

Big Data Storage and Parallel Analysis of Grid Equipment Monitoring System

Xiaoming Zhou^a, Anlong Su^a, Guanghan Li^a, Weiqi Gao^b, Chunhua Lin^b, Shidong Zhu^{c,*}
Zhenliu Zhou^c

^aState Grid Liaoning Electric Power Co., Ltd, Shenyang, 110004, China

^bState Grid Dalian Electric Power Co., Ltd, Dalian, 116001, China

^cShenyang Institute of Engineering, Shenyang, 110136, China

Abstract

With the analysis on data feature of grid equipment operation monitoring, this work focuses on discussing the big data storage scheme for grid equipment online monitoring data, and describes optimization measure of grid monitoring data analysis. Based on the characteristics of large data scale, multiple data types and low value density with the online monitoring data, we provide a big data storage scheme based on HDFS cloud platform using consistent hashing. Meanwhile, we also employ a multi-channel data acquisition system using multiscale multivariate entropy as the feature extraction algorithm of the multi-source power grid monitoring data. To validate the efficiency of the algorithm, we perform experiments using power grid equipment ledger data, chromatographic hydrocarbons data of transformer oil, microclimate data, and transformer vibration data for association analysis. The big data storage scheme and the feature extraction algorithm proved that it could reduce the communication overhead between storage nodes, efficiently improve system performance, and is suitable for the actual application of power grid monitoring system.

Keywords: big data storage; electric equipment monitoring data; hadoop distributed file system; consistent hashing

(Submitted on November 2, 2017; Revised on December 25, 2017; Accepted on January 15, 2018)

© 2018 Totem Publisher, Inc. All rights reserved.

1. Instruction

A very important means to ensure the quality and safety of electricity is to monitor and analyze the data of power transmission and transformation equipment. Currently, the online monitoring system of power grids include the power station monitoring system, insulators leakage current monitoring system, transformer partial discharge system, the system of monitoring the icing of transmission lines, etc. These systems have numbers of monitoring points and high sampling rates, so the monitoring data becomes very large. For example, in a province power grid company, every 10,000 sets of terminals would generate about 1920G data each day, and it becomes about 700TB per year. Furthermore, massive historical monitoring data needs to be preprocessed before data analysis, such as through data cleaning, data format conversion, signal denoising, feature extraction, pattern recognition, etc. These processes require the system to have high throughput processing, parallel processing and batching. Meanwhile, the multi-source monitoring data should have association analysis in the application of the power equipment monitoring big data storage, such as weather and environmental information. In general, the data storage and management in grid companies are built on the basis of large relational databases, so the methods of status monitoring analysis to power grid equipment use SQL queries that are based on relational database or data warehouse. This makes system performance degrade due to large size of data. Another problem in current status evaluation method of power grid equipment is that the results of malfunction analysis are based on abnormal data analysis, but the large amount of historical data and critical state data, which are very important for malfunction analysis in power grid equipment, are deprecated.

Given the above, the online monitoring data of power equipment has the important characteristics of large volume, large type, fast change and low value density, making it suitable for big data storage and analysis technology. In big data storage

* Corresponding author.

E-mail address: zhusd519@163.com

and calculation technologies applications corresponding to our scheme, we used HDFS (Hadoop distributed file system), MapReduce framework, and the Apache Hadoop project to be the cloud platform for grid equipment monitoring analysis processing for reasons below: First, HDFS provides high reliability and scalable storage capacity for storage of massive historical data. Second, Hadoop's MapReduce technology shields a large number of low-level communications details, making users focus on system business logic development. Third, HBase is a non-relational database. It is suitable not only for storing structured, semi-structured and unstructured data, but also for power equipment sampling data, timing waveform signal and other data storage. HBase provides low latency performance while making online query.

In this work, we describe the big data storage and distribution strategy of power equipment monitoring and propose using multiscale multivariate entropy as the feature extraction algorithm of the multi-source power grid monitoring data for the quality and stability of the power grid analysis.

