# GROUP PROJECT FINAL REPORT

# Research On Doctor-Patient Behavior And Medical Resource Allocation Of Internet Medical Platform

廖家瑞  (Jiarui Liao) 1801212884
刘德明  (Deming Liu) 1801212889
曹阳  (Yang Cao) 1801212825
王如思  (Rusi Wang) 1801212938
常鑫磊  (Xinlei Chang) 1801212779

## 1   Project Introduction

In recent years, China's medical and health undertakings have developed significantly, but this development also faces some challenges. For example, there is information asymmetry between doctors and patients, and the uneven distribution of medical resources between regions and between urban and rural areas. These phenomena have led to contradictions between doctors and patients and problems such as "inadequate and overly expensive medical service" in large hospitals. Specifically, compared with developed countries, China has been suffering from serious shortage of medical resources and unreasonable allocation of medical resources for a long time. According to the statistics of our ministry of health, our country has a huge population, accounting for more than one fifth of the world's population, but we only have two percent of the world's medical resources.

Internet medical treatment was created in order to solve these problems under the technical background of the time. Internet medical treatment will help patients search for information, and it will also facilitate the reorganization and sharing of resources such as doctors' knowledge and skills in space and time. Therefore, Internet medical treatment will greatly change the traditional medical model, and have a certain impact on patients' medical decision. How to use information and network technology to solve the imbalance of supply and demand and distribution of medical resources has become a hot issue in the field of research and practice.

This project will be carried out with multi-dimensional Internet data of the Haodf medical platform. **Established in 2006, Haodf medical platform is the leading Internet medical platform in China.** (https://www.haodf.com) We will crawl the data from 2008 to 2018 on Haodf website, and then carry out subsequent cleaning, analysis and visualization. The data set will contain all hospitals of Shanghai city department of the relevant data, including basic information on doctors and patients, the geographical position, participation, including the time of registration, login, and stay logged in time, number of posts, etc.), post (post, reply, questions), the content of the post (the length of the posts, posts, the theme of the emotion), evaluation data, etc. This project will examine the impact of geographic factors and willingness to consult from the perspective of doctor-patient interaction (reflected in doctor-patient counseling and response). The project will focus on the improvement effect of the Internet on the imbalance of medical resources, and also plan to visualize the geographical distribution of medical allocation on the Internet medical platform. Finally, this project hopes to put forward corresponding suggestions on the influence of Internet medical treatment on doctor-patient behavior and resource allocation.

# 2 Data Description

## 2.1 Data overview

On the website of Haodf, we collected all relevant data from most departments of all hospitals in some regions (Shanghai, Anhui and Gansu were sampled accordingly) The data includes basic information on doctors and patients, the geographical position, participation, (including the time of registration, login time, stay time, number of posts, etc), post (post, reply, ask questions), the content of the post (the theme of the posts, the length of the posts, the emotion), evaluation data, etc.

Generally speaking, the data we collected in the project can be divided into three parts.

The first part includes hospital information, such as location, grade, department name, etc (as shown in **Figure 1.**).

The second part includes the information of doctors, such as location, title and service status (as shown in **Figure 2.**).

The third part includes patient information, such as: location, consultation times, consultation length of time, etc (as shown in **Figure 3.**).

## 2.2 Intrinsic big data properties

The data of this project conforms to the inherent characteristics of big data.

1) **Variety:** both structured and unstructured. Structured data, such as the age of the patient, the title of the doctor, the number of consultants, etc. Unstructured data mainly includes the content of consulting details.

2) **Volume:** the project has collected relevant data of all Shanghai hospitals on the website of Haodf for nearly ten years, with the total number of tens of millions, and the storage capacity is TB, which is consistent with the characteristics of large data Volume.

3) **Velocity:** project data includes doctor-patient communication data all the time. The data changes quickly and must be processed quickly.

4) **Value:** the data of the project itself has a low Value density, but if appropriate processing is carried out, valuable information can be mined, such as the improvement effect of Internet medical model on the allocation of medical resources.

## 2.3 Relational databases

This project plans to use a **relational database: Microsoft Access** to store data. The Access database can connect the data tables used by the project to obtain the table of research objects, representing the relationship between the data. The connection relationship between the tables is shown in **Figure 4**.

