

ASHRAE Kaggle Competition

Alissa Stover, Sophia Skowronski, Ying Hua

Introduction

Purpose of competition¹: create a baseline model for which to compare improvements buildings make to their energy efficiency

Data available for competition: building characteristics, meter readings, weather data

¹<https://www.kaggle.com/c/ashrae-energy-prediction>



Data Journey

What are factors that impact energy use? Can you predict it for particular buildings?

Data cleaning

Subset data from 20,216,100 rows → 12,060,910 rows
Merge datasets and clean → 12,060,311 rows

Data exploration and feature engineering

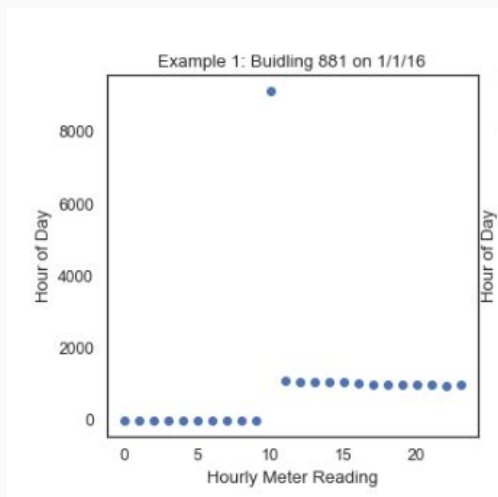
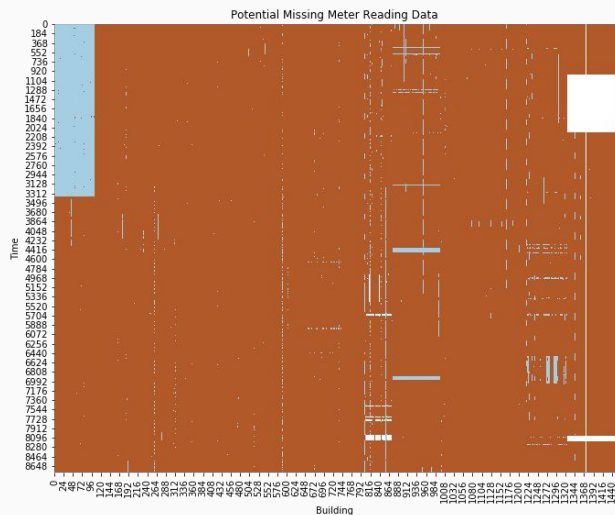
Univariate statistics, examine distributions, create new
variables (16 columns → 29 columns)

Data modeling

Examining linear relationships between variables &
evaluated models with Root Mean Squared
Logarithmic Error (RMSLE)

Data prep - Noticed data anomalies

Erroneous data: Some
disguised as 0 meter
reading; random spike in
energy reading



Data prep - Dealing with missing variables

Different
Treatments
for different
types of
missing data

Treatment of Missing Variables				
	Variable Name	% Missing	Treatment	Note
Outcome Variable	meter_reading	4%	Imputation using linear regression	<ul style="list-style-type: none">• No NA but we have reason to believe some data reported as 0 are erroneous• Also tried other methods of imputation including KNN, Naive Bayes. Linear regression were the most efficient to run and gave good results
Explanatory Variables	Building Variables			
	year_built	75%	Imputation using KNN	<ul style="list-style-type: none">• Also tried other methods of imputation including linear regression, Naive Bayes. KNN gave the best results
	floor_count	53%		
	Weather Variables			
	air_temperature	0.03%	Imputation using average of the values before and after; if NA, using backfill	<ul style="list-style-type: none">• Because temperature data is bounded by time and specific location, we think this method is most appropriate
	dew_temperature	0.08%		
	cloud_coverage	49%	Ignored	<ul style="list-style-type: none">• Most of the weather data did not have significant correlation with the response variable. As such, we priorities other variables to impute.• If we had more time we would impute these using a similar method for temperature
	precip_depth_1hr	36%		
	sea_level_pressure	7%		
	wind_direction	4%		
	wind_speed	0.2%		

