.695	68 s	1	3	65200	633 AL		
7	1003010500	2021	1.00301E+13	-2.450919095	-2.637889688	650.0833333	0.676666667	68 s	1	3	65200	633 AL		
8	1003010500	2022	1.00301E+13	2.479169337	2.561576355	666.2	0.694	68 s	1	3	65200	633 AL		
9	1003010600	2016	1.00301E+13	2.424568966	2.529761905	633.6666667	0.574166667	56 s	1	3	61800	691 AL		
10	1003010600	2017	1.00301E+13	1.407154129	1.45137881	642.5833333	0.5825	56 s	1	3	61800	691 AL		
11	1003010600	2018	1.00301E+13	1.997146933	2.00286123	655.4166667	0.594166667	56 s	1	3	61800	691 AL		
12	1003010600	2019	1.00301E+13	-0.152574698	-0.280504909	654.4166667	0.5925	56 s	1	3	61800	691 AL		
13	1003010600	2020	1.00301E+13	2.941551	2.953586498	673.6666667	0.61	56 s	1	3	61800	691 AL		
14	1003010600	2021	1.00301E+13	2.090549233	2.322404372	687.75	0.624166667	56 s	1	3	61800	691 AL		
15	1003010600	2022	1.00301E+13	7.888040712	7.343124166	742	0.67	56 s	1	3	61800	691 AL		
16	1003010704	2016	1.00301E+13	-0.542197077	-0.35318824	703.1666667	0.790833333	180 s	1	3	160288	1417 AL		
17	1003010704	2017	1.00301E+13	15.86868926	15.70073762	814.75	0.915	180 s	1	3	160288	1417 AL		
18	1003010704	2018	1.00301E+13	5.615219392	5.464480874	860.5	0.965	180 s	1	3	160288	1417 AL		
19	1003010704	2019	1.00301E+13	8.15417393	8.376511226	930.6666667	1.045833333	180 s	1	3	160288	1417 AL		
20	1003010704	2020	1.00301E+13	0.698424069	0.478087649	937.1666667	1.050833333	180 s	1	3	160288	1417 AL		
21	1003010704	2021	1.00301E+13	29.94842611	30.13481364	1217.833333	1.3675	180 s	1	3	160288	1417 AL		

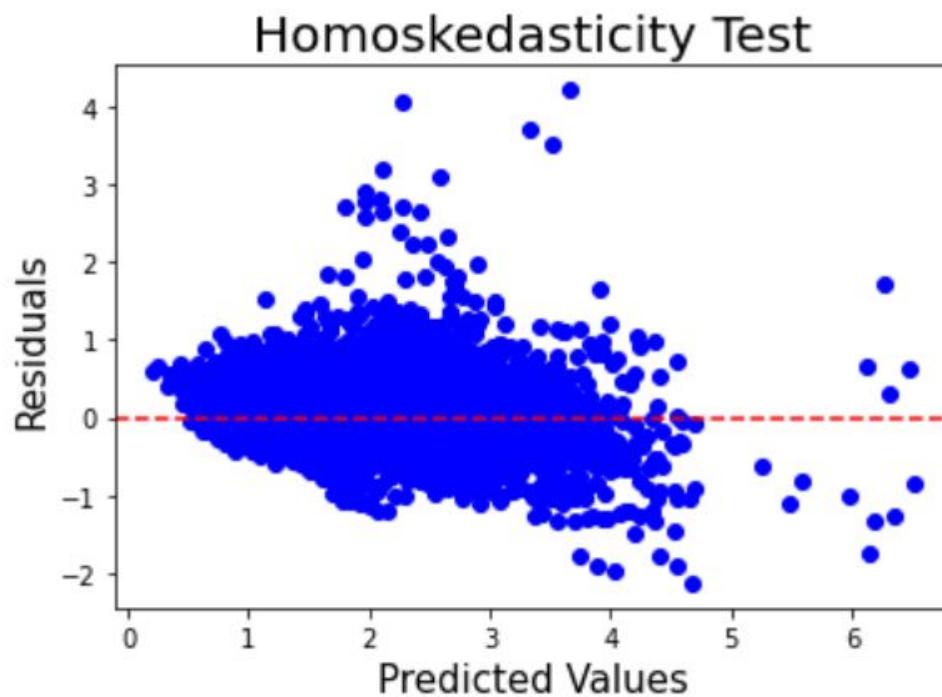
Time Series Models: Pooled OLS

Pooled OLS can be described as simple OLS (Ordinary Least Squared) model that is performed on panel data. It ignores time and individual characteristics and focuses only on dependencies **between** the individuals. It requires that there is no correlation between unobserved, independent variable(s) and the IVs (i.e. exogeneity).

$$y_{it} = x_{it}\beta + \alpha_i + u_{it} \quad cov(x_{it}, \alpha_i) = 0, \text{ for } t = 1, \dots, T, i = 1, \dots, N$$

Time Series Models: Pooled OLS

```
exog = sm.add_constant(dataset[vars])
mod = PooledOLS(dataset.y, exog, check_rank=False)
res = mod.fit(cov_type='clustered', cluster_entity=True)
```



- 3-year period data modeling 2017-2019 for 2020-2022
- Assumptions check (Homoskedasticity & Autocorrelation)
- Solutions: Random-Effects (RE) Model and Fixed-Effects (FE) Model

```
from statsmodels.stats.stattools import durbin_watson
pooled_OLS_dataset = pd.concat([dataset[vars], residuals_pooled_OLS], axis=1)
durbin_watson_test_results = durbin_watson(pooled_OLS_dataset['residual'])
print(durbin_watson_test_results)
```

✓ 0.4s

nan

- Panel data method for regression analysis
- All parameters being used are fixed, eliminating any randomness in the data
 - Uses the entity effect to determine model
- 4 Year Panel Period
 - Model trained using 2019-2022 data
- R-square as primary decision factor
- Feature selection was done by deleting the predictor variable with the highest p-value sequentially
- Does not include covariance variables (Drops these variables if present in dataset)

$$y_i = \sum_{j=1}^J \alpha_j z_{j[i]} + \beta x_i + \varepsilon_i; \quad \varepsilon_i \sim N(0, \sigma_y^2).$$