According to the data we collected in 2.1, we split the data into six tables: *Hospitals, Departments, Doctors, Consults, Patients, Disease_for_department* (as shown in **Figure 4** in the appendix). The database is designed based on 2NF (considering the difficulty of SQL query, we did not implement 3NF, because there might be more than 1 primary key in a table, which may bring

difficulty when designing SQL query statements.)

## 2.4 Data preprocessing

There is a lot of messy code and data redundancy in the data, so it is necessary to preprocess the data before the formal analysis (as shown in **Figure 5** and **Figure 6**).

# 3 Data Visualization

## 3.1 Main Idea

After data processing, we first perform data visualization. Good geographic information display can directly highlight the problem and also help to better predict the possible results. We first make a heat map of consulting Shanghai doctors across the country. And by collecting GDP per capita, the number of patients, and the number of netizens to intuitively test our hypothesis that the number of consultants is related to GDP and the number of Internet users. Due to the large amount of data, it is difficult to analyze the national data in a short time, so we chose to extract relevant data from Anhui Province and Shanghai for detailed analysis. Based on the results of data visualization, then we will mainly analyze the impact of Shanghai's medical level and other factors on Anhui patients' choice of local or Shanghai medical treatment. In the process of data visualization, we used the Power Map plug-in of Excel, the business intelligence software Tableau, and the Python drawing package Pyecharts.

## 3.2 Nationwide Analysis

Cluster chart (as shown in **Figure 7**): We combine the per capita GDP, the number of patients, and the number of netizens, and then normalize them in this Cluster diagram. It's to intuitively test our hypothesis that the number of consultants is related to GDP and the number of Internet users.

Heat map (as shown in **Figure 8**): This map can intuitively reflect the quantity and intensity of distribution through color and geographic location. It can be clearly seen that the distribution of consulting patients is closely related to distance and population density.

Combining the two graphs above, we can get several conclusions that: (1) The closer the area is to Shanghai, the greater the number of people seeking medical treatment in Shanghai. Obviously, Shanghai has the largest number. (2) The number of people consult Shanghai doctors is also related to GDP per capita. For example, Guangdong has a higher GDP per capita. Although it is far away from Shanghai, there are more people than other regions. (3) The choice of patients is also related to the Penetration of the network. The higher the proportion of Internet users, the corresponding number of people who use Haodf online is larger compared with the regions where distance and per capita GDP are close.

## 3.3 Province Level Analysis

We draw two bar charts on the doctor and patient level statistics of Anhui Province in consult list, and four regional maps of the region feature of Anhui Province.

The bar charts (as shown in **Figure 9** and **Figure 10**) show that:

(1) Although the departments are different, the main part of the consultations is to Grade A hospital.

(2) From the perspective of departments, the number of pediatric consultations occupies a

major advantage. Dermatology, surgery and internal medicine are for the main proportion of consultation. This conclusion is related to the characteristics of the disease.

(3) Patients prefer to consult Chief, Associate Chief, Attending Doctors online.

The regional maps (as shown in **Figure 11-14**) give an intuitive observation of the features of cities in Anhui Province. Through the analysis of the medical level scores, GDP per capita of cities in Anhui and distance and difference of medical level between Shanghai and cities in Anhui. Then we can have a unified understanding of the medical level and GDP of each city in Anhui Province. This information could help us in the model construction.

# 4 Model and Implementation

## 4.1 PCA

Principal component analysis is a statistical procedure that uses orthogonal transformation to convert a set of observations of possibly correlated variables into a set of values of linearly uncorrelated variables called principal components. The number of principal components is less than or equal to the number of original variables. This transformation is defined in such a way that the first principal component has the largest possible variance (that is, accounts for as much of the variability in the data as possible), and each succeeding component in turn has the highest variance possible under the constraint that it is orthogonal to (i.e., uncorrelated with) the preceding components. Principal components are guaranteed to be independent if the data set is jointly normally distributed. PCA is sensitive to the relative scaling of the original variables.

Since our purpose is to find the factors that influence the patient behavior, we use principal components analysis to construct the medical level indicator. Based on many medical references and data we collected, we decide constructing the medical level indicator by using eight features which include the elderly dependency ratio, Number of medical institutions, Number of hospital, Number of medical beds, Number of medical personnel, Number of practicing physician, Average hospitalization day and GDP.

