

# Aggregate Load Forecasting using Product of HMMs

Megha Gupta  
Department of Computer  
Science  
IIIT Delhi, India,  
meghag@iiitd.ac.in

Ullas Nambiar  
EMC, India  
Ullas.Nambiar@emc.com

Haimonti Dutta<sup>\*</sup>  
Department of Management  
Science and Systems  
State University of New York,  
Buffalo,  
New York, 14260  
haimonti@buffalo.edu

Amarjeet Singh  
Department of Computer  
Science  
IIIT Delhi, India  
amarjeet@iiitd.ac.in

## ABSTRACT

Real time spatio-temporal energy consumption data is captured by large scale deployment of smart meters. Data from these meters are usually sent to a base station (BS) where they are aggregated for analytics. Each BS aggregates the load derived from all the meters connected to that station. The readings received at the BS are adhoc and usually not synchronized in time. Different smart meters can send data points when they are collected resulting in inconsistent data including aggregating non-aligned time stamped readings, readings with missing values, repeated values, meter reset readings. We address the problem of learning from disparate data streams (with inconsistencies) by modelling streams as HMMs and the process of aggregating data at the BS as a Product of HMMs. This enables us to perform load forecasting using machine learning techniques. Empirical results are presented on two data sets - Reference Energy Disaggregation Data (REDD) and energy consumption data collected from faculty housing at IIIT-Delhi. The results show that this technique performs the best by combining via product, all the HMMs (corresponding to each data stream) with binary states (on, off or standby) and training time linear to the number of HMMs.

## Categories and Subject Descriptors

H.4 [Information Systems Applications]: Miscellaneous

## General Terms

Algorithm, Experimentation

<sup>\*</sup>The author is also affiliated to the Institute of Data Science and Engineering (IDSE), Columbia University and is an adjunct professor at IIIT-Delhi.

## Keywords

energy aggregation, ensemble learning, load forecasting, product of HMMs, aggregate forecasting

## 1. INTRODUCTION

Smart meters consisting of real time sensors, power outage notifications and power quality monitoring are widely used today. These meters provide a host of benefits like energy efficiency and savings, improved retail competition, better demand response actions, improved tariffs, lower bills due to better customer feedback, accurate billing, less environmental pollution, etc. [?] They generate huge amount of time series data which can be used for gaining meaningful insights through analytics. They can measure site specific information and also help agencies set different electricity prices for consumption based on the time of the day, seasons, holidays, etc. Based on the data collected from smart meters, a feedback sent to the customers by the utilities that can help consumers better manage their resources. McKerracher et al. [?] show that by providing real time feedback, consumers can reduce the consumption by 3-5%.

In recent years, machine learning has been applied to the problem of energy consumption and demand forecasting analysis. The role of the machine learning algorithm is to study the sensor data and provide alerts and warnings when anomalous behaviour occurs or to inform (and remind) customers when certain activities were performed, which rooms they occupied, and what appliances they used most frequently during that period. This information can be transmitted to customers in timely fashion via phone, email or the Internet. Chicco et al. [?] compared several clustering techniques (such as hierarchical, KMeans) and observed that the hierarchical clustering and modified follow-the-leader perform best among the rest K-Means, fuzzy K-Means to group customers with similar electrical behaviour [?]. Wijaya et al. [?] used classifiers like random forest, decision trees (J48), logistic and naive bayes to identify customers with similar electricity consumption profiles. Related problems involve study of trends of electricity consumption (steadily increasing, decreasing, cyclic, seasonal) and sudden anomalous behaviour (sudden peaks or drops on consumption) for indi-

vidual homes and across the community [?].

In this paper, we use Hidden Markov models (HMMs) to analyse the time series energy data. We model the data stream from each source as a HMM with its states represented as ON/OFF. For  $N$  sources, there are  $N$  HMMs and the total number of states collectively are  $2^N$ . The observations represent the energy consumed in a particular state. These observations are recorded at different time scales for different sources.

In order to aggregate the data from all the different sources, we build a machine learning model using products of HMMs (PoHMMs) and apply it to the energy aggregation problem. There are many reasons why the product model constructed from many HMMs is appropriate. First, in a high-dimensional space each model constraints a different subset of dimensions but their product constraints all of the dimensions. Second, HMMs alone are not efficient at capturing long range structure in time series [?] – in contrast to PoHMMs [?] allow each model to remember a different piece of information about the past. Two different proof of concepts are presented – first one on the REDD<sup>1</sup> data set and the other one on real data collected at the faculty housing in India.

**Organization:** This paper is organized as follows: Section ?? examines related work on data analytics on aggregated data of smart meters; Section ?? provides a review of products of Hidden Markov Models (HMMs) and how they relate to our application. The two proofs of concepts are introduced in Section ?? illustrating the effectiveness of the use of product of HMMs in the energy aggregation problem. Finally, Section ?? concludes the work.

## 2. RELATED WORK

In this section, we describe work that uses ensemble learning techniques and non-ensemble learning techniques to solve problems in energy domain.

### 2.1 Non-ensemble based learning techniques

**Energy Aggregation** is a method of combining data from different sources such that several unreliable data measurements combine to produce a more accurate signal by enhancing the common signal and reducing the uncorrelated noise. As the sensor network generates lot of data for the end user to process, there are automated methods employed to aggregate data. This data fusion is generally known as data aggregation which combines the data into a set of meaningful information [?]. The sensor nodes are organised in a tree structure, called aggregation tree. The leaves of this tree are the sensor devices, the internal nodes are the aggregator devices that takes the data from the leaves, aggregates it and sends it to its parent node which is the root of the tree. The main objective of data aggregation is to reduce the unnecessary information thereby reducing the network traffic and improving the privacy of the customers from internal and external entities by keeping only the necessary information [?].