2. The research

In general, most of the monitoring platforms of power transmission and transformation equipment in provincial Electric Power Corporations are using Web Service to integrate the monitoring data. This method can only receive the data generated from monitoring equipment, and the analysis should be processed by using relational database model and SQL query. This approach not only degrades system elasticity, which means the system should give a heavy cost when the cluster size grows or contracts, but also reduces fault-tolerant because the system only provides the transaction level of fault tolerance.

Some researchers have proposed solutions of using cloud computing technology for storage and analysis on power equipment monitoring data. D. W. Wang et al. [8] described a solution of intelligent grid condition monitoring cloud platform based on Hadoop technology, but they were unable to present an actual prototype system. S. Singh et al. [6] designed and implemented a cloud-based service architecture for power monitoring data, but they did not provide the methods of data analysis and distributed storage model. D. S. Yang et al. [10] discussed and analyzed the key technologies of high-speed storage and retrieval of the massive power data, but they did not give an actual model that has been implemented. G. Reeves et al. [4] first proposed the compression storage method of time series data, which had greatly improved the system storage capacity. But, the method had a low query performance. X. W. Wang et al. [9] described a method for reducing sampling rate to achieve the storage and processing of time series power data, but this method was only good for low sampling rate data analysis and processing. G. L. Qu et al. [3] proposed to save power quality data on PQDIF (power quality data interchange format) format, and designed and achieved reliable mass PQDIF storage based on HDFS. Q. L. Zhang et al. [11] designed a cloud storage scheme of massive lightning monitoring power data based on HDFS, but they did not give the actual methods on monitoring data storage mode and data analysis.

Some researchers also proposed many methods on monitoring data analysis. Z. Fadika et al. [2] implemented a parallel data analysis algorithm using MapReduce framework, which can effectively improve the query efficiency of massive power monitoring data. D. W. Wang et al. [7] achieved a parallel Bayesian classifier for transformer fault diagnosis based on MapReduce framework. This method provides higher diagnostic speed than that of a single-machine environment. Y. Zhang et al. [12] proposed a parallel data lossless compression algorithm for variable section measured data based on MapReduce framework for mass data of smart grids. S. Rusitschka et al. [5] designed a cloud model for real-time data stream management in smart power grids, and they also built an actual application of the cloud storage in smart power grids that realized intelligent online measurement and management.

Most of the above research results are based on the application of big data technology, but there is a lack of living examples against characteristics of power equipment monitoring data. In this work, we focus on the correlation between various modes of monitoring data transmission and transformation equipment. For the low value density of massive historical data, we propose a big data storage scheme based on Hadoop platform, on which we achieved common solutions about data analysis. The analysis method supported not only monitoring data that is generated from grid equipment, also but meteorological data and historical data, which will be very helpful for the fast recognition and fault diagnosis to guarantee the safety, reliability and the quality of power supply.

3. The storage scheme of power grid monitoring data

The monitoring data storage is the premise and basis for online monitoring data analysis in power grid, as most power transmission equipment monitoring data is currently stored on enterprise-level relational databases. In this way, it is difficult to satisfy the requirements of rapid analysis and judgment to monitoring data, for there are a large number of monitor index and data type in the monitoring system. The HDFS supports the flow-based data access mode and supports high concurrent,

high throughput of data access. In this respect, we propose the storage and analysis solution of power grid equipment monitoring data by using big data technologies. We will discuss the big data storage architecture, storage distribution strategy and HDFS storage performance optimization below.

3.1. The storage architecture

The important factors of building big data storage architecture should be the mode of accessing and cleaning data from data source, data interaction, data computing and data storage. The different kinds of monitoring data from grid operating systems, line monitoring systems, etc. are gathered to a monitoring center system. They should be cleaned first and then stored to HBase database or written directly to the HDFS file for permanent storage. This part can be called the data access layer. The underlying data storage of cloud platform is the Hadoop cluster based on NameNode management. On the basis of HDFS, we built a distributed column-oriented HBase database, which is used for mass monitoring data storage and management. So, the data that are stored directly in HDFS or HBase database can perform batch data analysis tasks using MapReduce.