The scree plot shows that three components can explain the 81.5% of total variance, which is enough for our analysis. We also find the intuitive meaning of these three components, we can see that the first components is highly correlated with GDP and medical facilities indicator, so the first component reflects the economic and medical resource level. The second components is highly corelated with the number of hospital and the number of medical institutions, which reflect the degree of medical competition. The third component is highly correlated with the elderly dependency ratio, which reflect the population and aging factor.

After contracting all the indicator to the three components, we construct the Medical level using these three components, the weight of each components can be calculated by the following formula:

$$Weight_i = \frac{Var(x_i)}{\sum_{i=1}^{3} Var(x_i)} = \frac{\lambda_i}{\sum_{i=1}^{3} \lambda_i}$$

$$Medical\ level = \sum_{i=1}^{3} Weight_i * Component_i$$

## 4.2  Logistic Regression

After using PCA to form 3 principle components, we add the weighted value of each component to form a new variable called Medical_Level. It means the comprehensive score of the city's medical level. And we combine Medical_Level and distance to Shanghai and GDP of this city into the patient's data. However, every patient now has the same city attributes, it is hard for Logistic Regression Model to discriminate the samples.

Therefore, we investigate the other information and invent a new variable called Utility Score. It measures the urgency or utility of a patient. For example, if a patient uses more words to describe his illness, the Utility Score will be higher. Also, if a patient makes a phone call to the doctor, the Utility Score will be higher too. In general, we define Utility Score as a comprehensive score of Length of Description, Numbers of Phone Calls, Total Time of Consulting.

Then we use Logistic Regression to process our data. The independent variables are city's distance to Shanghai, GDP, Medical_Level and the Utility Score of each patient. The dependent variable is whether the patient chooses Shanghai's hospitals. If it is true, it will be 1, otherwise 0. The distribution formula as follows:

$$P_{Y=1} = \frac{e^{\beta^T X}}{1 + e^{\beta^T X}}$$

The result shows that the accuracy, precision and recall are all about 60%, and the coefficients lead to our final results.

# 5  Summary and Outlook

The final question of our project is whether Internet medical platform is conducive to the allocation of medical resources? Our answer is yes! Taking Shanghai as an example, it can be seen from the visualization results that the consulting patients are all over the country, and people all over the country can enjoy the high-quality medical resources of Shanghai (especially Grade-A tertiary hospital) to a certain extent.

Regional factors and willingness to consult affect patients' decision-making. (1) The closer the distance is, (2) the higher the GDP is, and (3) the poor the medical level is, the more patients choose Shanghai, and (4) the stronger their willingness to consult, the more they choose Shanghai. That may because of the convenience of distance, the more abundant funds and the stronger willingness of better hospitals and doctors.

What else can we do in the future? The data can be combined with the patient's personal information (family income, duration of illness, age) to provide a more detailed analysis of the factors influencing the patient's decision. What's more, time variables can be added into the model to study the changes of patients' decision-making in different periods.

# Appendix



**Figure 1: Hospital information interface (Source: Haodf website)**



**Figure 2: Doctor information interface (Source: Haodf website)**

**Figure 3: Patients information interface (Source: Haodf website)**



**Figure 4: Relational Database based on 2NF**

**Figure 5: Using SQL query code shown on the left, we can see there exists messy code and redundancy in the data.**



**Figure 6: After data preprocessing, the data is cleaner and can be used in further analysis.**

**Figure 7: Cluster chart of province level features**



**Figure 8: Heat Map of online consult to Shanghai**

安徽咨询人数（按科室）_省

clusternam..

**Figure 9: Doctor level statistics of Anhui in consult list**



安徽咨询人数（按科室）_省

clusternam..

**Figure 10: Hospital level statistics of Anhui in consult list**

**Figure 11: The region feature of Anhui Province (I)**

# Medical level score of cities in Anhui Province

**Figure 12: The region feature of Anhui Province (II)**

GDP per capita of cities in Anhui Province

**Figure 13: The region feature of Anhui Province (III)**



Difference of medical level between Shanghai and cities in Anhui

**Figure 14: The region feature of Anhui Province (IV)**

**Figure 15: PCA Scree plot**

| Feature | Components_1 | Components_2 | Components_3 |
| --- | --- | --- | --- |
| elderly dependency ratio | -0.421 | -0.129 | 0.799 |
| Number of medical institutions | 0.520 | 0.582 | 0.499 |
| Number of hospitals | 0.349 | 0.758 | -0.276 |
| Number of medical beds | 0.861 | -0.405 | 0.054 |
| Number of medical personnel | 0.939 | -0.207 | 0.125 |
| Number of practicing physician | 0.936 | -0.238 | 0.157 |
| Average hospitalization day | 0.750 | -0.426 | -0.238 |
| GDP | 0.787 | 0.559 | 0.062 |

