

# Occupational Diseases Prediction Online Analysis Platform

(<http://predict.xjyg.net:666/>)

## 1. Introduction

How to model and forecast with limited data is a challenging task, especially in occupational health fields. In this online platform, we combined the grey systems theory and machine learning method to solve this issue. The GM models contain five models and they are even grey model (EGM), original difference grey model (ODGM), even difference grey model (EDGM), discrete grey model (DGM) and Verhulst which was the fitted value from occupational diseases as training data to train the machine learning models. They include five state-of-art models: K-Nearest Neighbor (KNN), Support Vector Regression (SVM), Random Forest (RF), Gradient Boosting Machine (GBM), and Artificial Neural Network (ANN). To the best of our knowledge, this is the first time that those five hybrid algorithm combining models are used to predict occupational diseases. The effectiveness and applicability of the models were assessed based on its ability to predict the incidence of occupational diseases in China.

## 2. Platform Function

### 2.1 Data Set

The data set of platform is the cases of occupational diseases from 2005 to 2017 which obtained from national health commission of the people's republic of China.

Data Set module will show the dataset to be used to analysis on the platform. Users can look through the fields of the data set and search the values from data set, that will let users fully understand the structure of data set here (Figure 1).

## Occupational Diseases Prediction Online Analysis Platform

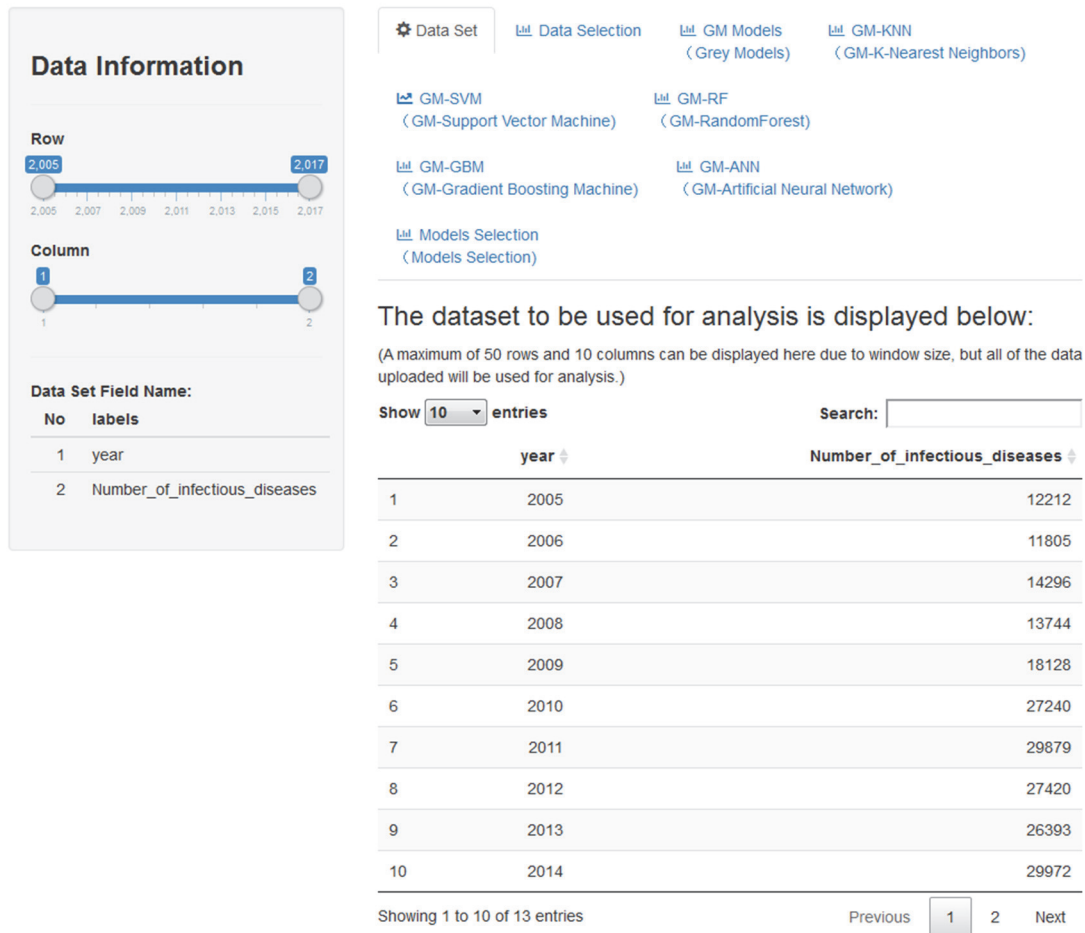
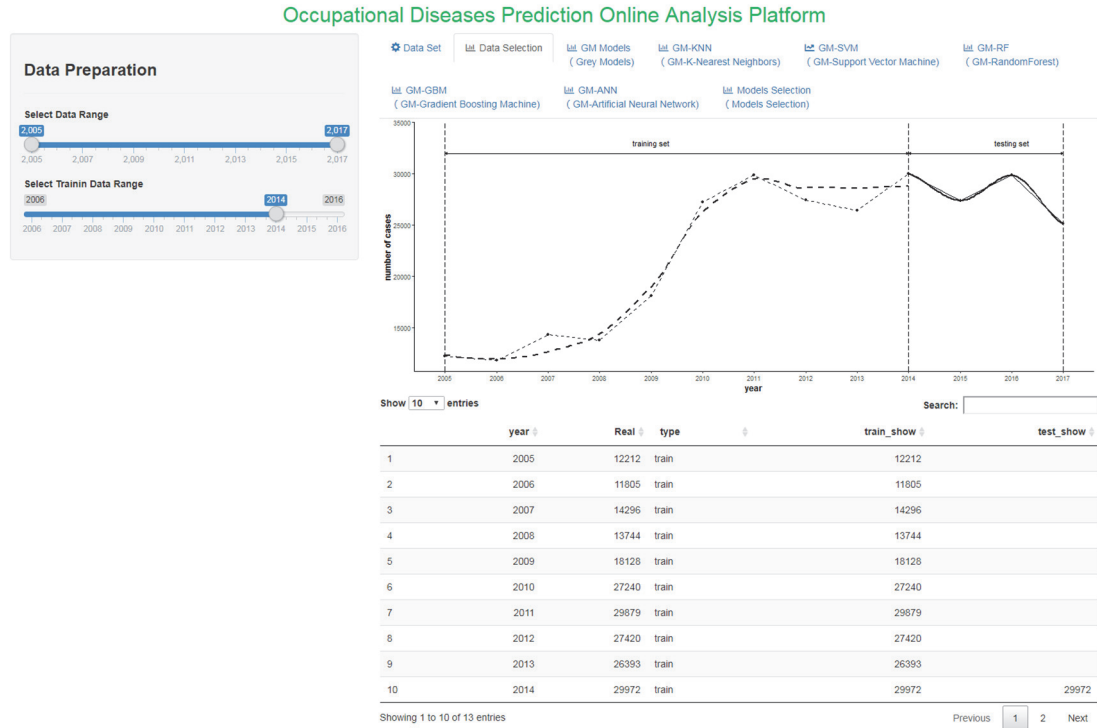


Figure 1. Data set Module

### 2.2 Data Selection

Data Selection module can select the data range for analysis, in order to give users a good experience, the default training data is from 2005 to 2014 that according to the paper of a comparative study on the prediction of occupational diseases in China with hybrid algorithm combining models. Users can see both plot of the dataset and the specific data set structure (Figure 2).



**Figure 2.** Data Selection Module

## 2.3 GM Models

In order to obtain the training set for the KNN, SVM, RF, GBM and ANN models, the five GM models, i.e., even grey model (EGM), original difference grey model (ODGM), even difference grey model (EDGM), discrete grey model (DGM) and Verhulst were used to fit the input of the five hybrid algorithm combing models with the training set of China occupational diseases data from 2005 to 2014.

