

# Long Term Probabilistic Load Forecasting and Normalization With Hourly Information

Tao Hong, Jason Wilson, *Member, IEEE*, and Jingrui Xie, *Associate Member, IEEE*

**Abstract**—The classical approach to long term load forecasting is often limited to the use of load and weather information occurring with monthly or annual frequency. This low resolution, infrequent data can sometimes lead to inaccurate forecasts. Load forecasters often have a hard time explaining the errors based on the limited information available through the low resolution data. The increasing usage of smart grid and advanced metering infrastructure (AMI) technologies provides the utility load forecasters with high resolution, layered information to improve the load forecasting process. In this paper, we propose a modern approach that takes advantage of hourly information to create more accurate and defensible forecasts. The proposed approach has been deployed across many U.S. utilities, including a recent implementation at North Carolina Electric Membership Corporation (NCEMC), which is used as the case study in this paper. Three key elements of long term load forecasting are being modernized: predictive modeling, scenario analysis, and weather normalization. We first show the superior accuracy of the predictive models attained from hourly data, over the classical methods of forecasting using monthly or annual peak data. We then develop probabilistic forecasts through cross scenario analysis. Finally, we illustrate the concept of load normalization and normalize the load using the proposed hourly models.

**Index Terms**—Load forecasting, load normalization, multiple linear regression models, weather normalization.

## I. NOMENCLATURE

<i>GSP</i> :	Gross state product.
<i>CDD</i> :	Cooling degree days.
<i>HDD</i> :	Heating degree days.
<i>Trend</i> :	A linear trend variable.
$T_{\max}$ :	Monthly peak temperature.
$T_t$ :	Current hour temperature.
$T_{t-k}$ :	Temperature of the previous $k$ th hour.
$T_a$ :	Average temperature of the past 24 hours.
<i>Month</i> :	Class variable, 12 months of the year.
<i>Weekday</i> :	Class variable, 7 days of a week.

Manuscript received February 14, 2013; revised May 03, 2013, June 18, 2013; accepted July 18, 2013. Date of publication September 10, 2013; date of current version December 24, 2013. Paper no. TSG-00120-2013.

T. Hong and J. Xie are with SAS Institute, Cary, NC 27513 USA (e-mail: hongtao01@gmail.com; rain.xie@sas.com).

J. Wilson is with North Carolina Electrical Membership Corporation, Raleigh, NC 27616 USA (e-mail: jason.wilson@ncemcs.com).

Color versions of one or more of the figures in this paper are available online at <http://ieeexplore.ieee.org>.

Digital Object Identifier 10.1109/TSG.2013.2274373

*Hour*: Class variable, 24 hours of a day.

*Day*: Class variable, code for days of a year.

## II. INTRODUCTION

LONG TERM LOAD forecasting (LTLF) provides peak demand and energy forecasts for one or more years, and can be expanded out to a horizon of a few decades. Utilities typically produce long term forecasts ranging from 20 to 50 years into the future. Such forecasts are often being used for planning by multiple departments in a utility, such as system planning, finance, demand side management, and power supply, etc. North Carolina Electric Membership Corporation (NCEMC), is one of the largest electric generation cooperatives in the U.S., and is comprised of a family of corporations formed to support 26 of North Carolina's electric distribution cooperatives. These cooperatives provide energy and related services to more than 950 000 households and businesses in 93 of North Carolina's 100 counties. At NCEMC, long term load forecasts serve as the important inputs to the power supply group to support decisions on electricity purchase contracts. Because NCEMC owns generation units, it is required to file Integrated Resource Planning (IRP) documents with the North Carolina Public Utilities Commission, thus expanding the scrutiny from the member cooperatives and NCEMC board, to the state regulatory commission.

In the regulatory environment, utility forecasters have to defend the long term forecasts internally to the utility's management and externally to the regulatory commission. Although forecasting by nature is a stochastic problem, most utilities today are still developing and using point forecasts instead of probabilistic forecasts. Due to the poor predictability of the climate, which is a main driver of electricity demand, it is unrealistic and unfair to judge a long term forecaster by comparing a few years of point forecasts with the corresponding actual values. Instead, there are two important questions that should be asked and answered properly when defending the long term forecasts: 1) *is the current scenario covered by the forecasts?* 2) *how accurate is the forecast given the current scenario?*

Most utilities today follow the LTLF practices similar to the ones established a few decades ago, when there was not high resolution data available. Since the type of low resolution data used in the traditional approach provides a limited number of observations for predictive modeling, the forecasters may not be able to use enough explanatory variables to capture all the salient features of electric load. When given the actual values of the weather and economy variables to re-forecast the loads under the current scenario, the model may still produce some significant errors, which can be hard to explain by the forecasters.

