

Express Supervision System Based on NodeJS and MongoDB

Li Liang

School of Computer Science
Communication University of China
Beijing Key Laboratory of China in
Security&Protection Industry
Beijing, China
18810361527@163.com

Ligu Zhu, Wenqian Shang, Dongyu Feng, Zida Xiao

School of Computer Science
Communication University of China
Beijing Key Laboratory of China in
Security&Protection Industry
Beijing, China

Abstract—Aiming at the functional requirements of the Express Supervision System, This paper discusses the advantages of using AngularJS to build the front-end framework, the advantages of using NodeJS to construct the back-end Web server, and the performance advantages of storing data based on MongoDB. This paper focuses on the storage solutions of using MongoDB to store large data and the statistical analysis solutions based on MapReduce. This paper argues on how to build Web services that meet the requirements of large data visualization based on NodeJs.

Keywords—Data storage; Statistic analysis; NodeJS; MongoDB;

MapReduce

I. INTRODUCTION (HEADING 1)

With the advent of the Internet+ era, online shopping has become an indispensable way of shopping in daily life. According to China online shopping market in 2015, China online shopping market in 2015 has reached 3.8 trillion yuan. As of 2016, the number of online shopping in China has reached 350 million. Hundreds of millions of online shoppers have generated a huge number of courier pieces. In 2015, the number of delivery has reached 20.5 billion. The average number of bills generated per day is 60 million. It is foreseeable that in the future, these figures will grow exponentially.

The data produced in Express is a very valuable digital resources, which includes an important indicator measuring the level of social economy in China. The data information from the Express have an important influence on the analysis of industry trends and the safety supervision. Therefore, the study of the data from Express has an important role in guiding the distribution of resources and standardizing the rules in Express. At present, there are a large number of off-line data in the security system accounting for about 400TB memory, which has reached 0.2% of the global production of printed materials(200PB). The data mainly includes the bill information, the state information and the branch information, which is the important sensitive data of the Express industry. Therefore, establishing a monitoring platform achieving the

management of the off-line data, has an important role in the prediction of the industry trend, the optimization of industry organization and safety supervision.

II. KEY TECHNOLOGIES

A. AngularJS

AngularJS is an open source JavaScript library maintained by Google to help a single page application run. Its goal is to enhance browser-based applications with MVC mode (MVC) capabilities, making development and testing easier.

B. MongoDB

MongoDB is a powerful, flexible, high performance, easy expanded way of data storage. It's a document-oriented database, not a relational database, is NoSQL.[1] The so-called document-oriented, use a more flexible "document" replace the "line" concept in original relational database. The document can be the value of array, documents and other complex data model. And the key of document are not pre-defined and will not be fixed. One of the main ideas of designing the MongoDB is that the operations that can be handed to the client are transferred from the server to the client, such as the generation of object id and other operational solutions. MongoDB as a general-purpose database, in addition to the ability to create, read, update, delete data, but also provides a series of unique features continue to expand. [2]

- MongoDB supports generic secondary indexes, allowing multiple quick queries, providing unique indexes, composite indexes, geospatial indexes, and full-text indexes.
- MongoDB supports "aggregated pipelines". Users can create complex aggregations from fragments, and automatically optimize them through the database.
- MongoDB supports a time limited set that applies to data that will expire at some point, such as a session. Similarly, MongoDB also supports fixed-size collections for storing recent data, such as logs.

- MongoDB supports a very easy-to-use protocol for storing file and file metadata.

MongoDB does not have some common features in relational databases, such as connecting queries and generating lines in order to get better scalability. Because these two functions in the distributed system are difficult to use. One of MongoDB's primary goals is to provide superior performance. MongoDB can dynamically fill the document, but also pre-allocation of data files to take advantage of additional space in exchange for stable performance. MongoDB uses as much memory as the cache, trying to automatically select the correct index for each query. In short, MongoDB in all aspects of the design is to ensure its stable performance. Although, MongoDB is very powerful and tries to preserve many of the features of a relational database, it does not pursue all the features of a relational database. Whenever possible, the database server will handle the generation and logic to the client to achieve. This compact design is one of the reasons MongoDB can achieve such high performance.