In order to verify the performance of selection by GM model's MAPE and RMSE, we can select the training data by minimum value of GM model's MAPE and RMSE, after verified by permutations and combination, we get the best model was built by all GM fitted values(Figure 3).

Users can see plot of GM models, prediction accuracy and fitted value here. It will update the data automatically according to user's choice.

## Occupational Diseases Prediction Online Analysis Platform

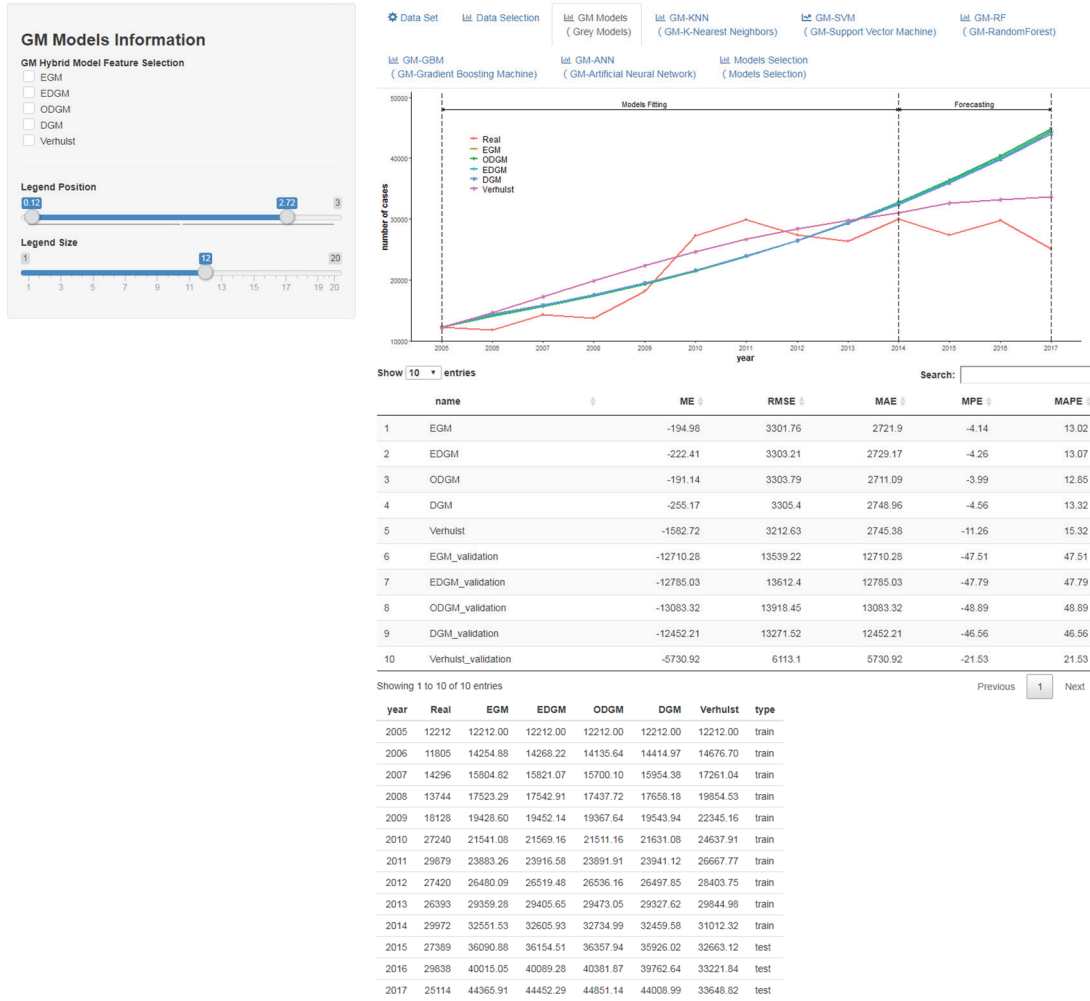


Figure 3. GM Models Module

### 2.4 GM-KNN (GM-K-Nearest Neighbor)

Platform used both KNN conventional method and weighted method to build the model respectively. In the conventional KNN method, we choose the most suitable parameter  $k=2$  for cross-validation. In KNN weighting method, we choose inversion weighting and  $k=2$  for cross-validation and grid scan.