Data prep - Timezone merge

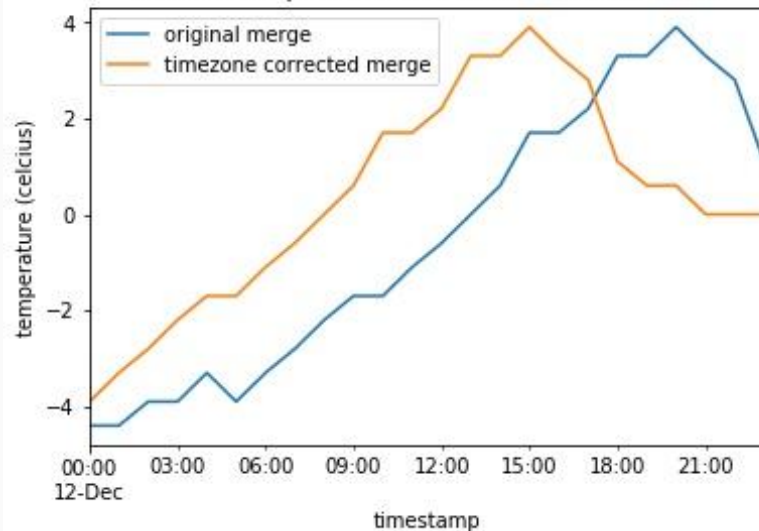
Meter data:

- Variables: building_id, meter, **timestamp**, meter_reading
- If merged w/ Building data: **site_id**

Weather data:

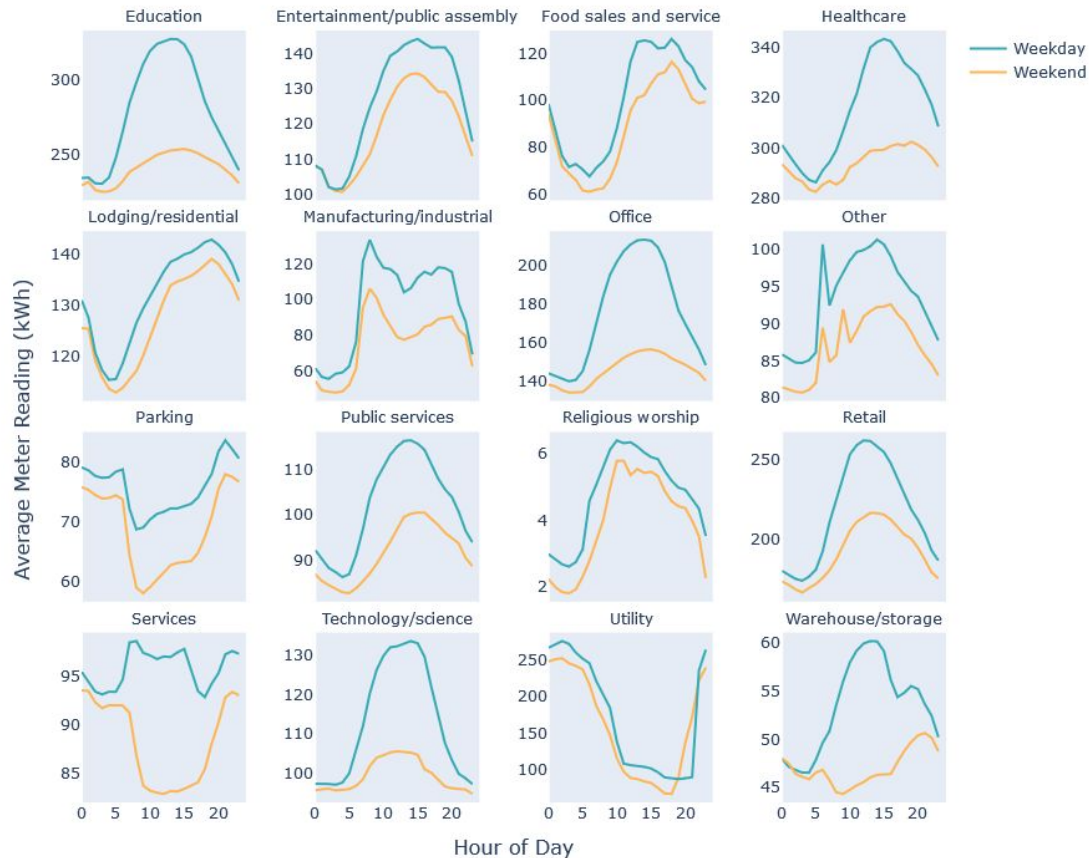
- Variables: **site_id**, **timestamp**, air_temperature, cloud_coverage, dew_temperature, precip_depth_1_hr, sea_level_pressure, wind_direction, wind_speed

Comparing daily air temperature patterns for different merges
Site 15 Temperatures on December 12th, 2016

