

Temporal Classification Method for Forecasting Power Load Patterns From AMR Data

Heon Gyu Lee*, Jin-Ho Shin**, Hong Kyu Park*, Young-il Kim**, Bong-Jae Lee**, and Keun Ho Ryu*[†]

*Database/Bioinformatics Lab., School of Electrical & Computer Engineering, Chungbuk National University, Korea

**Power Information Technology Group, Korea Electric Power Research Institute, Korea

Abstract : We present in this paper a novel power load prediction method using temporal pattern mining from AMR (Automatic Meter Reading) data. Since the power load patterns have time-varying characteristic and very different patterns according to the hour, time, day and week and so on, it gives rise to the uninformative results if only traditional data mining is used. Also, research on data mining for analyzing electric load patterns focused on cluster analysis and classification methods. However despite the usefulness of rules that include temporal dimension and the fact that the AMR data has temporal attribute, the above methods were limited in static pattern extraction and did not consider temporal attributes. Therefore, we propose a new classification method for predicting power load patterns. The main tasks include clustering method and temporal classification method. Cluster analysis is used to create load pattern classes and the representative load profiles for each class. Next, the classification method uses representative load profiles to build a classifier able to assign different load patterns to the existing classes. The proposed classification method is the Calendar-based temporal mining and it discovers electric load patterns in multiple time granularities. Lastly, we show that the proposed method used AMR data and discovered more interest patterns.

Key Words : Load Patterns, Temporal Pattern Mining, Load forecasting, Calendar-based temporal mining.

1. Introduction

Electricity load patterns prediction (or forecasting) has been an important issue in the power industry. Load patterns prediction deals with the discovery of power load patterns from load demand data. It attempts to identify existing load patterns and recognize new load forecasting methods, employing methods from sciences such as statistics (Huang, and

Shih, 2003) and data mining (Pitt, and Kirchen, 1999). In power system, data mining is the most commonly used methods to recognize and extract regularities in load data and thus has been the target of some investigations for its used in load pattern forecasting. In particular, it promises to help in the detection of previously unseen load patterns by establishing sets of observed regularities in load demand data. These sets can be compared to current

Received 2 October 2007; Accepted 13 October 2007.

[†] Corresponding Author: Keun Ho Ryu (khryu@dbl-lab.chungbuk.ac.kr)

load pattern for deviation analysis. Load patterns prediction using data mining is usually made by building models on relative information, weather, temperature and previous load demand data. Such prediction is aimed at short term prediction (Amjady, 2001), since mid and long term prediction may not be reliable because the results of prediction contain high forecasting errors. However, mid and long term (load patterns for longer period) forecasting on load demand is very useful and interest. Also, load demand data is temporal data which has timestamp. Previous researches such as clustering, classification and regression usually did not consider such time factor in temporal data or applied as static factor. Since the power load patterns have time-varying characteristic and very different patterns according to the hour, time, day and week and so on, it gives rise to the uninformative results if only traditional data mining is used. Therefore, if we consider time intervals under multiple time granularities in order to forecast load patterns to load demand data analysis with temporal dimension, we can discover useful load patterns during the given time interval. The purpose of this paper is to investigate the effectiveness and accuracy of classification method within load forecasting. To achieve this purpose, we attempt to apply clustering method and calendar-based temporal classification method and their use in load pattern forecasting. For the forecasting load patterns, the main tasks are the following: and a framework of our approaches is showed in Fig. 1.

1. Cluster analysis is performed to detect load pattern classes and the representative load profiles for each class.
2. Temporal associative classification method uses representative load profiles to build a classifier able to assign different load patterns to the existing classes.
 - ① Calendar pattern proposed in (Li, and Ning,

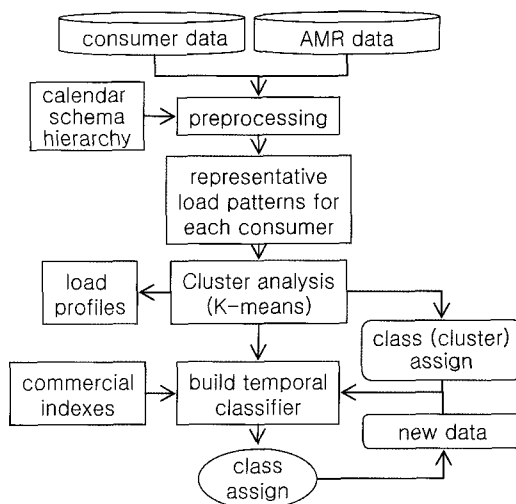


Fig. 1. A framework of temporal classification method.

