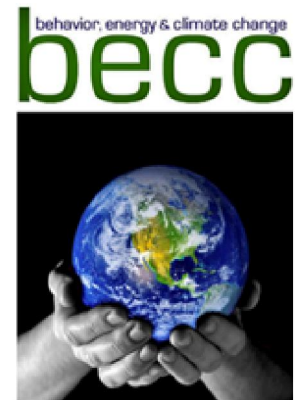


# Using Open Data to Predict Energy Usage

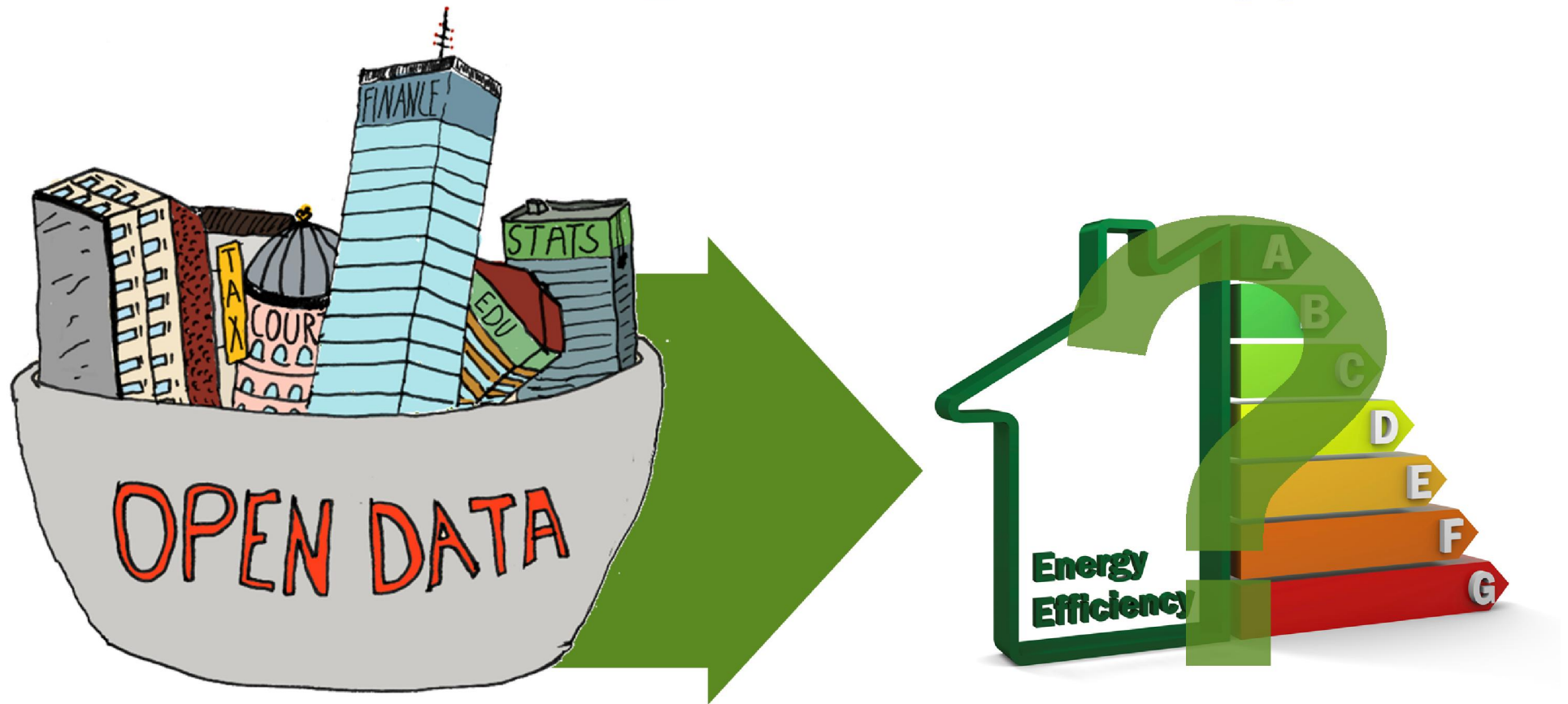
---

What tax lot data can tell us about energy usage intensity in New York City

*Theo Love*



# What can all this **data** tell us about how we use **energy**?



