

FORECASTING CUSTOMER'S ENERGY DEMAND USING MACHINE LEARNING

SAIFUL ABU

Department of Computer Science

APPROVED:

Christopher Kiekintveld, Chair, Ph.D.

M. Shahriar Hossain, Ph.D.

Paras Mandal, Ph.D.

Charles Ambler, Ph.D.
Dean of the Graduate School

©Copyright

by

Saiful Abu

2016

to my

MOTHER and FATHER

with love

Acknowledgements

I would like to express my deep-felt gratitude to my advisor, Dr. Vladik Kreinovich of the Computer Science Department at The University of Texas at El Paso, for his advice, encouragement, enduring patience and constant support. He was never ceasing in his belief in me (though I was often doubting in my own abilities), always providing clear explanations when I was (hopelessly) lost, constantly driving me with energy (*Where does he get it?!*) when I was tired, and always, *always* giving me his time, in spite of anything else that was going on. His response to my verbal thanks one day was a very modest, “It’s my job.” I wish all students the honor and opportunity to experience his ability to perform at that job.

I also wish to thank the other members of my committee, Dr. Luc Longpré of the Computer Science Department and Dr. Mohamed Amine Khamsi of the Mathematics Department, both at The University of Texas at El Paso. Their suggestions, comments and additional guidance were invaluable to the completion of this work. As a special note, Dr. Longpré graciously volunteered to act as my advisor while Dr. Kreinovich was working abroad in Europe. He was extremely helpful in providing the additional guidance and expertise I needed in order to complete this work, especially with regard to the chapter on NP-hard problems and the theory of NP-completeness.

Additionally, I want to thank The University of Texas at El Paso Computer Science Department professors and staff for all their hard work and dedication, providing me the means to complete my degree and prepare for a career as a computer scientist. This includes (but certainly is not limited to) the following individuals:

Dr. Andrew Bernat

He made it possible for me to have many wonderful experiences I enjoyed while a student, including the opportunity to teach beginning computer science students the basics of UNIX and OpenWindows (something I wish I had been taught when I first started), and the ability to present some of my work at the University of Puerto Rico, Mayagüez Campus.

Dr. Michael Gelfond

His influence, though unbeknownst to him, was one of the main reasons for my return to UTEP and computer science after my extended leave from school while island hopping in the navy. He taught me many things about computer science—and life. Among the many things he showed me was that there really is science in computer science.

And finally, I must thank my dear wife for putting up with me during the development of this work with continuing, loving support and no complaint. I do not have the words to express all my feelings here, only that I love you, Yulia!

NOTE: This thesis was submitted to my Supervising Committee on the May 31, 1996.

Abstract

Accurate electricity demand forecasting is an important problem as the failure to do so may be costly for both economic and environmental reasons. Power TAC simulation system provides a no risk platform to do research on smart grid based energy generation and distribution. Brokers are important components of the system. The brokers work as self-interested entities that try to maximize profits by trading electricity in various markets. To be successful, a broker has to forecast the electricity demand about its customers as accurately as possible, otherwise it will operate ineffectively. This proposed forecasting method uses a combination of cluster and classifiers. At first, the customers are clustered based on their weekly average usage. After that, energy usage history and related weather related information are combined together to train classifier for the cluster. To forecast for a new customer, the proposed method needs at least a week's energy usage history of the customer. The system assigns the new customer to one of the clusters based on its electricity usage history. The classifier for that cluster will be used to forecast the customer. This approach produced 13 % error compared to 31% relative absolute error observed against the moving average baseline predictor. The Power TAC system has six different types of customer such as customers with demand shifting capabilities, customer with no demand shifting capabilities, electric vehicles, thermal storage, wind production and solar production. Previous approaches to demand forecasting treated all types of customers equally. This work shows that a good forecasting system should treat customers of different type differently, otherwise the system will experience more error.

Table of Contents

	Page
Acknowledgements	iv
Abstract	vi
Table of Contents	vii
List of Figures	ix
Chapter	
1 Smart Grid and PowerTAC Competition	1
1.1 Traditional Electricity Distribution and Consumption System	1
1.2 Smart Grid	1
1.3 Smart Grid and Renewable Energy	2
1.4 Importance of accurate load forecasting	2
1.5 Power TAC System	2
1.5.1 Broker	2
1.5.2 Wholesale Market	2
1.5.3 Customers	3
1.5.4 Balancing Market	3
1.5.5 Weather Service	3
1.5.6 Tariff Market	3
2 Related Works	4
2.1 Variables in Electricity Demand	4
2.2 Electricity Load Forecasting Using Statistical Method	5
2.3 Load Forecasting using Machine Learning	5
2.4 Load Forecasting for a Specific Region	5
2.5 Load Forecasting using Clustering	5
2.6 Expert System based Load Forecasting	6
3 Customer Description	7
3.1 Customers	7
3.2 PowerTypes	7
3.2.1 consumption	7
3.2.2 Interruptible Consumption	7
3.2.3 Thermal Storage	8
3.2.4 Solar Production	8
3.2.5 Wind Production	9
3.2.6 Electric Vehicle	9
3.3 Statistics	9
3.3.1 Customer Vs PowerType	10
3.3.2 Population Vs PowerType	10
3.3.3 Total Energy Consumed Vs PowerType	10
4 Methodology and Result	13
4.1 The Baseline Electricity Forecasting Mechanisms	13

4.2	Proposed Electricity Demand Forecasting Mechanism	13
4.2.1	Demand Forecasting for Consumption Type Customer	14
4.3	Result	17
4.3.1	Finding number of clusters	17
4.3.2	Finding best predictor for each cluster	17
Appendix		
A	Some more stuff	24
	Curriculum Vitae	25

List of Figures

1.1	PowerTAC simulation environment.	3
3.1	Two days energy usage for the customer Brooksidehomes.	8
3.2	Two weeks energys usage of the downtown office customer.	8
3.3	Two days energys usage of the village 2 ns controllable customer.	9
3.4	A day's energys usage of the sf2 thermal storage customer.	9
3.5	Two week's energys usage of the sf2 thermal storage customer.	10
3.6	Two days energys usage of the SunnyHill solar production.	10
3.7	One week's energys usage of the SunnyHill solar production.	11
3.8	Number of customers vs Powertype.	12
3.9	Population vs Powertype	12
3.10	Energy vs PowerType.	12
3.11	Energy share for each power type.	12
4.1	cluster type vs average absolute error.	18
4.2	cluster type vs percent relative absolute error	18
4.3	cluster 0 average absolute error of 4 classifiers	19
4.4	cluster 0 average relative absolute error of 4 classifiers	19
4.5	cluster 1 average absolute error of 4 classifiers	19
4.6	cluster 1 average relative absolute error of 4 classifiers	19
4.7	cluster 2 average absolute error of 4 classifiers	19
4.8	cluster 2 average relative absolute error of 4 classifiers	19
4.9	cluster 3 average absolute error of 4 classifiers	19
4.10	cluster 3 average relative absolute error of 4 classifiers	19
4.11	Performance of the best classifier for each customer type. Customer Medical center was excluded as it was showing huge error.	21
4.12	average absolute error	21
4.13	average percent relative absolute error	21

Chapter 1

Smart Grid and PowerTAC Competition

In this chapter, I will describe the smart grid and Power TAC.

1.1 Traditional Electricity Distribution and Consumption System

In traditional power grids, there are three subsystems electricity generation, transmission, and distribution [3]. In electricity generation subsystem, the generator rotates a turbine in a magnetic field which generates electricity. The turbine rotates through the power of kinetic energy of water falling from a waterfall or a river with strong current, or from the energy of nuclear power plant, or energy received from burning coal or oil. Traditional energy generation system then transmits the electricity through transmission grid and electricity gets distributed through the distribution grid. This generation system is one way meaning a the electricity flow occurs from source node to consumption node only.