2001) is applied to AMR (Automatic Meter Reading) load data for the time expression of class association rules. This calendar pattern is based on user-specified calendar pattern and represents cyclic pattern.

② CARs (Class Association Rules) are discovered in given time. CARs are special subset of association rules with a consequent limited to class values only.

3. The generated temporal CARs are applied to build classifier for predicting load patterns in AMR load data.

2. Cluster Analysis

We describe clustering algorithms for generating the load profiles and class label which will be used classification module. The load pattern associated with any customer contains the information of commercial indexes such as contract assortment, industrial code and electricity use which recoded every 15 minutes. In order to perform clustering, we represent the load pattern for each consumer. The

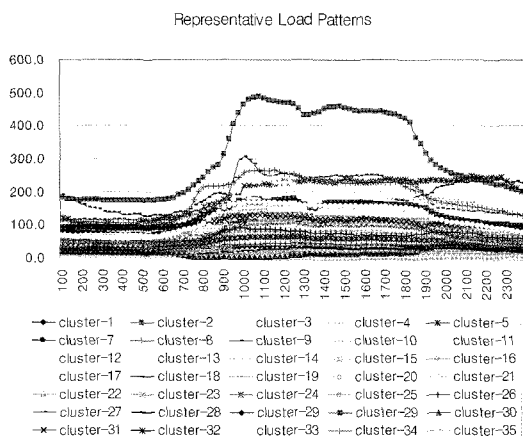


Fig. 2. 35 representative load patterns.

representative daily load pattern of the m th consumer is following:

$$V^{(m)} = \{V_{100}^{(m)}, \dots, V_h^{(m)}, \dots, V_H^{(m)}\} \quad (1)$$

where $h=100, \dots, H$ with $H=2345$, representing the 15 min. interval between the collected measurements.

In cluster analysis, K-means (Jain, and Dubes, 1998) is used to group the load patterns and the optimal clusters are obtained. The use of clustering in this step detects the number of classes as an input of the classification model. To define the number of classes, we performed the evaluation of the clusters compactness using the measure sum of the squared error (SSE). The obtained results are showed in Fig. 2. It is possible to see that 35 clusters would be good choice, considering the SSN (Park, et al, 2007).

3. Temporal Class Association Rules

In this section, we describe classification method and algorithms of Temporal CARs (Class Association Rules).

1) Calendar-based pattern

Let $I=\{a_1, a_2, \dots, a_m\}$ be a set of items, and a

transaction database $DB=\langle T_1, T_2, \dots, T_i \rangle$, where $T_i (i \in \{1 \dots i\})$ is a transaction which contains a set of items in I . The support (or occurrence frequency) of a pattern A , which is a set of items, is the number of transactions containing A in DB . A is a frequent pattern if A 's support is no less than a predefined minimum support threshold, ε .

We present a class of calendar related to temporal patterns called calendar patterns. Calendar pattern represents the sets of time intervals in terms of calendar schema specified by user. Calendar schema is a relational schema $R=(f_n:D_n, f_{n-1}:D_{n-1}, \dots, f_1:D_1)$, where each attribute f_i is a time granularity name such as year, month, day, and so on. Each D_i is a domain value corresponding to f_i . There is given a calendar schema (year: {1,2}, month: {1,2}, week: {1,2,3}, day: {1,2,3,4}). And a calendar pattern on the calendar schema R is a tuple on R of the form $\langle d_n, d_{n-1}, \dots, d_1 \rangle$. Each d_i is a positive integer in D_i or the symbol “*”. “*” denotes all the values corresponding to domain and means “every.” Exactly, it represents periodic cycles on the calendar pattern such as every week, every month, and so on.