There are a few ways to create weather scenarios for LTLF and weather normalization. A lot of utilities are using the average temperature profile (in hourly or daily interval for a year) from the previous few decades as the normal weather to derive the normal load, which is not a defensible approach: 1) an average temperature profile understates the peaks, so it can not accurately represent the normal weather; 2) a normal weather profile may not lead to a normal load profile due to the nonlinear relationship between the load and weather [1]. Another popular approach is to use normal or typical weather profiles created by third parties, such as National Oceanic and Atmospheric Administration (NOAA) and the U.S. Department of Energy (DOE). An advantage of this third-party profile approach is its simplicity. However, these weather profiles are not created specifically for utilities to calculate normalized load. Therefore, it is questionable that their best use is for normalizing the load profile in the utility industry. A more rigorous approach is based on Monte Carlo simulation, which is often adopted by the risk management teams in the utility industry. The quantitative risk analysts first analyze the distribution of temperatures on each hour of the year. They then create thousands of temperature profiles for scenario analysis. This simulation approach requires a lot of computational resources. The results, including thousands of load profiles, are sometimes too voluminous and become difficult to understand and be used by the system operators in practice.

Most literature in the load forecasting field has been devoted to short term load forecasting, of which the forecasting horizon is two weeks or less [1]–[7]. Not many papers have been devoted to LTLF, of which few papers present practical approaches verified through field implementations at utilities. An implementation of spatial load forecasting work at Madison Gas and Electric Company has been presented in [8]–[10]. A peak load forecasting methodology implemented at Australian Energy Market Operator (AEMO) has been reported in [11]. In this paper, we propose a probabilistic forecasting approach with hourly data, which is the continuation of Hong's load forecasting methodology presented in [1]. We dissect LTLF to three elements: predictive modeling, scenario analysis, and weather normalization. We then modernize each step with multiple linear regression (MLR) models and hourly data. The proposed approach has been deployed to many large and medium size utilities including NCEMC. The data required in the NCEMC case study includes hourly system load data at corporate level, which is available through NCEMC's Energy Management System, hourly weather data purchased from WeatherBank and annual economy data purchased from Moody's. Execution of the proposed approach on NCEMC data in automated mode can be finished within a day on a commodity server with an 8-core CPU and 32G RAM. This is well-acceptable for a once-per-year long term load forecasting task. In comparison with Fan's approach, which originated from a field implementation at an ISO, the approach proposed in this paper is more applicable to utilities operating within a regulatory environment, due to its relative simplicity and strong defensibility. The scope of this paper does not include forecasting under renewable penetration and demand response activities.

TABLE I  
MAIN AND CROSS EFFECTS OF THE STARTING MODELS

Model	Main Effects	Cross Effects
$C1$	$GSP, CDD, HDD, Month$	N/A
$C2$	$GSP, T_{max}, T_{max}^2, T_{max}^3, Month$	N/A
$B$	$Trend, T_t, T_t^2, T_t^3, Month, Weekday, Hour$	$T_t * Month, T_t^2 * Month, T_t^3 * Month, T_t * Hour, T_t^2 * Hour, T_t^3 * Hour, Weekday * Hour$
$S_{2010}$	$Trend, T_t, T_t^2, T_t^3, T_{t-1}, T_{t-1}^2, T_{t-1}^3, T_{t-2}, T_{t-2}^2, T_{t-2}^3, T_{t-3}, T_{t-3}^2, T_{t-3}^3, T_a, T_a^2, T_a^3, Month, Day, Hour$	$T_t * Month, T_t^2 * Month, T_t^3 * Month, T_t * Hour, T_t^2 * Hour, T_t^3 * Hour, T_{t-1} * Month, T_{t-1}^2 * Month, T_{t-1}^3 * Month, T_{t-1} * Hour, T_{t-1}^2 * Hour, T_{t-1}^3 * Hour, T_{t-2} * Month, T_{t-2}^2 * Month, T_{t-2}^3 * Month, T_{t-2} * Hour, T_{t-2}^2 * Hour, T_{t-2}^3 * Hour, T_{t-3} * Month, T_{t-3}^2 * Month, T_{t-3}^3 * Month, T_{t-3} * Hour, T_{t-3}^2 * Hour, T_{t-3}^3 * Hour, T_a * Month, T_a^2 * Month, T_a^3 * Month, T_a * Hour, T_a^2 * Hour, T_a^3 * Hour, Day * Hour,$

The rest of the paper is organized as follows: Section III reviews the fundamentals, including the models we start with; Section IV discusses the model selection approach and determines the length of historical data used in long term forecasting; Section V presents the long term probabilistic forecasts with cross weather and economy scenarios; Section VI introduces the methodology for load normalization; the paper is concluded in Section VI with discussions of potential future work.

### III. FUNDAMENTALS

#### A. Multiple Linear Regression

Multiple linear regression analysis has been widely used in the forecasting fields, including load forecasting. Detailed coverage on the theory of regression analysis and linear models is provided in [13]. Implementation of MLR in SAS is presented in [14]. A comprehensive guideline about how to apply MLR models to short term load forecasting is discussed in [1].

In this case study, we start with several MLR models: a classical model for monthly energy forecasting denoted as  $C1$ , a classical model for monthly peak forecasting denoted as  $C2$ , Tao's vanilla benchmark denoted as  $B$ , and a group of customized short term load forecasting models denoted as  $S$ . The models in  $S$  are derived using Hong's methodology documented in [1], where by default, 3 years of data are used for parameter estimation and the year after is used for variable selection. When using year 2010 for variable selection, we denote the resulting variable combination as  $S_{2010}$ . All of these starting models have the dependent variable *Load* and an intercept term. The main effects and cross effects are described in Table I, where each class variables consists of several 0–1 indicator variables.  $S_{2010}$  is used as an example of  $S$  models.