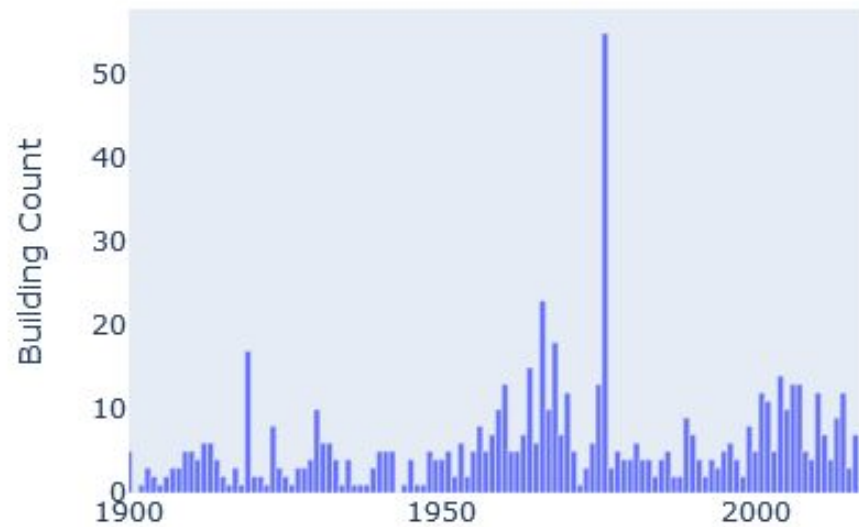
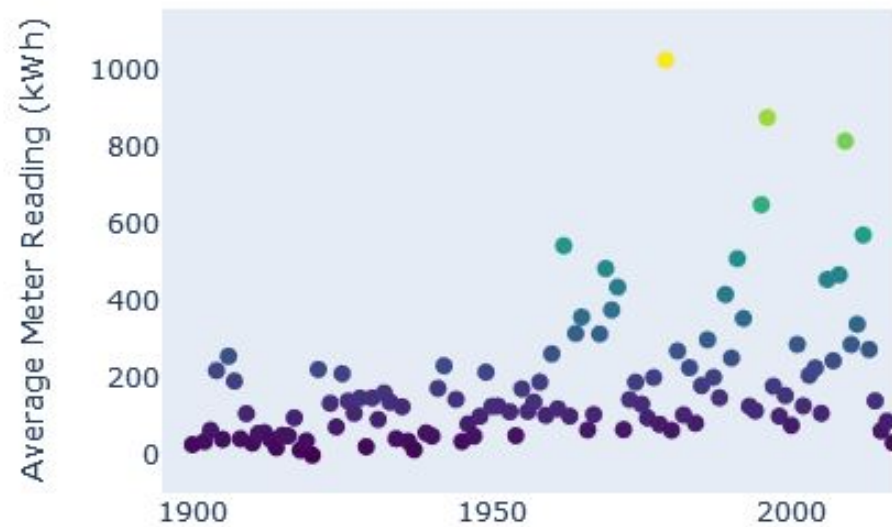