Time Series Models: Fixed Effects Regression

			Parameter Estimates						
			Parameter	Std. Err.	T-stat	P-value	Lower CI	Upper CI	
Dep. Variable:	mean_rent_per	R-squared:	0.6362						
Estimator:	PanelOLS	R-squared (Between):	0.7052						
No. Observations:	45732	R-squared (Within):	0.6362						
Date:	Thu, Jun 30 2022	R-squared (Overall):	0.7015						
Time:	15:54:39	Log-likelihood	7.167e+04						
Cov. Estimator:	Unadjusted	F-statistic:	2119.2						
Entities:	11767	P-value	0.0000						
Avg Obs:	3.8865	Distribution:	F(28, 33937)						
Min Obs:	1.0000	F-statistic (robust):	2119.2						
Max Obs:	4.0000	P-value	0.0000						
Time periods:	4	Distribution:	F(28, 33937)						
Avg Obs:	1.143e+04								
Min Obs:	1.1e+04								
Max Obs:	1.176e+04								
<hr/>									
PanelOLS Estimation Summary									
const		0.0042	0.0314	0.1325	0.8946	-0.0573	0.0656		
total_units		0.0035	0.0001	28.983	0.0000	0.0033	0.0038		
property_sqft		-3.566e-06	1.212e-07	-29.423	0.0000	-3.804e-06	-3.329e-06		
current_rent		0.0005	1.107e-05	46.318	0.0000	0.0005	0.0005		
property_value		1.568e-06	1.8e-08	87.116	0.0000	1.533e-06	1.603e-06		
median_income		1.961e-07	6.647e-08	2.9503	0.0032	6.583e-08	3.264e-07		
population		4.235e-05	7.621e-06	5.5568	0.0000	2.741e-05	5.729e-05		
white_pop		-1.633e-05	2.66e-06	-6.1376	0.0000	-2.154e-05	-1.111e-05		
black_pop		-7.713e-06	3.89e-06	-1.9826	0.0474	-1.534e-05	-8.787e-08		
asian_pop		-2.023e-05	5.47e-06	-3.6984	0.0002	-3.095e-05	-9.509e-06		
pacific_pop		6.475e-05	2.223e-05	2.9132	0.0036	2.118e-05	0.0001		
pop_0_19		-4.355e-05	8.339e-06	-5.2224	0.0000	-5.989e-05	-2.72e-05		
pop_20_29		-3.693e-05	8.612e-06	-4.2887	0.0000	-5.382e-05	-2.005e-05		
pop_30_39		-2.585e-05	9.065e-06	-2.8516	0.0044	-4.362e-05	-8.082e-06		
pop_40_49		-5.119e-05	9.326e-06	-5.4887	0.0000	-6.947e-05	-3.291e-05		
pop_50_59		-2.015e-05	9.62e-06	-2.0946	0.0362	-3.901e-05	-1.295e-06		
pop_60_69		-3.004e-05	1.022e-05	-2.9403	0.0033	-5.006e-05	-1.001e-05		
pop_70_over		-6.489e-05	1.073e-05	-6.0488	0.0000	-8.591e-05	-4.386e-05		
occupancy_rate		-0.0013	8.934e-05	-14.038	0.0000	-0.0014	-0.0011		
house_supply		3.31e-06	1.018e-06	3.2516	0.0011	1.315e-06	5.306e-06		
crime_rate		-0.0001	1.663e-05	-6.6664	0.0000	-0.0001	-7.829e-05		
ROI_SP500		-0.0800	0.0085	-9.3616	0.0000	-0.0968	-0.0633		
ROI_Tbill		1.5735	0.1290	12.200	0.0000	1.3207	1.8263		
ROI_Tbond		0.7152	0.0299	23.951	0.0000	0.6567	0.7738		
grocery_vists_per		0.1359	0.0360	3.7687	0.0002	0.0652	0.2065		
hospitals_per		0.0705	0.0200	3.5263	0.0004	0.0313	0.1097		
colleges_per		0.0569	0.0144	3.9412	0.0001	0.0286	0.0852		
private_schools_per		-0.0838	0.0315	-2.6576	0.0079	-0.1457	-0.0220		
employment_rate		0.1620	0.0194	8.3654	0.0000	0.1241	0.2000		

F-test for Poolability: 50.393

P-value: 0.0000

Distribution: F(11766, 33937)

Included effects: Entity

Fixed Effects Model: $y_i = \sum_{j=1}^J \alpha_j z_{j[i]} + \beta x_i + \varepsilon_i; \quad \varepsilon_i \sim N(0, \sigma_y^2).$

- Linear Regression of Y on X
- Adds to the specification a series of indicator variables Zj for each unit

Random Effects Model: $y_i = \alpha_j[i] + \beta x_i + \varepsilon_i; \quad \alpha_j \sim N(\mu_\alpha, \sigma_\alpha^2); \quad \varepsilon_i \sim N(0, \sigma_y^2).$

- α_j are not estimated directly
- Follow a specified probability distribution

- Similar to the Fixed Effect Regression
 - Panel Data methods for regression analysis.
 - Does not use entity effect to determine model

```
exog_vars = X_var_names
exog = sm.add_constant(df[exog_vars])
mod = PanelOLS(df.mean_rent_per, exog, entity_effects = False, check_rank=False, drop_absorbed=True)
RE = mod.fit()
print(RE)
```

- 4 Year Panel Period
 - Model trained from 2019-2022

```
Time periods: 4 Distribution:
Avg Obs: 1.145e+04
Min Obs: 1.1e+04
Max Obs: 1.176e+04
```