1.2 Smart Grid

In contrast to the traditional electricity generation system, Smart Grid are two-way [3]. So, any node in the distribution grid can produce electricity and push it to the distribution grid if necessary. The NIST report [3] states that the SG would make the electricity generation and supply robust against generator or distribution node failure, use renewable energy widely and efficiently, reduce greenhouse gas emission, reduce oil consumption by encouraging usage of electric vehicles, it will give customers more freedom to choose among energy sources. Smart grids will encourage usage of the electric vehicle as these vehicles have the ability to store power in a battery and transmit the power to the distribution grid if there is a necessity. The major challenge with the usage of renewable energy is it is uncertain. This uncertainty causes the ability to predict how much energy the SG can produce in a future time slot hard. The success of SG will need efficient methods to predict energy production [11].

1.3 Smart Grid and Renewable Energy

One of the major focus of Smart Grid will be using renewable energy. There are challenges involved with using this abundant source of energy [13]. People are already showing strong motivation to use renewable energy as indicated by the statistics that 20% of total energy is from the renewable sources which are second after coal 24%. Consumers are using renewable energy due to economic reward and environmental concern. A major challenge with renewable energy is the amount of the energy produced is greatly variable. Since the energy produced is volatile there must be a storage mechanism that balances out the surplus energy. The usage of rechargeable electric vehicles might serve the purpose of storage. Accurate prediction of the renewable energy might enable the electric car users to absorb surplus energy and push it back to the grid in peak hours if necessary.

1.4 Importance of accurate load forecasting

. Accurate load forecasting is important to ensure efficient fuel usage, reduce wastage of energy and planning proper operation of power generators [8]. Failure to do accurate load forecasting may lead to economic and environmental problems.

1.5 Power TAC System

Power Trading Agent Competition (Power TAC), is a low-risk system that simulates a smart grid based energy system. The power TAC simulation has several components such as wholesale market, brokers, customers and weather service. The system is trained on customers behavior of past years and uses real weather data. The following sections give a brief explanation of each component of the Power TAC.

1.5.1 Broker

Brokers represent the entities that buy energy from the wholesale market and sell to the customers. Contestants implement their own brokers. Each broker's objective is to maximize its profit. A successful broker has to buy and sell energy in a profitable way. The presence of several brokers in the system makes the environment competitive and every broker has to come up with a way to attract the customers.

1.5.2 Wholesale Market

The wholesale market is the bidding place for buying energy. Brokers submit their bids for a future timeslot in the wholesale market. If the bid was successful, the broker receives its desired amount by paying a certain amount of money.

1.5.3 Customers

A customer represents an entity that buys energy from the brokers. Customers subscribe to the tariffs that the brokers publish. The customers choose the most suited and affordable tariff for them by evaluating the existing tariffs in the market. They have to pay a certain amount of money to the brokers based on their tariff plans and energy usage.

1.5.4 Balancing Market

Balancing market represents the market from where the broker can buy energy in case of emergency. For example, if a broker has bought less amount of energy for a given timeslot and it finds it needs more energy then it can buy the necessary amount of energy from the balancing market. Usually, the balancing market transactions are costly for brokers than the wholesale market.

1.5.5 Weather Service

The weather service broadcasts weather forecast to the brokers. Many customer's energy usages vary based on the weather. The PowerTAC system uses the real weather data from the past.

1.5.6 Tariff Market

Brokers publish their tariff plans in the tariff market. A tariff holds information about the pricing of the energy. Customers, upon analyzing available tariffs, subscribe to their mostly suited tariff plan. Figure 1.1 shows a block diagram of the components of the powerTAC simulation environment.

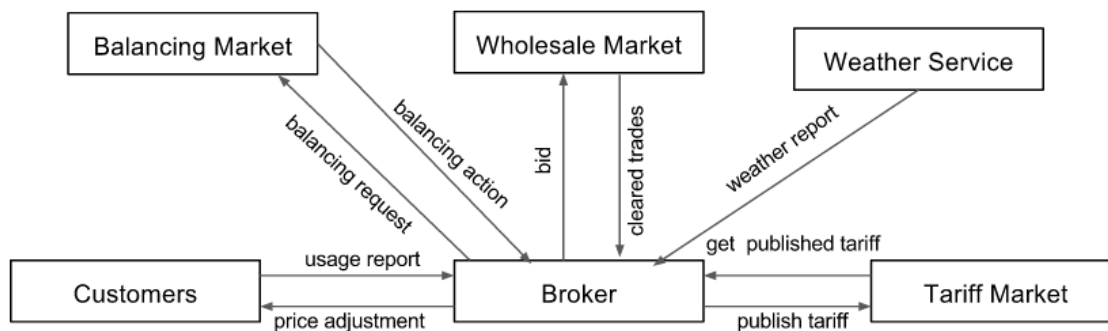


Figure 1.1: PowerTAC simulation environment.

Chapter 2

Related Works

There are mainly two types of load forecasting namely short term load forecasting and long term load forecasting. Short term load forecasting deals with forecasting up to a couple of weeks. Long term load forecasting may forecast customer's demand over month or year [2]. In this chapter, I have described different methods of energy load forecasting for long term and short term in the literature.

2.1 Variables in Electricity Demand

Studies such as [6], [5] and [2] have found that temperature has effect on electricity demand. The study in [5] was done in a region of Australia. It was found that, in a lower temperature the customers tend to use heaters and in a higher temperature they tend to use coolers. As a result, the increase or decrease of temperature from a certain point will cause the consumption of electricity to increase. In study [2], two demand forecasting models were proposed. One was univariate Auto Regressive Integrated Moving Average (ARIMA) and another one was univariate ARIMA model along with temperature depended transfer function. The model with temperature variable did better forecasting than the one without the temperature variable. On the other hand, the study in [1] showed that inclusion of temperature variable in forecasting model actually introduced more error in demand prediction. The aim of the study was to make forecasting about electricity usage of January based on past five years training data using a Support Vector Machine (SVM) forecaster. The reason behind of getting more error after including temperature variable may be because during January the temperature did not change much and the inclusion might have caused overfitting.

Weather variables such as wind speed and cloud cover has effect on electricity demand [6], [14]. As cloud cover increases, the demand for light increases too. The increased lighting demand causes increased electricity demand. The period where cloud cover was low, the electricity demand was also low [6]. High speed wind across wet walls help cool houses. High speed wind thus may cause reduced electricity demand due to reduced demand of air cooling [14].

In the survey article [4], the authors reported that the day of the week and the month of the year is highly correlated with customer's energy demand. The electricity load demand can be higher or lower based on the day of week. The weekends usually have different load demand pattern than the week days. Also, based on the hour of a given day, the load demand can be higher or lower too. The season also showed impact on electricity demand.

2.2 Electricity Load Forecasting Using Statistical Method

To make electricity load forecast, researchers have used statistical methods such as statistical average and Auto Regressive Integrated Moving Average (ARIMA). Agent TACTEX'13, the winner of the PowerTAC competition in 2013, used the statistical average to make electricity demand forecasting for an hour of a day of a week. In a week a customer has $24 * 7 = 168$ hours or slots where it can consume electricity. TACTEX'13 kept track of average usage of 168 weekly slots for each customer. To predict a future time slot, their agent would look at which weekly slot the future time slot would fall in. Then the agent used that weekly slot's average usage as the forecast of the future time slot. [2] have used ARIMA model for load forecasting. The ARIMA model uses both moving average and auto regression to forecast the demand. To make a forecast about a future time slot, the auto regression model uses some previously observed time slots values based on its degree. Moving average scheme would use the average of all the known time series data points to make a prediction about a future time slot.

2.3 Load Forecasting using Machine Learning

[10] the authors used various machine learning techniques to make 24 hours ahead load forecast for the Power TAC simulation. They found that hour of week, weather related features such as temperature cloud cover were influential to the electricity load demand. They created one machine learning electricity load forecasting module for each customer by extracting relevant features for the customers. The forecasting modules performed well for the customers that showed regularity in their energy consumption behavior. For the customers with load shifting capabilities to their favored hour and customers with irregular consumption patterns, the scheme did no perform well.