Interesting observations from EDA

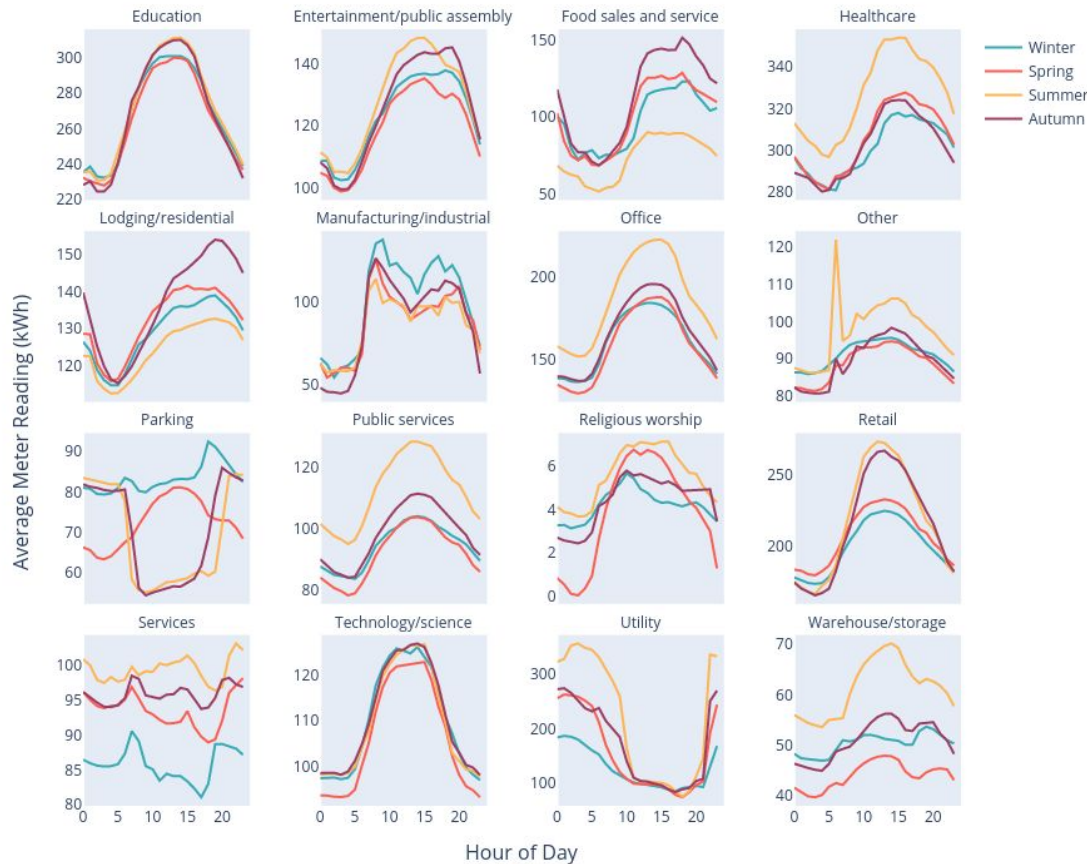
**Daily patterns vary
significantly between
primary uses**



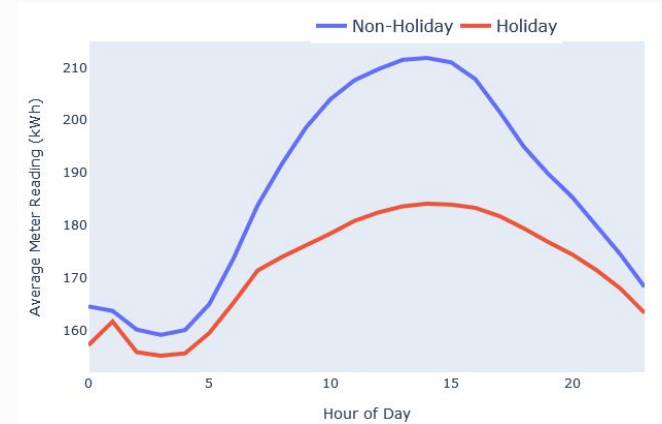
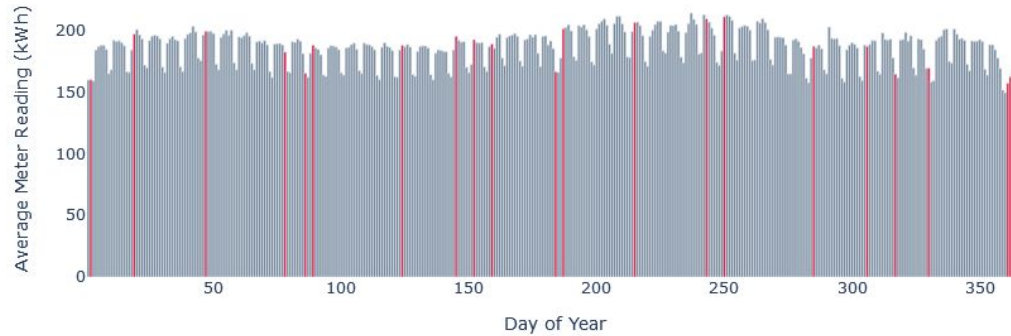
Interesting observations from EDA