The integral storage structure for power equipment monitoring data can be logically divided into four layers: data source, data interaction layer, computing layer and storage layer, as shown in Figure 1.

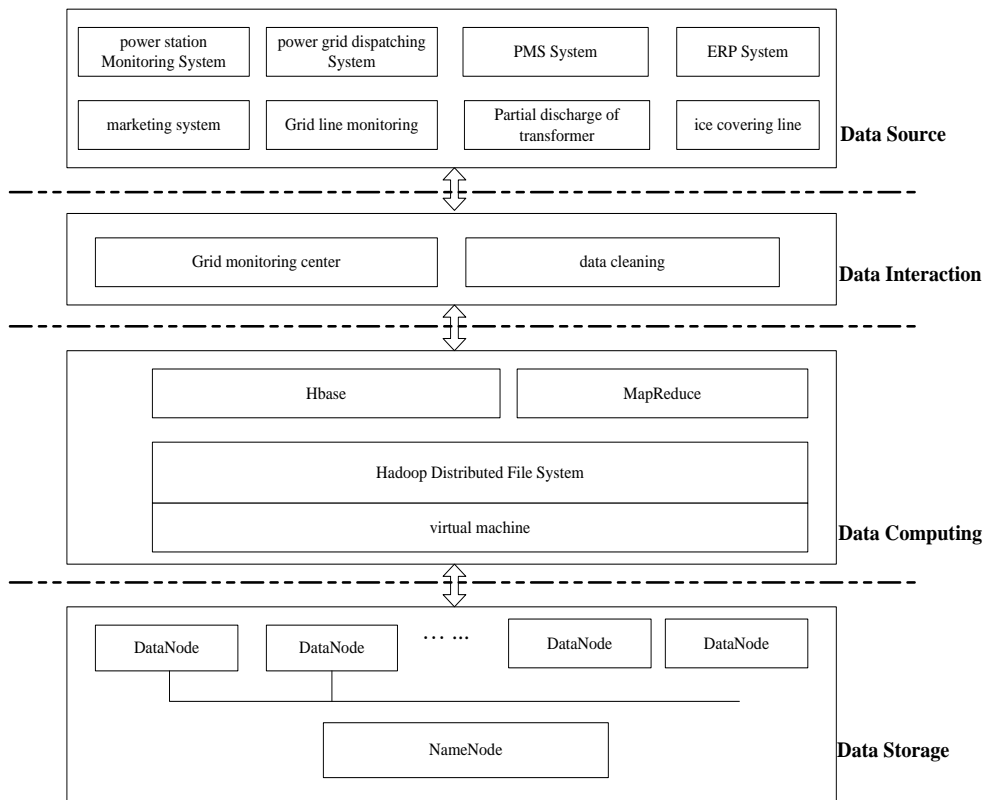


Figure 1. Storage architecture of power equipment monitoring data

3.2. The storage scheme based on consistence hashing

While the power equipment monitoring data are multi-source data and have huge data volume, the advantages of using distributed storage are to reduce network congestion and prevent system performance degradation. In many cases, different kinds of monitoring data have different effects to power supply quality, like the relationship between power equipment monitoring data and non-device factors (weather, temperature, etc.). They also have important effects to power supply quality. Meanwhile, the data distribution will also affect system performance during the parallel data processing by using MapReduce framework. Therefore, reducing the migration of related data is also key to the storage strategy.

For the purpose of improving efficiency of data analysis and processing, the multi-copy consistency hash algorithm is used widely among the different storage strategies. The main reasons are that it can achieve relational query of multiple-kind data and multi-channel data fusion. It also has obvious advantages of computational efficiency and data structure maintenance. The consistency hash algorithm can not only reduce the data migration when data nodes fail or increase, but can also

effectively solve the problems of load balancing and even distribution of the data on the condition of multi-copies data. When the data are uploaded to HDFS, they will be split into multiple blocks and stored to available nodes of the cluster according to a certain algorithm. The goal of data storage is to provide service for data query and analysis, and to improve the system performance and processing efficiency. Our storage strategy is analyzed as follows:

Firstly, the data should be distributed evenly to the nodes of the cloud storage cluster. As the data is processed in parallel on different nodes, load balancing can make high parallel computing on multiple nodes. Otherwise, unbalanced data distribution will lead to an idle state of some nodes, and other nodes are in the busy state. This will result in lower efficiency of the whole task. Secondly, the node failure in the cluster normal exists in the process of big data storage. The multiple copies strategy is the necessary means to solve the node failure, so the strategy should provide the solutions of data fault tolerance. Finally, the main factors that affect the system performance are I/O operations of storage media and data transmission in network during data storage. If the traffic between data nodes can be reduced, then the execution time of data analysis can be effectively reduced.