Users can see plot of GM-KNN models, prediction accuracy and fitted value here. It will update the data automatically according to user's choice (Figure 4).

## Occupational Diseases Prediction Online Analysis Platform

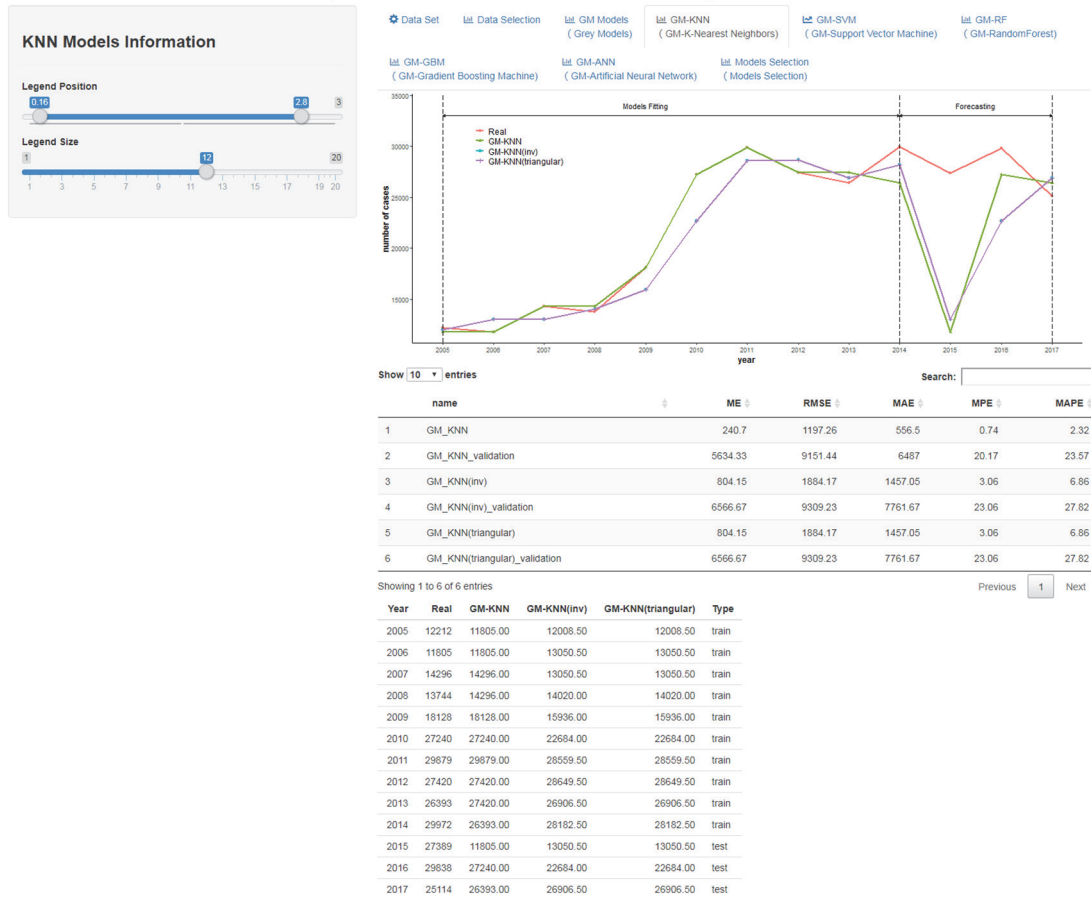


Figure 4. GM-KNN Module

## 2.5 GM-SVM (Support Vector Regression)

Platform built four SVM models with linear, polynomial, radial and sigmoid kernel respectively and with the cross-validation method.