TABLE II  
MODIFICATIONS TO THE DAYS OF A YEAR FOR MODEL  $S_{2010}$

Days of a Year (Original)	Day Code (Modified)
Wednesday (regular)	Tuesday
Day Before New Year's Day	Saturday
New Year's Day	Friday → Saturday; else → Sunday
Memorial Day	Saturday
Day After Memorial Day	Monday
Day Before Independence Day	Friday
Independence Day	Friday → Saturday; else → Sunday
Day Before Labor Day	Saturday
Labor Day	Saturday
Day After Labor Day	Thursday
Day Before Thanksgiving Day	Monday
Thanksgiving Day	Saturday
Day After Thanksgiving Day	Saturday
Day Before Christmas Day	Saturday
Christmas Day	Friday → Saturday; else → Sunday
Day After Christmas Day	Saturday

The *Day* variable is derived from the *Weekday* variable using rules described in Table II. We first group Tuesday and Wednesday together labeled as Tuesday. We then model some holidays and the surrounding days using weekdays and weekends [1]. For example, take New Year's Day: it is a fixed-date holiday. When it falls on a Friday, we modify the value of the *Day* variable to Saturday. Otherwise, we model it as a Sunday. We also model the day before New Year's Day as a Saturday.

#### B. Error Statistics

Despite of many criticisms, mean absolute percentage error (MAPE) is still a widely used error statistic in business forecasting. MAPE (%) can be calculated as follows:

$$MAPE = \frac{100}{N} \sum_{i=1}^N |(A_i - P_i)/A_i|, \quad (1)$$

where  $N$  is the number of observations,  $A_i$  represents the actual load, and  $P_i$  represents the predicted load.

Since the results of our case study are monthly energy forecasts and monthly peak forecasts, we also use the MAPE of monthly energy and MAPE of monthly peak to evaluate the forecasting accuracy. To calculate MAPE of monthly energy (or peak) based on hourly load forecasts, we have to first extract the actual and predicted monthly energy (or peak), and then apply (1) to the resulting series.

To properly answer the second question posted in Section II, the forecasts have to be evaluated based on ex post forecasting accuracy. Take one year ahead forecasting for example as covered in Section V. Assuming we are forecasting the monthly peaks of 2011, if we use the information available through the end of 2010 to forecast 2011, the resulting forecast is ex ante forecast, or "before the event" forecast. If we use the information available through the end of 2011 other than the loads of 2011 to forecast the loads of 2011, the resulting forecast is ex post forecast, or "after the event" forecast. At the beginning of 2012, instead of focusing on ex ante forecast of 2011, we should emphasize the ex post forecasting accuracy of 2011, which tells how the model behaves given the actual temperatures of 2011.

#### IV. PREDICTIVE MODELING

In this section, we first augment the  $S$  models to LTLF models, denoted as  $L$ , using the available data on and prior to 2006. We then determine the appropriate length of history for one year ahead load forecasting. At the end, we compare the ex post forecasting accuracy of the  $C1$ ,  $C2$ ,  $B$ ,  $S$ , and  $L$  models on a rolling basis using 2007 through 2010 [15].

##### A. Model Selection

The general health of the economy is what ultimately drives long term electricity consumption. We would like to extend the model group  $S$  for long term forecasting by adding a macroeconomic indicator,  $GSP$ . The same annual value of GSP is assigned to each hour of a year. We use  $GSP$  in this paper mainly due to two reasons: 1) the territory of NCEMC covers most of North Carolina, which makes  $GSP$  a good driver of the NCEMC's long term load; 2)  $GSP$  is easy to access and understand. If the utility's territory covers one or a few counties or cities,  $GDP$  (gross domestic product) by county or  $GMP$  (gross metropolitan product) can be used as the macroeconomic indicator. In practice, depending upon the drivers of the load, we can also use several other indicators and their combinations, such as housing stock, employment rate, number of jobs, etc. For the utilities, especially retail electricity providers, who provide services in deregulated environment, the total loads are highly impacted by customer churn. In those situations, we can use customer count as the macroeconomic indicator.

The augmentation to a long term forecasting model can be achieved in three ways:

- 1) Replace *Trend* by  $GSP$ . There is an inherent assumption in this approach: the loads sensitive to weather and calendar stay in the same profile over time, while there is part of a base load that growing linearly in proportion to the economic growth. If the forecasting horizon is within a few years, this approach can be a good approximation in practice. As the horizon becomes longer, there can be significantly more customers moving into the territory. Consequently, the weather and calendar sensitive loads should grow as well.
- 2) Divide *Load* by  $GSP$ . The inherent assumption for this approach is that the load is growing at exactly the same rate as the economic growth. In other words, there is no base load that stays constant while the economy is growing. Take a residential community as a counterexample. Before everyone moves in, the feeders, transformers, and street lights are already placed in the community, which lead to a small base load, including no-load loss of transformers, street lighting load, etc. As people are moving in during the next a few years, the total load of this system is growing. However, the small base load stays almost the same since day one. Several ways to extend this approach are to take the natural log or square root of the load or macroeconomic indicator, or both in some combination before performing the division, which allows load to grow faster or slower than the economy.
- 3) Replace *trend* by  $GSP$  and then add interactions between  $GSP$  and the existing main and cross effects. This approach

