

CUNY MSDS Capstone Project

---

**COMMERCIAL BUILDING ENERGY  
CONSUMPTION**

**ANALYSIS AND PREDICTION**

---

February 23, 2019

John Grando  
[john.grando@spsmail.cuny.edu](mailto:john.grando@spsmail.cuny.edu)

---

## Abstract

Commercial Building Energy Consumption accounts for approximately 25%<sup>1</sup> of the United States energy production profile. Many economical and sociological factors are pushing owners of these buildings to reduce energy consumption and optimize performance. However, it is difficult to say whether a building is operating efficiently or not. Using publicly available data, models can be constructed to predict major fuel consumption. Keywords: building energy consumption, predicted energy consumption, baseline energy model.

## Introduction

The concept of evaluating building performance typically requires two things; measuring the property in question's consumption, and comparing it to a standard practice equivalent. A baseline building consumption value is useful to inform people of expected operational performance; however, it is only useful if it is accurate. Additionally, there are benefits to using simple, available, predictors, which can make it easier to actually use the model.

Building owners, local governments, and utility providers are all looking for ways to reduce energy consumption. Reasons for doing so can vary all the way from social responsibility to economic gain. Some people want to show off an efficient building, others want to identify properties that are in need of improvement. However, in order to do this, a standard practice baseline value must be determined. Additionally, in order for the final product to be useful, it is important that the final set of predictors be parsimonious and be realistically available to users.

Comparing summary statistics between buildings, such as energy use per square foot, is not as simple as it seems because there are a multitude of factors that affect a building's energy consumption profile. The building use type can cause energy use to vary by a large amount; such as office buildings and refrigerated warehouses. Also, seemingly small factors, such as the hours of operation, may have significant impacts as well. This complexity of making similar comparisons creates a situation where it is difficult to determine whether a building is performing consistent with, or better than, other standard practice buildings.

Commercially, using the most popular example, ENERGY STAR<sup>2</sup> has implemented a benchmarking algorithm that scores buildings on a scale from 1 – 100 using market-available data. The output of this benchmarking algorithm is a unit-less score, as well as a reference 'baseline' building; however, the methodology is not released and it is unclear what factors are important to influence the energy consumption of the building. These barriers make it difficult to provide custom comparisons and nearly impossible to make batch predictions from a set of buildings, or variations of buildings.

Every few years, the U.S. Energy Information Administration (EIA) conducts a survey attempting to record pertinent features of these buildings, known officially as the Commercial Buildings Energy Consumption Survey (CBECS)<sup>3</sup>. While the survey is expansive

---

<sup>1</sup>EIA - [https://www.eia.gov/energyexplained/index.php?page=us\\_energy\\_commercial](https://www.eia.gov/energyexplained/index.php?page=us_energy_commercial)

<sup>2</sup>ENERGY STAR - <https://www.energystar.gov/>

<sup>3</sup>EIA Microdata

---

(i.e. more than 600 tracked features), it is essential to create a model that is usable and only requires predictors that can easily be attained by building operators. Therefore, series of models will be evaluated in order to determine the most important predictors which will then be used to train a final predictive model. Some iterations of specific building attributes (i.e. primary use) will be performed if there appears to be a disproportionate amount of error. After completion of the model, the predictors will be evaluated on how easy they are to attain, and will possibly be exchanged with simpler variables that are highly correlated.

## **Related Work**

The idea of determining building energy efficiency is not a novel concept in itself. As previously mentioned, ENERGY STAR has a building benchmarking tool<sup>4</sup>. Additionally, the United States Green Building Council has created the Arc Platform<sup>5</sup> which provides benchmarking and active monitoring features. While these platforms provide building comparisons in the form of an overall score, it is difficult to explore the space around the building attribute inputs themselves as well as compare consumption to randomly sampled buildings. With this functionality, a more direct comparison can be made and relative environmental impact can be measured.