2.4 Load Forecasting for a Specific Region

Regional load forecasting will enable the authority to know which regions need more energy. If they know which regions need more energy, they will know most suitable places to place electricity generator plants. [7] worked on load forecasting based on region. They divided electricity usage of Taiwan into 4 areas. For each region, they collected GDP, population, highest temperature and aggregated load. After that, they trained Artificial Neural Network (ANN) model for each region. For the baseline, they trained linear regression model for each region. The result showed that the ANN-based load forecasting methods performed better than the linear regression methods.

2.5 Load Forecasting using Clustering

[9] have used the clustering method to forecast customer's future electricity demand. They collected data from more than 4000 household customers in Ireland for about 6 months. Collected data included electricity usage at 30 minutes interval, appliances used in the

home and different socio-economic information about the people living in a particular house. They clustered each day's usage which they call load profiles. A customer's daily usage then can be assigned to one of those load profiles. The customer is then characterized by the mostly used load profile. The authors then trained a linear regression classifier that was built upon the socio-economic information of the customers, types of appliances used in the house and the description of the house to figure out the common load profile of the given household. The predicted load profile of the customer received from the linear regression model will be used to predict the demand of the customer for a given day. [2] noticed difference of behavior among customers. They manually clustered the population into four categories namely commercial, office, residential and industrial customers. In their paper [15], the authors proposed a novel demand prediction mechanism. In Power TAC competition, every broker is provided with past two weeks usage of all the customers or bootstrap usage. Their proposed broker clustered the customers based on the bootstrap data. For each cluster, the broker would make a linear regression model. The input variables included past average usage and weather related information. This approach of prediction clusters is based on the usage pattern of the customers. So this method may not be suitable for customers with irregular usage pattern such as customers with load shifting capabilities and electric vehicle customers.

2.6 Expert System based Load Forecasting

The authors in ref [12] have proposed an expert system based load forecasting method for the region, Virginia. The expert system would forecast the load of upcoming 24 hours. They observed the variables that are likely to affect the load. They came up with variables such as temperature, load of the previous hour, season, and day of the week have a strong correlation with the observed load. They implemented a computer program that mimicked how a human operator makes load forecast based on the independent variables. For a specific region's weather condition, their method worked well and required a limited amount of historical data.

From the review of the literature, the importance of weather related variables such as temperature, cloud cover and wind speed is evident. Also, the hour of the day and day of the week are highly correlated with the load demand. A combination of machine learning classifiers and clustering algorithms appears to be a better idea. For the methodology of [10] it will take a large number of predictors for the simulation system. Also, those predictors will not work if the name of the customer is changed or a new customer is introduced as each predictor is hard coded with a specific customer. It sounds reasonable to cluster the data first and then train machine learning classifier for each cluster. This approach will hold generality. Instead of training only on bootstrap data as the [15] have done, a wealth of data generated from the simulations can be used to train the cluster. Since the clustering is done offline, this approach will not suffer from the problem of having a time limit that the broker has to face if the cluster is trained during the competition. After the clustering is done, for each cluster, different machine learning classifiers can be trained to figure out which one performs the best. So, the broker will no longer stick to linear regression. This way, the training module will be able to deal with new customers.

Chapter 3

Customer Description

In this chapter, I will describe the customers present in the PowerTAC simulation system, some statistics about them and their attributes.

3.1 Customers

In PowerTAC simulation system the customers are the entities that buy and sell electricity. A customer subscribes to one of the tariffs published by the brokers and it pays or receives money according to the tariff plan. A customer can represent a population size of one to several thousand. For example, customers that represent an electric vehicle customer represent only one person and the customers that represent a village usually have several thousand populations. In PowerTAC environment there are 168 customers in total.

3.2 PowerTypes

In Power TAC, every customer has a power type. Power type determines the behavior of the customers. A customer that has power type related to production produces electricity. A customer that has a power type related to energy consumption consumes energy. In the following subsections, I describe different power types of the simulation.

3.2.1 consumption

A customer with power type consumption are the most common customers. They use the energy when they need it. They cannot shift their demand to a future timeslot. Usually, they have a regular pattern of their energy usage. Usually, they show a similar pattern for weekdays. They have similar kind of usage pattern for the weekends.

The figure 3.1 shows 2 days electricity usage of the BrookSideHomes customer. The pattern shows in a day, around at 10 am there is a growing need for electricity. During the night after 10 pm, the electricity consumption starts decreasing.

The figure 3.3 shows two weeks consumption of the downtown customer. The customer shows a similar pattern for all weekdays. It also shows distinguishable energy usage during the weekends.

3.2.2 Interruptible Consumption

Interruptible customers are smart enough to shift their energy demand in a timeslot where they can buy electricity at a reduced price. Because of this shifting capability, they don't



Figure 3.1: Two days energy usage for the customer Brooksidehomes.

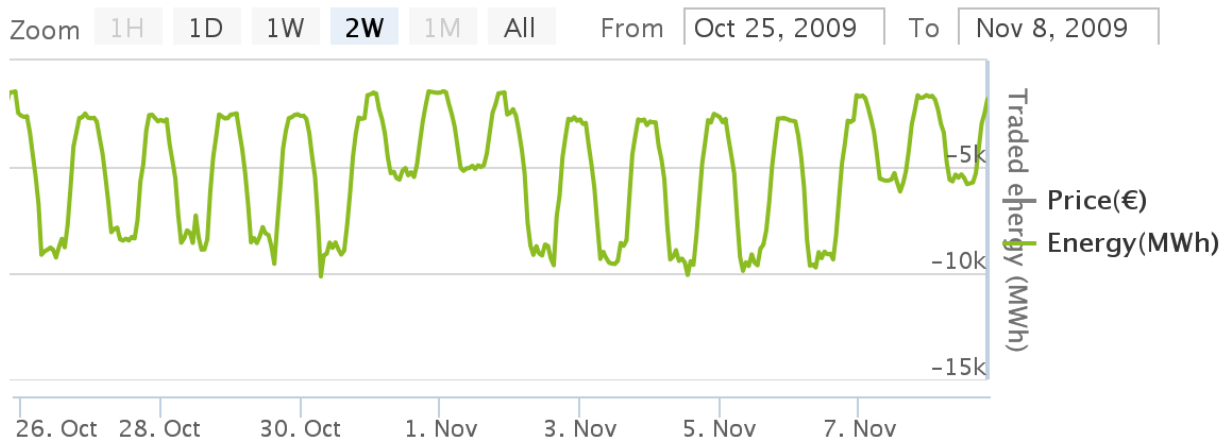


Figure 3.2: Two weeks energys usage of the downtown office customer.

show a regular usage pattern as the consumption customers do. Figure 3.3 shows a controllable customer's 2 days usage.

3.2.3 Thermal Storage

Thermal storage customers show a weekly pattern in their electricity usage. Also, during a day, their electricity usage in a day depends much on the energy they used in the last timeslot. Figure 3.4 and 3.5 shows a day and two week's energy usage of the thermal storage customer sf2.

3.2.4 Solar Production

Figure 3.6 shows two day's and figure ?? shows a week's energy production of the SunnyHill solar production customer. From the figures, we see that during the daytime the solar production customers usually produces electricity which is as expected.



Figure 3.3: Two days energies usage of the village 2 ns controllable customer.