TABLE III  
COMPARISON AMONG THREE WAYS TO ADD GSP ON MAPE OF HOURLY LOAD

Extension	2002	2003	2004	2005	2006	Average
1	4.3	4.2	4.9	6.3	3.7	<b>4.7</b>
2	3.6	3.6	3.7	7.5	8.3	5.3
3	4.2	4.7	4.9	6.8	4.3	5.0

TABLE IV  
COMPARISON AMONG DIFFERENT LENGTH OF HISTORICAL DATA

Length (yr)	2002	2003	2004	2005	2006	Average
1	4.9	5.5	4.6	4.1	4.7	4.8
2	3.6	4.2	4.8	5.1	3.5	<b>4.2</b>
3	4.3	4.2	4.9	6.3	3.7	4.7
4	4.2	4.8	4.9	7.4	4.6	5.2

TABLE V  
COMPARISON AMONG MODELS  $C1$ ,  $C2$ ,  $B$ ,  $S$  AND  $L$

MAPE	Model	2007	2008	2009	2010	Average
Annual Energy	$C1$	2.5	0.9	1.0	0.8	1.3
	$L$	0.2	0.1	0.1	0.1	0.1
Annual Peak	$C2$	1.1	7.5	13.8	8.6	7.8
	$L$	0.2	2.8	2.4	7.3	3.2
Monthly Energy	$C1$	3.4	2.0	2.7	3.4	2.9
	$L$	1.5	1.3	1.2	2.1	1.5
Monthly Peak	$C2$	3.8	4.2	9.3	6.9	6.1
	$L$	2.1	3.7	3.0	4.3	3.3
Hourly Load	$B$	4.8	5.1	5.0	4.9	5.0
	$S$	3.5	3.9	4.0	3.6	3.8
	$L$	3.5	3.3	3.4	3.7	3.5

assumes end-users' behavior changes as the economic environment changes. Since a significant amount of variables are being added through the additional interaction effects, the resulting model may be over-parameterized. Depending upon the forecasting horizon and the electricity usage pattern, this approach may not provide forecast results that are as accurate as the first two options.

Table III compares the MAPE of hourly loads of the three approaches discussed above for one year ahead forecasting. The MAPE values are generated on rolling basis with a history window fixed at 3 years. Take the 3.7% under 2004 for example. We used the second approach ("divide *Load* by *GSP*") mentioned above to augment the model  $S_{2004}$  to get the model for long term forecasting, denoted as  $L_{2004}$ . The parameters are then estimated using the load, temperature, and economy data from 2001 to 2003. Based on five years of validation results, we conclude that the "replace *Trend* by *GSP*" approach on average offers the lowest MAPE (4.7%) in this case study.

### B. Length of Training Data

The length of historical data for parameter estimation is another factor that impacts forecasting accuracy. Table IV lists the MAPE values generated on rolling basis with different length of history window. For example: observe 2005, with a MAPE of 7.4% in the last row. We use 4 years of history from 2001 to 2004 to estimate the parameters of the model  $L_{2005}$  ( $S_{2005}$  augmented by replacing *Trend* by *GSP*). Based on five years of validation results, we conclude that in this case study, using 2 years of historical data offers the lowest average MAPE (4.2%)

for forecasting one year ahead. While this rolling simulation approach can be used for determining multiple years ahead forecasting, we may not reach the same conclusion that 2 years of historical data is optimal for 5 years ahead forecasting.

### C. Comparison

We would like to compare the ex post forecasting accuracy of models  $C1$ ,  $C2$ ,  $B$ ,  $S$ , and  $L$ . Some of these models ( $C1$ ,  $C2$ , and  $B$ ) already have a pre-designated variable combination, while some ( $S$  and  $L$ ) require model identification. Some ( $C1$  and  $C2$ ) are based on monthly data, while some ( $B$ ,  $S$ , and  $L$ ) are based on hourly interval data. Due to the above characteristics, we have to apply different treatments to the models to calculate the MAPE values of ex post forecasts:

- 1) Classical models  $C1$  and  $C2$ : the variables are specified in Table I, while the parameters are estimated using the eight years of historical data prior to the year to be forecasted.
- 2) Tao's vanilla model  $B$ : the variables are specified in Table I, while the parameters are estimated using the three years of historical data prior to the year to be forecasted.
- 3) Customized short term forecasting model group  $S$ : to perform ex post forecasting for the loads of year  $y$ , we cannot use the loads of year  $y$  for model building, including the tasks of parameter estimation and variable selection. To avoid using the loads of year  $y$ , we first identify the model  $S_{y-1}$ , which is selected using the year  $y-1$  as the validation data and the three years  $y-4$  to  $y-2$  as the training data. Parameter estimation of  $S_{y-1}$  is based on the 3 years prior to year  $y$ , namely from year  $y-3$  to  $y-1$ .
- 4) Customized long term forecasting model group  $L$ : similar to the analogy above, we cannot use the loads of year  $y$  to build the model when ex post forecasting the same year. Therefore, we first identify model  $L_{y-1}$ , and then estimate the parameters based on two years of historical data,  $y-1$  and  $y-2$ .