The associated data can be mapped and stored to the same node by using consistent hashing, which can reduce the communication overhead that is generated from data transmission from Map node to Reduce node. Our approach used 3 copies for storage. Each copy uses different keywords as the calculation parameters in the hash mapping. The algorithm flow is described as follows:

- Suppose the number of storage nodes is n , each node calculates the hash and selects 32 bits as its values using the MD5 algorithm, by which the storage nodes can be mapped to the same hash ring whose value spaces are $0 \sim 2^{32} - 1$.
- The hash values are calculated by different parameters in the same hash ring. In our approach, the first copy uses the line ID as the parameter for hash operation for the purpose of monitoring and analyzing for a certain power line. The second copy uses timestamp as parameter for the purpose of Line status query and analysis in a specific time, and the third copy uses user-defined attributes as parameter, which users can define by the needs of application. Each copy will be stored in different nodes because of different parameters.
- The storage nodes and replica data are arranged by their hash values in the same hash ring, as shown in Figure 2. When mapping, the replica data starts clockwise in the hash ring, and is stored to the nearest storage node.
- When the storage node fails or is removed, the original copy of the data mapped to the failed node needs to be recalculated and stored to the next node in the clockwise turn. When a new data node is added, as the data affected is only associated with the new added node and the previous storage node, the data only need to be recalculated and mapped to the added node.

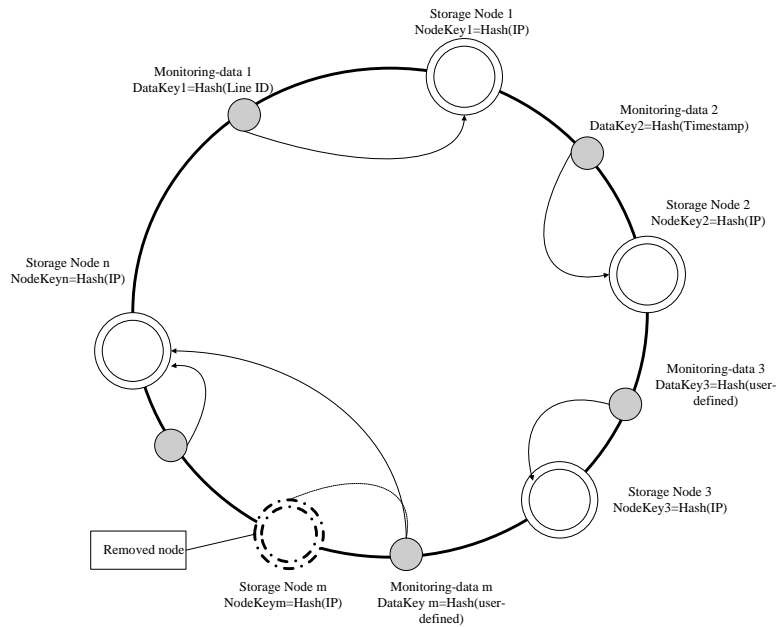


Figure 2. Consistent hashing storage mapping

4. Association analysis algorithm

The Multi-sensor measurement technology is referred to the data and are collected and stored by multichannel sequence. Based on this technology, the MMSE algorithm (Multiscale multivariate entropy) was proposed by M. U. Ahmed et al. [1] in

2011. This algorithm stores sample data to multiple files with a timestamp by multichannel. Before analysis, files are uploaded to HDFS and chunked to file block whose multiple copies are stored to different nodes. The signal segments with the same time stamp are assigned to the same MMSE calculation task. The data sieving and distribution are processed in Map stage. The MMSE calculation, whose result will be stored in HDFS, is not performed until the Reduce stage.