**Table1: Components matrix**

| 地区 | 到上海的距离（公里） | 人均地区生产总值（全市）/元 | 医疗水平综合得分 | 与上海医疗水平的差异 |
| --- | --- | --- | --- | --- |
| 合 肥 市 | 408.1 | 67689 | 13276.5175 | 5827.21859 |
| 淮 北 市 | 532.5 | 35324 | 6328.77547 | 12774.9606 |
| 亳 州 市 | 607.8 | 17769 | 3576.3282 | 15527.4079 |
| 宿 州 市 | 503.9 | 20895 | 4248.65928 | 14855.0768 |
| 蚌 埠 市 | 429.7 | 35542 | 6995.21678 | 12108.5193 |
| 阜 阳 市 | 563.5 | 15303 | 3242.62066 | 15861.1154 |
| 淮 南 市 | 451.5 | 33361 | 4861.05483 | 14242.6813 |
| 滁 州 市 | 315.5 | 30562 | 6026.79145 | 13076.9447 |
| 六 安 市 | 472.4 | 19211 | 4209.93965 | 14893.7965 |

| | | | | |
|---|---|---|---|---|
| 马鞍山市 | 286.1 | 60091 | 10923.092 | 8180.64415 |
| 芜　湖　市 | 289.6 | 64039 | 12188.9028 | 6914.8333 |
| 宣　城　市 | 259.3 | 35726 | 6864.70094 | 12239.0352 |
| 铜　陵　市 | 350.9 | 97193 | 10218.1817 | 8885.55439 |
| 池　州　市 | 387 | 36267 | 6887.73138 | 12216.0047 |
| 安　庆　市 | 424.6 | 28808 | 5877.07508 | 13226.661 |
| 黄　山　市 | 346.1 | 37306 | 7053.7353 | 12050.0008 |