Parameter Estimates

- R-square as primary decision factor
- Feature selection was done by deleting the predictor variable with the highest p-value sequentially

Time Series Models: Random Effects Regression

$$Y_{2025} = B_1 X_{2022} + B_2 X_{2022} + \dots + B_0$$

	tract	year		id	growth	growth_per_sqft	mean_rent	mean_rent_per	total_units	region	state_code	...	county_male_population
0	1003010500	2016	10030105002016	3.479098	3.968254	627.083333	0.655000	68.0	s	1	...		2705.50000
1	1003010500	2017	10030105002017	4.598007	4.071247	655.916667	0.681667	68.0	s	1	...		2728.00000
2	1003010500	2018	10030105002018	-5.107356	-4.889976	622.416667	0.648333	68.0	s	1	...		2701.50000
3	1003010500	2019	10030105002019	1.017539	1.285347	628.750000	0.656667	68.0	s	1	...		2991.00000
4	1003010500	2020	10030105002020	5.990722	5.837563	666.416667	0.695000	68.0	s	1	...		2801.00000

year		total_units	property_sqft	current_rent	property_value	income_tax_rate	corporate_tax_rate	state_sales_tax_rate	avg_local_sales_tax_rate	combined_rate	
2016	0	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	
2017	1	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	
2018	2	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	
2019	3	68.0	65200.0	X:2016	633.0	119546.416667	0.036667	0.065	0.04	0.052386	0.092386
2020	4	68.0	65200.0	X:2017	633.0	122164.416667	0.036667	0.065	0.04	0.052386	0.092386

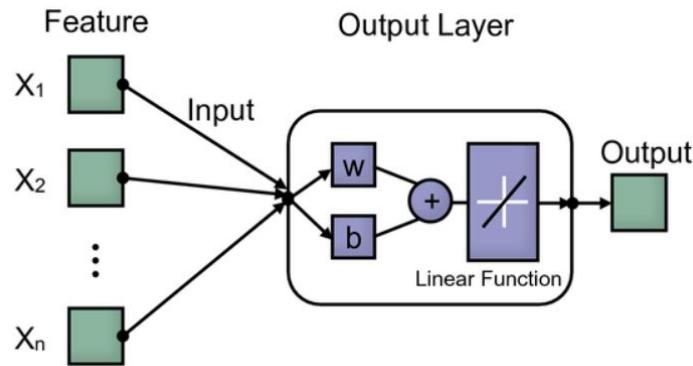
Random Effects Regression Results

RandomEffects Estimation Summary							predict	growth_25		
Dep. Variable:	mean_rent_per	R-squared:	0.6912							
Estimator:	RandomEffects	R-squared (Between):	0.7968							
No. Observations:	45732	R-squared (Within):	0.6099							
Date:	Thu, Jun 30 2022	R-squared (Overall):	0.7937							
Time:	15:54:41	Log-likelihood	6.302e+04							
Cov. Estimator:	Unadjusted									
Entities:	11767	F-statistic:	4263.5							
Avg Obs:	3.8865	P-value	0.0000							
Min Obs:	1.0000	Distribution:	F(24, 45707)							
Max Obs:	4.0000	F-statistic (robust):	4170.1							
Time periods:	4	P-value	0.0000							
Avg Obs:	1.143e+04	Distribution:	F(24, 45707)							
Min Obs:	1.1e+04			Parameter Estimates						
Max Obs:	1.176e+04			Parameter	Std. Err.	T-stat	P-value	Lower CI	Upper CI	
	const	-1.1344	0.0493	-23.014	0.0000	-1.2310	-1.0378			
	property_sqft	-3.386e-07	1.71e-08	-19.802	0.0000	-3.721e-07	-3.051e-07			
	current_rent	0.0004	5.238e-06	71.877	0.0000	0.0004	0.0004			
	property_value	1.034e-06	1.161e-08	89.061	0.0000	1.011e-06	1.057e-06			
	income_tax_rate	-0.2669	0.1281	-2.0833	0.0372	-0.5180	-0.0158			
	corporate_tax_rate	0.6044	0.1188	5.0851	0.0000	0.3714	0.8373			
	state_sales_tax_rate	1.2445	0.3724	3.3417	0.0008	0.5146	1.9744			
	max_local_sales_tax	0.7456	0.2168	3.4394	0.0006	0.3207	1.1705			
	avg_household_income	6.867e-07	1.747e-07	3.9303	0.0001	3.443e-07	1.029e-06			
	population	7.779e-06	3.204e-06	2.4276	0.0152	1.498e-06	1.406e-05			
	white_pop	-2.614e-05	2.369e-06	-11.034	0.0000	-3.079e-05	-2.15e-05			
	black_pop	-2.933e-05	3.056e-06	-9.5983	0.0000	-3.532e-05	-2.334e-05			
	asian_pop	-3.435e-05	4.52e-06	-7.6009	0.0000	-4.321e-05	-2.549e-05			
	pop_0_19	-1.369e-05	4.498e-06	-3.0434	0.0023	-2.25e-05	-4.873e-06			
	pop_20_29	3.146e-05	4.339e-06	7.2523	0.0000	2.296e-05	3.997e-05			
	pop_30_39	3.739e-05	5.221e-06	7.1604	0.0000	2.715e-05	4.762e-05			
	house_supply	1.059e-05	7.55e-07	14.025	0.0000	9.109e-06	1.207e-05			
	crime_rate	-7.62e-05	1.562e-05	-4.8771	0.0000	-0.0001	-4.558e-05			
	Tbond5	0.0001	1.825e-06	59.343	0.0000	0.0001	0.0001			
	hist_risk	13.408	0.9253	14.491	0.0000	11.594	15.221			
	SP500	-0.1674	0.0051	-32.833	0.0000	-0.1773	-0.1574			
	starbucks_vists_per	0.0376	0.0041	9.0829	0.0000	0.0295	0.0457			
	colleges_per	0.0686	0.0060	11.361	0.0000	0.0568	0.0804			
	public_schools_per	0.0489	0.0024	20.118	0.0000	0.0441	0.0537			
	private_schools_per	0.0411	0.0061	6.7698	0.0000	0.0292	0.0530			