The disadvantage of the above algorithm is that it will produce a large amount of communication overhead between the Map node and the Reduce node during data screening and distribution. In our approach, we used the timestamp as the key to calculate the Hash values, which causes the synchronous data to be mapped to the same data node. In this way, the Map task does not need to perform data screening by timestamps any more, but can perform the MMSE calculation task directly. Our approach has reduced the communication overhead of data transferred from Map nodes to Reduce nodes. The process of data distribution and feature extraction with our optimization algorithm is demonstrated in Figure 3.

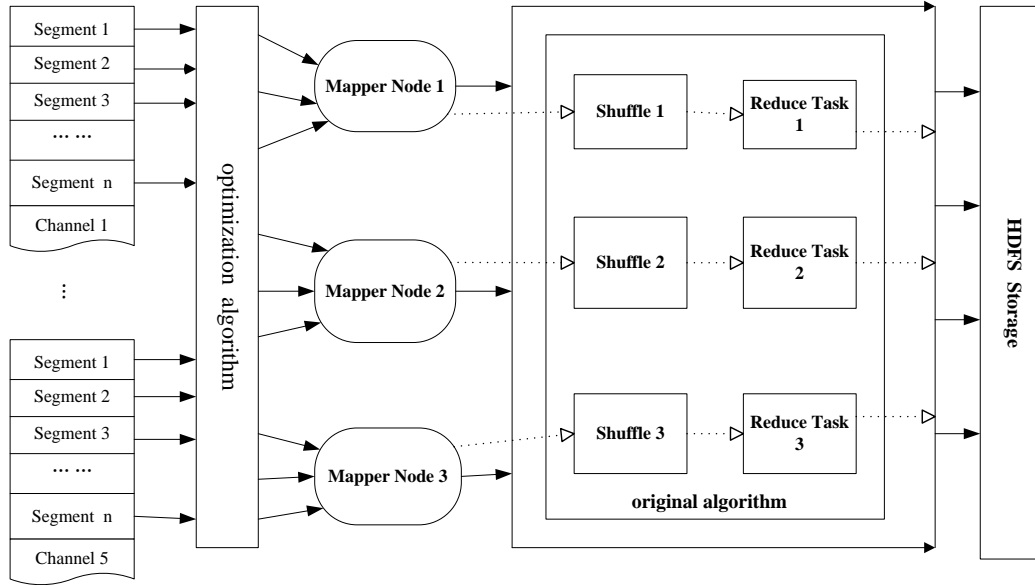


Figure 3. Optimized feature extraction process

The dotted square in Figure 3 describes the original MMSE algorithm, where the segmented data is firstly filtered and screened at the Map node, and then processed by shuffles. After that, the data will be sent to Reduce node for calculation. In our approach, the synchronous data with the same timestamp can be mapped to the same node directly, and the MMSE calculation can be directly executed in the local Map node. The optimized algorithm can save the data transmission overhead and accelerate the calculation process. The algorithm flow is described as follows:

- Suppose the time-series of original ‘ p ’ dimension channels are $\{x_{k,j}\}_{i=1}^N$, $K=1,2,\dots,P$, ‘ N ’ refers to the number of channel, ‘ ω ’ refers to scale factor, then multivariate time series can be expressed as Equation (1):

$$y_{k,j}^{\omega} = \frac{1}{\omega} \sum_{i=(j-1)\omega+1}^{j\omega} x_{k,i} \quad (1)$$

$1 \leq j \leq \frac{N}{\omega}$, $K=1,2,\dots,P$. when $\omega=1$, $\{y_{k,j}^{\omega}\}$ is the original time series.