**Table2: Logistic Regression Data with City Attributes**

| | patientId | score | city | distance | gdp | medical_level | medical_dif | patientId | prov | doctorId | name | doctorgrade |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 0 | 1013185098 | 8 | 阜阳 | 563.5 | 15303.0 | 3242.62066 | 15861.115440000000 | 1013185098 | 安徽 | 21587 | 杜永强 | 主任医师 教授 |
| 1 | 1013185098 | 8 | 阜阳 | 563.5 | 15303.0 | 3242.62066 | 15861.115440000000 | 1013185098 | 安徽 | 21587 | 杜永强 | 主任医师 教授 |
| 2 | 1013185098 | 8 | 阜阳 | 563.5 | 15303.0 | 3242.62066 | 15861.115440000000 | 1013185098 | 安徽 | 21587 | 杜永强 | 主任医师 教授 |
| 3 | 1013185098 | 8 | 阜阳 | 563.5 | 15303.0 | 3242.62066 | 15861.115440000000 | 1013185098 | 安徽 | 21587 | 杜永强 | 主任医师 教授 |
| 4 | 1013185098 | 8 | 阜阳 | 563.5 | 15303.0 | 3242.62066 | 15861.115440000000 | 1013185098 | 安徽 | 21587 | 杜永强 | 主任医师 教授 |
| 5 | 1013185098 | 8 | 阜阳 | 563.5 | 15303.0 | 3242.62066 | 15861.115440000000 | 1013185098 | 安徽 | 21587 | 杜永强 | 主任医师 教授 |
| 6 | 2255034500 | 5 | 亳州 | 607.8 | 17769.0 | 3576.3282 | 15527.4079 | 2255034500 | 安徽 | 16464 | 江山 | 副主任医师 |
| 7 | 2255034500 | 5 | 亳州 | 607.8 | 17769.0 | 3576.3282 | 15527.4079 | 2255034500 | 安徽 | 16464 | 江山 | 副主任医师 |
| 8 | 2255034500 | 5 | 亳州 | 607.8 | 17769.0 | 3576.3282 | 15527.4079 | 2255034500 | 安徽 | 16464 | 江山 | 副主任医师 |
| 9 | 2255034500 | 5 | 亳州 | 607.8 | 17769.0 | 3576.3282 | 15527.4079 | 2255034500 | 安徽 | 16464 | 江山 | 副主任医师 |
| 10 | 2255034500 | 5 | 亳州 | 607.8 | 17769.0 | 3576.3282 | 15527.4079 | 2255034500 | 安徽 | 16464 | 江山 | 副主任医师 |
| 11 | 2255034500 | 5 | 亳州 | 607.8 | 17769.0 | 3576.3282 | 15527.4079 | 2255034500 | 安徽 | 16464 | 江山 | 副主任医师 |
| 12 | 2255034500 | 5 | 亳州 | 607.8 | 17769.0 | 3576.3282 | 15527.4079 | 2255034500 | 安徽 | 16938 | 项平 | 主治医师 |
| 13 | 2255034500 | 5 | 亳州 | 607.8 | 17769.0 | 3576.3282 | 15527.4079 | 2255034500 | 安徽 | 16938 | 项平 | 主治医师 |
| 14 | 2255034500 | 5 | 亳州 | 607.8 | 17769.0 | 3576.3282 | 15527.4079 | 2255034500 | 安徽 | 16938 | 项平 | 主治医师 |
| 15 | 2255034500 | 5 | 亳州 | 607.8 | 17769.0 | 3576.3282 | 15527.4079 | 2255034500 | 安徽 | 16938 | 项平 | 主治医师 |
| 16 | 2255034500 | 5 | 亳州 | 607.8 | 17769.0 | 3576.3282 | 15527.4079 | 2255034500 | 安徽 | 16938 | 项平 | 主治医师 |
| 17 | 2255034500 | 5 | 亳州 | 607.8 | 17769.0 | 3576.3282 | 15527.4079 | 2255034500 | 安徽 | 16938 | 项平 | 主治医师 |
| 18 | 2255034500 | 5 | 亳州 | 607.8 | 17769.0 | 3576.3282 | 15527.4079 | 2255034500 | 安徽 | 24671 | 王静 | 主治医师 |

**Table3: Logistic Regression Data with Patient Attributes**

```
              precision    recall  f1-score   support

           0       0.58      0.69      0.63    410909
           1       0.64      0.52      0.58    433791

   micro avg       0.61      0.61      0.61    844700
   macro avg       0.61      0.61      0.60    844700
weighted avg       0.61      0.61      0.60    844700

[[-1.29484253e-03  1.76739044e-05 -1.20214732e-04   7.61322211e-02]]
('accuracy_score ', 0.6051177932993962)
('precision_score ', 0.6422969903463941)
('recall_score ', 0.5214884587278206)
[Finished in 33.0s]
```
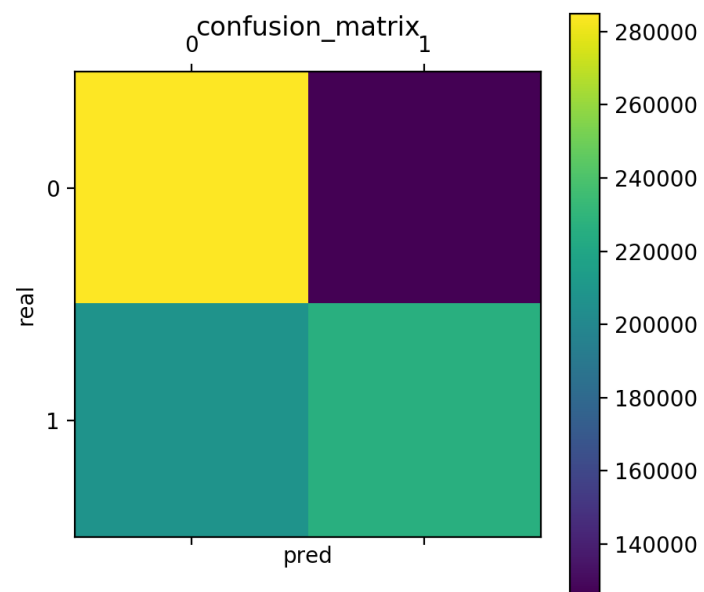
**Figure 16: Logistic Regression Results**

**Figure 17: Confusion Matrix**