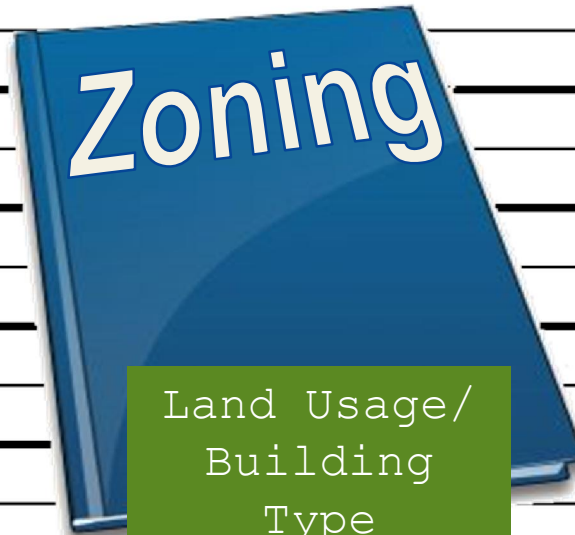
# It's the Usual Suspects

**WANTED**

For Explaining Energy Usage Intensity



Building  
Age



Land Usage/  
Building  
Type

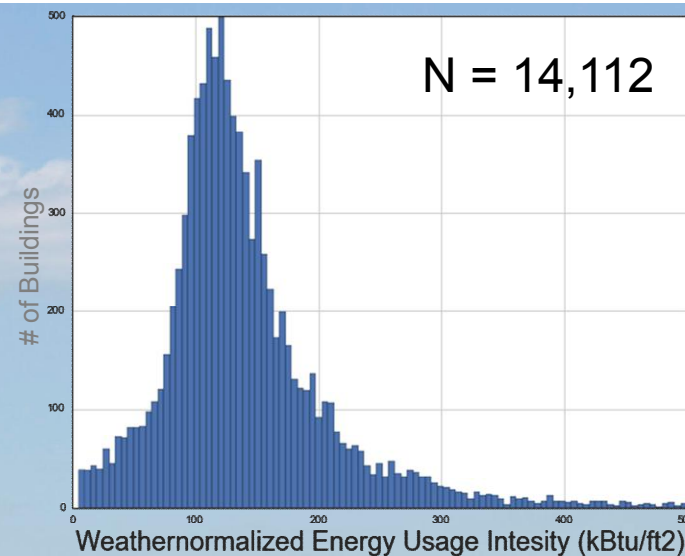


Building  
Value

But that's only **part** of the story...