For presentation of calendar pattern, we call a calendar pattern with k symbols a k -star calendar pattern (denoted e_k) and a calendar pattern with at least one symbol a star calendar pattern. In addition, we call a calendar pattern with no symbol a basic time interval (denoted e_0) According to the above example 1, the calendar pattern $\langle 1, 1, *, 3 \rangle$ represents time intervals which means the third day of every week of January in the first year.

2) Class association rules

CARs (Class Association Rules) are a combination of association rules mining and classification. CARs are a special subset of association rules whose antecedent is an itemsets and consequent are

restricted to the classification class label. Let $I=\{a_1, a_2, \dots, a_m\}$ be a set of all items in DB and Y to be a set of all class labels. CARs r is an implication of the form:

$$X \Rightarrow C \quad (2)$$

where $X \subset I, C \subset Y$. Antecedent of a CARs is also called itemset and a rule itself is called *rule_item*. *Rule_item* is large if the corresponding rule is frequent and accurate if the rule is confident. To find all strong association rules, CBA (Liu, 1998) uses an Apriori-like algorithm to generate large *rule_items*. And *k-rule_item* denotes a rule whose itemset has k items. In each pass, all large $(k-1)$ -*rule_items* are found. These $(k-1)$ -*rule_items* are used to generate candidate *k-rule_item* and selected candidates satisfy minimum support threshold, ε . For each large *rule_item*, the confidence of the corresponding rule is calculated and the rule is added to the set of all rules if the confidence satisfies minimum support threshold, δ .

Support and confidence of *rule_item* in a transaction DB is the following:

$$Support = \frac{ruleCount}{DB}, Confidence = \frac{ruleCount}{itemsetCount} \quad (3)$$

where *ruleCount* is the number of items in DB that contain the itemset and are labeled with class label and *itemsetCount* is the number of items in DB that contain the itemset.

3) Temporal CARs algorithm

Given a basic time interval $t(e_0)$ under a given calendar schema, we denote the set of transactions whose timestamps are covered by $t(e_0)$ as $DB[t]$. Temporal CARs over a calendar schema R is a pair (r, e) . Thus, temporal CARs (r, e) hold in DB if and only if the r satisfies ε and δ in $D[t]$ for each basic time interval e_0 covered by e . And CARs are an implication of the form:

$$\langle X \Rightarrow C, e \rangle \quad (4)$$

X : itemset, C : class label, e : star calendar pattern

We extend apriori (Agrawal, and Srikant, 1994) to discover large *rule_items*. On the data mining tasks, our algorithm produces the rules $TCAR_k(e)$ that satisfies ε and δ for all possible star calendar pattern on R . The *Temporal CARs* algorithm is given in Fig. 3. The algorithm generates all the large *rule_items* by making multiple passes over the data.

Let *k-rule_item* denote a *rule_item* whose *itemset* has k items. Let $L_k(e)$ denote the set of large *k-rule_item* for calendar pattern e . Each element of this set is of the following form:

$$\langle (itemset, itemsetCount), (class_label, ruleCount) \rangle$$

Let C_k be the set of candidate *k-rule_item*. The algorithm generates all the large *rule_items* by making multiple passes over the data. In each passes, the basic time intervals in the calendar schema are processed one by one. During the processing of basic time interval e_0 in pass k , the set of large *k-rule_items*,

```

Input transaction (Training data set) D
Output <TCARk(e), e> for all star calendar pattern e

forall basic time intervals e0 do
    L1(e0) = {large 1-rule_items in D[e0]};
    TCAR1(e0) = genRules(L1(e0))
end
forall (k=2; ∃ calendar pattern e such that Lk-1(e) ≠ ∅;
k++) do
    forall basic time intervals e0 do
        Ck(e0) = candidateGen(Lk-1(e0), min_sup);
        for each data case d ∈ D(e0) do
            Cd = subset(Ck(e0), d);
            for each candidate c ∈ Cd do
                c.itemsetCount ++;
                if d.class=c.class then c.ruleCount++;
            end
        end
        Lk(e0) = {c ∈ Ck | c.ruleCount ≥ min_sup};
        TCARk(e0) = genRules(Lk(e0), min_conf);
        forall calendar pattern e that cover e0 do
            update TCARk(e) using TCARk(e0);
        end
    end
end
    
```

Fig. 3. Temporal CARs algorithm.

