



# Modeling and Predicting Energy Consumption

Kevin Pham (kpham123),  
Vikram Shanker (vshanker),  
Vikas Yadav (vikasy)

Stanford  
CS221 P-Poster

## Problem

Identifying and reducing energy consumption is one of many necessary avenues in tackling climate change. A large source of energy consumption goes into the heating, cooling, and powering buildings. The goal of our project is to build a model that can be used to accurately **predict the pre-improvement energy consumption of buildings with energy improvements based on their use case, size, age, and local weather information**. The predicted usage can be compared with the actual usage to understand the value added by specific energy improvements, thereby creating a greater understanding of how these energy improvements impact energy usage.

## Dataset

- The data set consists of three types of data: **(1) building information, (2) local weather, and (3) energy usage**. The building data, for 1449 buildings, consists of its location, area, age, floor count, usage type. The weather and energy usage data have various time-series over one year sampled at every hour, total size ~1.5GB. The data source is a Kaggle competition scored by RMSLE [1].
- The weather data has air and dew temperatures, wind speed and direction, cloud coverage and precipitation.
- The two data sets (building and weather) have **missing values**. We addressed this using standard methods of dropping and imputing as follows:
  - cloud coverage and precipitation columns are dropped as they are missing more than 50% of values
  - missing floor counts are replaced by 1 as that is the mode (more than 90%) and missing age was replaced by the median.

## Choosing Features

We performed EDA by computing **cross-correlation** of all of the features with each other. This showed that weather data has a much lower correlation with meter reading and our current models are deprioritizing and/or removing those features during training. The meter reading is highly correlated with building area and slightly correlated with floor counts.



Correlation matrix also shows that air temperature is highly correlated with dew temperature, and wind direction is correlated with wind speed, thus, the feature dimension is reduced by dropping dew temperature and wind direction. We convert our single categorical var, “primary\_use”, with one-hot encoder. The timestamp field of hourly measurements across 2016 is decomposed into “month”, “weekday”, “hour” fields. However, the weather patterns in December are quite similar to January but are not captured. So we use **trig functions** to map time features to a unit circle and define 2 features per time feature, such as “month\_sin”, “month\_cos”. Also, we ran into general issues due to dataset size (basic manipulation and MemoryError) - after importing, we change 64-bit data types to the smallest possible size in numpy that can still represent the information.

Your most recent submission

Name	Submitted	Wait time	Execution time	Score
submission_02-Dec-2019-05-09-18.7z	12 hours ago	0 seconds	265 seconds	1.50

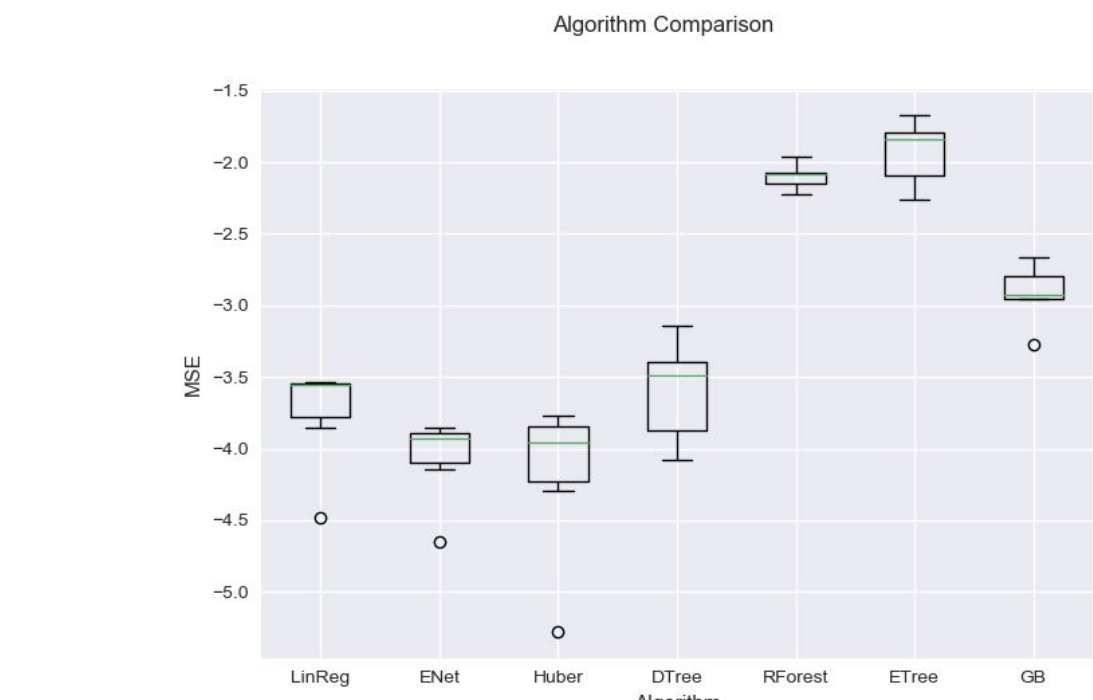
Complete

[Jump to your position on the leaderboard](#)