- The embedded vector $M=[m_1, m_2, m_3, \dots, m_p]$ is default parameter, $\tau=[\tau_1, \tau_2, \tau_3, \dots, \tau_p]$ refers to the time delay vector, the process of building N -n complex vector $Y_m(i)$ using $\{y_{k,j}^{\omega}\}$ is Equation (2):

$$\begin{aligned}
Y_m(i) &= y_{1,i}, y_{1,i+\tau_1}, \dots, y_{1,i+(m_1-1)\tau_1}, \\
&\quad y_{2,i}, y_{2,i+\tau_2}, \dots, y_{2,i+(m_2-1)\tau_2}, \dots, \\
&\quad y_{p,i}, y_{p,i+\tau_p}, \dots, y_{p,i+(m_p-1)\tau_p} \\
m &= \sum_{k=1}^p m_k, i=1,2,\dots,N-n, n=\max\{M\} \times \max\{\tau\}.
\end{aligned} \tag{2}$$

- The distance between $Y_m(i)$ and $Y_m(j)$ is defined as Equation (3):

$$d[Y_m(i), Y_m(j)] = \max_{i=1,2,\dots,m} \{|x(i+l-1) - x(j+l-1)|\} \tag{3}$$

- For the given threshold 'r', the probability of pi for each 'i' is calculated as Equation (4):

$$\begin{aligned}
B_i^m(r) &= \frac{1}{N-n-1} p_i \\
d[Y_m(i), Y_m(j)] &\leq r, i \neq j.
\end{aligned} \tag{4}$$

The probability refers to not only relational degree of $Y_m(i)$ and $Y_m(j)$ ($i \neq j$), but also degree of regularity of $Y_m(j)$.

- Then calculate the average probability of 'i' expressed as $B^m(r)$ is in Equation (5):

$$B^m(r) = \frac{1}{N-n} \sum_{j=1}^{N-n} B_j^m(r) \tag{5}$$

- Let 'm' to be 'm+1' in step2, execute step 3 to 5 repeatedly, then get $B^{m+1}(r)$.
- At last get the value of MMSE as Equation (6):

$$MMSE(M, \tau, r, N) = -\ln \left[\frac{B^{m+1}(r)}{B^m(r)} \right] \tag{6}$$

5. Experimental Results

5.1. The experimental data

The experimental dataset for storage and analysis are the real sample monitoring data from the Power Supply Company of Shenyang City. The associated data we selected for testing include equipment record information, oil chromatogram hydrocarbons (H_2 , CH_4 , C_2H_2 , C_2H_4 , C_2H_6 , etc.), sample data of transformer, micro-meteorology data and vibration value of transformer tank. Table 1 showed the associated sample data to be tested for storage and analysis with our approach. The sensor of transformer tank is located at the position of three quarters from the top of the tank and the detection frequency is 200 HZ.

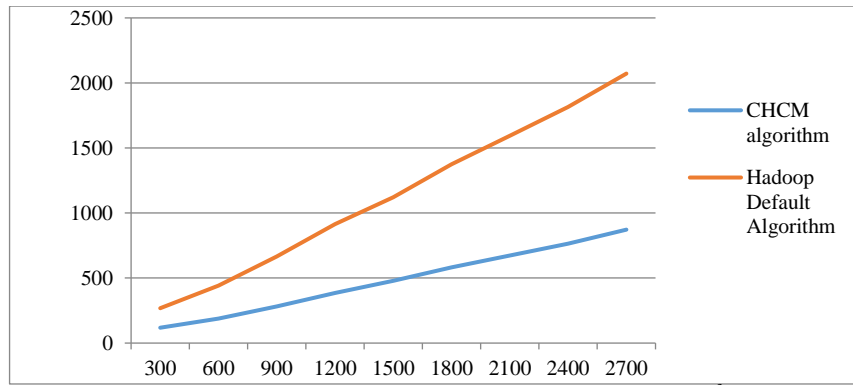
In the experimental dataset, we selected 1560 records from equipment ledger data, whose sizes are about 3.21 MB; the monitoring data of transformer THC are 1.586×10^6 records, size 1.68×10^6 MB; microclimate data are 3685 records, size 645 MB; and transformer tank vibration data are 1.165×10^6 records, size 1.17×10^6 MB. The system for testing are written by C++ language, which include data collecting interfaces, data cleaning module, relational query, multi-source data feature extracting algorithm, and analysis module. The testing system were built on Ubuntu operating system, which include one node named controller, one computing node, one network node and 9 storage nodes.

