



# Open Geospatial Machine Learning



Kevin Stofan



# Agenda

---

- Intro
- Spatial Data Models and Formats
- Spatial Data I/O
- Spatial Autocorrelation
- Spatial Weights Matrix
- Exploratory Spatial Data Analysis (ESDA)
- Rest of the Geospatial ML Workflow
- Related Studies and Competitions
- Discussion

# Goals

---

- Introduce geospatial machine learning workflow
- Exposure to Python spatial tools
- Spatially-explicit modeling
- External resources and further reading

# About Me

---

- Customer Facing Data Scientist at DataRobot
  - Pre- and Post- sales support for customers
  - Assist product and engineering teams with geospatial features
  - Consult customers with geospatial use cases
- Adjunct Professor at Penn State
  - Graduate level Geographic Information Systems (GEOG884)
  - Raster and vector data analysis
  - FOSS4G
- Applied Spatial Analysis
  - Point Pattern Analysis
  - Spatial Econometrics and geostatistics

# Repos and Contact

---

[https://github.com/TankofVines/data\\_intel](https://github.com/TankofVines/data_intel)

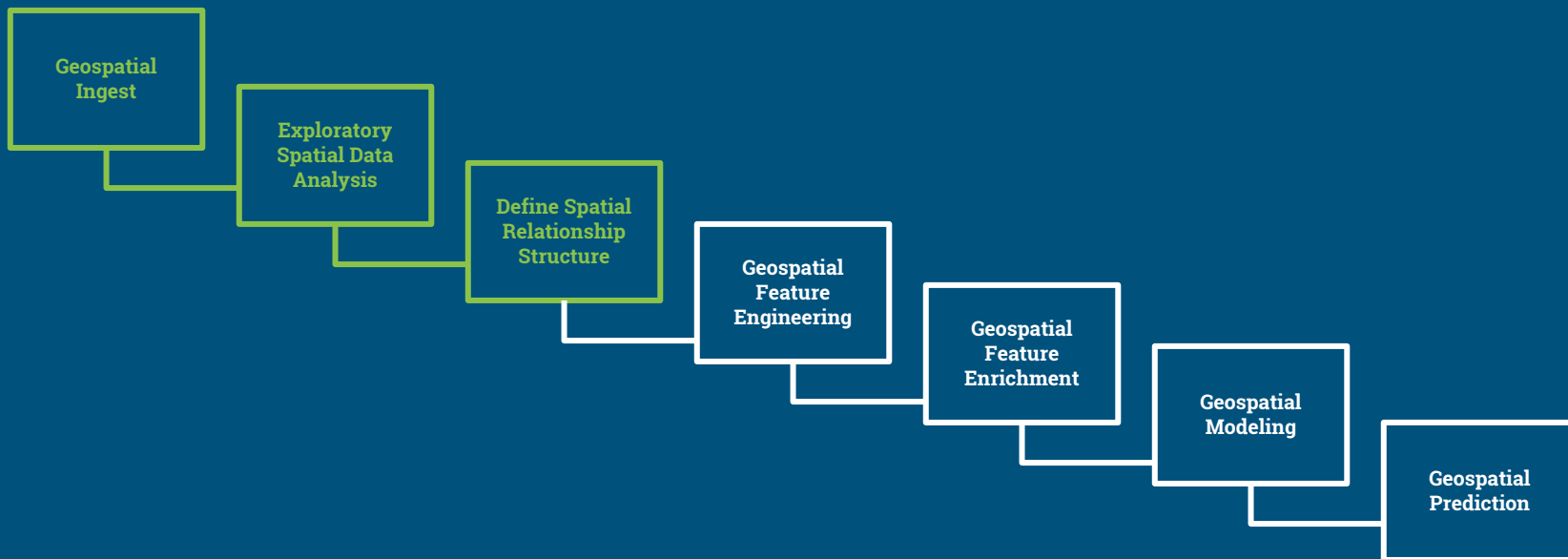
<https://github.com/TankofVines/odsc>

kevin.stofan@gmail.com

@tankofvines

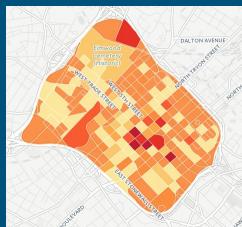
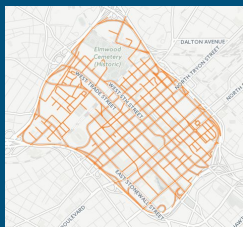
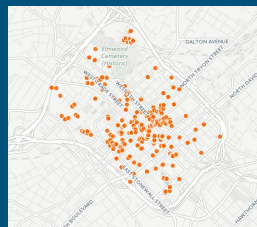
# Geospatial ML Workflow