C. NodeJS

Node.js written in C++ language, is a JavaScript operating environment. Node.js is a JavaScript runtime environment. Node.js uses the Google Chrome V8 engine for good performance, and also provides a lot of system-level APIs, such as file operations, web programming, and so on. The JavaScript code on the browser side is subject to various security restrictions at run time, and the operation of the client system is limited. Node.js uses event-driven, asynchronous programming, and designed for network services. Node.js design ideas take the event-driven as the core, it provides the vast majority of APIs that are event-based, asynchronous style. Take the Net module as an example, where the net Socket object has the following events: connect, data, end, timeout, drain, error, close, etc. The developer using Node.js needs to register the corresponding callback function according to its business logic. These callback functions are executed asynchronously, which means that although these functions appear to be registered sequentially in the code structure, they do not depend on the order in which they appear, but rather wait for the corresponding event to fire. The important advantage of Event-driven and asynchronous programming is that make full use of the system resources. The implementation of the code without waiting for a certain operation to complete, and the limited resources can be used for other tasks. This design is very suitable for back-end network service programming, which is the goal of Node.js. In server development, concurrency request processing is a big problem, and blocking functions can lead to the waste of resource and the delay of time. Through event registration, asynchronous function, developers can improve the utilization of resources, and performance will also improve. From the supported module provided by Node.js, we can see that many of the functions, including file operations, are executed asynchronously, which is different from traditional languages. In order to facilitate server development, Node.js' network modules are particularly large, including HTTP, DNS, NET, UDP, HTTPS, TLS, etc., developers can build a Web server on this basis.

III. SYSTEM ARCHITECTURE AND TECHNICAL SOLUTIONS

A. System Architecture

The Express Supervision platform is divided into four main levels: data integration, data storage, data analysis and data application layer at the vertical level. The figure below shows the specific hierarchical structure layout of Express Supervision platform.

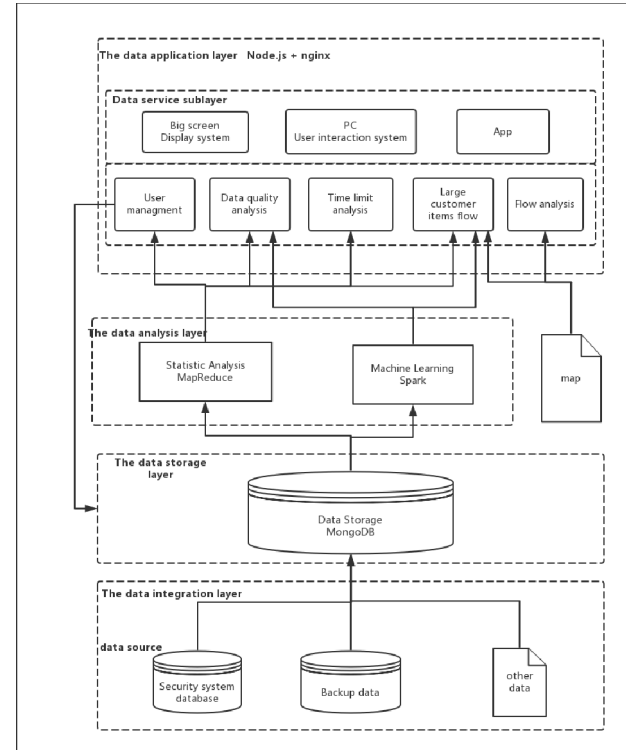


Fig. 1. System Architecture

In the data integration layer, we use real-time or dynamic method to export 11g data from the oracle database. At the same time, using the recovery method to export the backup data and using the data cleaning tools to achieve the data preprocessing. In the end, importing the data into MongoDB.

In the data storage layer, using the MongoDB to store data. Due to its fragmentation and weak consistency, MongoDB can ensure the speed of access. Coupled with the structure of its document storage, users can more easily access to the data.