**Energy Disaggregation** or Non-intrusive appliance load monitoring (NIALM) is a process in which the household's

aggregate electricity consumption is broken down into individual appliance's consumption. The motivation behind this task is twofolds; first, the reduction in energy consumption by using the information given to the household occupants about the individual appliance's electricity consumption; second, recommending the occupants to defer the use of appliance to a time of day when electricity is cheaper. Approaches for energy disaggregation from smart meter using ML techniques fall in two categories; first group is where sub-metered data is available to train appliance models prior to performing disaggregation; second group uses unsupervised disaggregation methods that require labelling of detected appliance, assuming the type and number of household appliance. Parson et al. [?] proposed an approach to NIALM that separates the energy consumption of individual appliance iteratively from the aggregated load using hidden Markov models. Several studies have been done in this regard, one of the unsupervised disaggregation method [?] that outperformed other unsupervised disaggregation methods is conditional factorial hidden semi-Markov model. This model when integrated with other features, accurately represents the individual appliance energy consumption. In one of the work, Kolter et al. [?] showed that by discriminately training the sparse coding algorithms using a method based on structured prediction, the performance of the algorithm significantly improved on the energy disaggregation task. In another study, Kolter et al. [?] exploited the additive structure of the FHMM to develop a convex formulation of approximate inference algorithm that achieves state-of-the-art performance in energy disaggregation problem. A survey paper [?] reviewed several NILM methods making use of steady state and transient load signature in addition to the state-of-the art load disaggregation algorithms.

**Load Forecasting** refers to the projection of electrical load required in a certain geographical area with the use of previous electrical load usage in the same area. It is extremely important for efficient power system planning and operation, energy purchasing and generation, load switching, infrastructure development. It encompasses various factors like, historical load, weather data, population, energy supply and price, time of the year, etc. It is usually divided into three categories, short-term forecasts (one hour to one week), medium-term forecasts (one week to one year) and long-term forecasts (more than a year). In short term load forecast, [?] and [?] used a three layer feed forward artificial neural network to predict daily load profiles. In a paper by [?], nonlinear autoregressive integrated neural network was used to predict daily load consumption. In medium term load forecasts, the author forecasts [?] the monthly load through knowledge based activities from the output of the ANN based stage providing yearly energy predictions. Similarly, in [?], time lagged feedforward neural network is used to do monthly forecasting on the basis of historical series of electrical load, economic and demographic variables. Also, the authors from Covenant University, [?] performed load forecasting of their own educational institute using the models based on linear, compound growth and cubic methods of regression analysis. In long term load forecasting, a study done by [?] compared the performance of support vector regression and multilayer perceptron neural networks. The results showed that the percentage error in SVM reduced to 15% when compared to neural networks trained with back

<sup>1</sup><http://redd.csail.mit.edu/>

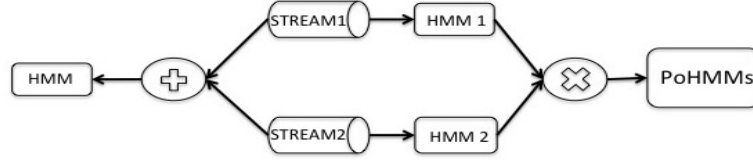


Figure 1: Problem Formulation

propagation algorithm. It was also observed that the parameter selection in SVM plays an important role in the performance of the model. Another work [?] uses support vector regression to derive non linear relationship between load and economic factors like GDP for long term forecasting in developing countries.

Load forecasting at the utility level can be done in three ways, that are completely aggregated, completely disaggregated and by using clustering for forecasting[CITE]. In completely aggregated method, the historical consumption data is combined first and then the prediction is performed. In completely disaggregated method, first the prediction at individual level data stream is performed and then combined together to compute the overall estimate at the utility level. In case of cluster based approach, individual data streams are clustered in specific segments, then segment level consumption is predicted which is finally combined to estimate the utility level consumption. In our approach, we have historical individual level data streams that are aggregated and hence used for next day load prediction using Product of HMMs as shown in figure ??

**Customer Segmentation** is the process of dividing a large homogeneous market into identifiable segments having similar demand characteristics. The appliance specific data has the potential to improve energy efficiency marketing by improving market segmentation, diversifying programs and transforming product development and evaluation [?]. This analysis is useful in various ways, like demand response system, intelligent distribution channel. The author [?] segments the customers based on contextual dimensions like location, seasons, weather patterns, holidays, etc which help with various higher level applications like usage-specific tariff structure, theft detection, etc. In [?], author proposes to infer occupancy states from the consumption data by using HMM framework. They investigate the effectiveness of HMM and model based cluster analysis in producing meaningful features of the classification.

## 2.2 Ensemble based learning techniques

Ensemble learning is a method where multiple learners are trained to solve the same problem. It constructs a set of hypothesis and combines them to generate the final result.

### 2.2.1 Prediction with expert advice

A study done by [?], proposed a Pattern Forecasting Ensemble Model (PFEM) comprising of five forecasting models using different clustering techniques, like k-means model, self-organising map model, hierarchical clustering model, k-medoids model and fuzzy c-means model. They have showed that on three real-world dataset, their proposed ensemble

model outperformed all the five individual model in case of day ahead electricity demand prediction. Another study [?] highlighted the importance of regularised negative correlation learning ensemble methodology on the problem of energy load hourly prediction. This method tried to overcome the problem of variability in neural network due to high sensitivtiness to the initial conditions. As this method combines the outputs of several neural networks, it achieves a marked reduction in error after introducing external data. An extension of HMMs, called Factorial Hidden Markov Model (FHMM) [?] is a class of ensemble based learning models that addresses the need for distributed hidden states in HMMs. The FHMM generalizes the HMM by representing the state using a collection, instead of single discrete variable. However, FHMMs being directed models, when conditioned on the observed sequence, the hidden state chains become independent making the inference easy but learning more complex. Thus, the exact inference becomes intractable, leaving one to resort to approximate inference techniques like Gibbs sampling or variational approximations.

