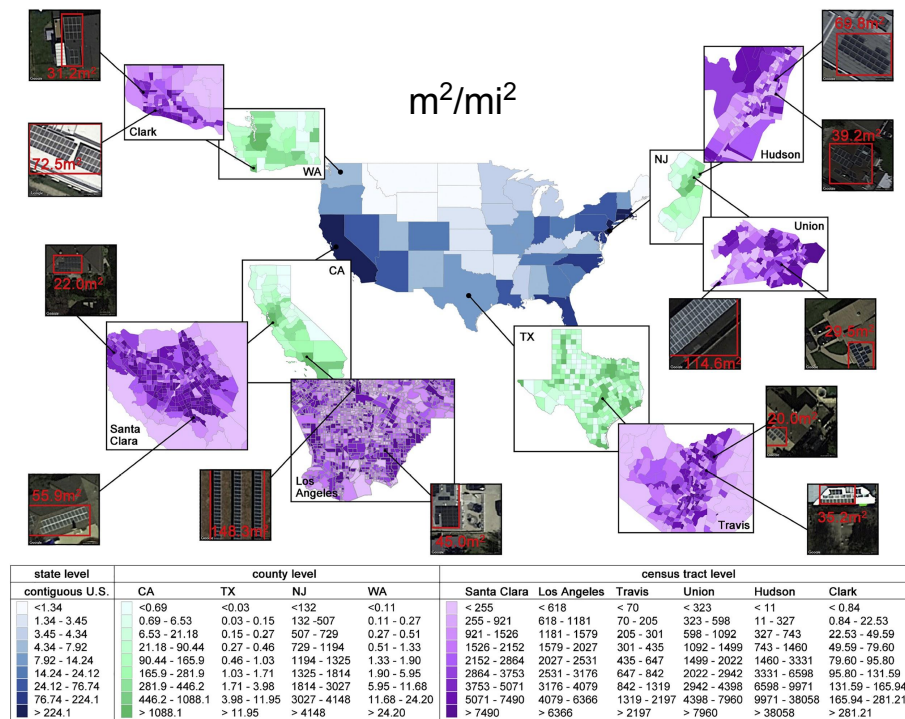


# DeepSolar

Jordan Landers | Thinkful Unit 3 Capstone

# Deep Solar Dataset

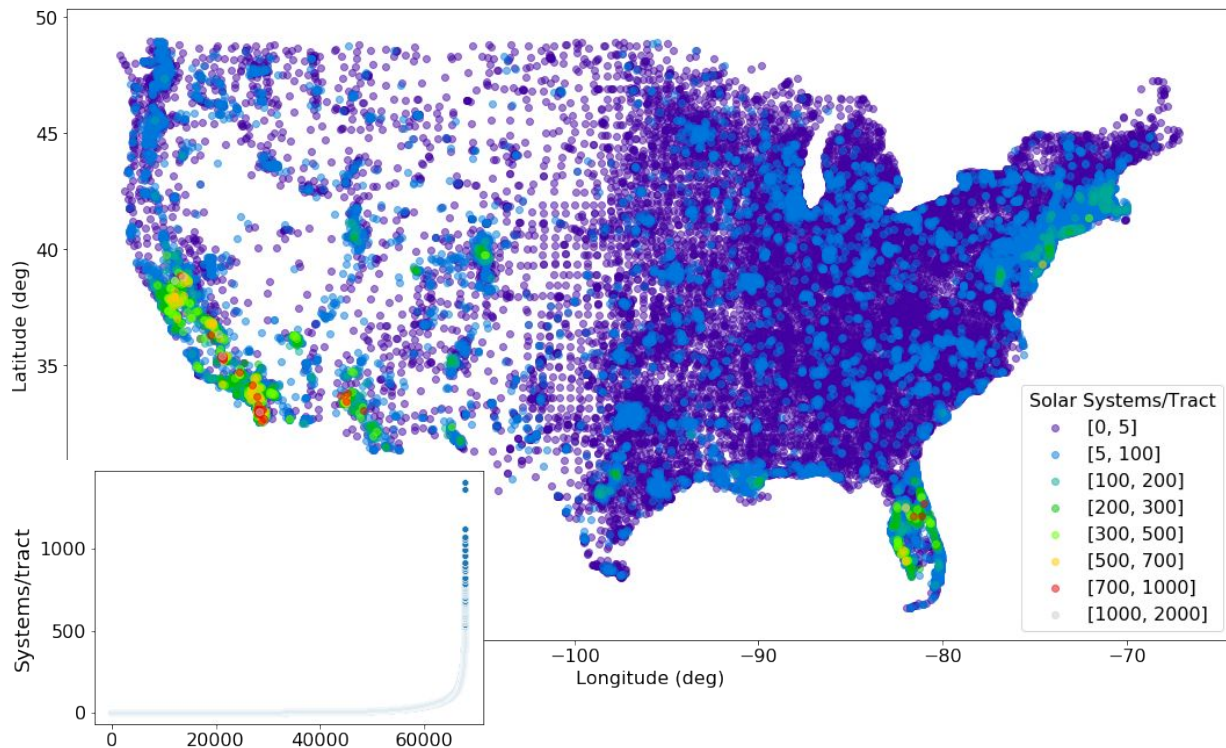


- Solar panel systems were counted from satellite images using a semi-supervised convolutional neural network
- Published dataset includes solar, geographic, and demographic data for 48 contiguous states at census tract resolution

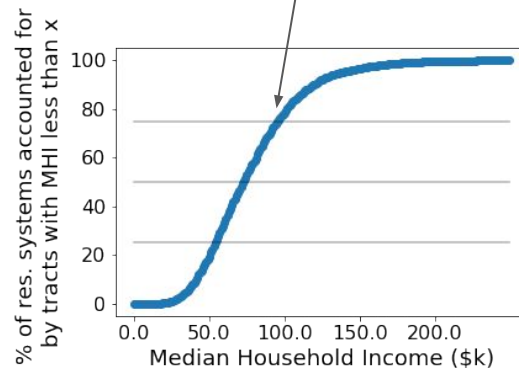
# Data Prep

- `voting_2012_dem_percentage`, `voting_2012_gop_percentage` **were** dropped because they were not reported in CO, CT, FL, GA, SC, UT, WY.