---



# Geospatial ML Techniques

## Geospatial Data Ingestion



KML

SHP

JSON

GPX

WKT

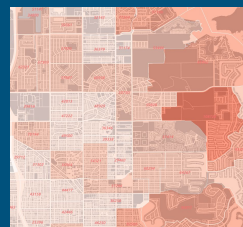
WBT

GML

## Spatial Feature Engineering/Enrichment



Spatial Imputation

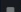






Feature Enrichment and Dasymetric Mapping



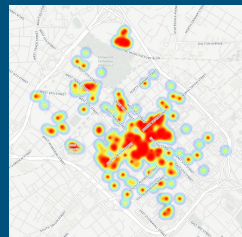
Derived features and transformations (e.g. area, centroid, contiguity, etc.)

## Spatially-explicit Models

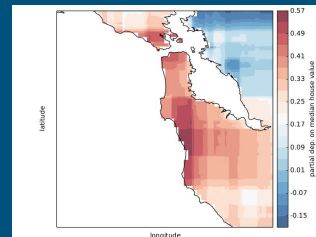
 Inverse Distance Weighting with Anisotropic Filtering <span>(BP145)</span> <span>(M106)</span>	Informative Features	64.0 %	0.6932	Run	🔒
<small>One-Hot Encoding   Matrix of word-grams occurrences   Missing Values Imputed</small>					
 Linear Regression with Spatial Lag <span>(BP147)</span> <span>(M110)</span>	Informative Features	64.0 %	0.6932	Run	🔒
<small>Converter for Text Mining   Auto-Tuned Word-In-Gram Text Modeler using token occurrences</small>					
 Empirical Bayesian Kriging <span>(BP148)</span> <span>(M111)</span>	Informative Features	64.0 %	0.6932	Run	🔒
<small>One-Hot Encoding   Missing Values Imputed   Standardizer</small>					
 Geographically Weighted Regression with Edge Corrections <span>(BP149)</span> <span>(M112)</span>	Informative Features	64.0 %	0.6932	Run	🔒
<small>Converter for Text Mining   Auto-Tuned Word-In-Gram Text Modeler using token occurrences</small>					
 Indicator Kriging (Aggressive Thresholds) <span>(BP144)</span> <span>(M113)</span>	Informative Features	64.0 %	0.6932	Run	🔒
<small>One-Hot Encoding   Missing Values Imputed   Standardizer</small>					

Spatial econometric, geostatistical, geographically weighted regression models

## Geospatial Visualization



Heatmaps, kernel density estimates, and hexagonal binning



Two-way partial dependence plots using coordinates

# A Note on Terminology

---



**courtney 8 claessens**

@sidewalkballet

 Follow

My friend is learning GIS and he keeps pronouncing it "jiss" and he won't listen to me and I want to burn this profession to the ground

12:30 AM - 5 Nov 2016

  5  30



# Tool: QuantumGIS

---

- Geographic Information Systems (GIS)
  - Collection - Maintenance - Storage - Analysis - Output - Distribution
  - Handles vector and raster data models
  - Historically dominated by ESRI and ArcGIS