## 3. REVIEW

A Hidden Markov Model (HMM) is a statistical Markov model that represents the probability distribution over a sequence of observations [?]. They are found useful in applications like speech [?], handwriting, gesture recognition, part-of-speech tagging, bioinformatics, etc. It has two properties, first, the observation at time  $t$ ,  $y_t$  is generated by a process whose state at time  $t$ ,  $s_t$  is hidden from the observer and second, is that this hidden state process satisfies Markov property which states that given the value at state  $s_{t-1}$ , the value at current state  $s_t$  is independent of all the states prior to  $t - 1$ . The subscripts  $i$  and superscripts  $j$  indicate the model at  $i$ th time and the  $j$ th HMM. The state space of the HMM is discrete, that is a state can take 2 values denoted by ON and OFF. The observed values represent the aggregated load/energy collected from different data streams at time  $t$ . In order to define probability distribution over the sequence of observation, it is important to define probability distribution over the initial state  $P(s_1)$ , the transition probability  $P(s_t|s_{t-1})$  and the observed probability  $P(y_t|s_t)$  where  $y_t$  is the observation at time  $t$ . Following a notation in [?], HMM is composed of a 3-tuple  $\{A, B, \pi\}$  where  $A$  is the transition probability,  $B$  is the observed probability and  $\pi$  is the initial state probability. HMMs solve three fundamental problems: 1. Given the model  $\lambda = \{A, B, \pi\}$ , and observation sequence  $Y = \{y_1, \dots, y_T\}$ , how do we efficiently compute the probability of the sequence of observations given the model, that is  $P(Y|\lambda)$ . 2. Given the model  $\lambda$  and observation sequence  $Y$ , what is the

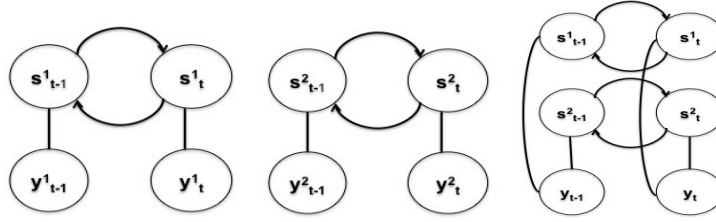


Figure 2: HMM  $S^1$  and  $S^2$  and PoHMMs,  $P = S^1 \times S^2$

Floor 0-5	Floor6 – 11	Total
00:01:31+05:30, 786.51	00:01:30+05:30, 1867.34	00:01:29+05:30, 3369.30
00:02:01+05:30, 787.24	00:02:00+05:30, 1850.54	00:01:59+05:30, 3343.39
00:02:32+05:30, 787.54	00:02:31+05:30, 1832.85	00:02:30+05:30, 3339.27

Table 1: Timestamp readings of energy consumption for faculty housing building

underlying state sequence  $\{s_1, \dots, s_T\}$  that best explains the observations. 3. How do we adjust the model parameters  $(A, B, \lambda)$  so as to maximize the probability of observation sequence given the model  $P(Y|\lambda)$ . There is no optimal way to estimate the model parameters given any infinite observation sequence. However, using an iterative procedure like Baum-Welch or gradient techniques, model parameters  $\lambda$  can be chosen such that  $P(O|\lambda)$  is locally maximized. In this paper, we deal with the third problem as it involves learning parameters by training the model with the past energy consumption data and then using these parameters to perform load forecasting.

**An Example:** Figure ?? shows the HMMs  $S^1$  and  $S^2$  generated by two data streams. The energy consumed by first six floors and the top six floors of the faculty housing building is represented by the two data streams,  $S^1$  and  $S^2$  individually. In order to know the energy consumption of all the 12 floors, simple addition of two data streams will not work as the readings in both the data streams are not synchronized in time. Table ?? shows the timestamp readings of the energy consumed by the two data streams and the entire building. As can be seen in table ??, the recording in both the data streams is shifted by 1 second. Also, the recording for each data stream is sampled every 30 seconds.

### 3.1 Product of HMMs

PoHMM is a model that combines several HMMs by multiplying their individual distribution together and then renormalizing them as can be seen in equation ???. Here,  $d$  is a vector in discrete space,  $\theta_m$  is all parameters of individual model  $m$ ,  $f_m(d|\theta_m)$  is the probability of  $d$  under model  $m$ , and  $c$  indexes all the possible vectors in the data space. Its representation includes both directed and undirected links where the hidden states are causally connected to the other hidden states but non causally related to the visible states. This causes different conditional independence relationships among the variables in graphical model. The figure ?? is a product of two HMMs  $P = S^1 \times S^2$  where the superscript in  $S^1$  indicates the  $k$ th HMM. The number of states in the PoHMM is the product of states in  $S^1$  and  $S^2$  which is 4 in our case. The connections formed in the  $P$  depend on the links in the multiplying HMMs.

$$p(d|\theta_1 \dots \theta_n) = \frac{\prod_m f_m(d|\theta_m)}{\sum_c \prod_m f_m(c|\theta_m)} \quad (1)$$

### 3.2 Training the model by minimising contrastive divergence

To fit the model to the data, we need to maximize the likelihood of the dataset or minimise the Kullback-Liebler (KL) divergence between the data distribution,  $P^0$  (distribution at time 0) and  $P_\theta^\infty$  (also written as  $p(d|\theta_1 \dots \theta_n)$ ) which is the equilibrium distribution over the visible variables. KL divergence is defined as the non-symmetric measure of the difference between two probability distributions  $P^\infty$  and  $P^0$ . It measures the information lost when  $P^\infty$  is used to approximate  $P^0$  as shown mathematically in equation ??.