## **Problem Statement**

Commercial buildings are complex and encompass a wide variety of purposes. However, they all must be powered, and require a considerable amount of energy to operate properly. There are a variety of texts that are dedicated to the analysis of building energy consumption, and determining operating efficiency, such as ASHRAE Guideline 14, etc. etc..... Particularly, there is a well thought out process for auditing commercial buildings, known as ASHRAE Audits, which start at the lowest level (I) and progress to the highest level (III) as the opportunity for energy and cost savings becomes more apparent (<http://aea.us.org/3143-2.html>). As part of the initial audit process, an assessment of the building's overall operational efficiency is gauged. Typically, an auditor will walk through the building, analyze utility bills, and make the closest comparisons they can based on experience. This takes years of experience and sometimes requires highly tuned spreadsheets that have been developed over years. It can take a surprising amount of time just to determine if a building is operating efficiently or not.

The CBECS data set provide some insight as to what building attributes most greatly affect building energy consumption. Over XXX survey questions are recorded and coupled with major fuel consumption. These fuel sources are Electricity, Natural Gas, District Heat, Fuel Oil, and ... However, it would not be useful to construct a model with a large number of predictors, as it would require a large amount of time and effort to compile the necessary information in order to provide a prediction. Therefore, one of the main focuses for this study will be to extract the fewest amount of predictors necessary in order to make accurate predictions.

---

<sup>4</sup><https://www.energystar.gov/buildings/about-us/how-can-we-help-you/benchmark-energy-use/benchmarking>

<sup>5</sup><https://arcskoru.com/>

---

Given the complex nature, it is unlikely a linear regression will provide the best prediction accuracy. This point is especially highlighted by the fact the the goal of this study is produce a parsimonious set of predictors, which means a small subset must be selected. Therefore, an investigation into more complex, nonlinear, algorithms will be performed in order to keep the number of necessary predictors as low as possible while still capturing complex interactions.

## **General Process**

Given the large number of features in the survey responses, it is not possible to analyze each one individually. Therefore, the first steps in the process will be centered around selecting a smaller subset. A few algorithms will be used in order to try and reduce bias. First, a partial least squares model will be fit to the response. Second, a random forest regression tree will be used in order to try and extract any nonlinear relationships. Third, an attempt to construct a lasso regression model will be made. Finally, a backward selection linear model will also be fit in order to see if an automated approach can be taken. The magnitude and contribution percentage of each variable will be considered in selecting features from this model. Also, the various error rates from each preliminary model will be used as a benchmark for the final model performance.

After the preliminary set of models have been run and summarized, the extracted variables will be analyzed in order to verify their importance, gauge their potential predictive power, and to check whether they are easily attainable for a building operator/owner. This step is very important because it is essential worthwhile variables are used to predict the outcome. Selecting a variable that, for one reason or another, is erroneous may lead to reduced predictive power in the final model. If a variable did pass our initial analysis but doesn't actually have much predictive power (i.e. it only changes values by a slight amount) then it may not be worthwhile to select it at all. All selected variables increase the complexity of the model; therefore, we wish to only select those that will matter. Finally, the predictor must be usable, and 'knowable'. Ultimately, this tool will not be usable if a very difficult and hard to understand, and/or attain, variable value is used. These three concepts will be used in the analyzation of the candidate variables from the the preliminary analysis.

Finally, a neural network model will be built to take the verified subset of features and make predictions for the selected major fuel use. A variety of hyperparameters will be tested, using cross-validation, and compared on a common error metric. This step will reveal the optimal hyperparameter combination to use for the model. The prospective model will then be retrained on the entire training and validation data. This model's selected error metrics will then be compared to the preliminary models, which should be considered a floor. Next, the model's metrics will compared to a similarly trained neural network model that uses every available predictor, which should be considered a competitor. Once this is done, the model can then be re-assessed for feature selection as well as analyzed for the value/tradeoff of adding/removing certain features.

---

## **Data Pre-Processing**

Given

### **Electricity**

**Variable Selection**