In the data analysis layer, the system use the MapReduce to achieve the statistics analysis and Machine Learning to complete large data processing and analysis.

In the data application layer, Using Node.js + Nginx cluster, based on the standard restful cloud service interface to provide data services, the big data through the visualization has been presented in the front end.

B. Data Storage Solutions

The goal of the data storage solution for the express data is to meet the needs of storage, backup and high-performance read and write of large data (including historical data and real-time data) in the express delivery industry. It is designed to provide efficient data for data analysis and the basis for the protection of the performance.

a) Overview of the Overall Architecture: For the huge and complex data in Express industry, combined with the hardware performance and technical route requirements. The server cluster is divided into: analysis and application server group, Mongo DB storage server group. Analysis and application server group consists of two SX204-12 large data all-in-one machine, and each machine is equipped with four dual server nodes, which achieve the data analysis and processing, storage node routing and configuration control. Mongo DB storage server group consists of eight SX206-12 large data all-in-one machine, each large data one machine configuration 6 single server node, which achieve the data distributed storage[3]. It is shown in the below pictures:

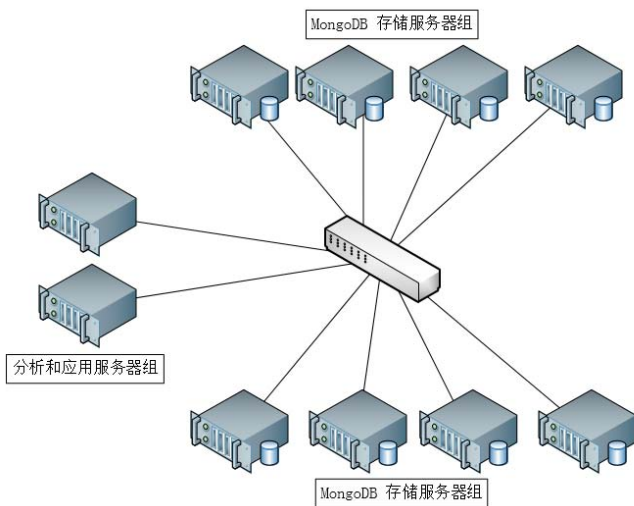


Fig. 2. System Architecture

b) MongoDB Distributed Storage Architecture:

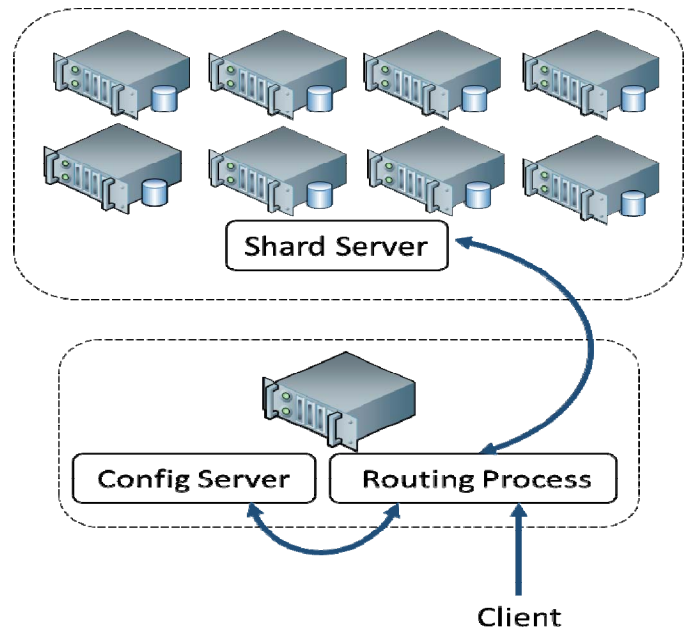


Fig. 3. MongoDB Distribute Storage Architecture Model

Figure 3 shows the MongoDB distributed storage architecture model, the specific components are as follows:

Shard Server: For the characteristics of express data, as well as robustness and efficiency of MongoDB distributed storage, the 48 nodes of storage server group are divided into 16 slices, with 3 copies in each slices(master, slave and arbitration node). Using the above configuration can be better to avoid a single point of failure of the host, but also automatically switch from the master and slave copy, taking into account the performance and efficiency[4].