Table 1. The experimental dataset

Equipment Record				Sample Data of Hydrocarbons ($\mu\text{L/L}$)	micro-meteorology data			Vibration Value of Transformer Tank (mV) Frequency: 200HZ
ID	Name (220KV)	Place (city)	Time		Temperature ($^{\circ}\text{C}$)	Humidity (%RH)	temperature exchange value ($^{\circ}\text{C}$)	
51	LingBei Line	Shenyang	Nov-16	73.225,71.228,73.995,78.364,... ...,75.353	-8	66	0.78	189,199,205,185,,196
52	Kaifa Line	Fushun	Dec-16	78.362,77.985,73.668,81.263,... ...,74.235	-11	70	1.1	196,185,184,193,,188
53	Shanchan Line	Shenyang	Jan-07	79.412,73.894,76.952,83.223,... ...,76.369	-18	68	0.65	192,196,180,191,,179
54	Yuxiong Line	Yingkou	Jan-17	78.523,71.962,77.695,75.543,... ...,77.492	-16	72	0.96	187,206,179,187,,196
55	Dehua Line	Shenyang	Feb-17	81.885,75.662,74.387,78.989,... ...,82.556	-9	65	1.22	203,210,203,185,,188

5.2. The results of query testing

As the four kinds of sample data shown in Table 1 are able to reflect the working situations of electric transmission line, the performance of relational query has the important reference value for storage testing. We made a comparison of query run time with the data stored in the cloud using consistent Hashing and the data that is processed by standard algorithms of Mapeduce. The results are showed in Figure 4 below.

Figure 4. Time cost comparison of relational data query $\times 10^3$ records

From the experimental result, it can be seen that the time cost of relational data query using consistent Hashing approach is about 42 percent of the time cost of data query using standard algorithms of Hadoop platform. The main reason is that the data were distributed by consistent Hashing approach before it was stored to the cloud, so the data connection of relational query can be executed in the local server and avoid the data transmission between Map node and Reduce node. And another advantage is that it reduced the time cost of disk I/O operations and time cost of Reduce task start.

5.3. The results of association analysis testing

For performance testing of association analysis, the objects of transformer tank vibration data we selected are three-phase winding transformer (model *SFPSZ-120000/220*), which include five ICP-type acceleration vibration sensors (100 mV/g) and its sampling frequency is 1000 Hz. The sensors are located at three quarters from the top. The sample data were stored in five files according to tunnel. As the size of each file is about 21.5MB, we replicated every file to the size of 0.78 GB for the purpose of testing the performance of large-scale dataset, whose size is 3.75GB in total. The data settings are: sampling site is 4315 points, multi-scale factors are 7 and 15, embedding dimension is set as ' $M = [2, 2, 2, 2, 2]$ ', time delay vector is set as ' $f = [1, 1, 1, 1, 1]$ ', the threshold parameter is set as ' $r = 0.48$ '. The dataset included 2700 records in total, and the time cost comparison with different multi-scale factor is shown in Figure 5 below.

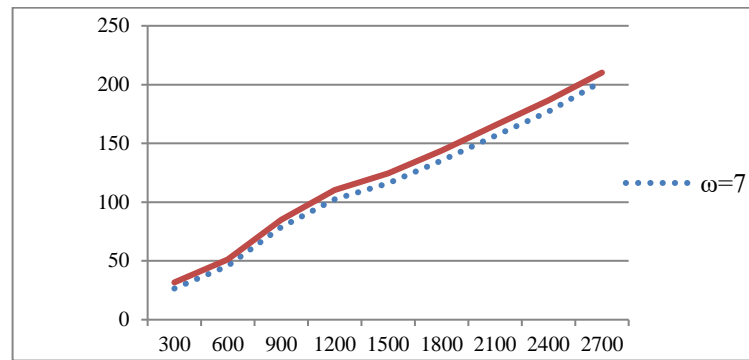


Figure 5. Time cost comparison with different multi-scale factor

It can be seen from the figure that the time cost of analysis using MMSE algorithm are at flat-lining growth with the increase of data. As the computing process is operated in Map processing, the time cost of MMSE computing avoid the influence of communication bandwidth. It proved that the approach had stable system performance and was suitable for large-scale data storage and analysis.