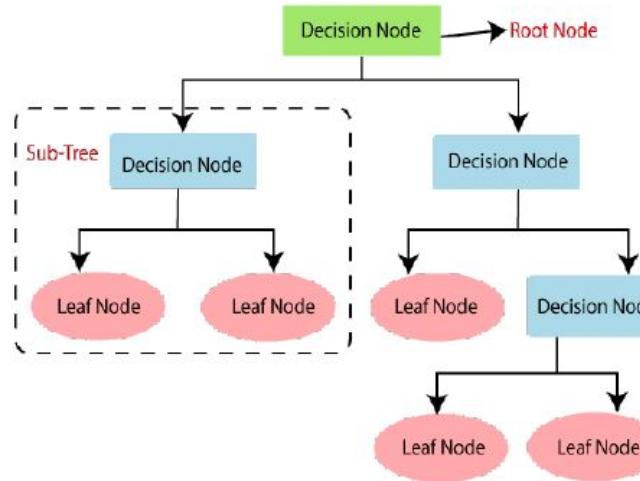
Random Effect Regressions Results (FE Variables)

Traditional Machine Learning (ML) Models

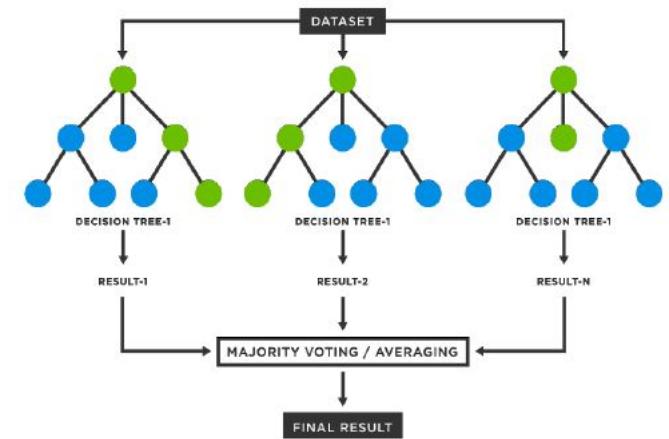
Linear Regression



Decision Tree

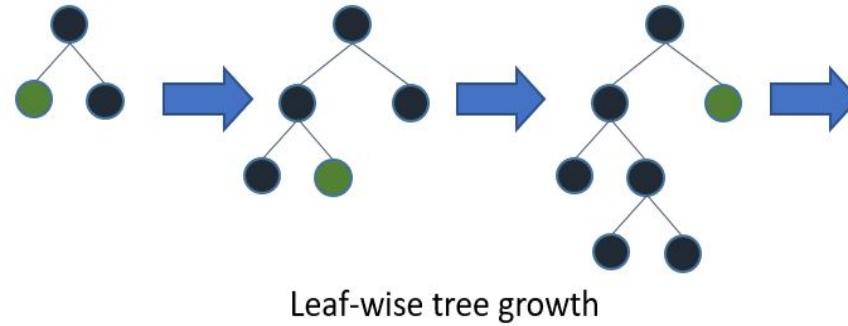


Random Forest

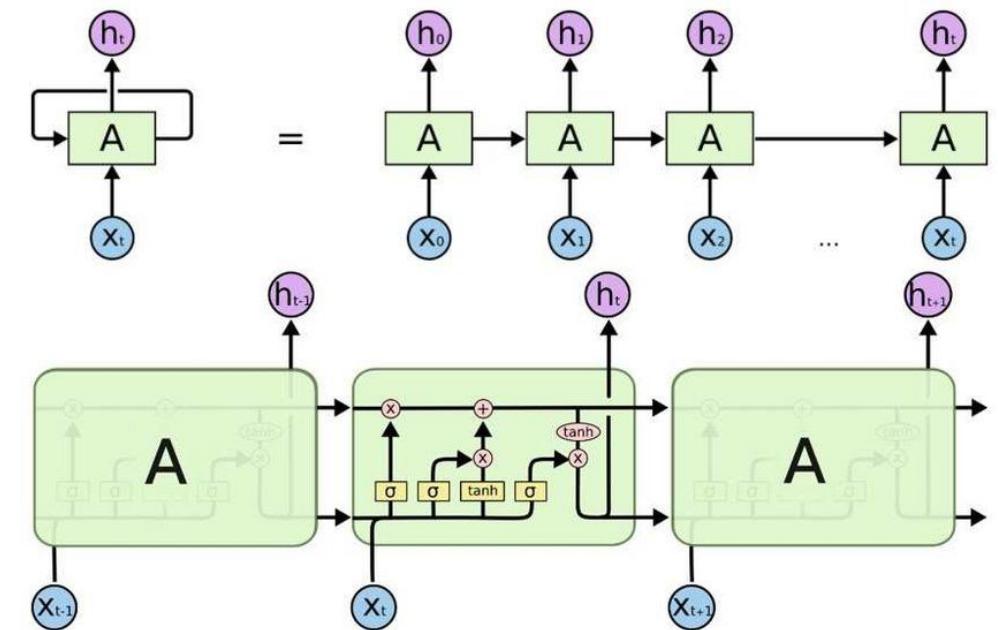


Traditional ML Models