Routing Process: As the routing service on the entire database storage speed and read and write performance will have a greater impact, it will be allocated for the three high-performance server nodes, responsible for managing the fragment. The client is routed through the front-end route so that the entire cluster looks like a single database. The client application can be transparently used. The Routing Process does not store data from the Configuration Server.

Configuration Server: these 3 nodes reused by routing server are responsible for storing the entire cluster configuration information, that is, data and fragmentation of the corresponding relationship between slices.

c) Configuration Settings:

Chunks: In order to avoid Chunk is too small, resulting in frequent movement of stored data, combined with the actual size of the data, we divided by the size of 256MB.

Data Block: According to the daily generation of 100 million data, each data is about 0.5KB. We calculated that a single set need to set 200 data blocks.

Data Organization: The data is stored in the physical storage node according to the size of the data block and

Chunks by using the date key as the key value and the label item is mapped. The efficiency of the data operation is improved evenly.

Index: In order to improve the efficiency of data manipulation and management, the data is indexed with key values that uniquely label data.

d) Database Structure

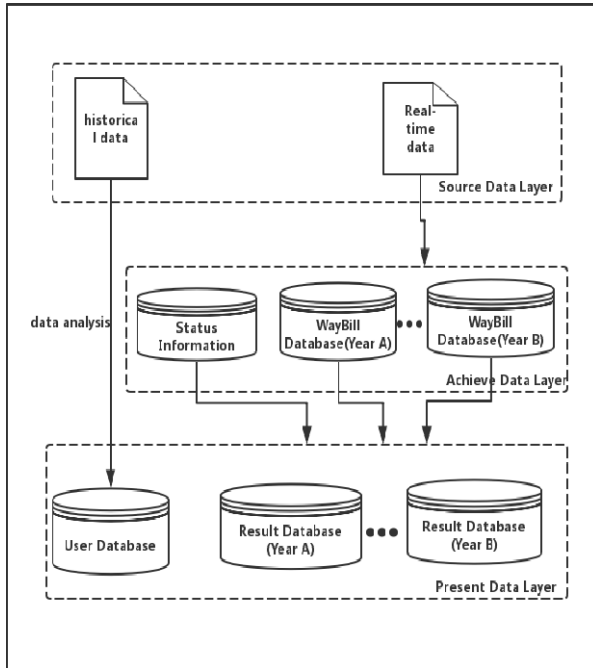


Fig. 4. Database Structure

Source Data Layer: All historical information is stored in the historical database, limiting a single set to a maximum of 100 million documents. Real-time data will be imported into the database, after its preprocessing, generating state information and waybills data temporary tables, which are into stored in a specific set after preprocessing.

Achieve Data Layer: The information in the state information database is only used as the auxiliary information of the next stage data analysis. Therefore, it is not necessary to do preprocessing in the information storage process, and the storage method of the waybill database can be stored. In accordance with the year of the waybills, a waybill database is maintained every year.

Present Data Layer: The data of covering all time range(Including historical data, real-time data collection), and the user-related information are processed and stored in the database, limiting the individual collection to a maximum of 100 million (100 million) documents.

C. Statistics Analysis Solutions

Mongo DB is one of the most popular large data storage platforms, and it can be used to provide source data as the

bottom storage layer of cloud computing technology, such as Spark, Hadoop, Pig, Hive, Drill and other computing framework. At the same time it provides the aggregation pipeline, MapReduce features to support the calculation of large data, statistics, classification and other needs. Aggregation pipeline and MapReduce are collectively referred to as aggregation operations in MongoDB.

a) Aggregation Pipeline: Aggregation Pipeline is a newly introduced aggregation framework for MongoDB 2.2, which follows the Unix-style data processing pipeline. The output is a result of the aggregation after that the documents in the collection enter the multi-level pipeline[5]. The first step in the processing pipeline may contain a plurality of components, including filtering, projecting, grouping, sorting, limiting, skipping and so on.