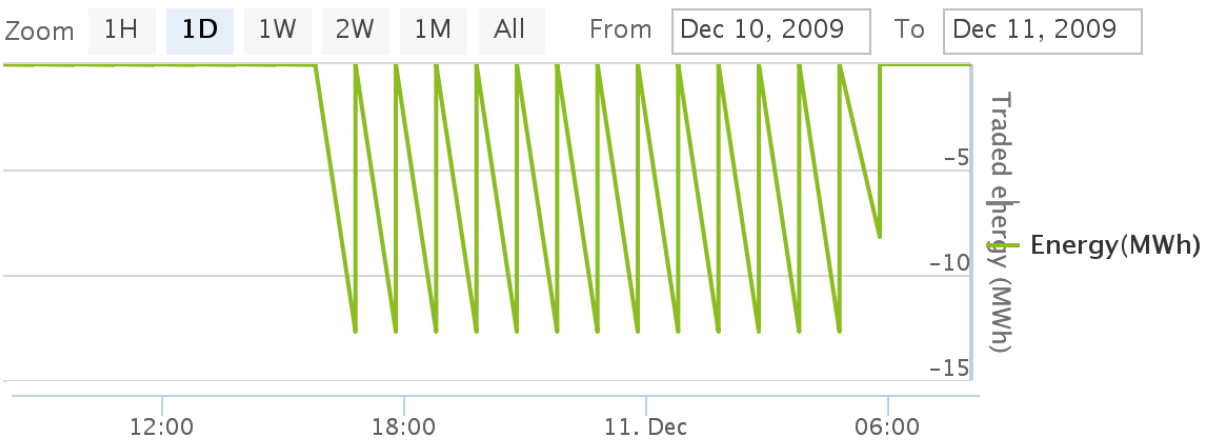


Figure 3.4: A day's energies usage of the sf2 thermal storage customer.

3.2.5 Wind Production

Wind production customers generate energy from the wind.

3.2.6 Electric Vehicle

An electric vehicle customer represents one electric vehicle. Their usage of energy is quite irregular and hard to predict.

3.3 Statistics

In this section, I present some statistics on the customers available in the system.

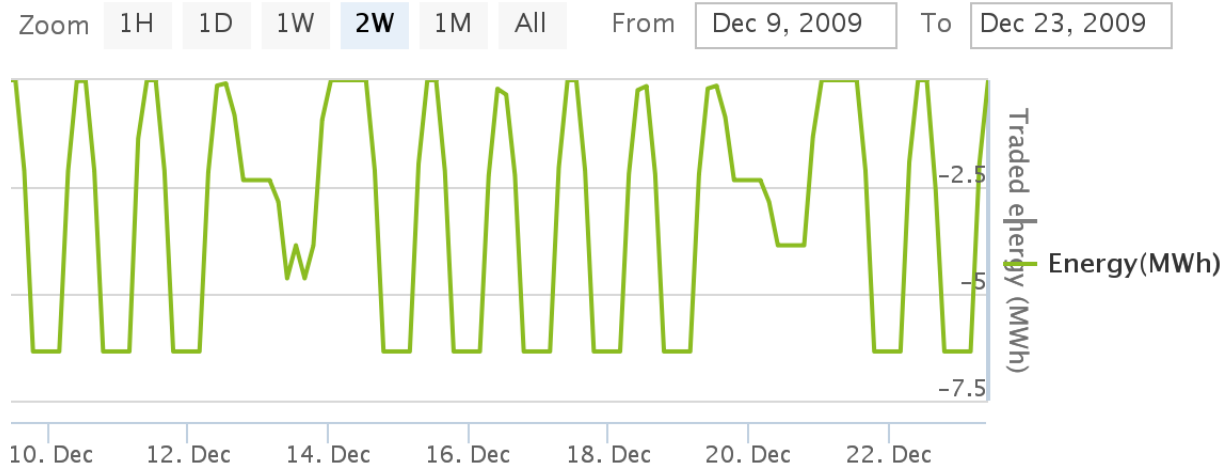


Figure 3.5: Two week's energys usage of the sf2 thermal storage customer.

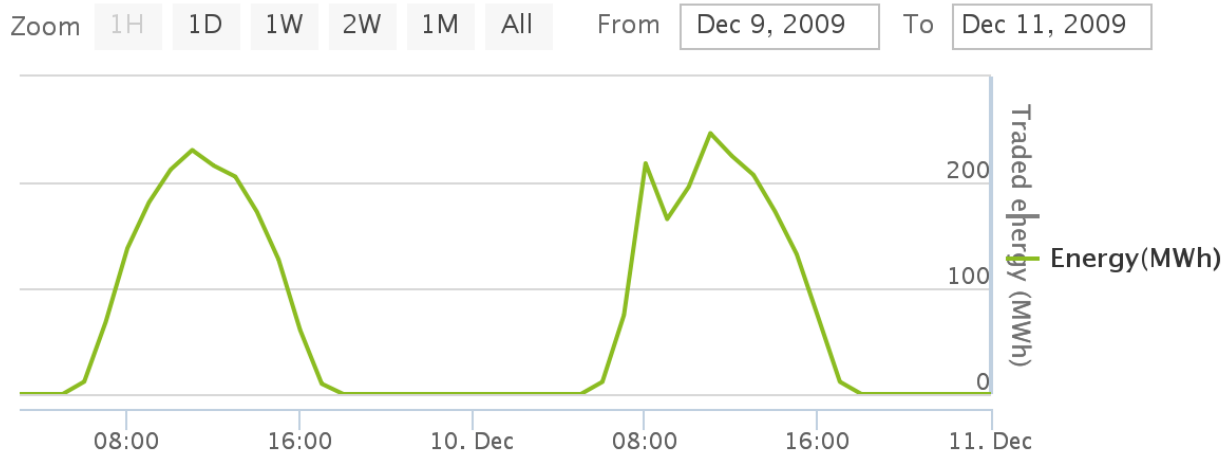


Figure 3.6: Two days energys usage of the SunnyHill solar production.

3.3.1 Customer Vs PowerType

In the figure 3.8 we can see the system has more customer with the power type electric vehicle than any other power types. This is because the electric vehicle represents a population of size 1.

3.3.2 Population Vs PowerType

From figure 3.9 by far the powertype of consumption has the most number of population.

3.3.3 Total Energy Consumed Vs PowerType

From figure 3.10 we can see that the consumption type customers uses the most amount of the electricity.

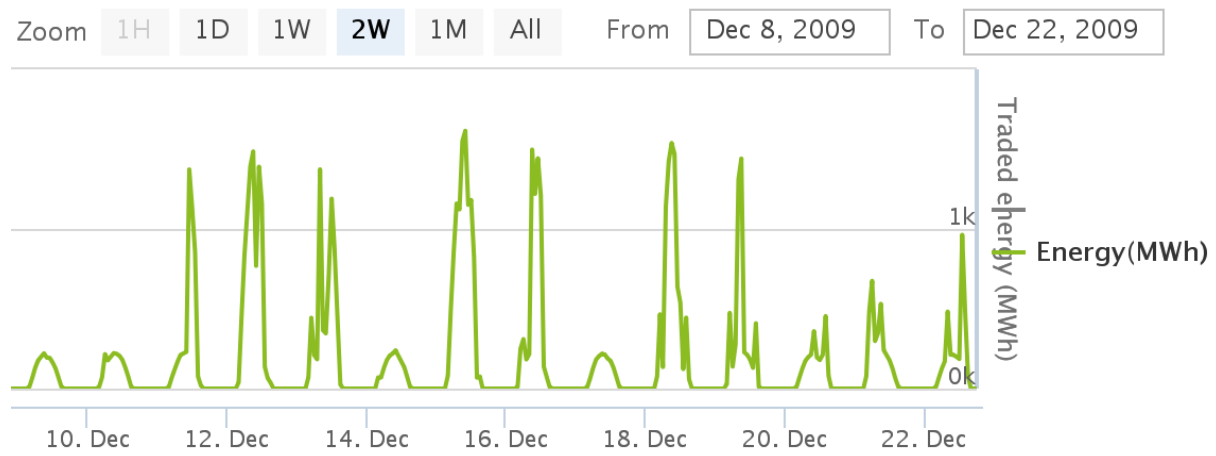


Figure 3.7: One week's energys usage of the SunnyHill solar production.