# Spatial Data Models

---

Real World Entities

Vector Objects

Pixel-based Raster

# Spatial Data Models

---

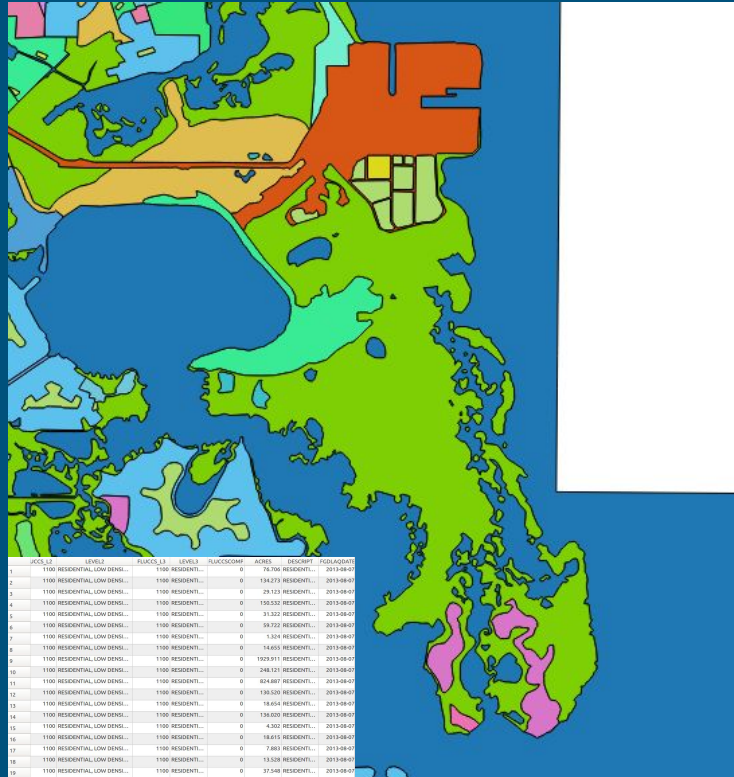


# Spatial Data Models

---



# Spatial Data Models



# Vector Data Types

---



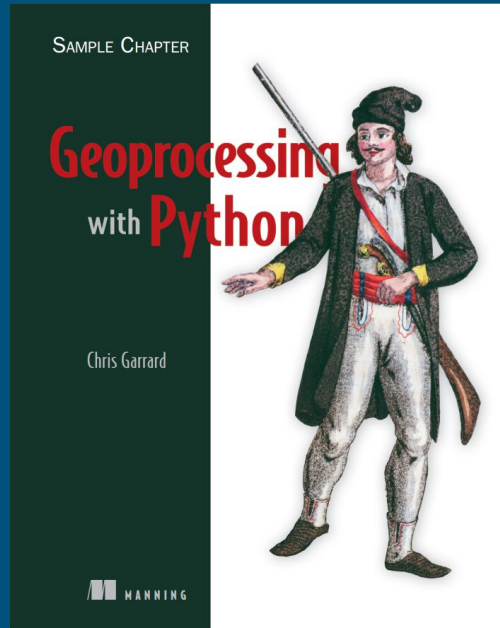
# Vector Data Formats

---

- File-based
  - ESRI Shapefiles
  - GeoJSON
  - File Geodatabase
- Database
  - PostGIS Table
  - SpatialLite
  - Various proprietary DBs (Oracle, Mongo, MS SQL Server)
- Binary/Text
  - WKT/WKB

# Tool: OGR/GDAL

---



<https://www.manning.com/books/geoprocessing-with-python#downloads>



# Tool: Geopandas

---

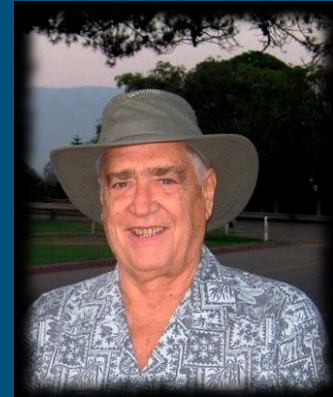
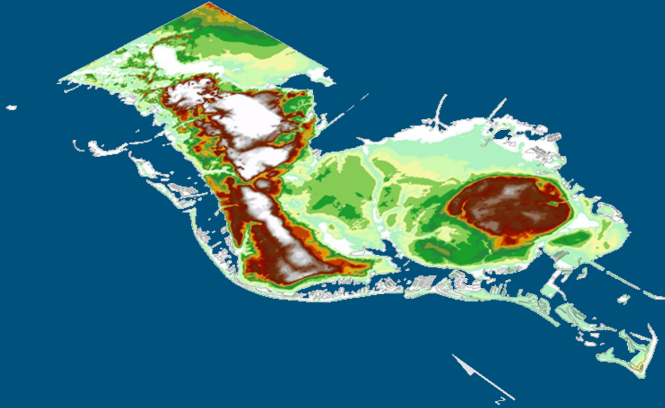
- Python Geospatial Library
  - I/O
  - Feature Engineering
  - Visualization