Objectid("56dd7f120a5288c0c40dee28")	Objectid
210670000052	String
陈刚	String
湖北省荆门市东宝区解放路95号	String
42	String
13681172939	String
13681172939	String
2015-04-16 11:24:45.619Z	Date
朱文海	String
广州市天河区桃园中308号1404房	String
62	String
18922260503	String
18922260503	String
2015-04-16 22:06:47.632Z	Date

Fig. 5. Waybill Information

Figure 5 shows a waybill document that stores fields such as Receiver_name, Receiver_mobile, Receiver_date, and so on. The goal is to count the number of different users (differentiated by the Receiver_name + Receiver_mobile combination field) and the first use of time. The processing flow is divided into three steps, preprocessing, statistics and result output.

Preprocessing: Preprocessing: First, according to the Receiver date field, using the sort command \$sort of MongoDB to sort the data in the collection at the first stage of the pipeline. During this period, you can use indexes to sort quickly, and the result of this step is to provide an ordered data set for calculating the earliest time[6]. Otherwise it will be slower to sort after other processing steps and take up a lot of memory.

Statistics: Using the mongo DB grouping command \$group to group the documents according to the specific fields(id). Specifically, specifying the Receiver_name + Receiver_mobile combination field as the _id field, which will identify different users. In this case, you also need to count the user's earliest time of using and the number. The earliest time can be queried using the \$ first command, and \$ first is more efficient in the sorted data set in the corresponding field. The number of users is passed through the \$sum operator, which adds the value to the calculated result for each document in the group[7]. In the application case, setting value to 1 which can be used as a statistic value.

The output: the final step is to output the result, and the result document can be written to the target collection via the \$ out operation, and the operation can also bypass the Mongo DB document size limit. The problem with the target set is

described as follows: When the target set does not exist, the collection is created, which is visible after the end of the aggregation; when the target collection already exists, the \$ out operation automatically overwrites the original collection.

v _id	{2 fields}	Object
name	刘睿源	String
mobile	15924701976	String
useDate	2015-01-08 04:29:40.893Z	Date
count	12	Double

Fig. 6. User Information

Where _id field records the user ID field, userDate records the user's earliest time, and count is the total number of times the user's use.

b) *MapReduce*: MapReduce is software architecture model proposed by Google for large-scale data sets (greater than 1TB) of the parallel operation. Concepts "Map" and "Reduce", and their main ideas are borrowed from the functional programming language, as well as from the vector programming language borrowed features. MongoDB implements MapReduce using the same processing paradigm. In the map function, the values in the current document are mapped to the key-value pair and emit; the reduce function takes the arguments and merges the values of the same key into a document. By dividing the calculations into parallel add operations, MapReduce can be to achieve large-scale data computing functions[8].

For example, we assume that the goal is to obtain the total amount of shipments of different users and the time of the first use. First implementing the map part, the Receiver_name, Receiver_mobile are taken as a key to group, and setting the value of {'count': 1, 'dates': [this.Receiver_date]}. The count field is the number of users, and the dates array stores the user time data. The format of the data is used to ensure the consistency of the values and reduce the time array.

In the reduce phase, the value of the current key is counted. Count field can be a simple summation. After the dates array merges the time data from the current values into an array, it need to be reassigned the value of the smallest time in the array. The other optimizations are sorted using the sort key of the map, which reduces the number of operations[9]. Finally, our calculations are shown in Figure 7.

v _id	{2 fields}	Object
name	刘睿源	String
mobile	15924701976	String
value	{2 fields}	Object
count	12	Int32
dates	2015-01-08 04:29:40.893Z	Date

Fig. 7. User information

IV. SYSTEM DISPLAY

A. User Query

The supervisor can search for the relevant information and data in the user database. The user database information is mainly divided into individual users and institutional users.