6. Conclusions

In this work, we first analyzed the monitoring data characteristic of transmission and transformation equipment in power grid, due to the data features of large scale, unstructured format, low value density etc., we described a big data storage approach using the multi-copy consistency hash algorithm, which could distribute data to the nodes of cloud storage cluster evenly based on HDFS platform. By this scheme, the duplicates were mapped by different feature parameter, and it is proven that the approach could reduce the communication-load of storage nodes and increase system efficiency. For association analysis, we proposed that features of multi-source data were extracted by multi-channel data collection and storage and analyzed by using multiscale multivariate entropy algorithm. It has been proven that this approach could improve system efficiency in the application of large scale data analysis in our experiment.

Acknowledgements

The authors acknowledge the State Grid Corporation Science and technology project (Contract number: 2017YF-34), the Liaoning Natural Science Foundation of China (Grant: 2015020020).

References

1. M. U. Ahmed and D. P. Mandic, "Multivariate Multiscale Entropy Analysis," *Signal Processing Letters*, vol. 9, no. 2, pp. 91-94, 2012
2. Z. Fadika, M. Govindaraju, and R. Canon, "Evaluating Hadoop for Data-intensive Scientific Operations," in *Proceedings of 2012 IEEE 5th International Conference on Cloud Computing (CLOUD)*, pp. 67-74, Hawaii, USA, June 2012
3. G. L. Qu, H. S. Yang, and Y. Zhang, "Quick Resolution about Mass Data of Power Quality Data Exchange File Based on Map-Reduce Module," *Power System Technology*, vol. 38, no. 6, pp. 1705-1711, 2014
4. G. Reeves, J. Liu, and S. Nath, "Managing Massive Time Series Streams with Multi-scale Compressed Trickle," *Proceedings of the VLDB Endowment*, vol. 2, no. 1, pp. 97-108, 2009
5. S. Rusitschka, K. Eger, and C. Gerdes, "Smart Grid Data Cloud: A model for Utilizing Cloud Computing in the Smart Grid Domain," in *Proceedings of 2010 First IEEE International Conference on Smart Grid Communications (SmartGridComm)*, pp. 483-488, Mayland, UK, Jan 2010
6. S. Singh and Y. Liu, "A Cloud Service Architecture for Analyzing Big Monitoring Data," *Tsinghua Science & Technology*, vol. 21, no. 1, pp. 50-70, 2016
7. D. W. Wang and X. J. Liu, "Parallel Fault Diagnosis Method of Power Equipment Based on MapReduce," *Electric Power Automation Equipment*, vol. 34, no. 10, pp. 116-120, 2014
8. D. W. Wang, Y. Q. Song, and Y. L. Zhu, "The Information Platform of Smart Power Grids Based on Cloud Computing," *Power System Automation*, vol. 34, no. 22, pp. 7-12, 2010
9. X. W. Wang, X. M. Zhai, and X. L. Jiang, "Self-adaption SPIHT Data Compression of Insulator Leakage Current," *Transactions of China Electro Technical Society*, vol. 26, no. 12, pp. 190-196, 2011
10. D. S. Yang, J. J. Chen, and M. Zhang, "Research and Application on Key Technology of High-speed Storage and Retrieval for Big Data," *Electronic Testing*, vol. 3, pp. 62-63, 2014
11. Q. L. Zhang, Z. Q. Zhou, and S. Q. Gu, "Application for Mass Data of Monitor Lightning in Cloud," *Automatic Control on Electrical Power System*, vol. 36, no. 24, pp. 58-63, 2012
12. Y. Zhang, H. G. Yang, and M. Q. Ye, "Management Program of Mass Electric Energy Quality Monitoring Data Based on Distributed File System," *Power System Automation*, vol. 38, no. 2, pp. 116-120, August 2014