Daily Trends by Season

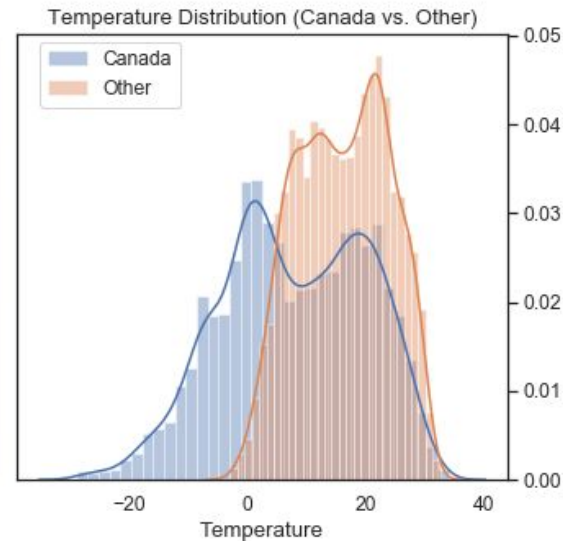
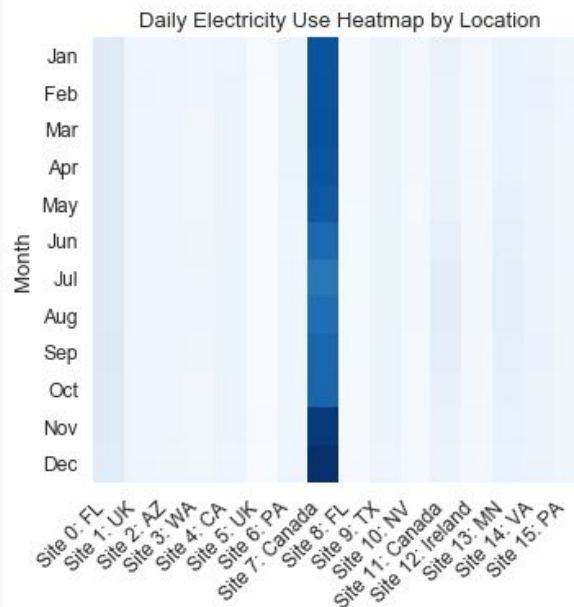