$$P^0 || P_\theta^\infty = \sum_d P^0(d) \log P^0(d) - \sum_d P^0(d) \log P_\theta^\infty(d) \quad (2)$$

$$= -H(P^0) - \langle \log P_\theta^\infty \rangle_{P^0}$$

where  $||$  represents Kullback-Leibler divergence,  $d$  is the data vector in discrete space,  $\theta_m$  is all the parameters of individual model  $m$ ,  $P^0$  is the data distribution at time 0,  $P_\theta^\infty$  (fantasy data) is the equilibrium distribution obtained after prolonged Gibbs sampling (figure ??),  $H(P^0)$  represents the entropy which is ignored during optimisation as  $P^0$  does not depend on the parameters of the model, angle brackets denote the expectation over the distribution specified by the subscript.

In Gibbs sampling, each variable draws a sample from its posterior distribution given the current states of the other variables. The hidden states of all the models are conditionally independent given the data and hence can be parallel updated as shown in Figure ??. This contributes to an important consequence of product formulation. At time  $t=0$ , the observed variables represent a data vector,  $d$  and the hidden variables,  $s$  of all the models are updated in parallel with samples from their posterior distribution given the

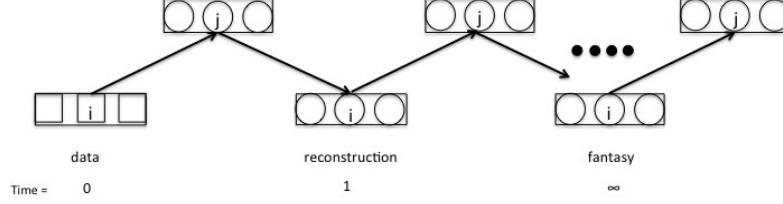


Figure 3: Visualization of Gibbs sampling

observed variables,  $y$ . At time 1, the visible variables are updated to generate a reconstruction of the original data vector from the hidden variables and the hidden variables are again updated simultaneously. This prolonged sampling helps the Markov chain to converge to the equilibrium distribution which helps to attain the unbiased estimate of the gradient of the PoHMMs (equation ??) where  $D$  corresponds to the data,  $\log f_{\theta_m}$  is also be written as  $f_m(D|\theta_m)$ .

$$\left\langle \frac{\partial \log P^\infty(D)}{\partial \theta_m} \right\rangle_{P^0} = \left\langle \frac{\partial \log f_{\theta_m}}{\partial \theta_m} \right\rangle_{P^0} - \left\langle \frac{\partial \log f_{\theta_m}}{\partial \theta_m} \right\rangle_{P_\theta^1} \quad (3)$$

Since the samples from the equilibrium state require computation and also have high variance as they come from the entire model's distribution, it poses a difficulty in determining the estimate the derivative. Therefore, the optimisation is performed on the different objective function called contrastive divergence, defined in equation ?. Contrastive divergence is the difference between  $P^0||P_\theta^\infty$  and  $P_\theta^1||P_\theta^\infty$  where  $P_\theta^1$  is the distribution over the one-step reconstruction of the data vectors generated by one full step of Gibbs sampling. The intuition behind using contrastive divergence is to leave the initial distribution  $P^0$  over the visible variables unaltered and also the intractable expectation over  $P_\theta^\infty$  gets cancelled out. Instead of comparing the initial and final derivatives,  $P^0$  and  $P_\theta^\infty$ , the Markov chain is run for one full step and the parameters are updated to avoid the chain to wander away from the initial distribution on the first step. As  $P^1$  is a step closer to  $P^\infty$  which guarantees that  $P^0||P_\theta^\infty$  will always exceed  $P_\theta^1||P_\theta^\infty$  ensuring a non negative value unless  $P^0 = P_\theta^1$ . If  $P^0 = P_\theta^1$ , then it implies that the chain is already in an equilibrium state, that is  $P^0 = P_\theta^\infty$  hence making the value of contrastive divergence as 0.

$$-\frac{\partial}{\partial \theta_m} (P^0||P_\theta^\infty - P_\theta^1||P_\theta^\infty) = \left\langle \frac{\partial \log f_{\theta_m}}{\partial \theta_m} \right\rangle_{P^0} - \left\langle \frac{\partial \log f_{\theta_m}}{\partial \theta_m} \right\rangle_{P_\theta^1} + \frac{\partial P_\theta^1}{\partial \theta_m} \frac{\partial (P_\theta^1||P_\theta^\infty)}{\partial P_\theta^1} \quad (4)$$

In equation ??, the first two terms on the right hand side are tractable as it is easy to sample from  $P^0$  and  $P_\theta^1$  but the third term represents the effect on  $P_\theta^1||P_\theta^\infty$  of the change of the step reconstruction caused by the change in the  $\theta_m$ . Extensive simulations show that it is small and rarely differs from the result of other two terms, hence can be safely ignored. Therefore in contrastive divergence, the parameters are learned according to the equation ??.

$$\Delta \theta_m \propto \left\langle \frac{\partial \log f_{\theta_m}}{\partial \theta_m} \right\rangle_{P^0} - \left\langle \frac{\partial \log f_{\theta_m}}{\partial \theta_m} \right\rangle_{P_\theta^1} \quad (5)$$

The contrastive divergence algorithm for training the PoHMM has the following steps:

Each model's gradient  $\frac{\partial}{\partial \theta_m} P(Y|\theta_m)$  ( $Y = \{y_t\}_{t=1}^T$  is a visible variable) is calculated on a data point using forward backward algorithm. For each model, a sample is taken from the posterior distribution of paths through state space. At each time step, the distributions are multiplied and renormalized together to get the reconstruction distribution. A sample from the reconstruction distribution is drawn at each time step to get a reconstructed sequence. Each model's gradient is computed on the new sequence  $P(\hat{Y}|\theta_m)$ . Parameters are updated as per equation ??