# Spatial Autocorrelation

---

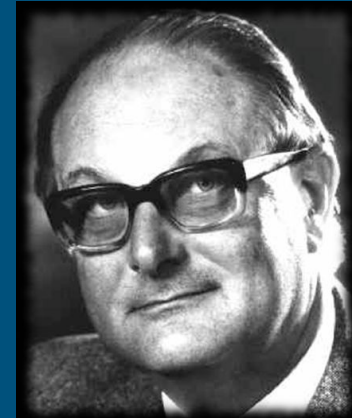
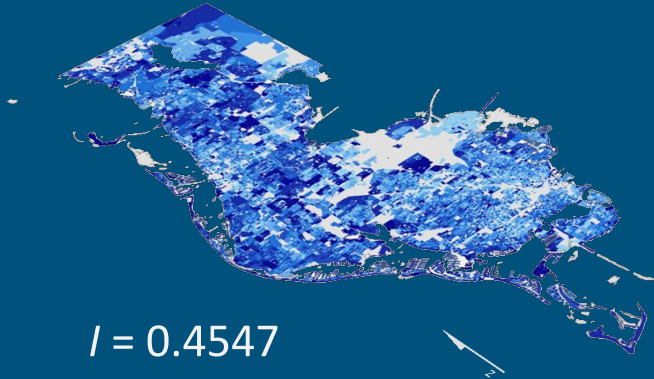
“...everything is related to everything else, but near things are more related than distant things.”

- Waldo Tobler



# Global Spatial Autocorrelation

$$I = \left[ \frac{\text{observations}}{\text{normalize for overall dataset variance}} \right] \times \left[ \frac{\text{weighted covariance}}{\text{observation total map spatial neighbor weight}} \right]$$



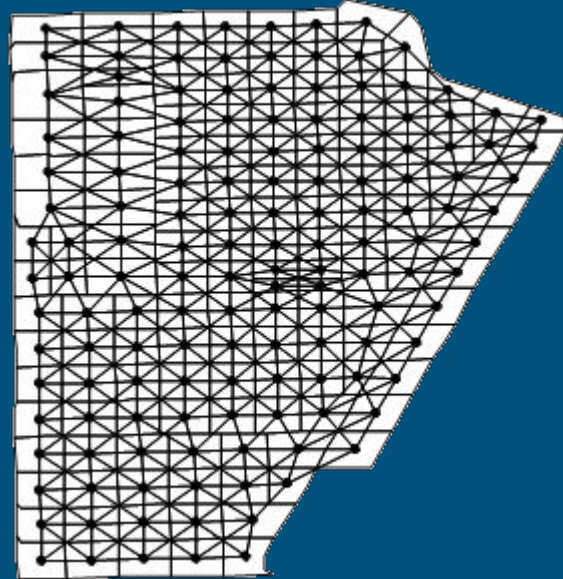
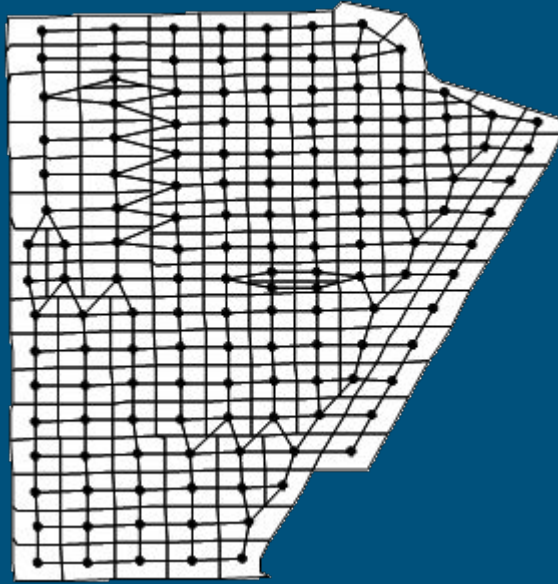
# Local Spatial Autocorrelation