$L_k(e_0)$ is first computed, and then $TCAR_k(e_0)$ is used to update the $TCAR_k(e)$ for all the calendar patterns that cover e_0 . Assume that we have calendar schema (year, month, day). In the algorithm, we may need to consider, e.g., calendar pattern $\langle 2000, *, 1 \rangle$ as well as $\langle *, 1, * \rangle$. These two patterns have an overlapping basic time interval, $\langle 2000, 1, 1 \rangle$. In our algorithm, we use the large *rule_items* for $\langle 2000, 1, 1 \rangle$ to derive the large *rule_items* for $\langle 2000, *, 1 \rangle$ and $\langle *, 1, * \rangle$ to avoid duplicate tasks.

In the first pass, we compute the large *1-rule_items* for each basic time interval by counting the supports. These *1-rule_items* are used to generate candidate *2-rule_items* by *candidateGen*($L_{k-1}(e_0)$). The next scan of the data is performed to count which of the candidate *2-rule_items* satisfy ε for each basic time interval. And then, the algorithm produces the rules $TCAR_k(e_0)$ using *genRules* function. The algorithm iterates in this fashion, starting each subsequent pass with the seed set of rules found to be large in the previous pass. Finally, large *k-rule_items* in $TCAR_k(e_0)$ is used to update the $TCAR_k(e)$ for each star calendar pattern that covers the basic time interval. After the basic time interval e_0 is processed in pass k , the large $TCAR_k$ for all the calendar patterns e that cover e_0 are updated as follows. We associate a counter count with each candidate $TCAR_k$ for each star calendar pattern. The counters are initially set to 1.

When $TCAR_k(e_0)$ is used to update $TCAR_k(e_0)$, the counters of the *k-rule_items* in $TCAR_k(e_0)$ that are also in $TCAR_k(e_0)$ are incremented by 1, and the *k-rule_items* that are in $TCAR_k(e_0)$ but not in $TCAR_k(e_0)$ are added to $TCAR_k(e)$ with the counter set to 1.

4. Constructing Classifier

In this section, we describe to build the efficient

classifier using Temporal CARs. Generated classifier for each calendar pattern is following format:

$$(r_1, r_2, \dots, r_n, \text{default_class}, e) \quad (5)$$

where, r_i is the generated Temporal CARs and e is the calendar pattern. However, the number of rules generated by algorithm can be huge. To make the efficient classifier, we need to prune rules generated. To reduce the number of rules generated, we perform two types of rule prune. First rule prune is the pruning from calendar patterns of Temporal CARs. After we discover all large *rule_items*, we remove all the (r, e) if we have other (r, e') and e is covered by e' . Second, we use general and high-confidence rule to prune more specific and lower confidence ones. Before the pruning, all rules are ranked according to the following criteria.

Given two rules r_i and r_j , $r_i > r_j$ (or r_i is ranked higher than r_j) if

- 1) $\text{conf}(r_i) > \text{conf}(r_j)$ or
- 2) $\text{conf}(r_i) = \text{conf}(r_j)$, $\text{sup}(r_i) > \text{sup}(r_j)$ or
- 3) $\text{conf}(r_i) = \text{conf}(r_j)$ and $\text{sup}(r_i) = \text{sup}(r_j)$, but r_i is generated before r_j .

A rule $r_1: X \Rightarrow C$ is said a general rule w.r.t. rule $r_2: X' \Rightarrow C'$, if only if A is a subset of A' . First, we need to sort the set of generated rules for each calendar pattern. This sorting guarantees that only the highest rank rules will be selected into the classifier. And then, given two rules r_1 and r_2 , where r_1 is a general rule w.r.t. r_2 . We prune r_2 if r_1 also has higher rank than r_2 .

After a set of rules is selected for classifier, we classify new data. First, we select the rules whose calendar patterns can cover basic time intervals of cases. Then, we discover the rules matching case in the selected one and classify class labels of the discovered rules. If all the rules matching the new data have same class label, the new case is assigned to that

label. Otherwise, we classify the new case as class label of the rule with higher confidence and support.