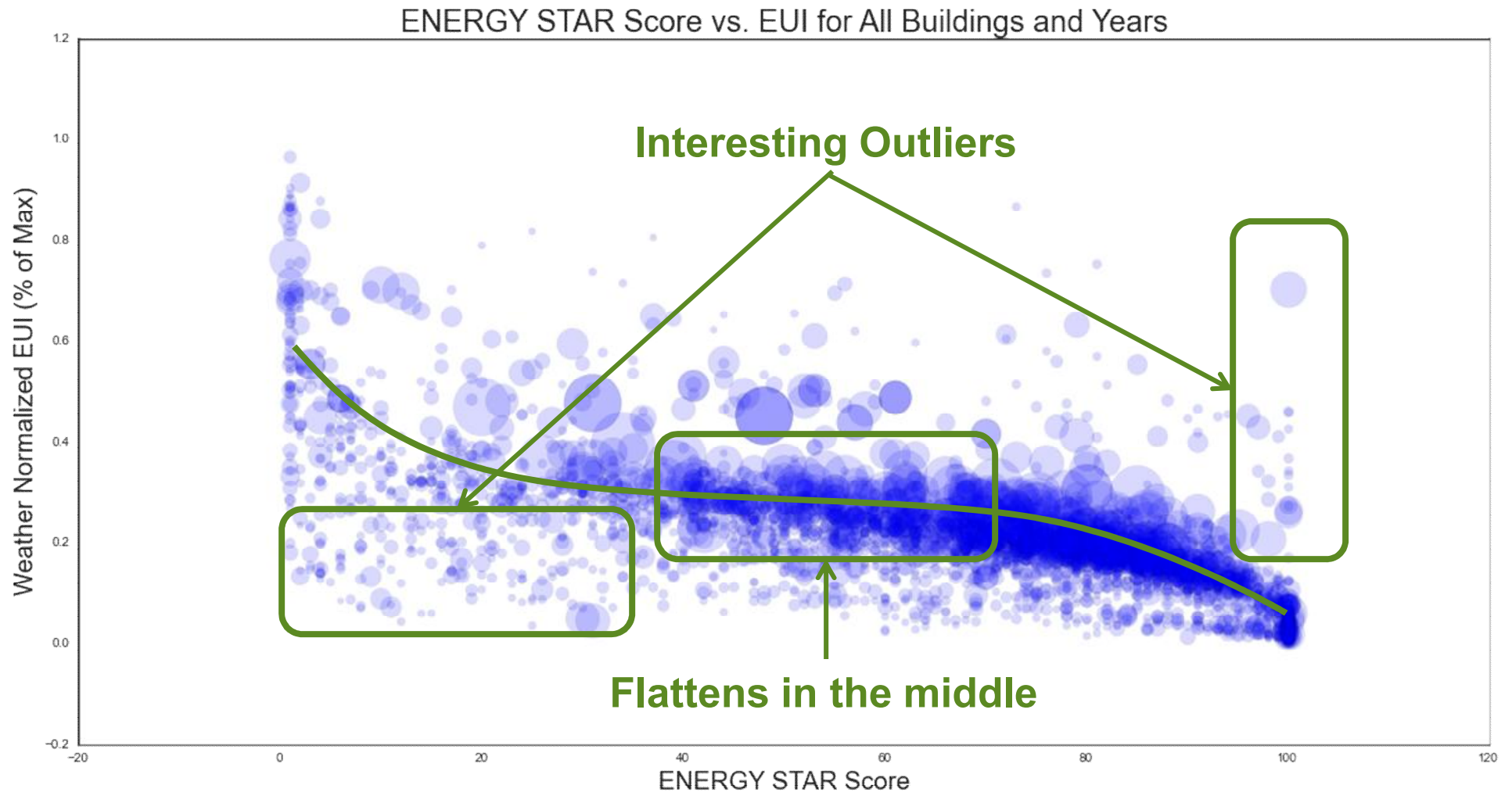
# NYC<sup>TM</sup>

## Local Law 84

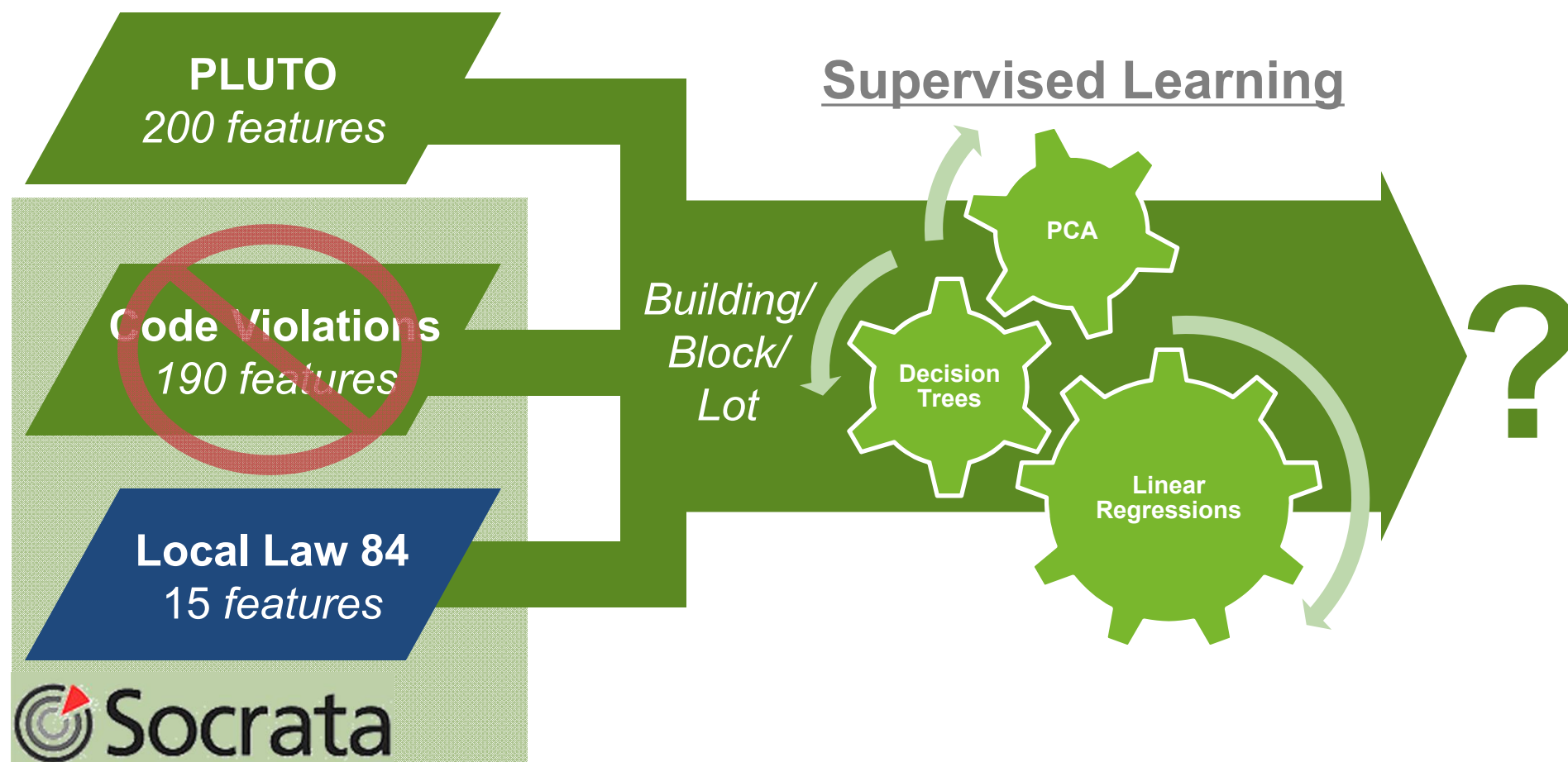




# Existing Benchmarks