- Geographic/climate related values missing from 5208 tracts, affecting as many as 25% of tracts in certain states
  - For sets of missing tracts accounting for less than a threshold percentage of the respective county land area, missing values were replaced with county averages. 0%, 25% and 30% were investigated as threshold percentages.
- Dropped `null`, `infinite`, and `NaN` values
- Dropped categorical variables

# Target Variable: solar\_systems\_count\_residential



~75% of solar installations accounted for by median household income < \$100k

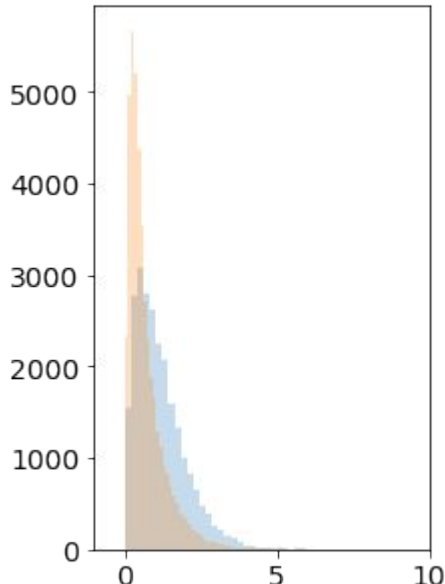
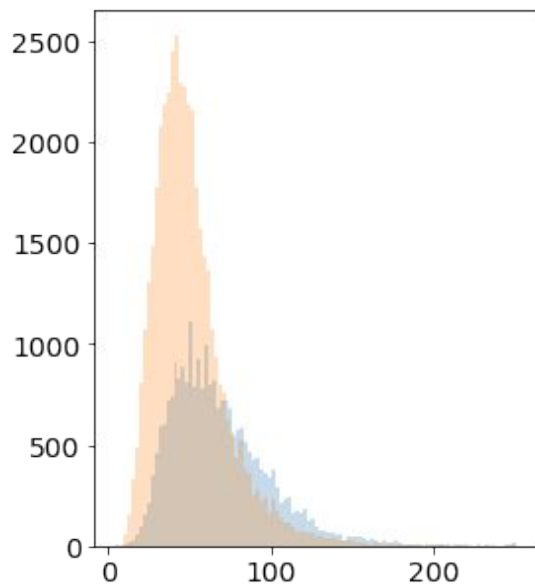