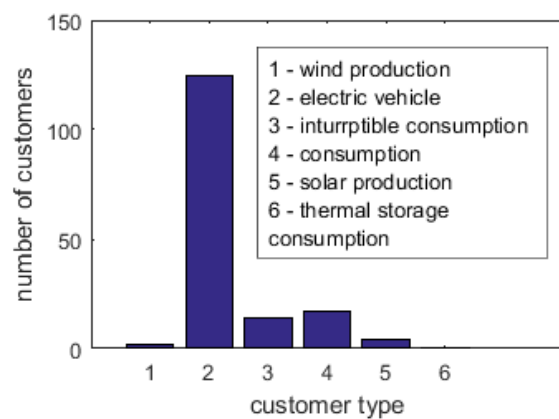


Figure 3.8:
Number of customers vs Powertype.

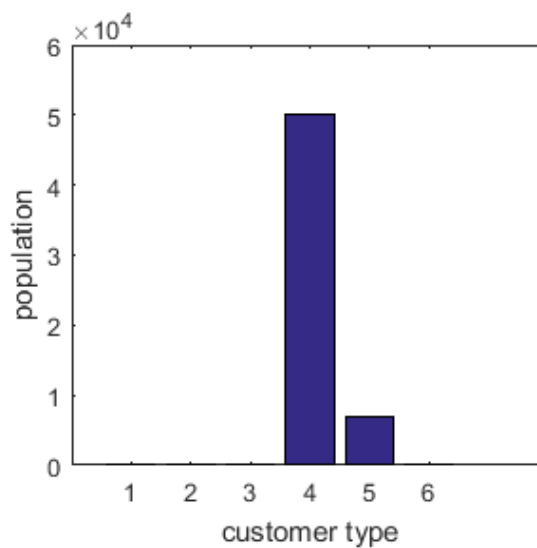


Figure 3.9: Population vs Powertype

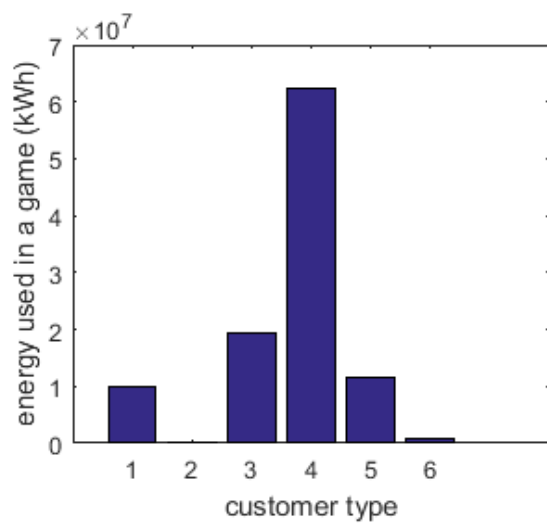


Figure 3.10: Energy vs PowerType.

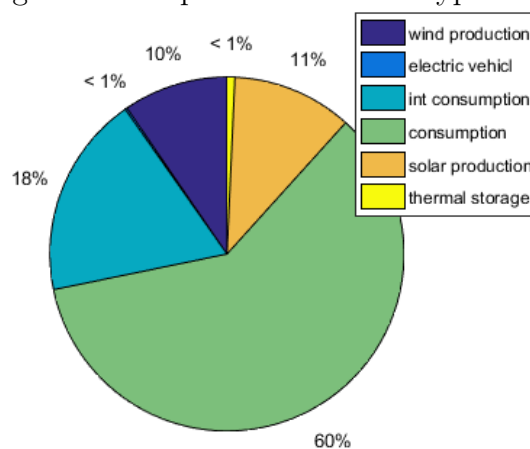


Figure 3.11:
Energy share for each power type.

Chapter 4

Methodology and Result

Traditionally, a single type of predictor served to predict the energy demand of all power type customers. Since each power type customers acts differently, I have attempted to attack each type of customer separately to make a prediction mechanism that performs better than the baseline predictor.

4.1 The Baseline Electricity Forecasting Mechanisms

The first baseline energy forecasting mechanism is the default prediction mechanism provided by the PowerTAC system. It exploits the fact that usage of a timeslot of a customer in a specific date is highly correlated with the day of the week and the time slot. To make a prediction it stores the average energy usage of an hour of a week. So, for each customer, it uses $24 * 7 = 168$ memory to remember average usages. As soon as it learns about a new usage information of an hour of a week, it updates old average using the following algorithm.

Algorithm 1 Update average usage for $customer_i$ for day d and timeslot t , $newUsage$

- 1: $avgUsage = \text{get average usage of } customer_i \text{ at day } d \text{ and time slot } t$
 - 2: $avgUsage = 0.7 * avgUsage + 0.3 * newUsage$
-

Algorithm 2 forecast usage for day d and timeslot t for $customer_i$

- 1: $avgUsage = \text{get average usage of } customer_i \text{ at day } d \text{ and time slot } t$
 - 2: return $avgUsage$
-

The second baseline forecasting mechanism is designed to make energy forecasts for a single customer. In general, if there are n customers in the system, we will need n energy forecasters each one trained on the data of a single customer. I went further by checking different machine learning algorithms such as M5Tree, Linear Regression, M5P rules and REP tree for each customer and picked the best performing one for each customer.

4.2 Proposed Electricity Demand Forecasting Mechanism

In this section, I will describe how I attempted to make energy demand forecaster for consumption power type customers.

4.2.1 Demand Forecasting for Consumption Type Customer

For the consumption type customers, the algorithm 3 describes the proposed method of forecasting energy demand and how it was compared to the baseline methods.

Algorithm 3 Make electricity demand forecasting for consumption type customer

- 1: Extract features for each time slot for each customer [algorithm 4, 5 and 6]
 - 2: train kmeans cluster for different sizes of k [algorithm 7]
 - 3: train linear regression classifier for each cluster and compute error [algorithm 8]
 - 4: pick suitable value for k by observing the errors
 - 5: for each cluster, find the best performing predictor for that cluster [algorithm 9]
 - 6: train individual classifier for each customer to make the second baseline [algorithm 10]
 - 7: evaluate performance using test data [algorithm 11]
-

The algorithm 3 begins with extracting information from the game log files. All the activities that occurred in a game can be found in a game log. In power TAC the shortest time unit is an hour and it is called time slot. Activities such as buying or selling electricity occur during a time slot. At the beginning of a time slot, the system notifies the broker that a new time slot is about to begin. The system also notifies the brokers with weather forecast about the future time slots. As a time slot ends, the broker receives information about its customer's energy usage which is called tariff transaction report. Algorithm 4 refers how the extraction program retrieves necessary information from tariff transaction report. As the broker gets notification of the beginning of a new time slot, the extraction program has all the information related to energy usage and weather data of the previous time slot available by this time. Algorithm 5 shows the procedure of writing the information of the known time slot's information in training instance file. Once the simulation ends, the extraction program knows the average energy usage of all the customers during a week. In a week there are $24 * 7 = 168$ hours or time slots. The extraction program writes all the 168 hourly averages of a week to a file. This is explained in algorithm 6.

Algorithm 4 extract information from transactionReport sent to broker after each time slot through TariffTransactionHandler call back method

- 1: timeSlot = get time slot from transactionReport
 - 2: customerName = get customer name from transactionReport
 - 3: energyUsed = get energy used from transactionReport
 - 4: addUsage(customerName, timeSlot, energyUsed)
-

Next, all the average weekly usages are combined together to make training set for the clustering algorithm. I have used kMeans clustering algorithm to cluster the training set. I have trained clusters of size = 4, 5, 6, 7, 8, 9, 10 and 11. The algorithm 7 describes the procedure of making clusters from the training instances. Once a kMeans of cluster size k is made, a program groups the hourly usages of the customers in the same cluster and combines them to make training set for machine learning classifier. This training set is

Algorithm 5 write extracted data after timeSlot update message received from TimeSlotUpdateHandler call back method

```
1: knownTimeSlot = timeSlot - 1
2: for each customer do
3:   day = get day of knownTimeSlot
4:   hour = get hour of knownTimeSlot
5:   statisticalData = get statistics of the customer of day and hour
6:   weatherData = get weather data of knownTimeSlot
7:   trueUsage = get true usage of customer in knownTimeSlot
8:   trainingInstance = create training instance by combining statistical data, weather
   data and true usage
9:   writeToFile(trainingInstance)
10: end for
```

Algorithm 6 write average electricity usage of the customers of each hour of the week

Require: information of all timeslots has been received

```
1: for each customer do
2:   trainingInstance = create empty training instance
3:   for each day of week do
4:     for each hour of day do
5:       averageUsage = get average usage of day and hour of customer
6:       append averageUsage to the trainingInstance
7:     end for
8:   end for
9:   writeToAvgUsageFile(trainingInstance)
10: end for
```

used to train linear regression classifier. To test the performance of the classifiers, I have separated five game logs and they were not used for training purposes. The algorithm 8 describes how the cluster based predictor's performance was evaluated.