### 3.3 Inference in PoHMM

The main feature of PoHMMs is its undirected graphical modelling with no direct connection among the latent variables ( $S_t^1$  and  $S_t^2$ ) as they only interact indirectly via observed variables ( $Y_t$ ). The hidden variables all the experts are rendered independent when conditioned on visible variables. So, if the inference in each of the constituent model is tractable then the inference in the product is also tractable. To generate a data point in this model, all the models in PoHMMs generate an observation and if they all generated the same point then it is accepted else they again generate an observation until all the models agree to it. Therefore all the models have some influence over the generated data. So, the inference determines the the probability that all the models would have taken in order to generate the given observation.

## 4. APPLICATIONS USING PRODUCT OF HMMS

**Aim:** In this section we demonstrate how PoHMMs can be used to model data streams and perform load forecasting. Proof-of-concepts are provided on two data sets - REDD<sup>2</sup> and on the energy data collected from faculty housing at IIIT Delhi.

### 4.1 Data Description

- **Dataset 1:** The Reference Energy Disaggregation Data Set (REDD) contains power consumption data from real homes, for the whole house as well as for each individual circuit in the house (labeled by the main type

<sup>2</sup><http://redd.csail.mit.edu/>

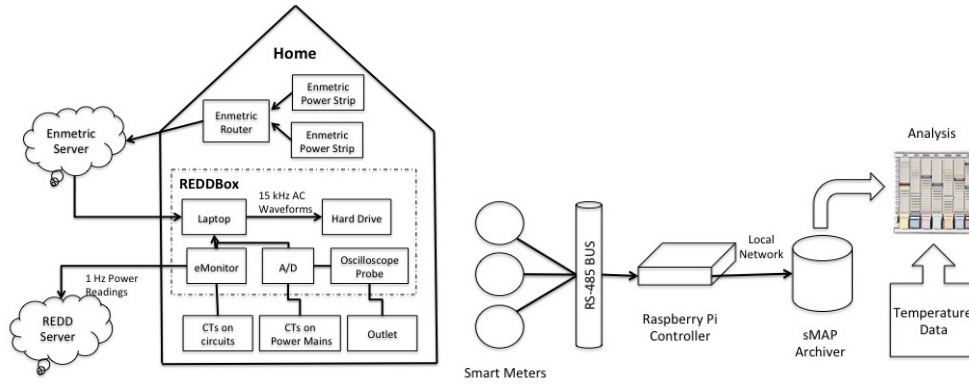


Figure 4: Schematic diagram of REDD and Faculty housing building

of appliance on that circuit). The REDD data set contains two main types of home electricity data: high-frequency current/voltage waveform data of the two power mains (as well as the voltage signal for a single phase), and lower-frequency power data including the mains and individual, labeled circuits in the house. Its schematic diagram is shown in figure ?? . Experiments reported here use ‘house 2’ data from REDD. This dataset has 318759 records and 2 columns. We randomly sample 300 records for our initial experiment.

- **Dataset 2:** This data represents the energy consumed by the IIT Delhi’s faculty housing building. As a part of research, a team from the institute has installed various temperature, light and motion sensors to perform real world studies and to analyse user preferences for energy conservation. Its schematic diagram is shown in figure ?? . For our analysis, we selected 1 month’s historical data ranging from 01-01-2014, 00:01 hours to 31-01-2014, 23:59 hours. The two smart meters installed, captured data from all the 12 floors. The first meter records readings from ground to 5th floor generating one stream of data and the second meter generates a stream from 6th to 11th. Both these streams are aggregated to obtain the aggregated load of the housing building. The dataset includes timestamp and power consumed in watts. There are 84133 records in this dataset. We also have the total energy consumed by the faculty housing building which would serve as the ground truth to compare the aggregated load using PoHMMs.
- **Dataset 3:** The third dataset that we have used is the Enernoc data set<sup>3</sup>, which consists of power consumption data from different industrial customers. This contains energy consumption data of 100 industrial consumers from January 2012 to December 2012. In our analysis, we have used 66% of data for training and 33% for testing. We are comparing the load forecast by using PoHMMs vs single HMM. The metric used to quantify the similarity between the distribution of both the methods is KL divergence (refer equation??). The comparisons using both the above approaches are done to observe the behaviour at different time scales,

number of consumers and type of consumers. Figure ?? shows how KL Divergence gets affected with varying data samples and experts for 1 month, 4 months and 6 months data. We can see that KL is

## 4.2 Problem Formulation

The disparate energy data streams collect readings at different time scales. Each of the data stream is modelled as a HMM with cardinality 2, that is either ON or OFF. The process of aggregating the energy data from different data streams is modelled through PoHMMs.

Each energy data stream is used to train the model, till the time the objective function, that is contrastive divergence reaches a threshold value. Once the model is trained from a randomly sampled data stream, the parameters learned by the model are provided to the randomly sampled test set (total energy consumed data) to obtain the conditional probability distribution of the gaussians given the data. Similarly, all the data streams are used for training the model, and the parameters learned are then applied on the test set to obtain the conditional probability of the gaussians given the data. After all the data streams are used to obtain the probability distribution, we use the data stream that correspond to the total energy consumed from the house/ building to train the model and hence obtain the probability distribution  $P$  of the gaussians. These probability distributions are then compared with the product of the probability distributions  $Q$  obtained from the individual data streams. The evaluation of how well the learning has taken place is done by using a Kullback-Leibler divergence. KL divergence of two probability distributions  $P$  and  $Q$ ,  $D_{KL}(P||Q)$  is the measure of information lost when  $Q$  is used to approximate  $P$ . Figure ?? shows a flowchart that depicts the problem formulation.