92% of households have income < \$100k

# Challenges

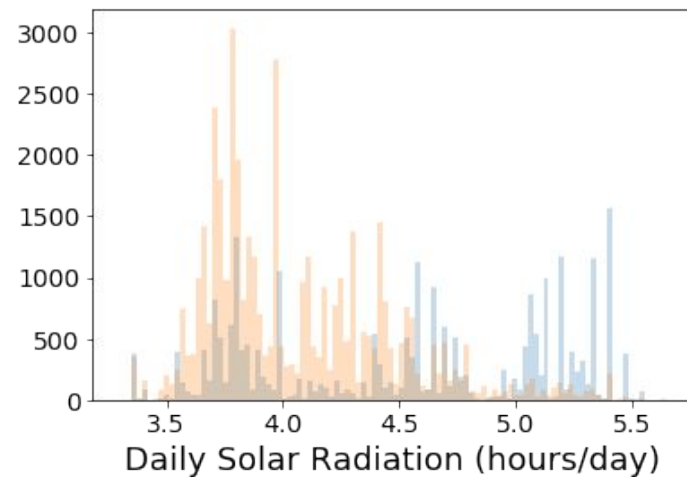
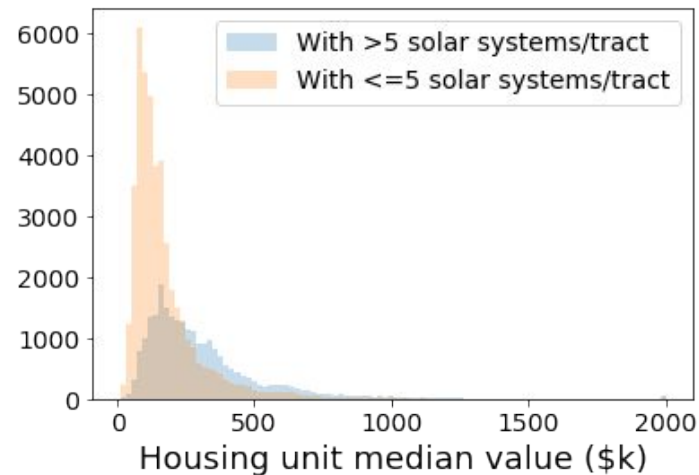
- Data reflect cumulative effect of various economic conditions and incentive programs, demographic data is a snapshot that may/may not reflect the characteristics of the population that installed solar systems
- Cost of living varies so money related variables are locally calibrated
  - \$67k (average `median_household_income` in CA) represents less buying power in CA than it does in TN, but to the model \$67k is \$67k
- Class imbalance (1.13% of households in dataset have solar installations)
- Inconsistent distribution
  - 62% of residential solar systems are located in CA (skew of residential solar systems = 7.4)
  - CA represents 11% of households
- Data do not reflect a system at saturation
  - Solar is being adopted at an increasing rate, so a household without solar is not necessarily a “non-solar household” as much as “not *yet*-solar household”

# Exploring the Data



Median Household Income (\$k)    People with higher ed (k)