In Table V, we list the MAPE (and absolute percentage error for annual interval summary) values of annual energy, annual peak, monthly energy, monthly peak, and hourly load from the five model groups. Table V first shows that the LTLF models ( $L$ ) derived based on the proposed approach have much lower MAPE values than the classical models  $C1$  and  $C2$  on one year ahead ex post forecasting. On monthly energy and peak forecasting, the proposed approach reduces the MAPE by over 45%. Table V also shows that the performance of  $L$  improves on both model  $B$  and model group  $S$ .

Figs. 1 and 2 show the line plots of monthly energy and monthly peak profiles from 2007 to 2010, which confirms that the proposed approach leads to more accurate forecasts than does the counterpart.

The classical approach based on monthly data leads to significantly higher error than the proposed approach. This is because the monthly data (peak temperature, HDD, and CDD) cannot tell: 1) which hour of the day and which day of the week the high/low temperatures fall into; 2) the variation of the temperatures throughout a day; 3) the temperature profiles for modeling recency effect [1]. In addition, the HDD and CDD require the forecasters to specify the threshold or comfortable zone, which may not be very defensible.

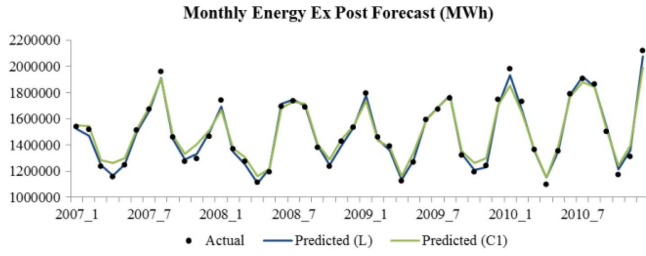


Fig. 1. Comparison on ex post forecasts of monthly energy (2007–2009).

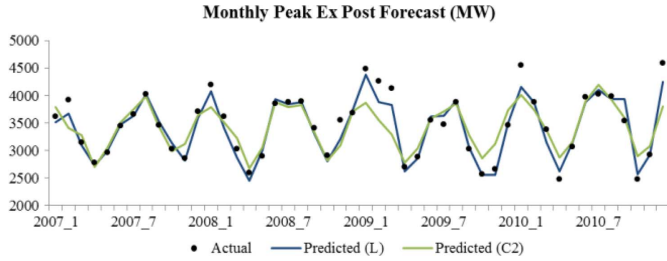


Fig. 2. Comparison on ex post forecasts of monthly peak (2007–2009).

## V. SCENARIO ANALYSIS

Forecasting is, by nature, a stochastic problem. Due to the uncertainty in climate and economic forecasts, long term load forecasters are encouraged to provide multiple forecasts based on different scenarios. This section discusses how to create weather and economic scenarios. Since 2011 is a year that many U.S. utilities had trouble forecasting, we use 2011 as an example to illustrate the proposed methodology.

### A. Weather and Economic Scenarios

The pros and cons of several existing means to create weather scenarios have been discussed in Section II. In this paper, we use actual temperature profiles from the history to create weather scenarios. There are three components that should be clearly specified in the one year ahead load forecasting process for a given year  $y$ ; or multiple years ahead load forecasting process for a given horizon starting from year  $y$ :

- 1) How to model the system, such as combination of weather and calendar variables, incorporation of macroeconomic indicator(s), and length of load, weather, and economy history for parameter estimation. Since the load of year  $y$  should be excluded from model building, we can use model  $L_{y-1}$ , which is identified using the most recent years of information.
- 2) How many years of temperature history to use. Different organizations may adopt different practices when selecting the length of temperature history, which ranges from 20 years to 50 years. NOAA, for instance, uses 30 years of history to create and update the typical meteorological year (TMY). In this paper, we also use 30 years of temperature history, from  $y-30$  to  $y-1$ , to create 30 weather scenarios for year  $y$ . If the year  $y$  is a leap year, i.e., 2008, and the year of weather scenario is based on a non-leap year, i.e., 1991, we fill in 02/29/2008 with 02/28/1991's temperatures. If the year  $y$  is a non-leap year, i.e., 2011, and the year of weather scenario is based on a leap year, i.e., 2000, we can remove the temperatures of 02/29/2000. Based on each

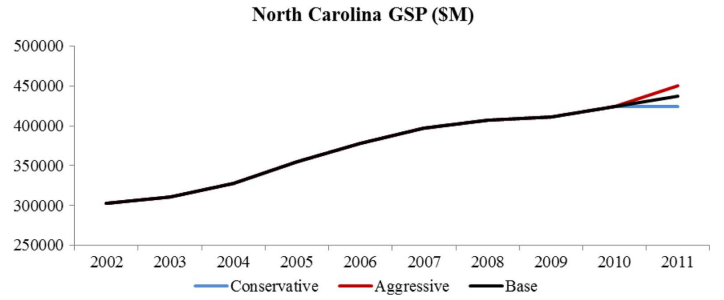


Fig. 3. History (2002–2010) and forecast (2011, 3 scenarios) of GSP.

weather scenario, we can generate an hourly load profile for the year  $y$  using the model  $L_{y-1}$ .