### 4.2.1 Implementation Details

The implementation of the product of hidden Markov model is obtained from Iain Murray’s website<sup>4</sup>. It implements the technique described in paper [?].

## 4.3 Empirical Results

In both the datasets, experiments are performed by tuning some parameters and keeping some fixed. In case of REDD,

<sup>3</sup><http://open.enernoc.com/data/>

<sup>4</sup><http://homepages.inf.ed.ac.uk/imurray2/code/>

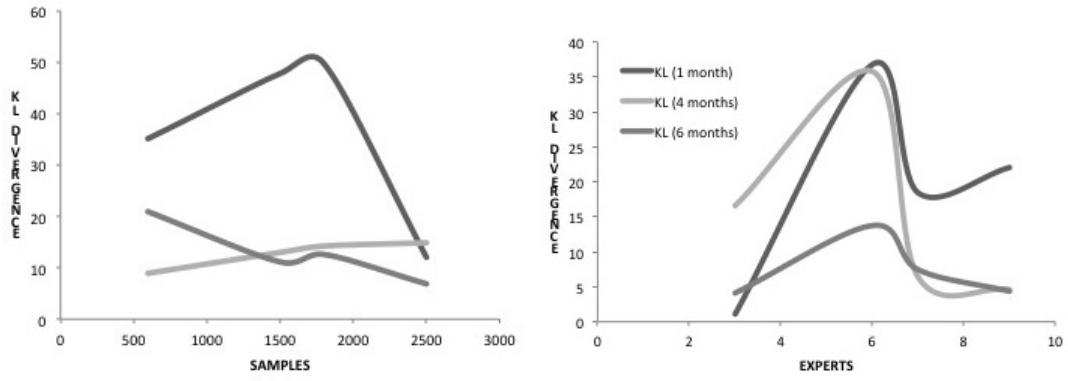


Figure 5: Plots using Enernoc data set

Samples	$KLDiv$	$Iterations$	$T(sec)$
300	2.4864	18600	$186.212 \pm 9.087$
500	0.6761	10200	$106.564 \pm 10.046$
1000	1.1088	11200	$158.521 \pm 1.97$
1500	3.8829	5300	$92.896 \pm 8.075$
2000	1.8686	6900	$130.98 \pm 1.932$
2500	0.4733	9900	$215.563 \pm 2.471$
3000	2.8204	11000	$258.213 \pm 1.918$
3500	1.2332	7900	$204.661 \pm 1.713$
4000	0.8959	10400	$292.666 \pm 0.619$
4500	1.1118	7200	$222.558 \pm 1.967$
8000	6.392	8100	$381.635 \pm 2.952$
10000	8.276	10500	$887.932 \pm 13.824$
15000	0.7201	9400	$1368.514 \pm 13.605$

Table 2: Effect of varying the samples (REDD)

Samples	$KLDiv$	$Iterations$	$T(sec)$
100	0.26219	45100	257
300	0.19753	43200	222
500	0.55493	44800	260
700	0.32847	44000	249
900	3.9486	42600	221
1100	4.9274	44700	317
1300	3.0425	43100	276
1500	3.1128	44400	303
2000	1.9192	44400	306
2500	1.7122	44100	370
3000	1.4686	43300	331
3500	1.2663	43200	370
4000	1.0793	43200	403

Table 3: Effect of varying the samples (housing data), KL Div in e-04

Samples	$T_P$	$T_H$	$Iter_P$	$Iter_H$	$KLDiv$
600	31	24	5600	13100	35.01
1200	41	15	5300	5500	28.85
1500	68	16	7700	5400	47.78
1800	52	112	5400	34300	49.72
2100	69	40	6100	11200	35.75
2500	83	25	6900	5500	12.0329

Table 4: Effect of varying the samples (Enernoc)

Experts	$T_P$	$T_H$	$Iter_P$	$Iter_H$	$KLDiv$
3	26	41	5700	28200	1.04
4	49	153	10233	98200	4.35
5	158	36	32000	22300	20.72
6	76	10	15000	6000	36.74
7	117	2798	21000	1430000	18.445
9	455	420	70500	189000	22.04

Table 5: Effect of varying the HMMs (Enernoc)

Tables ??, ?? and ?? show the effect of varying data samples, threshold and no. of HMMs on KL divergence. In this experiment, data samples are randomly selected from the dataset, KL divergence is chosen as the performance metric, iterations are performed to obtain the samples from the equilibrium distribution by using Gibbs sampler on the hidden and visible variables and the average time taken (sec) with the standard deviation on three such iterations is noted.

In case of faculty housing dataset, Tables ??, ?? and ?? show the effect of varying data samples, threshold and no. of HMMs on the KL divergence. As per table ?? and ??, there is no linear relationship between no. of samples and KL divergence. The best performance was attained when 2500 and 300 samples were randomly chosen from the REDD and faculty housing dataset respectively. In table ?? and ??, there is an indirect relationship between threshold and KL divergence value, that is lower the threshold, higher the KL divergence. In table ?? and ??, it can be seen that the KL divergence attains its minimum value when the no. of HMMs multiplied are same as the number of data streams aggregated. Hence, the error is minimum when the number of HMMs reaches the number of aggregated data streams, that is when all the appliances in ‘house 2’ (total appliances are 9) of REDD are used in PoHMMs and when both the data streams of faculty housing building are aggregated together.