Algorithm 7 create kmeans cluster of size k from weekly usage training instance file

```
1: data = load weekly average usage file
2: kmeansCluster = build kmeans cluster of size k
3: save kmeansCluster
```

Based on the errors observed from different kMeansCluster based predictions, I fixed the number of clusters. Once the number of the clusters was fixed, a program creates several machine learning predictors to see which one performs best for a given cluster. The machine learning classifiers that were tried out are linear regression, M5P rules, M5 Tree, REP tree. In the runtime, a customer will be grouped in a cluster based on its weekly usage. Once the program knows the cluster assigned to a customer, the program will load the corresponding demand forecaster to make electricity demand forecast about the customers.

At this phase, I have the proposed cluster-based customer's demand forecaster. Next,

Algorithm 8 find error of kmeans clusters of different size

```
1: for each cluster size k do
2:   get the kMeansCluster of size k
3:   for cluster in KMeansCluster do
4:     combine slot based training instances of that cluster
5:     train linear regression classifier based on the combined data
6:     save the classifier for cluster
7:   end for
8: end for
9: for each training instance do
10:  compute error of the instance using each kMeansCluster
11: end for
```

Algorithm 9 find best classifiers of each cluster of kmeans cluster of size k

```
1: for each cluster in kMeansCluster do
2:   combine slot based data of the all the customers in cluster
3:   train available classifiers on the combined data using 10 fold cross validation
4:   choose the classifier with minimum error
5:   save the classifier for making demand forecasting for cluster
6: end for
```

the baseline predictor that needs a machine learning classifier for each customer is built. At first, the training instances are combined based on the name of the customer. This means for n customers n training set is constructed, each of the training set has only the information of a single customer. A training set related to a customer is used to create machine learning classifiers for that customer. Several classifiers had been tried out to figure out which classifier performs the best for a customer. The best performing classifier was chosen to predict about a customer. Algorithm 10 explains the procedure of getting the best classifier.

Algorithm 10 find best classifiers created for each individual customer

```
1: for each customer do
2:   combine all slot based training instance of the customer
3:   train available classifiers on the combined data using 10 fold cross validation
4:   choose the classifier with minimum error
5:   save the classifier for making prediction about the customer
6: end for
```

The next phase is testing the performance of the proposed and baseline methods. For testing, I had used five game logs that were not used for training purposes. For each test instance, all three methods output was observed to figure out the performance. The algorithm 11 shows the mechanism of testing.

Algorithm 11 performance evaluation of each method

```
1: for each test instance do
2:   classify the test instance using moving average usage [algorithm 2]
3:   classify the test instance using individual prediction mechanism
4:   classify the test instance using cluster based predictor
5:   calculate and accumulate errors of each mechanism [algorithm12]
6:   update moving average baseline predictor based on the information from the test
     instance [algorithm 1]
7: end for
8: find average error from the accumulated errors for each forecasting mechanism
```

Algorithm 12 calculate error from the predicted value and the true value

```
1: absoluteError = abs(predictedValue - trueValue)
2: relativeAbsoluteError = (absoluteError / trueValue ) * 100 %
```

4.3 Result

4.3.1 Finding number of clusters

At first, I have segmented the customer using KMeans clustering algorithm with cluster sizes = 4, 5, 6, 7, 8, 9, 10 and 11. For KMeans with size k, we will have k clusters. For each of the k clusters, I had a linear regression predictor. I observed the relative percentage error and absolute average the above cluster sizes. Figure 4.1 shows the result. From, the figure it is clear that the size of the cluster does not have a big impact on the prediction performance.

To keep things simple, I have decided to choose Kmeans cluster of size 4. When $k = 4$ was chosen, table 4.1 shows the cluster assignment for each customer. It can be seen that, cluster-0 held most of the offices, cluster 2 held most of the village types, cluster 3 held the medical center, cluster 1 held large housing such as brooksidehomes, centerville homes and large offices such as downtown offices and centerville offices.

4.3.2 Finding best predictor for each cluster

I have used the following features for a given timeslot to train prediction models.

- Temperature
- Cloud Cover
- Wind Speed
- Average of the Slot
- Standard Deviation of the Slot

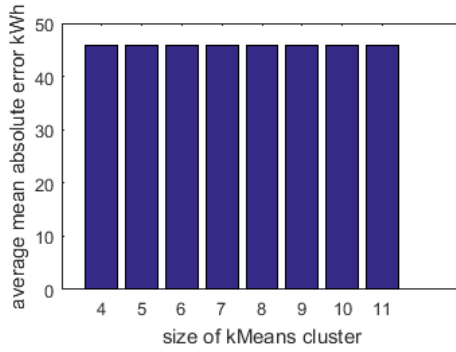


Figure 4.1:
cluster type vs average absolute error.

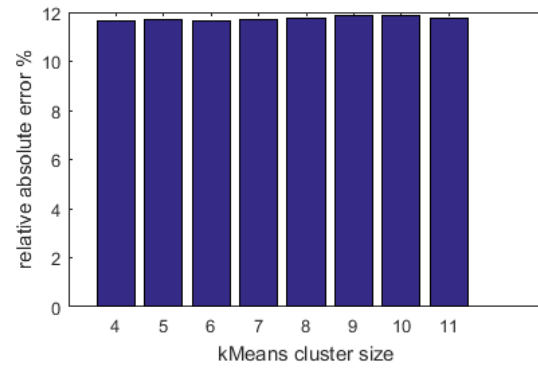


Figure 4.2:
cluster type vs percent relative absolute error

Customer Name	Assigned Cluster Number
BrooksideHomes	0
CentervilleHomes	0
DowntownOffices	1
EastsideOffices	1
OfficeComplex 1 NS Base	0
OfficeComplex 1 SS Base	0
OfficeComplex 2 NS Base	0
OfficeComplex 2 SS Base	0
Village 1 NS Base	2
Village 1 RaS Base	2
Village 1 ReS Base	2
Village 1 SS Base	2
Village 2 NS Base	2
Village 2 RaS Base	2
Village 2 ReS Base	2
Village 2 SS Base	2
MedicalCenter@1	3

Table 4.1: Assigned cluster for each customer

Next, I have tried out M5Tree, Linear Regression, M5P rules and REP tree machine learning classifiers to see which one performs the best for each of the 4 clusters. Figure 4.3, 4.5, 4.7, 4.9 show that M5P, M5P, REPTree and M5RULES are the best predictors for cluster 0, 1, 2 and 3 respectively.

The next step is to find the best classifiers for each of the customers. Based on the data from each of the customers, the four types of classifiers described in previously were tried out. For each customer, the following classifiers performed the best.