- 3) How to extract normalized peak and energy. From each hourly load profile, we first derive monthly peak (or energy) profiles. We then find the median of the monthly peaks (or energy) for each month. The results are the normalized monthly peak (or energy) forecast. Many organizations also require the forecasts at the 10th and 90th percentiles to support the decision making processes.

Most utilities purchase economic forecasts from third parties for LTLF. The economic forecasts usually come with multiple scenarios. In this paper, we use three macroeconomic scenarios: base, aggressive, and conservative scenarios for the year of 2011 as shown in Fig. 3. For each macroeconomic scenario, we can have the same 30 weather scenarios as mentioned above. In total, we can create 90 cross scenarios.

### B. Probabilistic Forecasts

Figs. 4 and 5 show one year ahead forecasting of 2011's monthly peak and energy respectively. There are 30 dashed lines representing the forecasts obtained using the 30 weather scenarios combined with the base economic scenario. In addition, we plot the 5 scenarios extracted from the 90 cross scenarios, including 10th (gray) 50th (black) and 90th (green) percentiles of the load with base economic scenario, and median load with conservative (blue) and aggressive (red) economic scenarios. The actual monthly peaks and energy of 2011 are labeled as black dots.

In practice, the 90th percentile is often used to represent a severe scenario that may happen one out of ten times. It does not mean that the load will never exceed this bound. Among the 12 monthly peaks shown in Fig. 5, the actual peak of May 2011 does exceed the 90th percentile line, which is reasonable considering the definition of the 90th percentile.

Sometimes the extreme estimates are unrealistic, because the given temperature scenario can be out of range of the training data. For instance, in Fig. 5, the extreme scenario of Jan 2011 exceeds 6000 MW, which is driven by an extremely cold year in the 1980s. Since the 90th percentile derived from the 30 scenarios is not sensitive to the extreme value, it is still reliable and practical to use such a 90th percentile curve for planning purposes.

## VI. LOAD NORMALIZATION

Due to the variation in climate from year to year, most utilities conduct some form of weather normalization processes to esti-



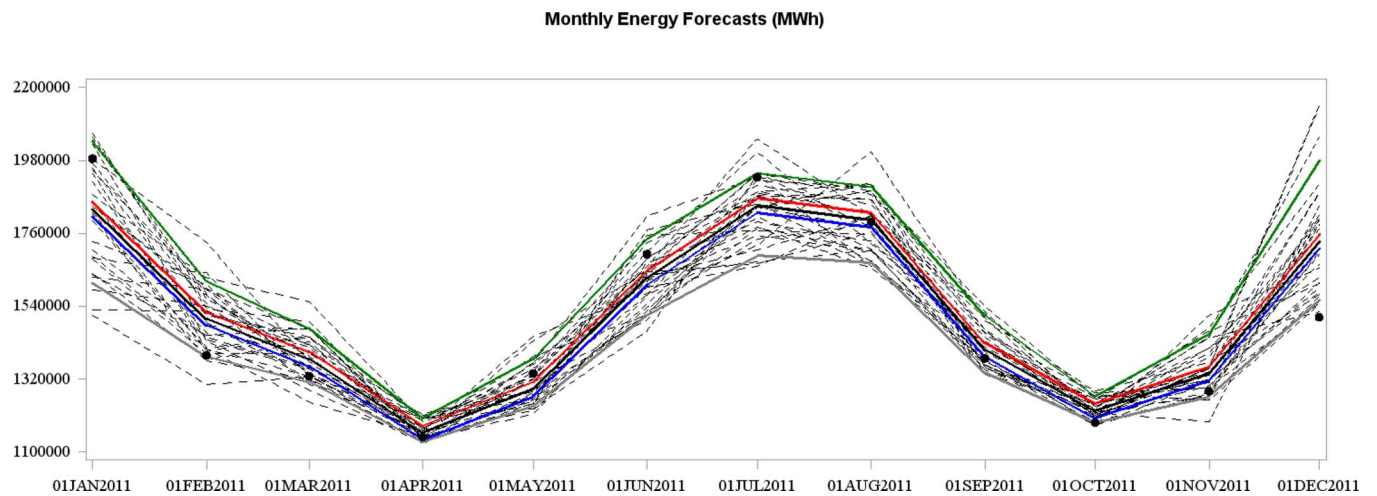


Fig. 4. Ex ante forecasts of 2011 monthly energy.