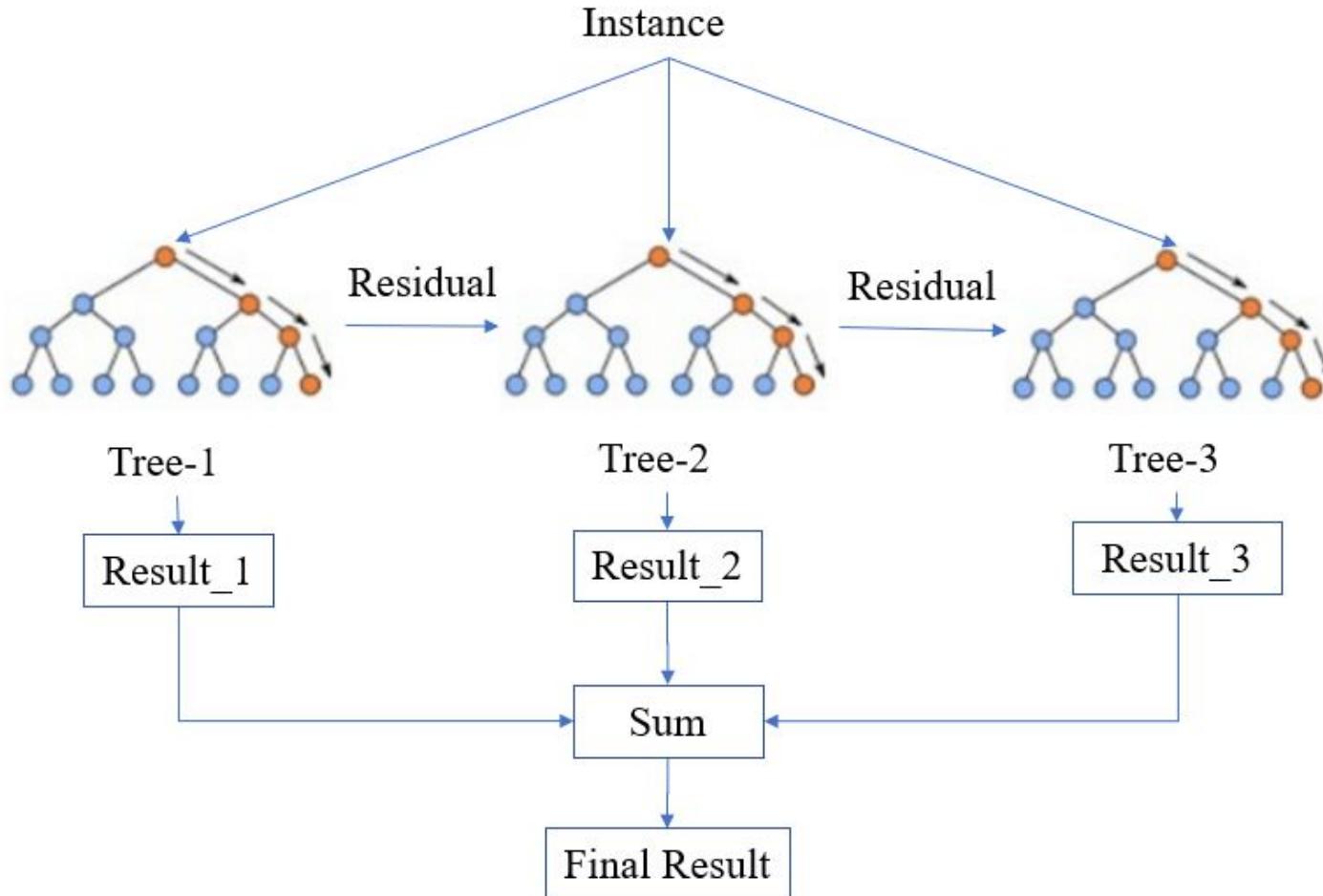
LightGBM



Long short-term memory

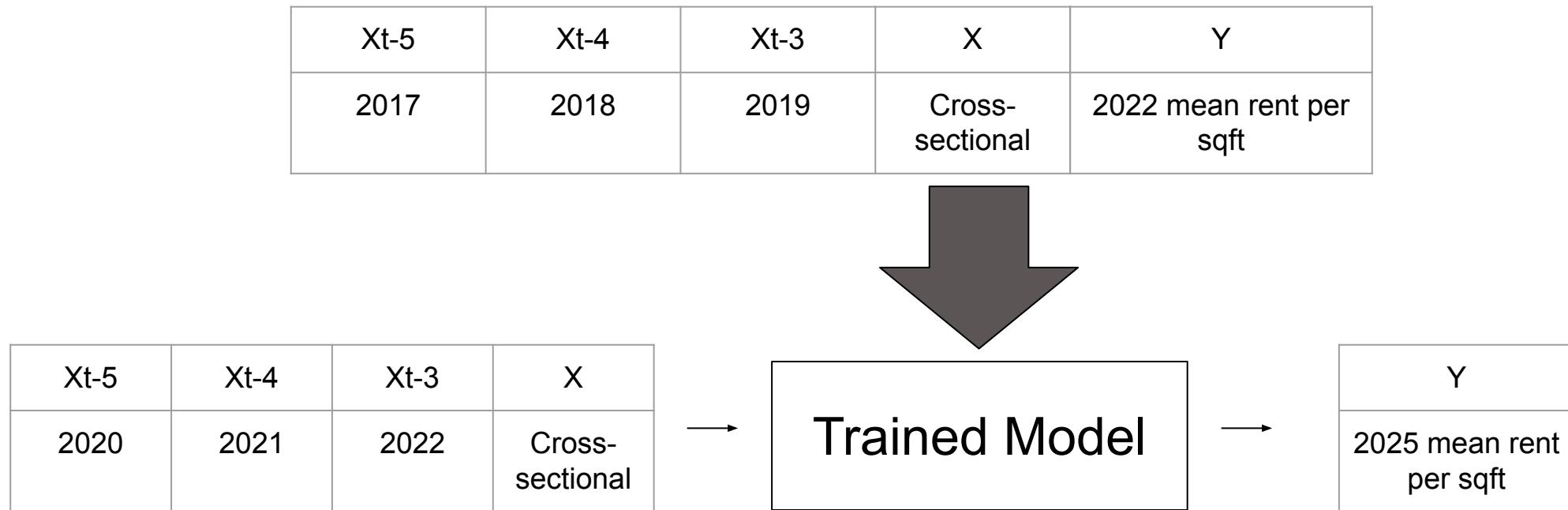


Traditional ML Model: XGBoost



- Decision tree based ensemble model
- Designed for speed and performance
- Suit with small-to-medium structured/tabular data
- Easy for overfitting but overcome by tuning properly
- Training time is pretty high