Holidays



Interesting observations from EDA

Cold weather
driving high energy
consumption



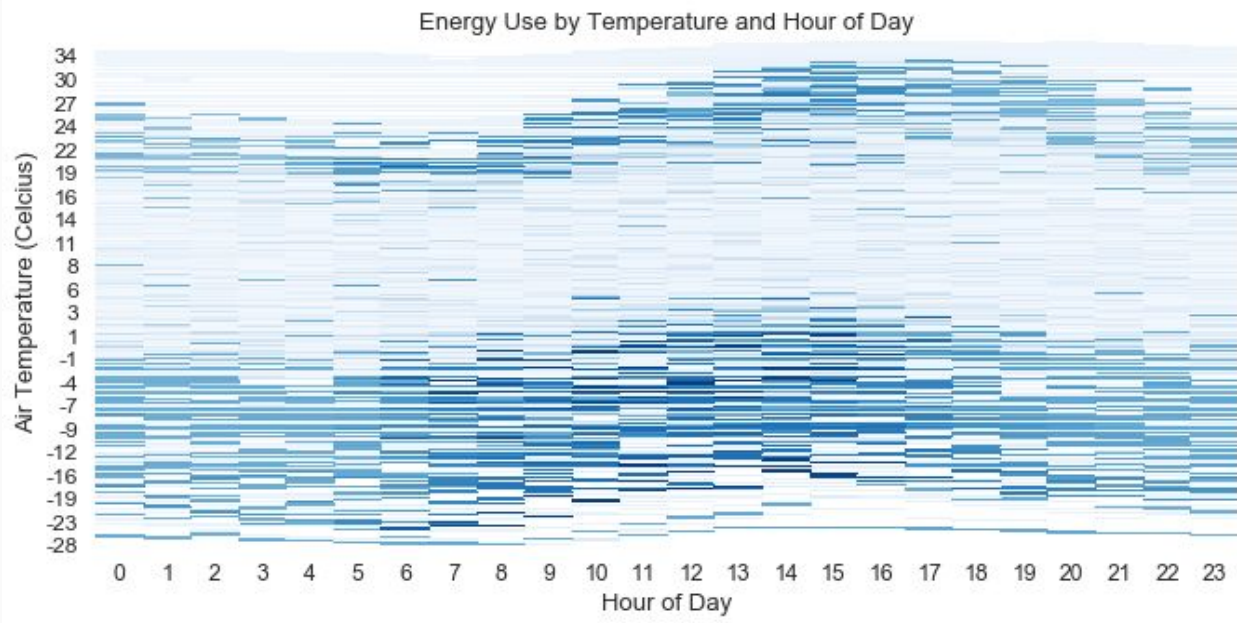
Interesting observations from EDA

Interesting patterns

when we look at

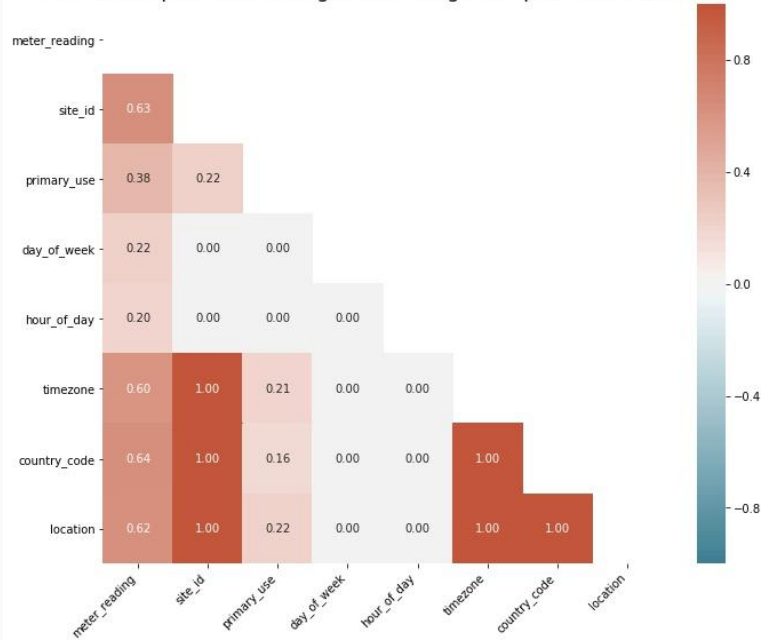
electricity use by

temperature & time

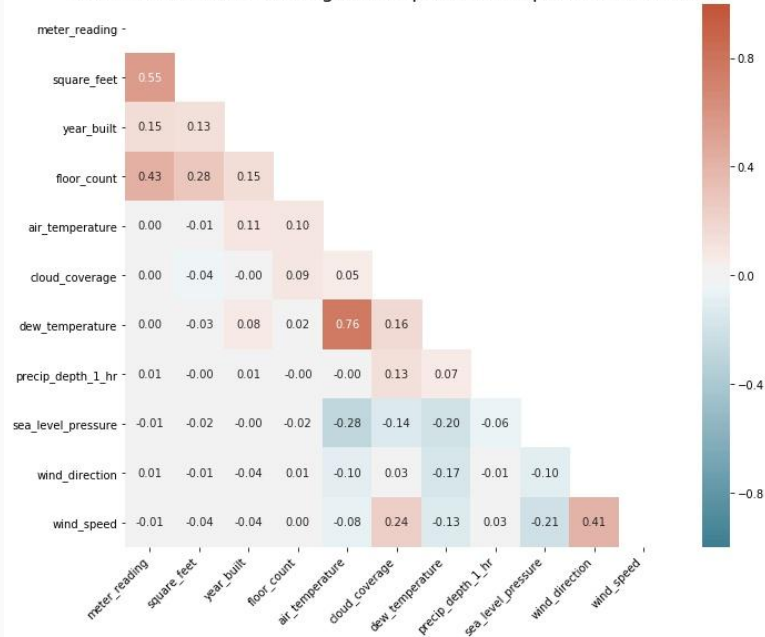


Model Building

Relationships between target and categorical predictor variables



Correlations between target and quantitative predictor variables



Model Building

3 predictive models and
good RMSLE scores vs.
Kaggle leaderboard
(0.93)

Type	Variables	Model 1	Model 2	Model 3
Meta	Square Ft	✓	✓	✓
	Floor Count	✓	✓	✓
	Primary Use	✓	✓	✓
	Site ID	✓	✓	✓
	Year Built			✓
Temporal	Day of the Week	✓	✓	✓
	Hour of Day	✓	✓	✓
Weather	Air Temperature		✓	✓
	Dew Temperature			
Results	RMSLE on test data	1.0133	1.0124	0.9824

Key takeaways

- What they say is true: data preparation accounts for 80% of work
- Datetime methods for time series
- Different plotting tools
 - Seaborn is pretty awesome!
 - Plotly is a good intro to interactive visualization
 - Matplotlib(a classic)
- Dabbled a bit into sklearn and ML
 - For imputing missing variables (KNN, Naive Bayes, Linear regression)
 - Building models (tried random forest but ran out of memory)

If we had more time...

- Examine and model all meter types (not just electricity)
- Examine in more detail weather data we did not include in our model (cloud coverage, sea level pressure, wind, etc)
- Try other model building techniques (higher order, ML models)
 - Attempted random tree but ran out of memory