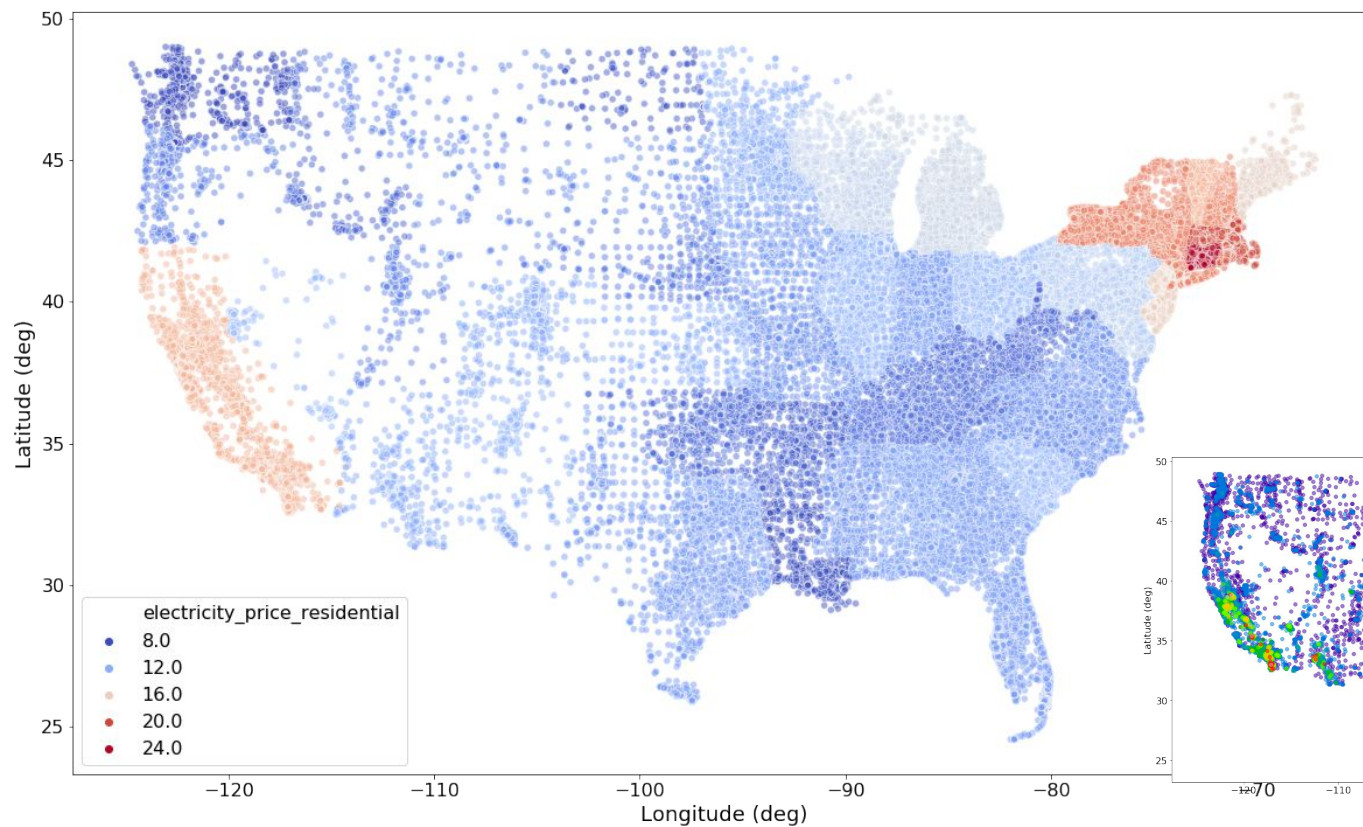
Populations are statistically different ( $p=0$ ) but distributions clearly overlap.



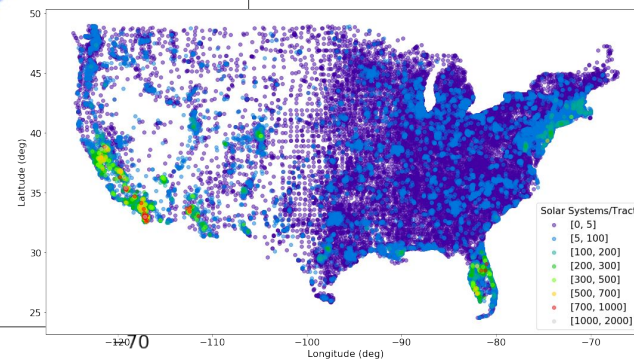
	Median Household Income	Count Higher Ed	Housing Unit Median Value	Daily Solar Radiation
--	----------------------------	-----------------	------------------------------	--------------------------

> 5 installations per tract				
count	24533	24533	24533	24533
mean	68653	1176	311485	4.54
Standard dev.	31159	884	233107	.63
min	9100	6	9999	3.35
max	250001	20482	2000000	5.65

<= 5 installations per tract				
count	43403	43403	43403	43403
mean	50779	702	173116	4.07
Standard dev.	23723	648	147227	.4
min	4147	0	9999	3.3
max	250001	10471	2000001	5.65