---



# Spatial Weights Matrix

---



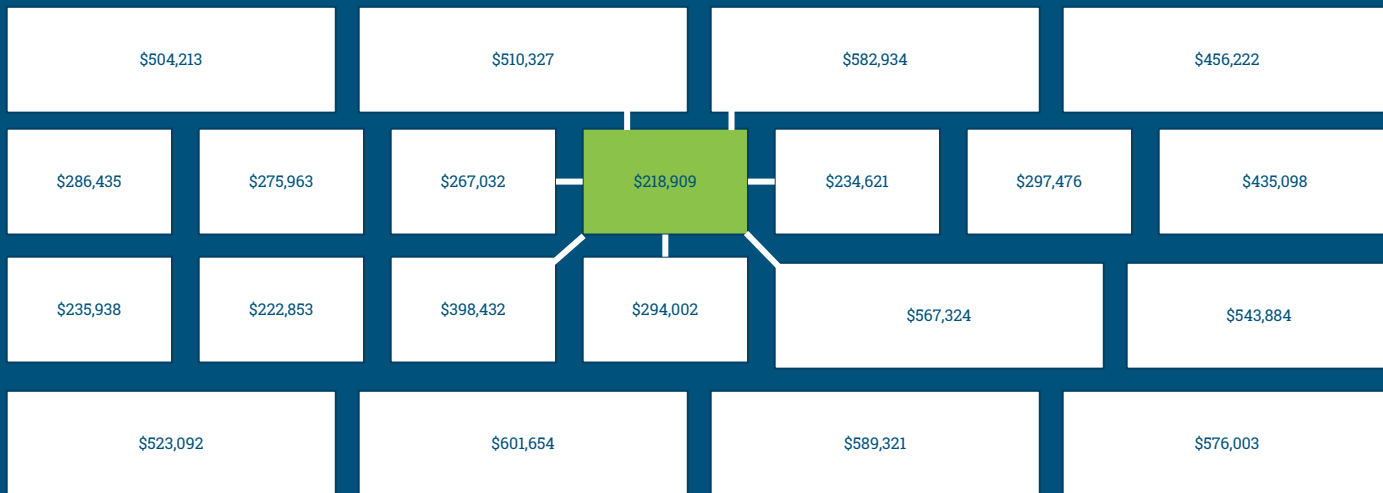
# Tool: GeoDa

---

- Exploratory Spatial Data Analysis (ESDA)
  - Desktop
  - Spatial Autocorrelation
  - Spatial Regression

# Spatial Lag

---



# Higher Order Weights

---

\$504,213	\$510,327	\$582,934	\$456,222			
\$286,435	\$275,963	\$267,032	\$218,909	\$234,621	\$297,476	\$435,098
\$235,938	\$222,853	\$398,432	\$294,002	\$567,324		\$543,884
\$523,092	\$601,654	\$589,321		\$576,003		