The figure 4.11 shows error percentage of each of the predictors type for each of the

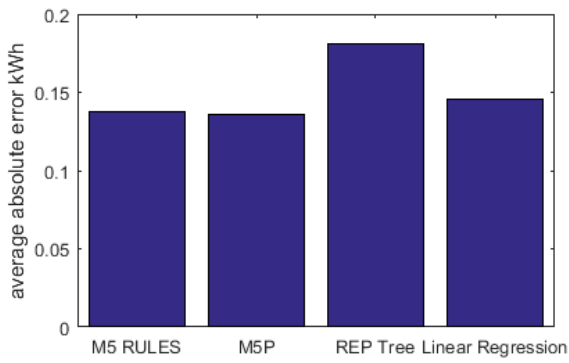


Figure 4.3:
cluster 0 average absolute error of 4 classifiers

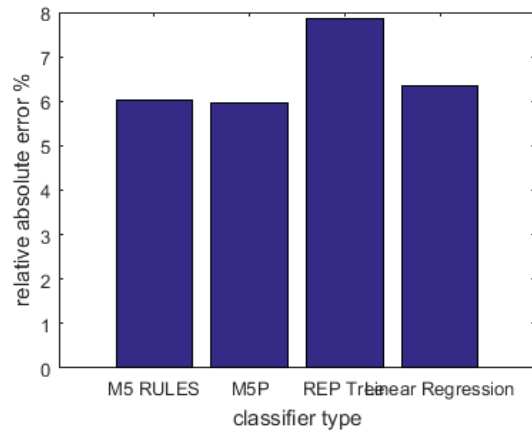


Figure 4.4:
cluster 0 average relative absolute error of 4 classifiers

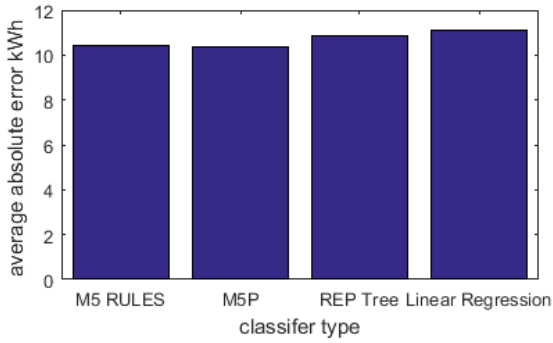


Figure 4.5:
cluster 1 average absolute error of 4 classifiers

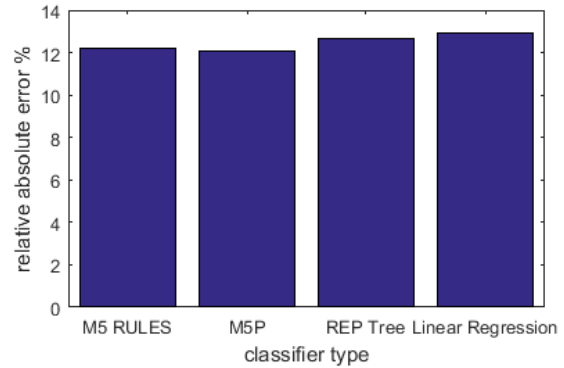


Figure 4.6:
cluster 1 average relative absolute error of 4 classifiers

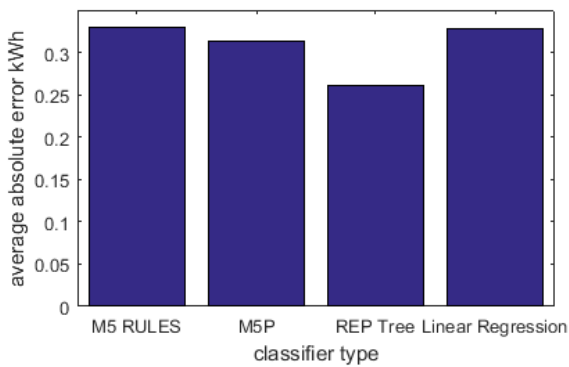


Figure 4.7:
cluster 2 average absolute error of 4 classifiers

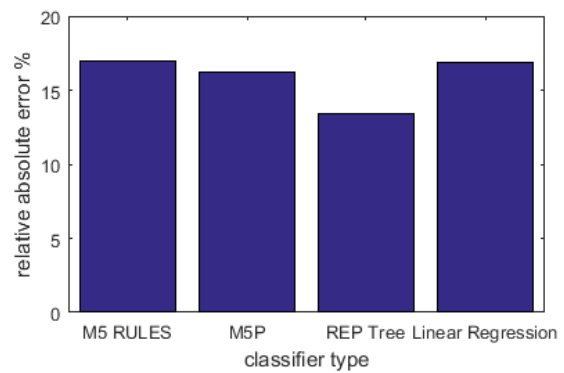
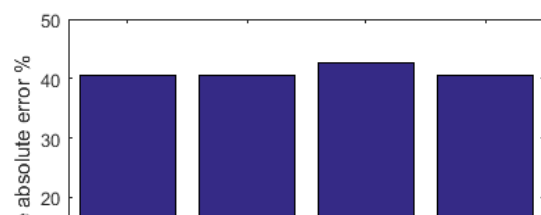
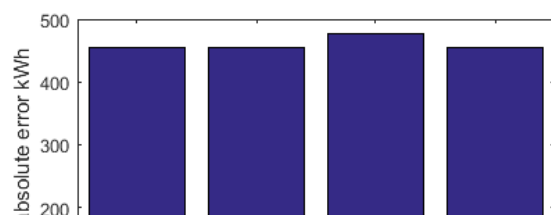


Figure 4.8:
cluster 2 average relative absolute error of 4 classifiers



Customer Name	Best Predictor Type
BrooksideHomes	M5P
CentervilleHomes	M5P
DowntownOffices	M5P
EastsideOffices	M5P
OfficeComplex 1 NS Base	LinearRegression
OfficeComplex 1 SS Base	LinearRegression
OfficeComplex 2 NS Base	LinearRegression
OfficeComplex 2 SS Base	LinearRegression
Village 1 NS Base	M5P
Village 1 RaS Base	LinearRegression
Village 1 ReS Base	M5P
Village 1 SS Base	M5P
Village 2 NS Base	LinearRegression
Village 2 RaS Base	M5P
Village 2 ReS Base	M5P
Village 2 SS Base	M5P
MedicalCenter@1	M5P

Table 4.2: Best individual predictor for each customer

customer types.

Finally, the cluster based forecasting and the two baselines were tested with data extracted from 5 test files that were not used for training. From Figure 4.12 we can see that cluster based prediction mechanism performed almost as good as the mechanism where n predictors are needed for n customers. And it did well than the default moving average prediction scheme.