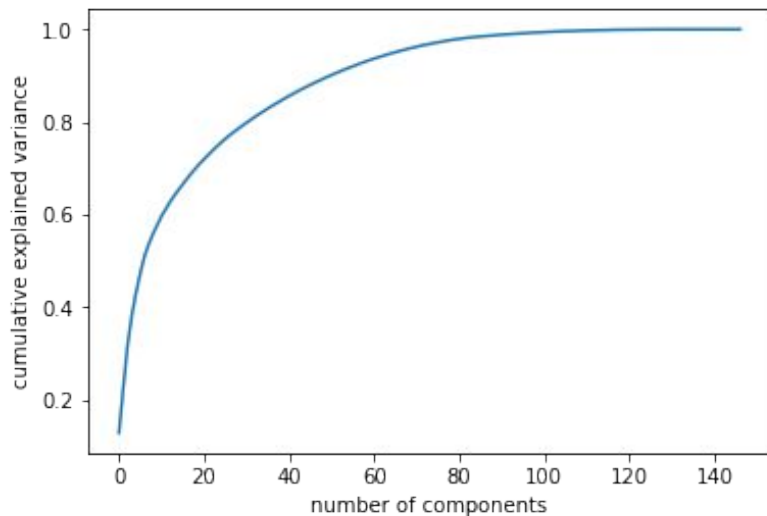
Indirect relationship  
between solar  
installations and  
electricity price





# Feature Set

140 variables in published data set pertaining to residential Heating, Weather, Income, Education, Employment, Race, Age, Occupation, Commute, Political, Energy consumption/Pricing, Incentives

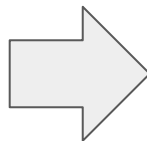


- PCA suggests that over 99% of variance can be accounted for using 92 principle components
- No strong feature loading in first 15 components

# Modeling

## Considerations

- Regression
- Complicated relationships between variables
- 100 - 140 variables
- 50,000 training examples

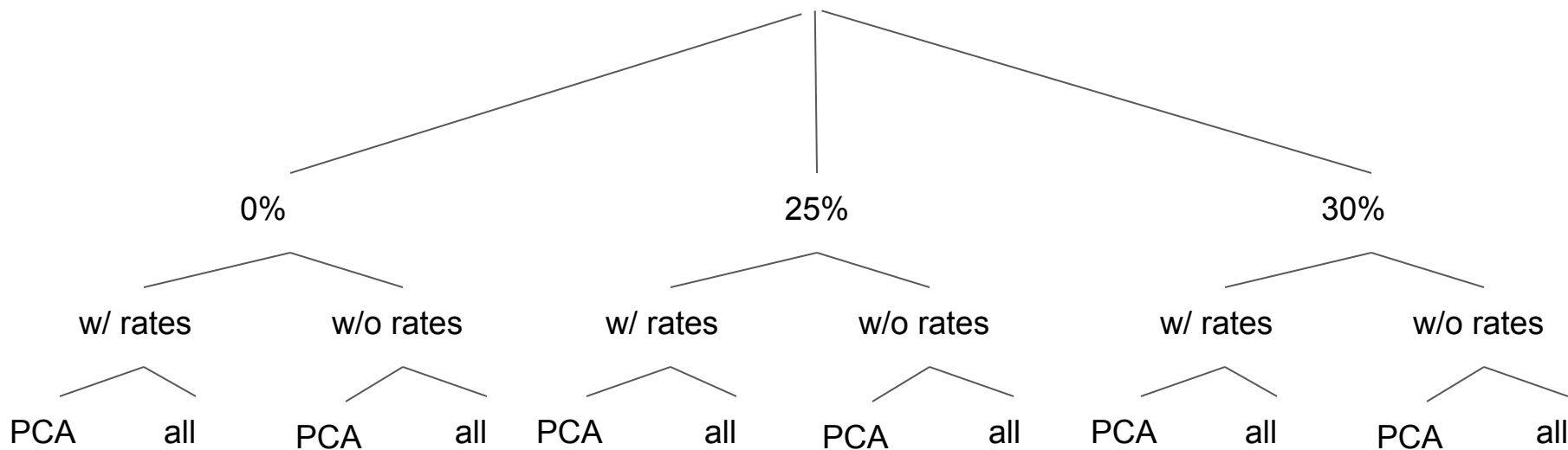


## Model Candidates

- Ridge
- Random Forests
- Support Vector Machines
- Gradient Boosting

# Model Sets

- 0%, 25%, 30% = threshold percent of county land requiring imputed geo values
- w/ rates, w/o rates = including/excluding rate variables from feature set
- PCA, all = reduce feature set using PCA, or retain all variables