The user's statistical data can be obtained by the three steps of the above data aggregation processing. The result is shown in the following figure

You can browse and access information about individual users and institutional users. Personal user database information items include: name, address, telephone, ID number, code number, etc. Institutional user database information items are divided into: organization name, address, telephone, organization code, code number and other information.

ID	姓名	地址	电话	备注
1	李华	北京市海淀区中关村大街10号	1342-463-8841	...
2	王明	上海市浦东新区陆家嘴环路1000号	15016-8956-1120	...
3	张三	广州市天河区珠江新城华夏路10号	1724-806-3887	...
4	赵四	深圳市南山区科技园科兴科学园	1412-488-8725	...
5	王五	成都市高新区天府大道北段1700号	1849-545-8777	...
6	陈六	武汉市武昌区中南路100号	1470-329-9527	...
7	李七	杭州市西湖区文三路100号	1388-191-2346	...
8	李八	深圳市南山区科技园科兴科学园	1663-555-8904	...
9	李九	上海市浦东新区陆家嘴环路1000号	1598-224-7837	...
10	李十	广州市天河区珠江新城华夏路10号	1212-365-8881	...

ID	机构名称	地址	电话	备注
1	中国邮政集团有限公司	北京市东城区东直门内大街100号	1342-463-8841	...
2	上海浦东发展银行股份有限公司	上海市浦东新区陆家嘴环路1000号	15016-8956-1120	...
3	深圳农村商业银行股份有限公司	深圳市福田区福田街道福安社区	1724-806-3887	...
4	四川省农村信用社联合社	四川省成都市高新区天府大道北段1700号	1412-488-8725	...
5	山东农村商业银行股份有限公司	山东省济南市经二路纬三路	1849-545-8777	...
6	湖南农村商业银行股份有限公司	湖南省长沙市岳麓区岳麓大道	1470-329-9527	...
7	江西农村信用社联合社	江西省南昌市红谷滩新区红谷大道	1388-191-2346	...
8	广州农村商业银行股份有限公司	广东省广州市天河区珠江新城华夏路10号	1663-555-8904	...

Fig. 8. User Query Page

B. Statistic Chart

a) *Regional User Statistics*: Users can browse the real number of users in the provinces dynamically. Users can also see the number of national, provincial and municipal users and the corresponding trends. On the histogram, the user understands the number of users in each province by the depth of the color and the situation. Through the realization of the visualization tool, you can also choose to select a province to browse to the province each city each county's new user statistics.

V. CONCLUSION

This paper introduces the construction of the Express Supervision System. First it introduced the key technology used in system construction. Then it mainly introduces the system architecture, the data storage solutions based on MongoDB database and the statistical analysis solutions based on MapReduce. Finally, this paper presents user query and statistical analysis modules for the Express Supervision System. In the future, we will use the method of machine learning to optimize the statistical analysis method and get a more accurate analysis result.

ACKNOWLEDGMENT

This paper is partly supported by “Key Cultivation Engineering Project of Communication University of China (Project number: 3132016XNG1606 and 3132016XNG1608)”, “Cultural technological innovation project of Ministry of Culture of P.R.China (Project number: 2014-12)”, and partly supported by “The comprehensive reform project of computer science and technology, department of science and Engineering”. The research work was also supported by “Chaoyang District Science and Technology Project (CYXC1504)”.