Traditional ML Models: Steps



- Y: Mean Rent Per SQFT in year t
- X: Cross- sectional data ; Time-series data in year t-3, t-4, t-5
- Applied the 10-fold cross validation to find the best hyperparameters
- R-squared as criteria

Traditional ML Models: Steps

Linear Regression Models

Region	R-Squared Value
All Regions	.952
Florida	.923
SouthEast	.926
South	.931
SouthWest	.920
Western	.867
Southern California	.877

Traditional ML Models: Results

Model	Parameters				Avg R squared
Linear Regression					0.9517
Decision Tree	max_depth	min_leaf	min_split	splitter	0.9291
	10	10	50	best	
Random Forest	max_depth	min_leaf	min_split	n_estimators	0.9466
	10	5	10	100	
XGBoost	learning_rate	max_depth	n_estimators		0.9644
	0.1	4	800		
LightGBM	learning_rate	max_depth	n_estimators		0.9618
	0.1	4	500		
RNN - LSTM	epochs	batch_size	verbose		0.8640
	30	50	1		

Best Model Results and Evaluation

Sort by Fixed Effect

	tract	FE_rent_growth_2025	XG_rent_growth_2025	LGB_rent_growth_2025
0	04013217001	1.684731	0.161208	0.171603
1	06059062645	1.525346	0.136040	0.141538
2	06059062658	1.473987	0.120287	0.123380
3	06037264000	1.325111	0.066382	0.054277
4	48201333204	1.297981	0.083040	0.071199
5	18043070801	1.276547	0.561004	0.605929
6	06059063010	1.217841	0.142234	0.144315
7	06059062649	1.181896	0.100328	0.108008
8	48201421402	1.153177	0.050812	0.068860
9	22051020517	1.034263	0.095851	0.103653

Sort by XGBoost

	tract	FE_rent_growth_2025	XG_rent_growth_2025	LGB_rent_growth_2025
0	18043070801	1.276547	0.561004	0.605929
1	06037800205	0.586632	0.466287	0.452792
2	13021011100	0.369543	0.435106	0.423584
3	40131050105	0.097562	0.390880	0.380768
4	37119005200	-0.255410	0.368388	0.349797
5	22113950600	-0.153783	0.352597	0.374065
6	40109102500	-0.134369	0.344746	0.340528
7	12081002025	0.192506	0.342930	0.330297
8	13135050747	0.019494	0.342184	0.329722
9	40109101300	0.009879	0.339300	0.354523

	tract	FE_rent_growth_2025	XG_rent_growth_2025	LGB_rent_growth_2025
0	18043070801	1.276547	0.561004	0.605929
1	06037800205	0.586632	0.466287	0.452792
2	13021011100	0.369543	0.435106	0.423584
3	40131050105	0.097562	0.390880	0.380768
4	04013050620	0.364423	0.326314	0.379444
5	22113950600	-0.153783	0.352597	0.374065
6	12091020800	-0.213282	0.328451	0.367889
7	40109101300	0.009879	0.339300	0.354523
8	37119005200	-0.255410	0.368388	0.349797
9	40109102500	-0.134369	0.344746	0.340528