# Results: PCA (n=100)

\* Same Training/Test sets used for all tests except where noted

	Ridge	Random Forests	SVM-R	Gradient Boosting-Regression
Training/Test 1	0.435 $\mp$ 0.028	.474 $\mp$ 0.045	0.608 $\mp$ 0.029	0.660 $\mp$ 0.035
Training/Test 2	0.434 $\mp$ 0.012	.471 $\mp$ 0.015	0.623 $\mp$ 0.042	0.640 $\mp$ 0.008
Training/Test 3	0.430 $\mp$ 0.031	.485 $\mp$ 0.018	0.627 $\mp$ 0.030	0.654 $\mp$ 0.034
CV score	0.433 $\mp$ 0.026	.477 $\mp$ 0.030	0.621 $\mp$ 0.034	0.652 $\mp$ 0.031
Total score	0.445 $\mp$ 0.004	.468 $\mp$ 0.019	0.594 $\mp$ 0.022	0.647 $\mp$ 0.012
Time to fit (s)	0.154 $\mp$ 0.006	448.2 $\mp$ 8.37	1961 $\mp$ 25.6	669.2 $\mp$ 6.6

# Results: Full Feature Set

\* Same Training/Test sets used for all tests except where noted

	Ridge	Random Forests	SVM-R	Gradient Boosting-Regression
Training/Test 1	.47 $\pm$ .031	.637 $\pm$ .029	.62 $\pm$ .025	.777 $\pm$ .02
Training/Test 2	.463 $\pm$ .015	.643 $\pm$ .02	.632 $\pm$ .039	.775 $\pm$ .028
Training/Test 3	.469 $\pm$ .011	.632 $\pm$ .02	.630 $\pm$ .01	.772 $\pm$ .015
CV score	.468 $\pm$ .021	.637 $\pm$ .02	.627 $\pm$ .028	.775 $\pm$ .022
Total score	.48 $\pm$ .011	.643 $\pm$ .023	.61 $\pm$ .017	.771 $\pm$ .01
Time to fit (s)	.111 $\pm$ .005	390.8 $\pm$ 1.26	1452 $\pm$ 25.5	569.7 $\pm$ 2.04

# Model Results: feature set comparison

\* Same Training/Test sets used for all tests except where noted

	With rates	Without rates	30% threshold	No imputed data	No dropped columns
Training/Test 1	0.777 $\pm$ 0.02	0.775 $\pm$ 0.02	0.775 $\pm$ 0.027	0.77 $\pm$ 0.026	0.773 $\pm$ 0.019
Training/Test 2	0.775 $\pm$ 0.028	0.771 $\pm$ 0.026	0.776 $\pm$ 0.017	0.776 $\pm$ 0.014	0.781 $\pm$ 0.027
Training/Test 3	0.772 $\pm$ 0.015	0.767 $\pm$ 0.015	0.773 $\pm$ 0.036	0.774 $\pm$ 0.013	0.79 $\pm$ 0.011
CV score	0.775 $\pm$ 0.022	0.771 $\pm$ 0.021	0.775 $\pm$ 0.023	0.772 $\pm$ 0.019	0.781 $\pm$ 0.021
score	0.771 $\pm$ 0.01	0.77 $\pm$ 0.003	0.773 $\pm$ 0.005	0.781 $\pm$ 0.004	0.783 $\pm$ 0.006
Time to fit (s)	569.7 $\pm$ 2.04	584.2 $\pm$ 69.4	746 $\pm$ 45.4	1955.9 $\pm$ 1210.4	1573.5 $\pm$ 1215.0
			Different training/test sets used		