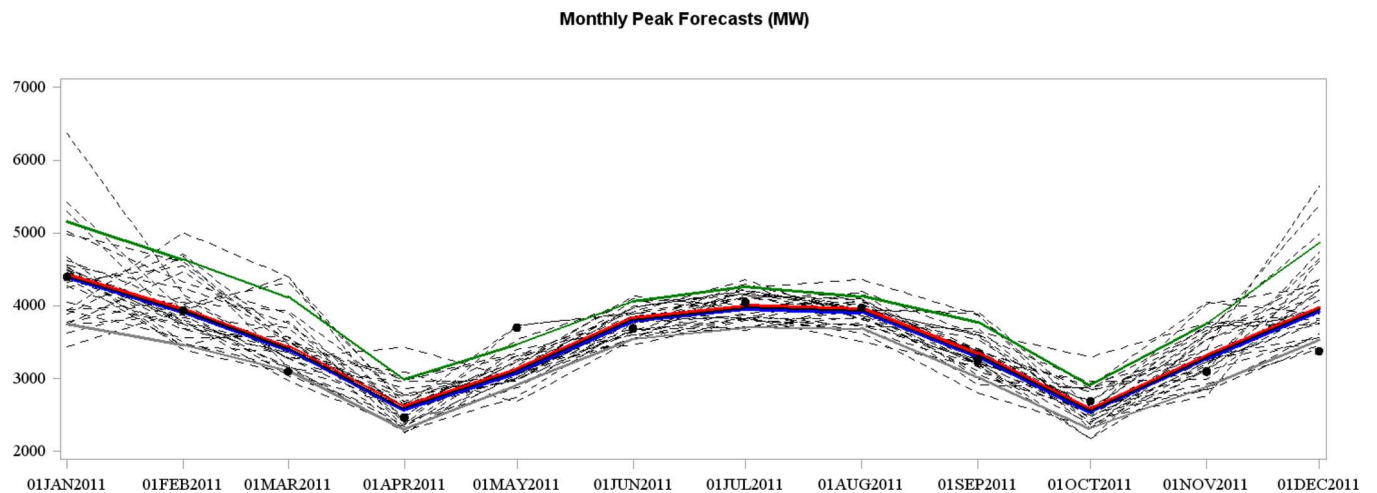


Fig. 5. Ex ante forecasts of 2011 monthly peak (30 + 5 scenarios).

mate the normalized load profile. There are two business needs for such processes: 1) understanding the load growth *without* the impact of climate change; 2) understanding the variation of the load *with* the impact of climate change. Due to the nonlinear relationship between load and weather [1], a normal weather profile usually does not lead to a normal load profile. Comparing with the conventional term *weather normalization*, a more accurate description to the process of estimating the load profile without impact of climate change should have been *load normalization against weather*.

Similar to creating weather scenarios for LTLF as discussed in Section V-A, there are three components that should be clearly specified in the load normalization process for a given year  $y$ : 1) how to model the system; 2) how many years of temperature history to use; 3) how to extract normalized peak and energy. The second and third components can be treated the same way as discussed in Section V-A, while the first one is slightly different.

When normalizing the historical load of a given year  $y$ , we should identify a model that concurrently best represents the system status in the year  $y$ , and has strong predictive power to answer the “what-if” questions. Since all the information in-

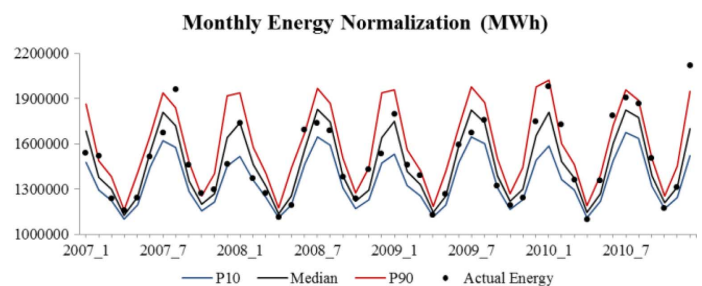


Fig. 6. Monthly energy normalization (2007–2010).

cluding load, temperature, and economy of the year  $y$  is available for load normalization, we can use the model  $L_y$ , which is identified using the data through the end of year  $y$ .

Figs. 6 and 7 present the load normalization results for monthly energy and peaks from 2007 to 2010, where the 10th percentile, median, and 90th percentile load profiles are colored in blue, black, and red respectively. The actual peaks are labeled as black dots. As shown in Fig. 6, the actual monthly energy of December 2010 is above the 90th percentile line. This is due to 3 consecutive very cold weeks, which rarely happened in the past several decades.

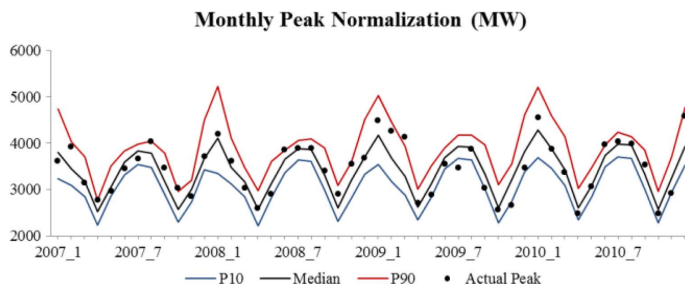


Fig. 7. Monthly peak normalization (2007–2010).

## VII. CONCLUSION

In this paper, we presented a practical approach to LTLF. We modernized predictive modeling, weather normalization, and probabilistic forecasting with MLR models and hourly information. Through a case study at NCEMC, we showed how this method can create superior accuracy and defensibility of the forecast results over the classical approach based on monthly data. In particular, we proposed the concept of load normalization, and demonstrated a simulation approach to normalizing the load against weather.

In future work, as an expansion of the proposed methodology, we would like to further explore the following directions: 1) incorporation of high resolution spatial information; 2) how data cleansing could help improve long term load forecasts; 3) understanding how the forecast errors of explanatory variables contribute to the error of ex ante forecasts.

## ACKNOWLEDGMENT

The authors gratefully acknowledge the insights, support, and encouragement from Tom Laing, Director of Market Research at NCEMC, during the project implementation and paper preparation.