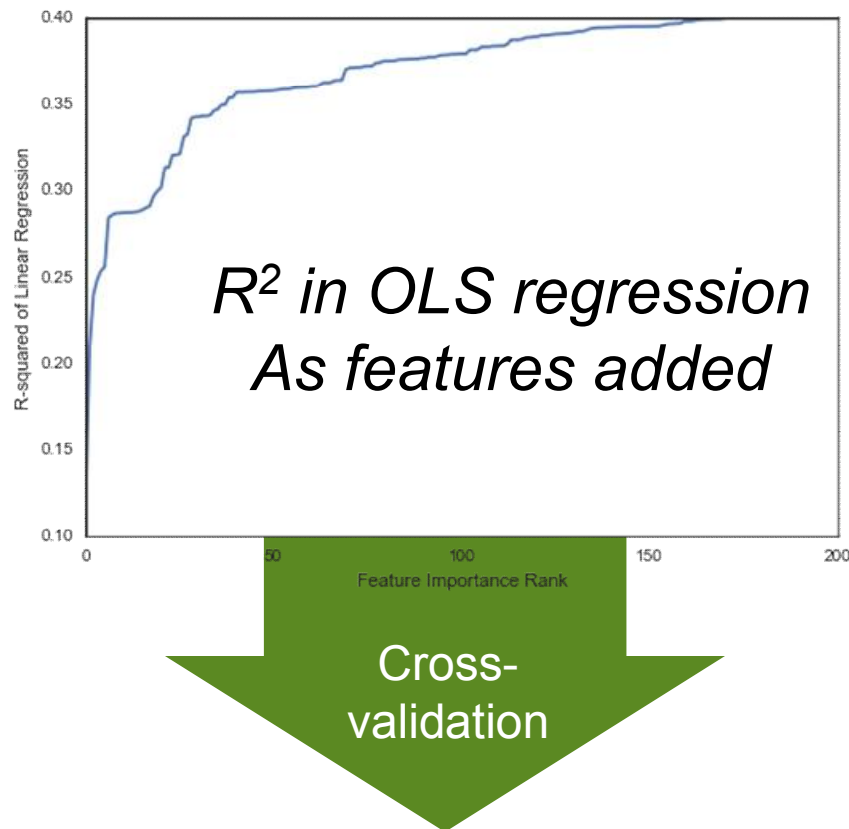


# Going Further with New Data



 **Socrata**

# What Was Found?



Gini Importance

## Tier 1 Features

- Is it a hospital?