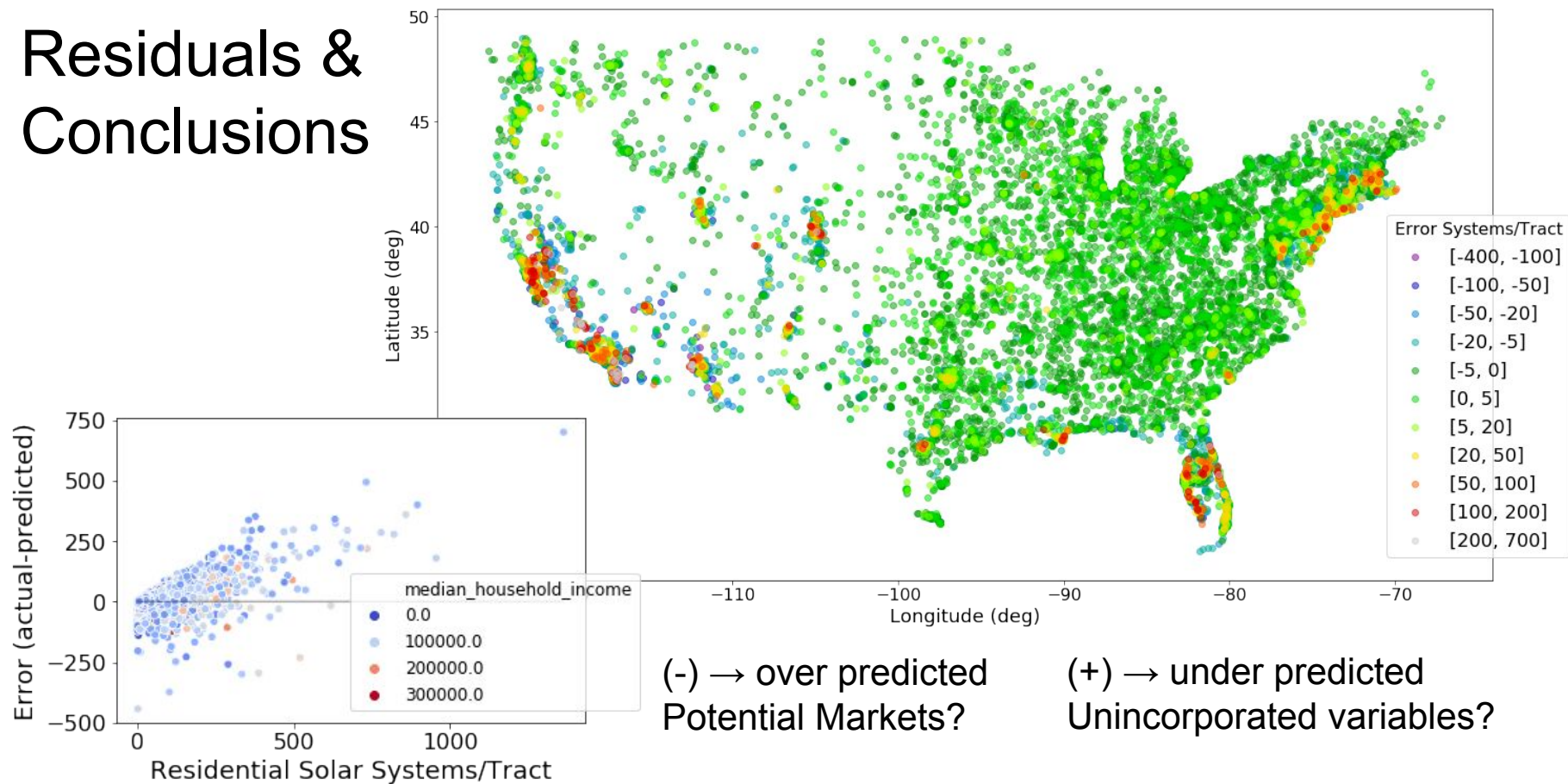
# Model Optimization & Performance

Used GridSearchCV to optimize performance of Gradient Boosting (~9 hrs)

	Option 1	Option 2
criterion	friedman_mse	mae
loss	ls	huber
max_depth	3	6
n_estimators	400	500
learning_rate	.01	.1
min_samples_split	2	4

Test score:  $0.77 \pm 0.007$   
CV score (n=6):  $0.775 \pm .023$   
Mean time to fit:  $723.3 \pm 113.3$

# Residuals & Conclusions



(-) → over predicted  
Potential Markets?

(+) → under predicted  
Unincorporated variables?



# What's next?

## **Solar companies & public utilities have stakes in residential solar predictions**

- Where to market solar?
  - Affinity Propagation to investigate groups of tracts that are similar to those that were overpredicted
- Where to prepare the grid for upcoming shifts in usage patterns?

## **Potential next steps:**

- Clustering data and training separate models for each cluster
  - Separating data based on thresholds of single variables did not yield interesting results, but clusters based on unsupervised learning might do better