2444	Kevin Pham		1.50	1	12h
Your First Entry					

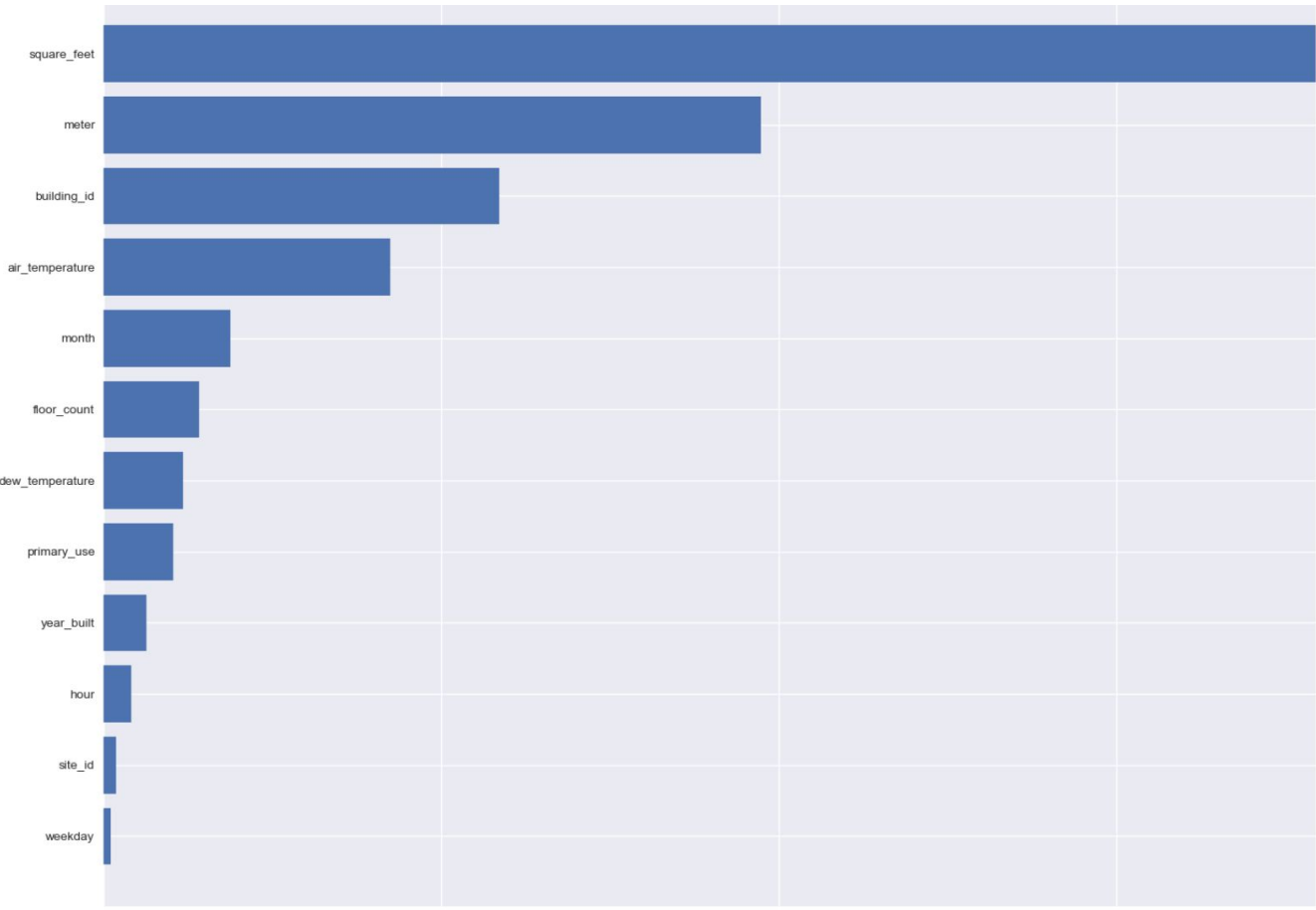
## Hybrid RF + Gradient Boosting

To explore models in SKLearn, a variety of linear and decision-based approaches were trained and fitted using default hyperparameters on a subset of the data and the different values of MSE after 5-fold cross-validation is shown in box-plots.



Support vector regressor was not considered due to its poor runtime performance on large number of data points. The basic linear regression, ElasticNet (combo of Ridge and Lasso by reducing model complexity and reducing number of features), and Huber (less susceptible to outliers) performs the lowest. A single decision tree is also mediocre but we see that the innovations of random forest, extra trees, and gradient boosting algorithms are promising. Random forests (and its cousin Extra Trees) are easier to train but are generally outperformed by modern boosting algos that are carefully tuned to avoid overfitting. To balance this tradeoff, the chosen approach is to train two concurrent models 1) Extra trees, 2) LightGBM and average their prediction.