5. Experimental Result

In this section, a case study concerning a database with load patterns from 231 consumers is considered and this information has been collected by KEPRI (Korea Electric Power Research Institute). The collected load patterns were made during a period of three month (January, February and March) in 2007. The instant power consumption for each consumer was collected with a cadence of 15 min. The commercial indexes related with contract assortment power, industrial classification code, and electricity use code are also applied.

1) Data preprocessing

To preprocessing the 3 month's training data, we define calendar schema and domain by a hierarchy of calendar concepts.

$$R = (\text{Month: } \{1, \dots, 3\} \text{ week: } \{1, \dots, 4\}, \text{ working day: } \{1, \dots, 7\}, \text{ hours: } \{1, \dots, 5\})$$

where domain values of hours represent early morning(0-7hour), morning(7-9hour), daytime(9-17hour), evening(17-20hour), and late night(20-24hour).

This calendar schema considers mid and long periods to predict mid and long term load patterns. The class labels are also assigned representative load patterns of cluster in training data. To compare the load patterns, we use feature of load shape (Ernould, 1982), able to capture relevant information about the consumption behavior, must be create the classifier. These features must contain information about the daily load curve shape of each consumer and presented in Table 1.

Table 1. Load curve shape features.

Feature	Definition	Period
Load Factor	$s_1 = \frac{\text{Pattern}_{Avg, for \text{ day}}}{\text{Pattern}_{Max, for \text{ day}}}$	1 day
Night Impact	$s_2 = \frac{1}{3} \frac{\text{Pattern}_{Avg, for \text{ day}}}{\text{Pattern}_{Max, for \text{ day}}}$	8 hours (23 ~ 07)
Lunch Impact	$s_3 = \frac{1}{8} \frac{\text{Pattern}_{Avg, for \text{ day}}}{\text{Pattern}_{Max, for \text{ day}}}$	3 hours (12 ~ 15)

Since the extracted features contain continuous variables, those variables also must be made discrete. Therefore, entropy-based discretization (Fayyad, and Irani, 1993) has been used because the intervals are selected according to the information they contribute target variable.

Fig. 4 shows the preprocessing results from load pattern data of AMR.

Feature	Type	Description
Date	datetime	YYYYMMDD
Contract power	nominal	Different 29 value
Industrial code	nominal	Different 158 value
Electricity use code	nominal	Different 21 value
AMR Capacity (15min. Interval)	100	continuous Min.: 0.3 ~ Max.: 490
	115	continuous Min.: 0.3 ~ Max.: 490
	...	continuous Min.: 0.3 ~ Max.: 490
	2345	continuous Min.: 0.3 ~ Max.: 490
class	cluster	nominal {cluster1, ... cluster 35}



Data preprocessing

Feature	Type	Description
Calendar pattern	nominal	Calendar expression
Contract power	nominal	Different 29 value
Industrial code	nominal	Different 158 value
Electricity use code	nominal	Different 21 value
AMR Capacity (15min. Interval)	S1	nominal. Discrete value
	S2	nominal. Discrete value
	S3	nominal. Discrete value
class	cluster	nominal {cluster1, ... cluster 35}

Fig. 4. Result of preprocessing of AMR data.

2) Classifier evaluation

In our experiment, first, we build a classifier based on Temporal CARs from the preprocessed AMR training data. The accuracy was obtained by using the methodology of stratified 10-fold cross-validation (CV-10). One of the criteria for evaluating classifier is the accuracy of the classification results. We would like to be able access how well the classifier can classify. For this purpose, the sensitivity and specificity measures were used and accuracy is defined as:

$$accuracy = sens. \cdot \frac{Positive}{Positive + Negative} + spec. \cdot \frac{Negative}{Positive + Negative} \quad (6)$$

$$\left[sens. = \frac{True_Positive}{Positive} \right] \quad (7)$$

$$\left[spec. = \frac{True_Negative}{Negative} \right] \quad (8)$$

We have two important thresholds for the classifier performance (min. support and min. confidence). These thresholds control the number of patterns selected for constructing classifier so we used these thresholds.