## 5. CONCLUSION

In this paper, we are solving the problem of load aggregation (energy consumption) for disparate energy data sources using the ensemble based learning technique called Product of HMMs. Several challenges are faced while computing the aggregated energy, few of which include non-aligned timestamp readings, missing values, meter reset readings. This technique produces the combined probability distributions of several simpler distributions and then renormalizes this output. The optimisation problem is to minimise the contrastive divergence between the two probability distributions of the data at time 0 and the data at time 1 (one-step reconstruction using gibbs sampling). The evaluation of the algorithm is done by computing the KL divergence between the product of the energy data distributions and the total energy data distribution. From the results shown in table ?? and ??, we can see that the algorithm has performed best when the number of appliances (REDD) reached the actual number of appliances (9) in the house and the number of data streams in case of housing data reached the actual number of streams (2) used in the 12 floor faculty housing data. Therefore, this signifies that this algorithm is applicable for this kind of application and can be used for load forecasting.

## Acknowledgment

The authors would like to thank *EMC<sup>2</sup>* for supporting this research through grant number JRA092014-Q3-5.

## 6. REFERENCES

- [1] R. Achata. Long term electric load forecasting using neural networks and support vector machines. *IJCST*, 3(1), 2012.
- [2] A. Albert and R. Rajagopal. Smart meter driven segmentation: What your consumption says about you. *Power Systems, IEEE Transactions on*, 28(4):4019–4030, Nov 2013.
- [3] K. C. Armel, A. Gupta, G. Shrimali, and A. Albert. Is disaggregation the holy grail of energy efficiency? the case of electricity. *Energy Policy*, 52(0):213 – 234, 2013. Special Section: Transition Pathways to a Low Carbon Economy.
- [4] A. Bakirtzis, V. Petridis, S. Kiartzis, and M. Alexiadis. A neural network short term load forecasting model for the greek power system. *Power Systems, IEEE Transactions on*, 11(2):858–863, May 1996.
- [5] D. Bassi and O. Olivares. Medium term electric load forecasting using tlfn neural networks. *International Journal of Computers, Communications & Control*, 1(2):23–32, 2006.
- [6] A. Brown and G. Hinton. Proceedings of artificial intelligence and statistics 2001. In *Products of Hidden Markov Models*, number GCNU TR 2000-008, 2001.
- [7] A. D. Brown. Product model for sequences. *Gatsby Unit*, 2001.
- [8] H. Chen, C. Canizares, and A. Singh. Ann-based short-term load forecasting in electricity markets. In *Power Engineering Society Winter Meeting, 2001. IEEE*, volume 2, pages 411–415 vol.2, 2001.
- [9] G. Chicco, R. Napoli, and F. Piglion. Comparisons among clustering techniques for electricity customer classification. *Power Systems, IEEE Transactions on*, 21(2):933–940, May 2006.
- [10] T. Chow and C. Leung. Nonlinear autoregressive integrated neural network model for short-term load forecasting. *Generation, Transmission and Distribution, IEE Proceedings-*, 143(5):500–506, Sep 1996.
- [11] D. J. Cook and C. Chen. Teasing detailed home habits from aggregate energy consumption data. *IEEE Smart Grid*, February 2012.
- [12] M. Falvo, R. Lamedica, S. Pierazzo, and A. Prudenzi. A knowledge based system for medium term load forecasting. In *Transmission and Distribution Conference and Exhibition, 2005/2006 IEEE PES*, pages 1291–1295, May 2006.
- [13] M. D. Felice and X. Yao. Short-term load forecasting with neural network ensembles: A comparative study [application notes]. *IEEE Comp. Int. Mag.*, 6(3):47–56, 2011.
- [14] Z. Ghahramani. Hidden markov models. In *World Scientific Series In Machine Perception And Artificial Intelligence Series*, chapter An Introduction to Hidden Markov Models and Bayesian Networks, pages 9–42. World Scientific Publishing Co., Inc., River Edge, NJ, USA, 2002.
- [15] Z. Ghahramani and M. I. Jordan. Factorial hidden markov models. *Mach. Learn.*, 29(2-3):245–273, Nov. 1997.
- [16] W. R. Heinzelman, A. Ch, and H. Balakrishnan. Energy-efficient communication protocol for wireless microsensor networks. In *System Sciences*, pages 3005–3014, 2000.
- [17] G. E. Hinton. Training products of experts by minimizing contrastive divergence. *Neural Computation*, 14(8):1771–1800, 2002.



Threshold	<i>KLDiv</i>	<i>Iterations</i>	<i>T(sec)</i>
.1	0.473	9900	210.6 $\pm$ 1.493
.05	0.443	10900	240.607 $\pm$ 2.436
.01	0.454	18000	431.536 $\pm$ 14.509
.005	0.509	49800	1167.243 $\pm$ 43.412

**Table 6: Effect of varying threshold (REDD)**

HMMs	<i>KLDiv</i>	<i>Iterations</i>	<i>T(sec)</i>
3	5.559	10700	233.664 $\pm$ 0.579
4	0.188	19900	465.634 $\pm$ 5.275
5	0.432	13400	338.416 $\pm$ 3.988
6	8.736	28100	606.062 $\pm$ 7.534
7	5.054	17300	411.457 $\pm$ 10.051
8	0.436	10700	260.544 $\pm$ 27.862
9	0.15	20600	474.579 $\pm$ 14.619

**Table 8: Effect of varying HMMs (REDD)**

Threshold	<i>KLDiv</i>	<i>Iterations</i>	<i>T(sec)</i>
0.9	0.768	49300	153.67 $\pm$ 13.57
0.8	0.768	49300	190.3 $\pm$ 34.35
0.7	0.769	49400	193.67 $\pm$ 31.64
0.6	0.792	49500	218.3 $\pm$ 43.93
0.5	0.793	49600	250.67 $\pm$ 3.51
0.4	0.854	49600	191 $\pm$ 35.55