Sort by LightGBM

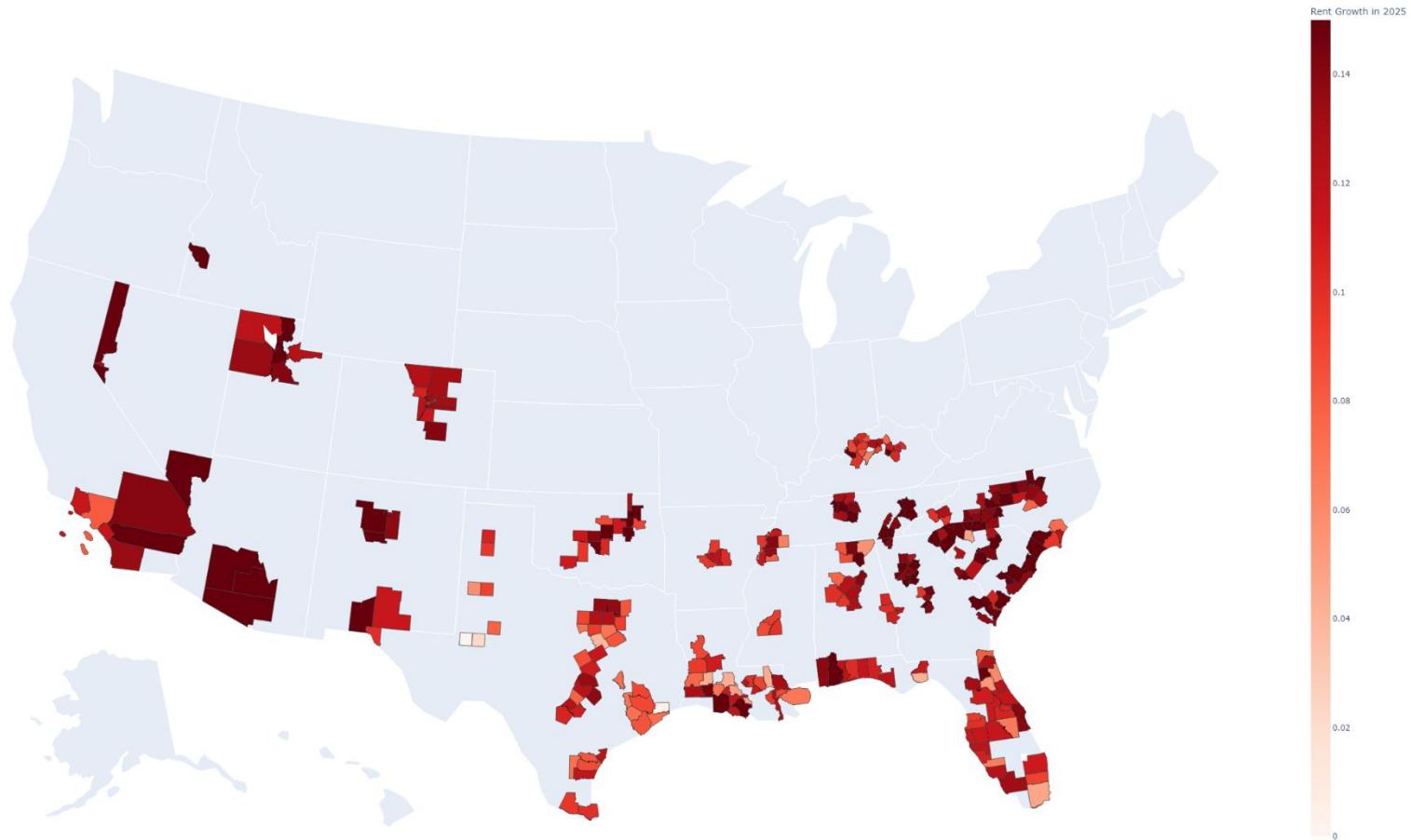
Final Model Selection and Analysis

- Final model: **XGBoost**
- R-Squared Value: **0.9644**
- Expected 3 Year ROI* (\$10m): **\$9.17 Million Profit**
- Causes of Time/Series Model Limitations:
 - Short Time Period?
 - Linearity?

* Assuming 4% cap rate, 50% expense rate, and \$1M invested in each of the top 10 tracts.

Overview of Results

1. Tract: 18043070801 Expected Growth: 0.561
2. Tract: 06037800205 Expected Growth: 0.466
3. Tract: 13021011100 Expected Growth: 0.435
4. Tract: 40131050105 Expected Growth: 0.391
5. Tract: 37119005200 Expected Growth: 0.368
6. Tract: 22113950600 Expected Growth: 0.353
7. Tract: 40109102500 Expected Growth: 0.345
8. Tract: 12081002025 Expected Growth: 0.343
9. Tract: 13135050747 Expected Growth: 0.342
10. Tract: 40109101300 Expected Growth: 0.339



Qualitative Significance of Tracts

1. Tract: 18043070801
Location: (Floyd County, Indiana)
Expected Growth: 0.561
Summary: Hispanic population increased, access to many public schools, near Indiana University Southeast and museums

2. Tract: 22113950600
Location: (Vermilion County, Louisiana)
Expected Growth: 0.353
Summary: Total housing units increased by 10%, 9.5% population growth, hospitals are closely accessible

3. Tract: 13021011100
Location: (Bibb County, Georgia)
Expected Growth: 0.435
Summary: Near attractions, access to many public schools, university, banks, restaurants

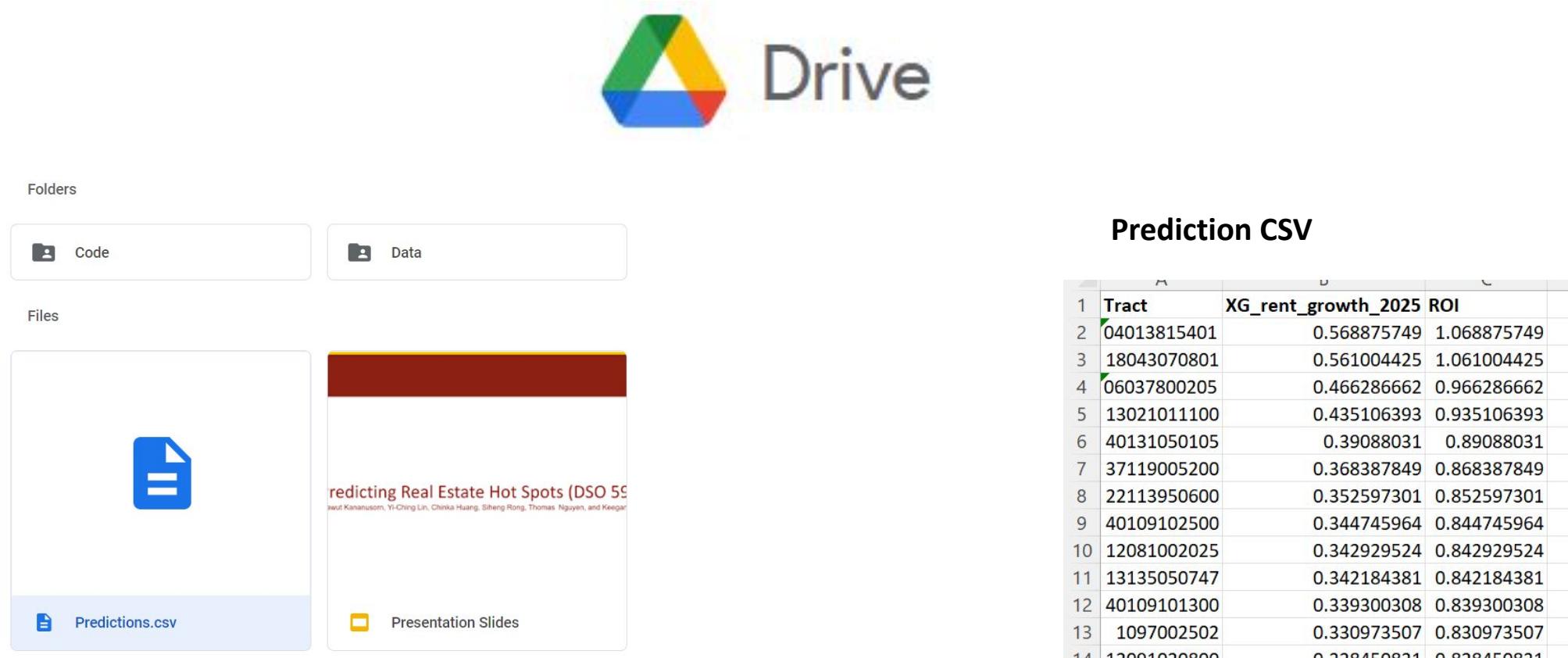
4. Tract: 13135050747
Location: (Gwinnett County, Georgia)
Expected Growth: 0.342
Summary: Total housing units increased by 31.8%, population growth 55.5%, near attractions, access to many schools and hospitals