Fig. 5 is the result of testing accuracy according to different minimum confidence and minimum support. Fig. 6 shows the example of extracted classification

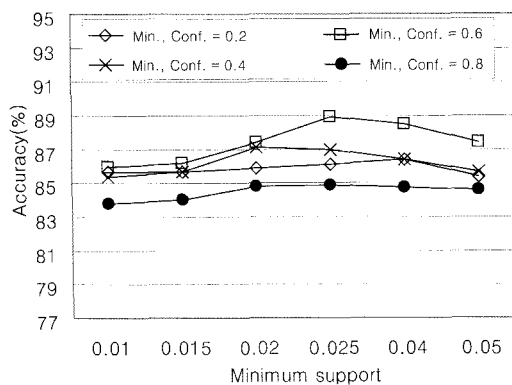


Fig. 5. Effect of support and confidence on accuracy.

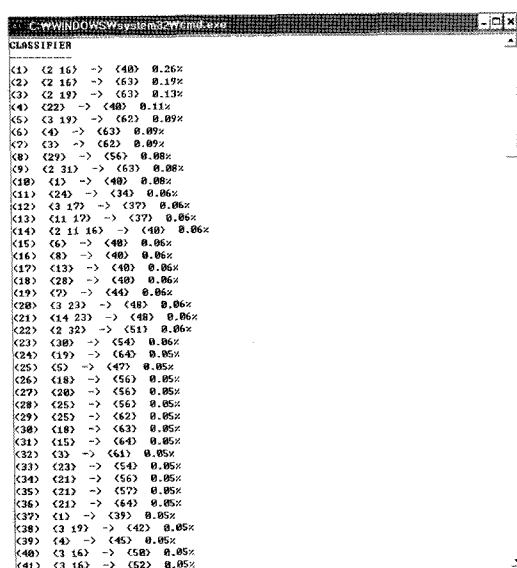


Fig. 6. the example of extracted rules for <3, *, 6, 3>

rules for calendar pattern, <3, *, 6, 3> that represents time intervals which means 9-17hour of the Saturday of every week of in March.

6. Conclusion

In this paper, we proposed a novel temporal pattern mining to predict power load patterns.

The proposed main mining tasks include clustering method and temporal classification method. Cluster analysis is used to define load pattern classes and the representative load profiles for each class. Classification method uses representative load profiles to build a classifier able to assign different load patterns to the existing classes. The proposed classification method is the Calendar-based temporal mining and it discovers electric load patterns in multiple time granularities. In experiment, the applied K-means and temporal classifier tested KEPRI AMR data and discovered interest load patterns.

Acknowledgements

This research was supported by development of AMR system interfacing model on internet GIS environment project of the Korea Electric Power Research Institute (KEPRI)

References

- Huang S. J. and Shih K. R., 2003. Short-term load forecasting via ARMA model identification including non-Gaussian process considerations. *IEEE Trans. Power System*, 18(2): 673-679.
- Pitt B. and Kirchen D., 1999. applications of data mining techniques to load profiling. *In Proc. IEEE PICA*, pp. 131-136.
- Amjady N., 2001. Short-term hourly load forecasting using time-series modeling with peak load capability. *IEEE Trans. Power System*, 16(3): 498-505.
- Li Y. and Ning P., 2001. Discovering Calendar-based Temporal Association Rules, *In Proc. of the 8th Int'l Symposium on Temporal Representation and Reasoning*.
- Jain A. K. and Dubes R. C., 1988. *Algorithms for Clustering Data*. Englewood Cliffs, NJ: Prentice-Hall.
- Park H. K., Kim Y., Park J. H., and Ryu K. H., 2007. Cluster Analysis for Region Electric Load Forecasting System. *to be appeared in Int'l Symposium Remote Sensing 2007*.
- Liu B. and Ma Y., 1998, Integrating classification and association rule mining. *In Proc. of the 4th Int'l Conf. Knowledge Discovery and Data Mining*.
- Agrawal R. and Srikant R., 1994. Fast algorithms for mining association rules. *In Proc. of the 20th VLDB Conference*, pp 487-499.
- Ernault M., 1982. Analysis and forecast of electrical energy demand. *Revue General d'Electricite*, No. 4.
- Fayyad U. M. and Irani K. B., 1993. Multi-Interval discretization of continuous-valued attributes for classification learning. *In Proc. of the 13th Int'l Joint Conf. on Artificial Intelligence*, pp. 1022-1027.