## Tier 2 Features

- Is it an office?
- Year built

## Tier 3 Features

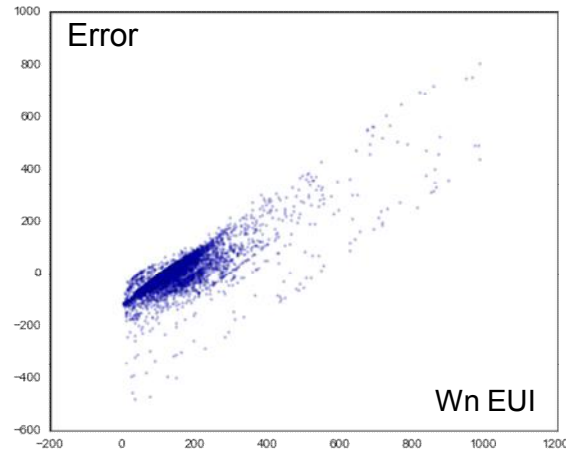
- Value of building/land
- Year of last renovation
- Building size/usage type

## Best Model: Random Forest Regression

$R^2$  from 23% to 29% with standard deviation ~3%

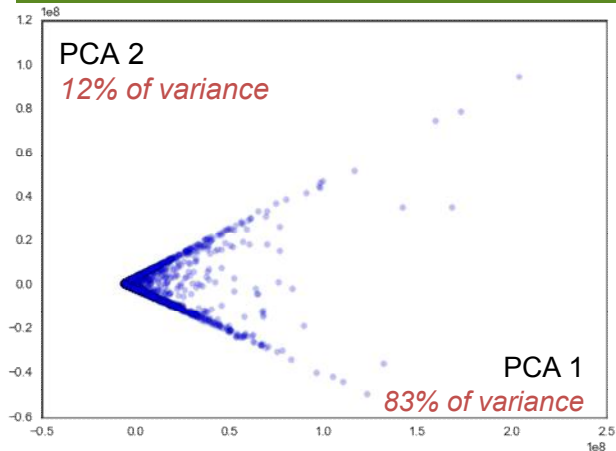
# Hints in the Data

## Residuals



- » Clear linear relationship
- » Further classification striations

## Primary Component Analysis



- » Collapses to two main features
- » First component explains vast majority of variance



# What Does it Mean?

- » Around **30%** of variance explained by features in PLUTO
- » **Sub-sector** analysis **crucial**
- » Type > Age > Value > Size

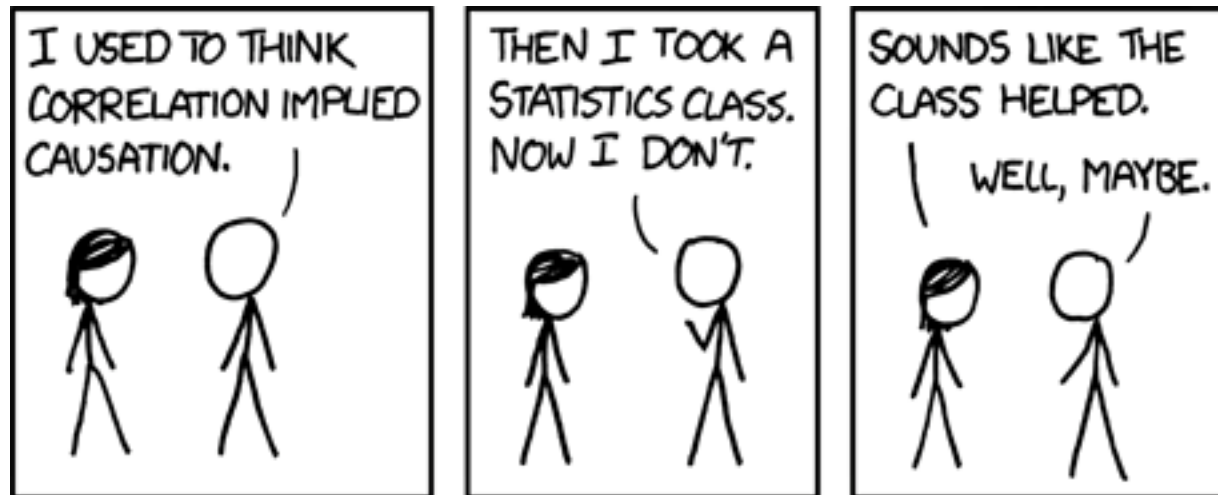
# What Now?

- » **White paper** with additional details
- » **Predicting** out of sample (new years)
- » Bring in **more data** to improve model
- » Applying and examining in **new areas**

***“Large repositories of **public data** can help **improve** our understanding of **energy usage**.”***

» Digging for **efficiency opportunities**

- Significant features could show hidden usage drivers
- Outliers may mean opportunity



© xkcd.com

---

# Theo Love

» [tlove@greenenergyeconomics.com](mailto:tlove@greenenergyeconomics.com)



behavior, energy & climate change  
**becc**