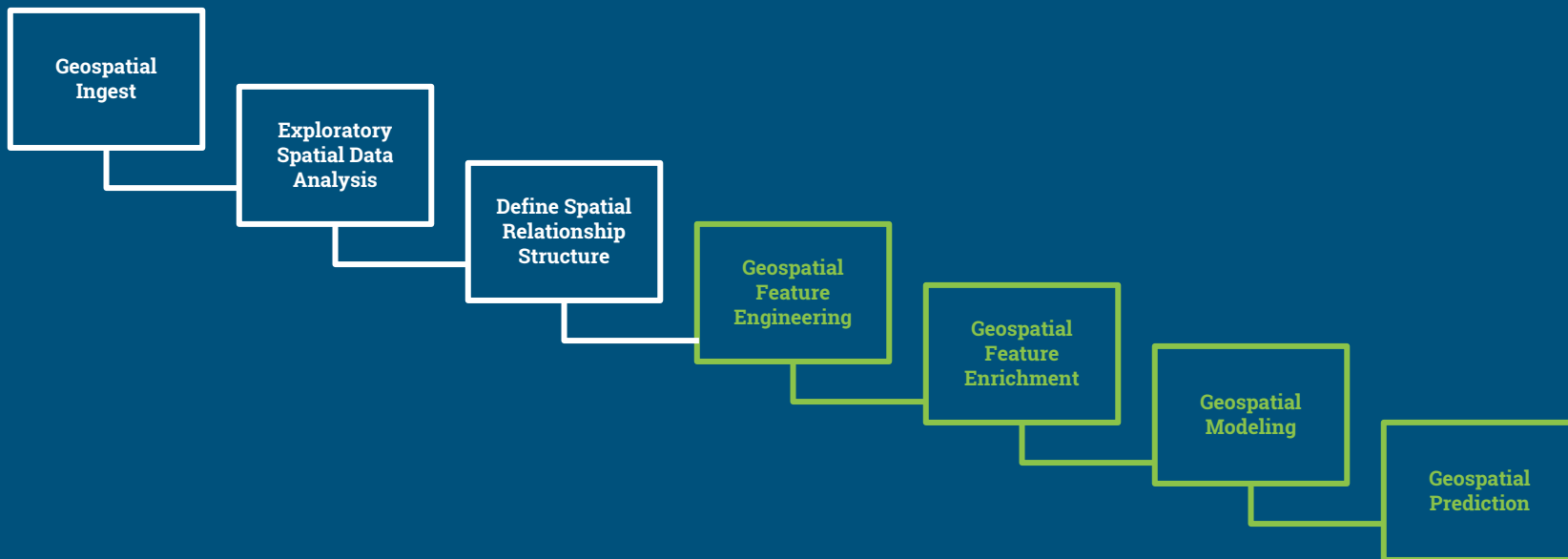
# Tool: PySAL

---

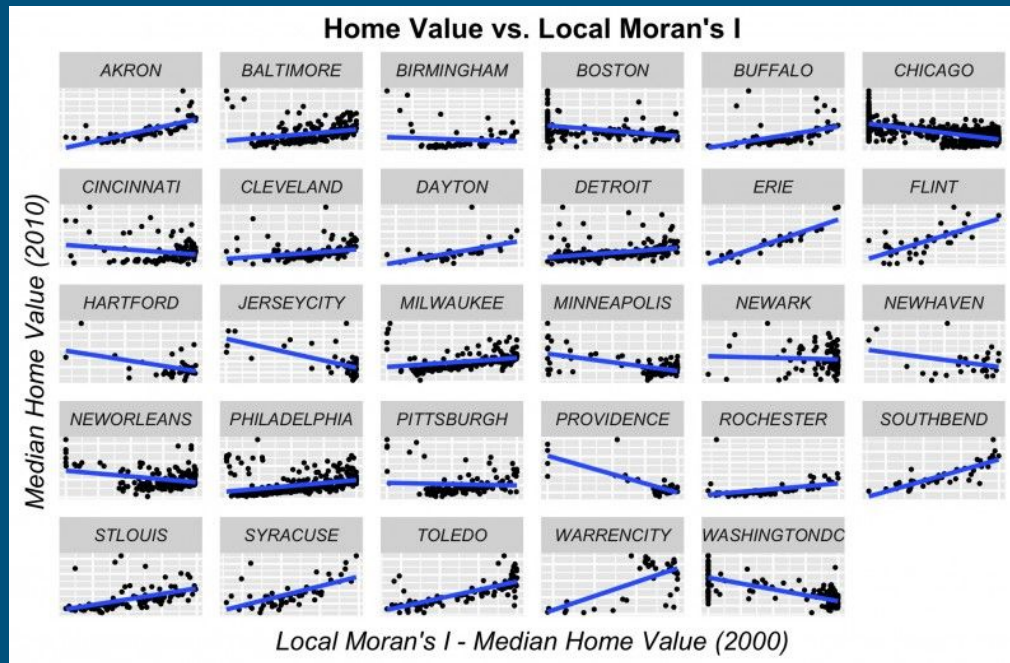
- Python Library for ESDA
  - GeoDa Equivalent
  - Spatial weights
  - Spatial lag, smoothing, regionalization, and more

<http://pysal.readthedocs.io/en/v1.11.0/#>

# Geospatial ML Workflow




# Further Reading: Urban Spatial




# Further Reading: Kaggle Zillow Zestimate

---

 Featured Prediction Competition

## Zillow Prize: Zillow's Home Value Prediction (Zestimate)

Can you improve the algorithm that changed the world of real estate?

 Zillow · 1,042 teams · 7 months to go (3 months to go until merger deadline)

OverviewDataKernelsDiscussionLeaderboardMore

My Submissions

Submit Predictions

**\$1,200,000**

Prize Money

<https://www.kaggle.com/c/zillow-prize-1>

# Discussion

---