**Table 7: Effect of varying threshold (housing data), KL Div in e-04**

HMMs	<i>KLDiv</i>	<i>Iterations</i>	<i>T(sec)</i>
2	0	49200	233.33 $\pm$ 7.23
3	0.81	49200	237.67 $\pm$ 10.59
5	0.77	49300	242.33 $\pm$ 3
10	2.12	49300	228.33 $\pm$ 46.23
15	2.55	49300	269.33 $\pm$ 17.61
20	2.23	49300	305.66 $\pm$ 2.08
25	2.19	49300	302.33 $\pm$ 8
30	2.24	49300	336.33 $\pm$ 9.5
35	1.94	49300	331.33 $\pm$ 9

**Table 9: Effect of varying HMMs (housing data), KL Div in e-04**

- [18] R. Kawamoto, A. Nazir, A. Kameyama, T. Ichinomiya, K. Yamamoto, S. Tamura, M. Yamamoto, S. Hayamizu, and Y. Kinosada. Hidden markov model for analyzing time-series health checkup data. In *MedInfo*, volume 192 of *Studies in Health Technology and Informatics*, pages 491–495. IOS Press, 2013.
- [19] H. Kim, M. Marwah, M. F. Arlitt, G. Lyon, and J. Han. Unsupervised disaggregation of low frequency power measurements. In *SIAM*, pages 747–758. SIAM / Omnipress, 2011.
- [20] F. Klopfer and G. Wallenborn. Empowering consumers through smart metering. *BEUC, The European Consumer Organisation*, 2011.
- [21] J. Z. Kolter, S. Batra, and A. Y. Ng. Energy disaggregation via discriminative sparse coding. In J. Lafferty, C. Williams, J. Shawe-Taylor, R. Zemel, and A. Culotta, editors, *Advances in Neural Information Processing Systems 23*, pages 1153–1161. Curran Associates, Inc., 2010.
- [22] J. Z. Kolter and J. Ferreira. A large-scale study on predicting and contextualizing building energy usage. In *AAAI*. AAAI Press, 2011.
- [23] J. Z. Kolter and T. Jaakkola. Approximate inference in additive factorial hmms with application to energy disaggregation. In *AISTATS*, volume 22 of *JMLR Proceedings*, pages 1472–1482. JMLR.org, 2012.
- [24] F. Martinez Alvarez, A. Troncoso, J. Riquelme, and J. Aguilar Ruiz. Energy time series forecasting based on pattern sequence similarity. *Knowledge and Data Engineering, IEEE Transactions on*, 23(8):1230–1243, Aug 2011.
- [25] C. McKerracher and J. Torriti. Energy consumption feedback in perspective: integrating australian data to meta-analyses on in-home displays. *Energy Efficiency*, 6(2):387–405, 2013.
- [26] P. Mirowski, S. Chen, T. Kam Ho, and C.-N. Yu. Demand forecasting in smart grids. *Bell Labs Technical Journal*, 18(4):135–158, 2014.
- [27] O. Parson, S. Ghosh, M. Weal, and A. Rogers. Using hidden markov models for iterative non-intrusive appliance monitoring. In *Neural Information Processing Systems workshop on Machine Learning for Sustainability*, December 2011. Event Dates: 17 December 2011.
- [28] L. R. Rabiner. Readings in speech recognition. In A. Waibel and K.-F. Lee, editors, *Readings in speech recognition*, chapter A Tutorial on Hidden Markov Models and Selected Applications in Speech Recognition, pages 267–296. Morgan Kaufmann Publishers Inc., San Francisco, CA, USA, 1990.
- [29] I. A. Samuel, A. A. Awelewa, et al. Medium-term load forecasting of covenant university using the regression analysis methods. *Journal of Energy Technologies and Policy*, 4(4):10–16, 2014.
- [30] W. Shen, V. Babushkin, Z. Aung, and W. L. Woon. An ensemble model for day-ahead electricity demand time series forecasting. In *Proceedings of the Fourth International Conference on Future Energy Systems, e-Energy '13*, pages 51–62, New York, NY, USA, 2013. ACM.
- [31] S. Shi and A. Weigend. Taking time seriously: Hidden Markov experts applied to financial engineering. In *CIFER '97: Proc. of the Conf. on Computational Intelligence for Financial Engineering*, pages 244–252. IEEE, 1997.
- [32] D. I. Stern. The role of energy in economic growth. *Annals of the New York Academy of Sciences*, 1219(1):26–51, 2011.
- [33] G. Taban and A. A. C?rdenas. Data aggregation as a method of protecting privacy in smart grid networks. *Power Systems, IEEE Transactions on*, March 2012.
- [34] G. W. Taylor and G. E. Hinton. Products of hidden markov models: It takes n >1 to tango. In *Proceedings of the Twenty-Fifth Conference on Uncertainty in Artificial Intelligence*, UAI '09, pages 522–529,

Arlington, Virginia, United States, 2009. AUAI Press.

- [35] T. K. Wijaya, J. Eberle, and K. Aberer. Symbolic representation of smart meter data. In *Proceedings of the Joint EDBT/ICDT 2013 Workshops*, EDBT '13, pages 242–248, New York, NY, USA, 2013. ACM.
- [36] T. K. Wijaya, T. Ganu, D. Chakraborty, K. Aberer, and D. P. Seetharam. Consumer segmentation and knowledge extraction from smart meter and survey data. In *SIAM International Conference on Data Mining (SDM14)*, number EPFL-CONF-196276, pages 226–234, 2014.
- [37] Z. Zhang and S. Ye. Long term load forecasting and recommendations for china based on support vector regression. In *Information Management, Innovation Management and Industrial Engineering (ICIII), 2011 International Conference on*, volume 3, pages 597–602, Nov 2011.
- [38] A. Zoha, E. Gluhak, M. A. Imran, and S. Rajasegarar. Article non-intrusive load monitoring approaches for disaggregated energy sensing: A survey, 2012.