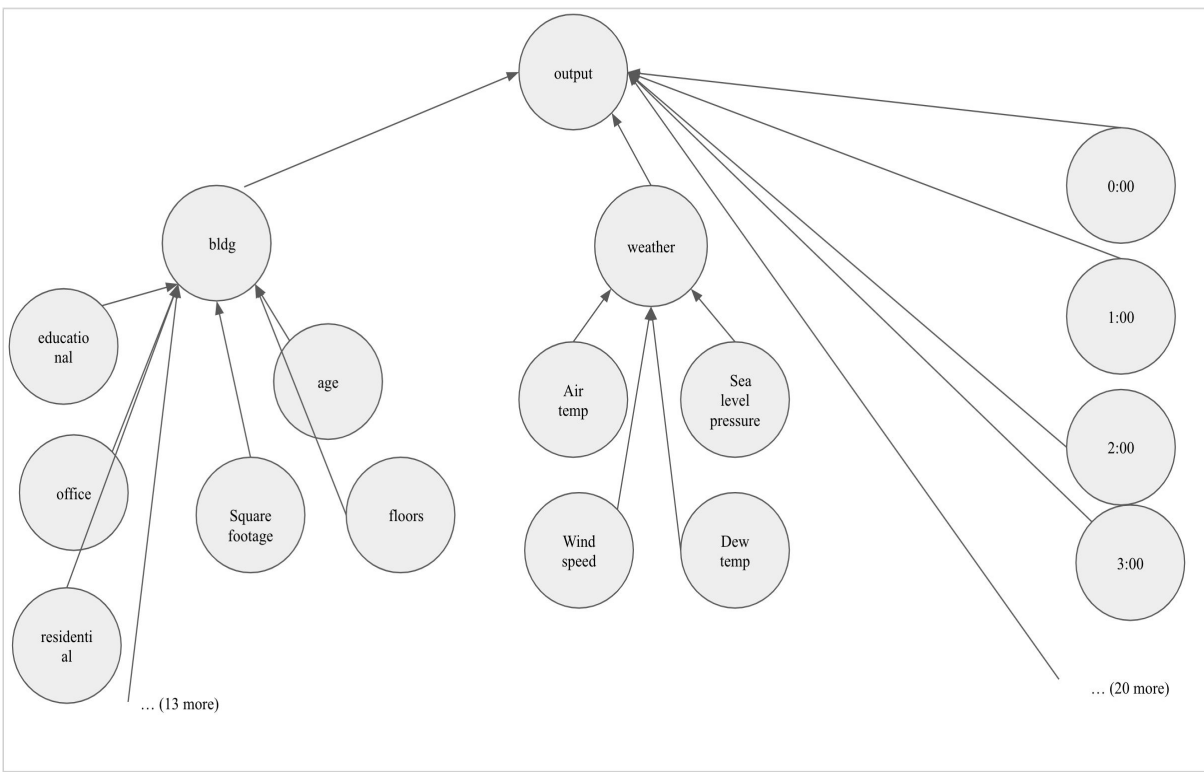
## (continued)



Gradient boost algos have seen recent progression and top proposed algorithms are XGBoost, LightGBM, and recently CatBoost. LightGBM is chosen due to fast training time as CatBoost seems to only outperform with very high numbers of categorical variable. Above, we graph the feature importances of the trained model, showing that meter type, building\_id, air\_temperature are prioritized. Of interest is square\_feet and building\_id are 1:1 uniquely mapped, so square\_feet possibly has high importance due to its large comparative magnitude and is a candidate for normalization. By using default parameters of ExtraTrees and LightGBM, the model was able to achieve RMSLE of 1.5 and at the time of writing, 2444/2933 on leaderboard. Training on Surface Book 2 took 2.25 hrs with n\_jobs=2, or using 2 cores that are hyperthreaded into 4 threads. Further work is necessary to explore parallelization and Google Cloud for training speedup.

## Neural Net

The neural network (still a work in progress) will use two topographies. The first is the following network topography. The fundamental idea behind it is that the features of the building that affect energy consumption are independent of the weather features. Finally, the time of day is also independent. Therefore, they are each grouped independently in their own hidden layer.



The second topography will be a fully connected topography. We will explore performance with different numbers of hidden nodes.

## Link to Video

[shorturl.at/euxlZ](https://shorturl.at/euxlZ)

## References

[1] <https://www.kaggle.com/c/ashrae-energy-prediction/data>  
[2] <https://www.energy.gov/eere/buildings/about-building-energy-modeling>  
[3] <https://arxiv.org/pdf/1607.06332.pdf>  
[4] [https://www.researchgate.net/publication/316653386\\_Modeling\\_energy\\_consumption\\_in\\_residential\\_buildings\\_A\\_bottom-up\\_analysis\\_based\\_on\\_occupant\\_behavior\\_pattern\\_clustering\\_and\\_stochastic\\_simulation](https://www.researchgate.net/publication/316653386_Modeling_energy_consumption_in_residential_buildings_A_bottom-up_analysis_based_on_occupant_behavior_pattern_clustering_and_stochastic_simulation)



Stanford  
University



# Next Steps

## More Data Exploration

- Identify outliers and/or trends that could give more features

## Gradient Boosting Model

- Explore sensitivity to: data scaling/normalization
- Hyperparameter search (built-in, RandomizedCV, 'hyperopt')
- Model training speedup (parallelization?)

## Neural Net Model

- finish implementation
- explore topographies (number of hidden nodes, whether it is fully connected)