## REFERENCES

- [1] T. Hong, "Short term electric load forecasting," Ph.D. dissertation, Graduate Program of Operation Research and Dept. Electrical and Computer Engineering, North Carolina State Univ., Raleigh, NC, USA, 2010.
- [2] T. Hong, P. Wang, and H. L. Willis, "A naive multiple linear regression benchmark for short term load forecasting," in *Proc. 2011 Power Energy Soc. Gen. Meet.*, Detroit, MI, USA, Jul. 24–29, 2011.
- [3] S. Fan and R. J. Hyndman, "Short-term load forecasting based on a semi-parametric additive model," *IEEE Trans. Power Syst.*, vol. 27, no. 1, pp. 134–141, Feb. 2012.
- [4] S. Fan, K. Methaprayoon, and W. J. Lee, "Multiregion load forecasting for system with large geographical area," *IEEE Trans. Ind. Appl.*, vol. 45, no. 4, pp. 1452–1459, Jul.–Aug. 2009.
- [5] S. Fan, L. Chen, and W. J. Lee, "Short-term load forecasting using comprehensive combination based on multimeteorological information," *IEEE Trans. Ind. Appl.*, vol. 45, no. 4, pp. 1460–1466, Jul.–Aug. 2009.
- [6] A. Motamedi, H. Zareipour, and W. D. Rosehart, "Electricity price and demand forecasting in smart grids," *IEEE Trans. Smart Grid*, vol. 3, no. 2, pp. 664–674, Jun. 2012.
- [7] N. Amjadi, F. Keynia, and H. Zareipour, "Short-term load forecast of microgrids by a new bilevel prediction strategy," *IEEE Trans. Smart Grid*, vol. 1, no. 3, pp. 286–294, Dec. 2010.
- [8] T. Hong, "Spatial load forecasting using human machine co-construct intelligence framework," Master's thesis, Graduate Program of Operation Research and Dept. Industrial and Systems Engineering, North Carolina State Univ., Raleigh, NC, USA, 2008.
- [9] T. Hong, S. M. Hsiang, and L. Xu, "Human-machine co-construct intelligence on horizon year load in long term spatial load forecasting," in *Proc. 2009 Power Energy Soc. Gen. Meet.*, Calgary, AB, Canada, Jul. 26–30, 2009.
- [10] D. J. Barger, "MGE experience with INSITE spatial load forecasting," in *Proc. 2011 Power Energy Soc. Gen. Meet.*, Detroit, MI, USA, Jul. 24–29, 2011.
- [11] R. J. Hyndman and S. Fan, "Density forecasting for long-term peak electricity demand," *IEEE Trans. Power Syst.*, vol. 25, no. 2, pp. 1142–1153, May 2010.
- [12] T. Hong, *Electric Load Forecasting: Fundamentals and Best Practices*. Cary, NC, USA: SAS, May 2012 [Online]. Available: <http://courses.drhongtao.com/sasbelf>
- [13] M. Kutner, C. Nachtersheim, J. Neter, and W. Li, *Applied Linear Statistical Models*. New York: McGraw-Hill/Irwin, 2004.
- [14] R. Littell, W. Stroup, and R. Freund, *SAS for Linear Models*. Cary, NC, USA: SAS, 2002.
- [15] L. J. Tashman, "Out-of-sample tests of forecasting accuracy: An analysis and review," *Int. J. Forecasting*, vol. 16, no. 4, pp. 437–450, 2010.

**Tao Hong** received his B.Eng. in automation from Tsinghua University, Beijing, China, a M.S. in EE, a M.S. with co-majors in OR and IE, and a Ph.D. with co-majors in OR and EE from North Carolina State University, Raleigh, NC, USA.

He is Senior Industry Consultant at SAS Institute, Cary, NC, USA, where he leads the forecasting vertical of the Energy Business Unit. The long term spatial load forecasting methodology implemented in his M.S. thesis and the short term forecasting methodology proposed in his Ph.D. dissertation have been commercialized and deployed to many utilities worldwide.

Dr. Hong serves as the Chair of the IEEE Working Group on Energy Forecasting, Guest Editor in Chief of IEEE Transactions on Smart Grid Special Issue on Analytics for Energy Forecasting with Applications to Smart Grid, and General Chair of Global Energy Forecasting Competition.

**Jason Wilson** (M'11) received the B.A. degree in economics from the University of Alaska, Anchorage, AK, USA. He is currently completing his M.S. thesis in resource and applied economics from the University of Alaska Fairbanks, AK, USA.

He is Load Research Analyst at North Carolina Electric Membership Corporation (NCEMC), Raleigh, NC, USA. He has been working in vertically integrated electric cooperatives, and G&T cooperatives since 2008. His research interests include power and energy, and the societal implications of technology.

**Jingrui Xie** (A'12) graduated as the 1st ranked student with her B.S. degree in finance from Sun-Yat Sen University, Guangzhou, China, and received her M.S. degree in economics from Duke University, Durham, NC, USA.

She is Senior Associate Analytical Consultant at SAS Institute, Cary, NC, USA, with expertise in statistical analysis and forecasting. She is the primary statistician developer for SAS Energy Forecasting solution and has been working on energy forecasting projects with several U.S. utilities since she joined SAS. Her research interests are in the field of power and energy.