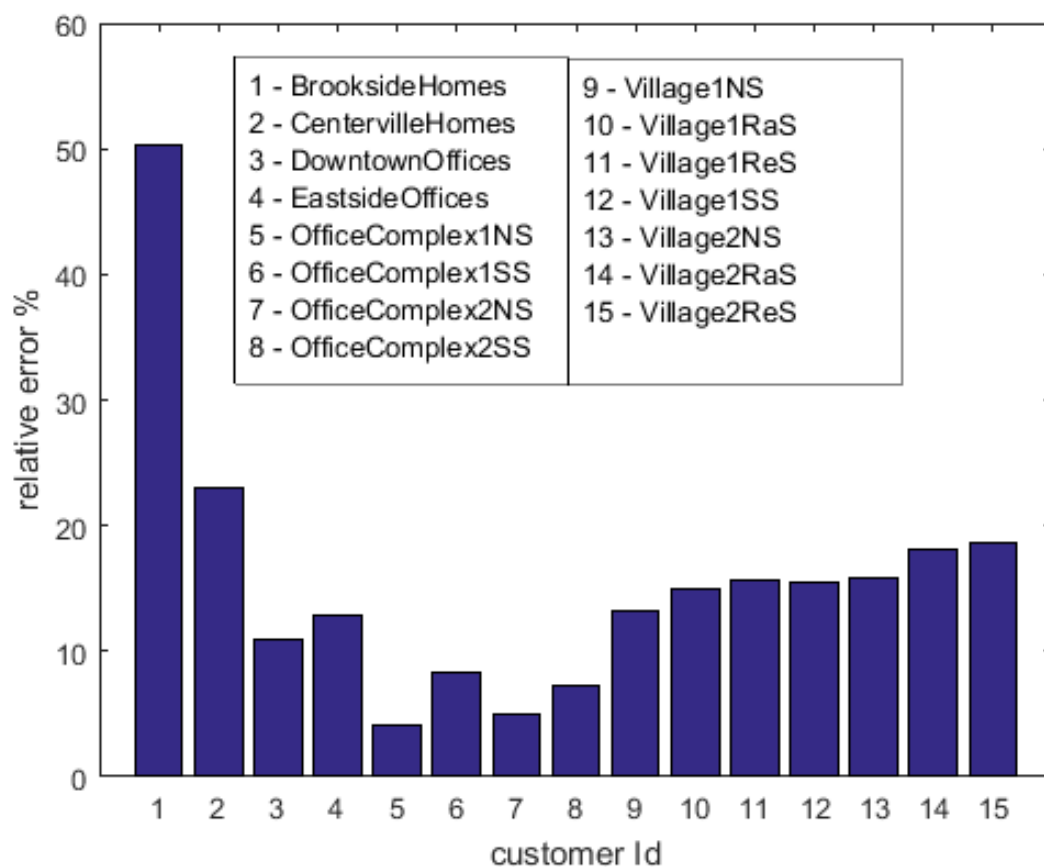


Figure 4.11: Performance of the best classifier for each customer type. Customer Medical center was excluded as it was showing huge error.

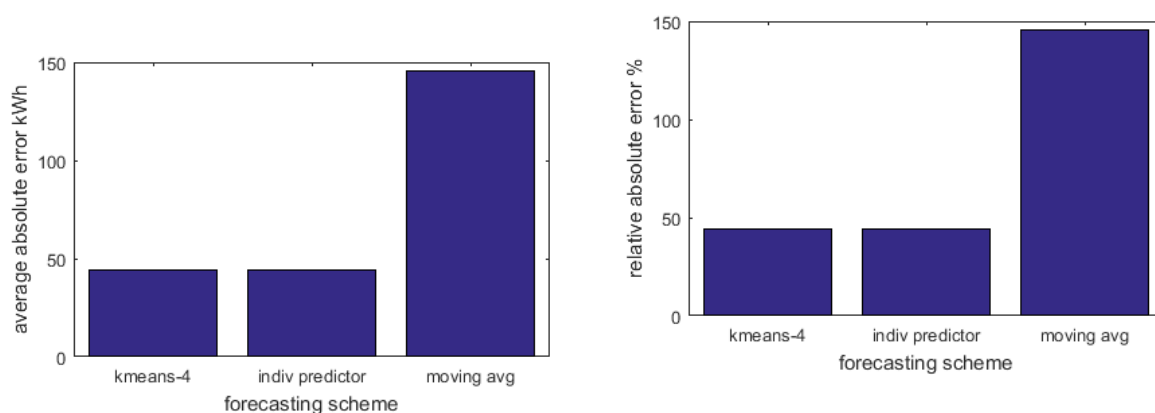


Figure 4.12: average absolute error

Figure 4.13: average percent relative absolute error

References

- [1] Bo-Juen Chen, Ming-Wei Chang, and Chih-Jen Lin. Load forecasting using support vector machines: A study on eunite competition 2001. *Power Systems, IEEE Transactions on*, 19(4):1821–1830, 2004.
- [2] MY Cho, JC Hwang, and CS Chen. Customer short term load forecasting by using arima transfer function model. In *Energy Management and Power Delivery, 1995. Proceedings of EMPD'95., 1995 International Conference on*, volume 1, pages 317–322. IEEE, 1995.
- [3] Xi Fang, Satyajayant Misra, Guoliang Xue, and Dejun Yang. Smart gridthe new and improved power grid: A survey. *Communications Surveys & Tutorials, IEEE*, 14(4):944–980, 2012.
- [4] Heiko Hahn, Silja Meyer-Nieberg, and Stefan Pickl. Electric load forecasting methods: Tools for decision making. *European Journal of Operational Research*, 199(3):902–907, 2009.
- [5] Melissa Hart and Richard de Dear. Weather sensitivity in household appliance energy end-use. *Energy and Buildings*, 36(2):161–174, 2004.
- [6] Ching-Lai Hor, Simon J Watson, and Shanti Majithia. Analyzing the impact of weather variables on monthly electricity demand. *IEEE transactions on power systems*, 20(4):2078–2085, 2005.
- [7] Che-Chiang Hsu and Chia-Yon Chen. Regional load forecasting in taiwan—applications of artificial neural networks. *Energy conversion and Management*, 44(12):1941–1949, 2003.
- [8] Fang Liu, Qiang Song, and Raymond D Findlay. Accurate 24-hour-ahead load forecasting using similar hourly loads. In *Computational Intelligence for Modelling, Control and Automation, 2006 and International Conference on Intelligent Agents, Web Technologies and Internet Commerce, International Conference on*, pages 249–249. IEEE, 2006.
- [9] Fintan McLoughlin, Aidan Duffy, and Michael Conlon. A clustering approach to domestic electricity load profile characterisation using smart metering data. *Applied energy*, 141:190–199, 2015.
- [10] Jaime Parra Jr and Christopher Kiekintveld. Initial exploration of machine learning to predict customer demand in an energy market simulation. In *Workshops at the Twenty-Seventh AAAI Conference on Artificial Intelligence*, 2013.

- [11] Cameron W Potter, Allison Archambault, and Kenneth Westrick. Building a smarter smart grid through better renewable energy information. In *Power Systems Conference and Exposition, 2009. PSCE'09. IEEE/PES*, pages 1–5. IEEE, 2009.
- [12] Saifur Rahman and Rahul Bhatnagar. An expert system based algorithm for short term load forecast. *Power Systems, IEEE Transactions on*, 3(2):392–399, 1988.
- [13] Andre Richter, Erwin van der Laan, Wolfgang Ketter, and Konstantina Valogianni. Transitioning from the traditional to the smart grid: Lessons learned from closed-loop supply chains. In *Smart Grid Technology, Economics and Policies (SG-TEP), 2012 International Conference on*, pages 1–7. IEEE, 2012.
- [14] Ina Rüdenauer and Carl-Otto Gensch. Energy demand of tumble driers with respect to differences in technology and ambient conditions. *Öko-institut, Freiburg*, 2004.
- [15] Xishun Wang, Minjie Zhang, Fenghui Ren, and Takayuki Ito. Gongbroker: A broker model for power trading in smart grid markets. In *2015 IEEE/WIC/ACM International Conference on Web Intelligence and Intelligent Agent Technology (WI-IAT)*, volume 2, pages 21–24. IEEE, 2015.

Appendix A

Some more stuff

This is an example of how to add an appendix.

Curriculum Vitae

Patrick Thor Kahl was born on July 12, 1961. The first son of Ulf Thor Gustav Kahl and Carolyn Kahl, he graduated from Coronado High School, El Paso, Texas, in the spring of 1979. He entered Auburn University in the fall of 1979, and, in the spring of 1982, The University of Texas at El Paso. In 1985 he joined the United States Navy where he served for eight years, most of it aboard the submarine USS Narwhal (SSN671). In the fall of 1993, after being honorably discharged from the navy, Patrick resumed his studies at The University of Texas at El Paso. While pursuing his bachelor's degree in Computer Science he worked as a Teaching Assistant, and as a programmer at the National Solar Observatory at Sunspot, New Mexico. He received his bachelor's degree in Computer Science in the summer of 1994.

In the fall of 1994, he entered the Graduate School of The University of Texas at El Paso. While pursuing a master's degree in Computer Science he worked as a Teaching and Research Assistant, and as the Laboratory Instructor for the 1995 Real-Time Programming Seminar at the University of Puerto Rico, Mayagüez Campus. He was a member of the Knowledge Representation Group and the Rio Grande Chapter of the Association for Computing Machinery.

Permanent address: 6216 Sylvania Way
El Paso, Texas 79912-4927