REFERENCES

- [1] Y., G., et al. Analysis of data storage mechanism in NoSQL database MongoDB. in 2015 IEEE International Conference on Consumer Electronics - Taiwan. 2015
- [2] V., A. and M.R. C. MongoDB and Oracle NoSQL: A technical critique for design decisions. in 2016 International Conference on Emerging Trends in Engineering, Technology and Science (ICETETS). 2016.
- [3] Zhu, Q. and P. Zhou, The System Architecture for the Basic Information of Science and Technology Experts Based on Distributed Storage and Web Mining. 2012. p. 527-530.
- [4] Kookarinrat, P. and Y. Temtanapat, Analysis of Range-Based Key Properties for Sharded Cluster of MongoDB. 2015. p. 1-4.
- [5] Bonnet, L., et al., Reduce, You Say: What NoSQL Can Do for Data Aggregation and BI in Large Repositories. 2011. p. 483-488.
- [6] Hong, T.P., et al., Efficient data preprocessing for genetic-fuzzy mining with MapReduce. 2015. p. 88-89.
- [7] Arora, S. and I. Chana, A survey of clustering techniques for big data analysis. 2014. p. 59-65.
- [8] Nabavinejad, S.M. and M. Goudarzi, Faster MapReduce Computation on Clouds through Better Performance Estimation. IEEE Transactions on Cloud Computing, 2017. PP(99): p. 1-1.
- [9] Pol, V.V. and S.M. Patil, Implementation of on-process aggregation for efficient big data processing in Hadoop MapReduce environment. 2016. p. 1-5.

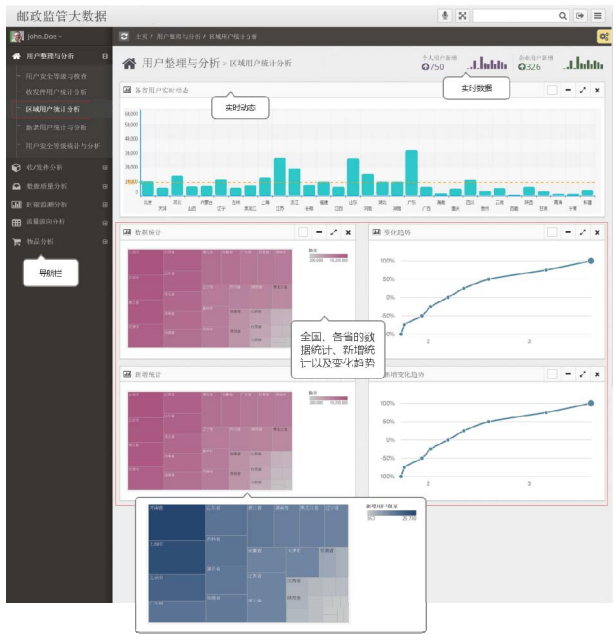


Fig. 9. Regional User Statistics

b) User Safety Level Statistics: With the addition of dynamic users can understand the users increased by real time, and you can choose to browse to different levels of security by the number of users statistics and changes in the situation, so that different levels of user data statistics can become more intuitive and clear. Users can also see the growth of different security levels of users, as well as the proportion of all different levels of users.

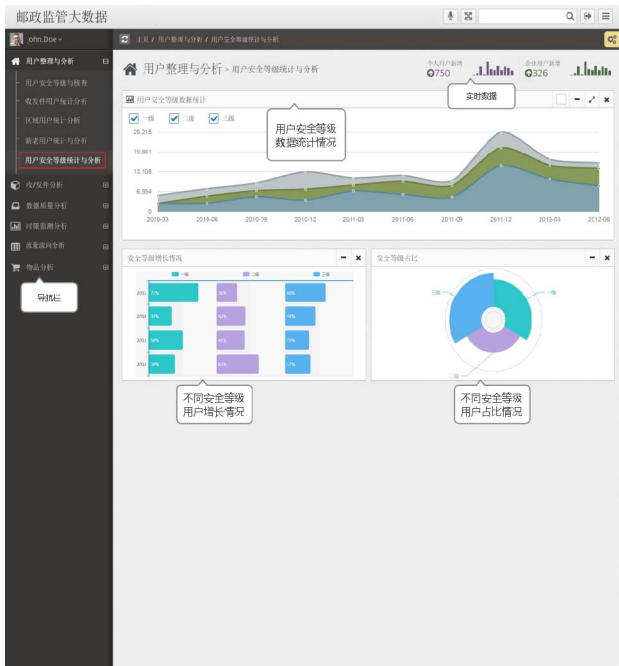


Fig. 10. User Safety Level Statistics