5. Tract: 37119005200
Location: (Mecklenburg County, North Carolina)
Expected Growth: 0.368
Summary: white population increased by 90%, Hospitals, clinics or other healthcare facilities are closely accessible, neighborhood walkability is good

Qualitative Significance of Tracts

6. Tract: 06037800205
Location: (Los Angeles County, California)
Expected Growth: 0.466
Summary: Malibu area, close to city, high income people, almost 10% WFH and 40% migration from last year
7. Tract: 40131050105
Location: (Rogers County, Oklahoma)
Expected Growth: 0.391
Summary: In small city near Rogers State University, and many hospitals,
8. Tract: 40109102500
Location: (Oklahoma County, Oklahoma)
Expected Growth: 0.345
Summary: High population growth in Oklahoma (~17%), new jobs, low unemployment rate(~1.6%), ~30% movie-in last year
9. Tract: 40109101300
Location: (Oklahoma County, Oklahoma)
Expected Growth: 0.339
Summary: High population growth in Oklahoma (~17%), new jobs, low unemployment rate(~1.6%)
10. Tract: 12081002025
Location: (Manatee County, Florida)
Expected Growth: 0.343
Summary: In the golf area, near many attractions, recently government investment in infrastructure, education, communities

Prediction and Code Submission



The image shows a Google Drive interface. At the top is the Google Drive logo. Below it, under 'Folders', are 'Code' and 'Data'. Under 'Files', there is a 'Predictions.csv' file (blue icon) and a presentation titled 'Predicting Real Estate Hot Spots (DSO 55)' (yellow icon). A preview of the presentation slide is visible, showing the title and some text. To the right, a table titled 'Prediction CSV' displays data from the 'Predictions.csv' file.

	Tract	XG_rent_growth_2025	ROI
1	04013815401	0.568875749	1.068875749
2	18043070801	0.561004425	1.061004425
3	06037800205	0.466286662	0.966286662
4	13021011100	0.435106393	0.935106393
5	40131050105	0.39088031	0.89088031
6	37119005200	0.368387849	0.868387849
7	22113950600	0.352597301	0.852597301
8	40109102500	0.344745964	0.844745964
9	12081002025	0.342929524	0.842929524
10	13135050747	0.342184381	0.842184381
11	40109101300	0.339300308	0.839300308
12	1097002502	0.330973507	0.830973507
13	12001000000	0.328450021	0.828450021
14			

<https://drive.google.com/drive/folders/1Fj6biog-kUq5GLSRYcDUvkIzCTc-RIKj?usp=sharing>

Data Sources

Data	Factors	Sources
Economic Data	<ul style="list-style-type: none"> - Tax Rates (Sales, Income, Corporate, Combined) - Debt - Bond - National Indicators (SP 500) 	<ul style="list-style-type: none"> - Individual Tax Rate: State Individual Income Tax Rates - Corporate Tax Rate: State Corporate Income Tax Rates - Sales Tax Rate: State Sales Tax Rate - Debt: Debt in America - US Bond: US Bond and Returns
Location Data	<ul style="list-style-type: none"> - Starbucks, Banks, Public School, Private Schools, Hospitals (Normalized by Location Density) - Zip code area 	<ul style="list-style-type: none"> - Starbucks dataset: Starbucks Locations Worldwide Kaggle - Public schools: Public Schools Public Schools HIFLD Open Data (arcgis.com) - Private schools: Private Schools Private Schools HIFLD Open Data (arcgis.com) - Colleges: Colleges and Universities Colleges and Universities HIFLD Open Data (arcgis.com) - Banks: FDIC Insured Banks FDIC Insured Banks HIFLD Open Data (arcgis.com) - Hospitals: Hospitals Hospitals HIFLD Open Data (arcgis.com)
Demographic Data	<ul style="list-style-type: none"> - Population (gender, age, race) - Median Household Income - Employment Rate 	<ul style="list-style-type: none"> - ACS - TIGER
Housing Data	<ul style="list-style-type: none"> - Occupancy Rate - Housing Supply - Property Values 	<ul style="list-style-type: none"> - Yardi Property - FRED New Private Housing Units Authorized by Building Permits: 1-Unit Structures - Zillow Rent Index https://www.zillow.com/research/data/
Crime Data	<ul style="list-style-type: none"> - Crime Rate 	Crime Data Explorer

Research for modeling:

- Filling missing values
 - <https://www.machinelearningplus.com/time-series/time-series-analysis-python/>
- Fixed Effects, Random Effects, Pooled OLS
 - <https://www.ucl.ac.uk/~uctqiax/PUBLG100/2016/week8/seminar8.html>
 - <https://towardsdatascience.com/a-guide-to-panel-data-regression-theoretics-and-implementation-with-python-4c84c5055cf8>
- LSTM
 - https://github.com/melanieshi0120/COVID-19_global_time_series_panel_data/blob/master/Time_series_panel_data_model_part_confirmed_cases.ipynb
 - <https://melaniesoek0120.medium.com/covid-19-global-data-time-series-prediction-with-lstm-recurrent-neural-networks-f7825c4a1f6f>
- Feature Selection
 - <https://towardsdatascience.com/application-of-feature-selection-techniques-in-a-regression-problem-4278e2efd503>

Questions?