Users can see plot of GM- SVM models, prediction accuracy and fitted value here. It will update the data automatically according to user's choice (Figure 5).

## Occupational Diseases Prediction Online Analysis Platform

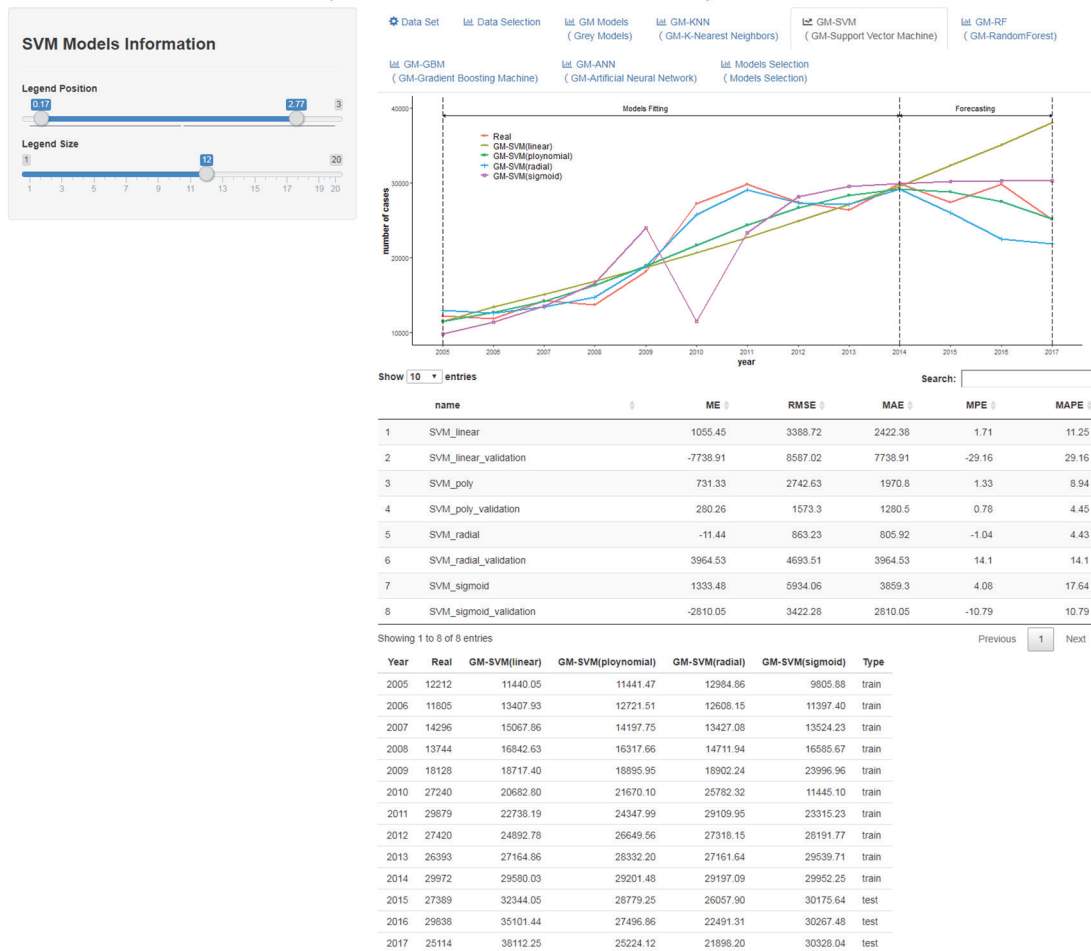


Figure 5. GM-SVM Module

## 2.6 GM-RF (Random Forest)

GM-RF model is built with the optimum parameters of  $mtry=1$  and  $ntree=30$  after selecting from 500 trees.

Users can see plot of GM-RF models, prediction accuracy and fitted value here. It will update the data automatically according to user's choice (Figure 6).



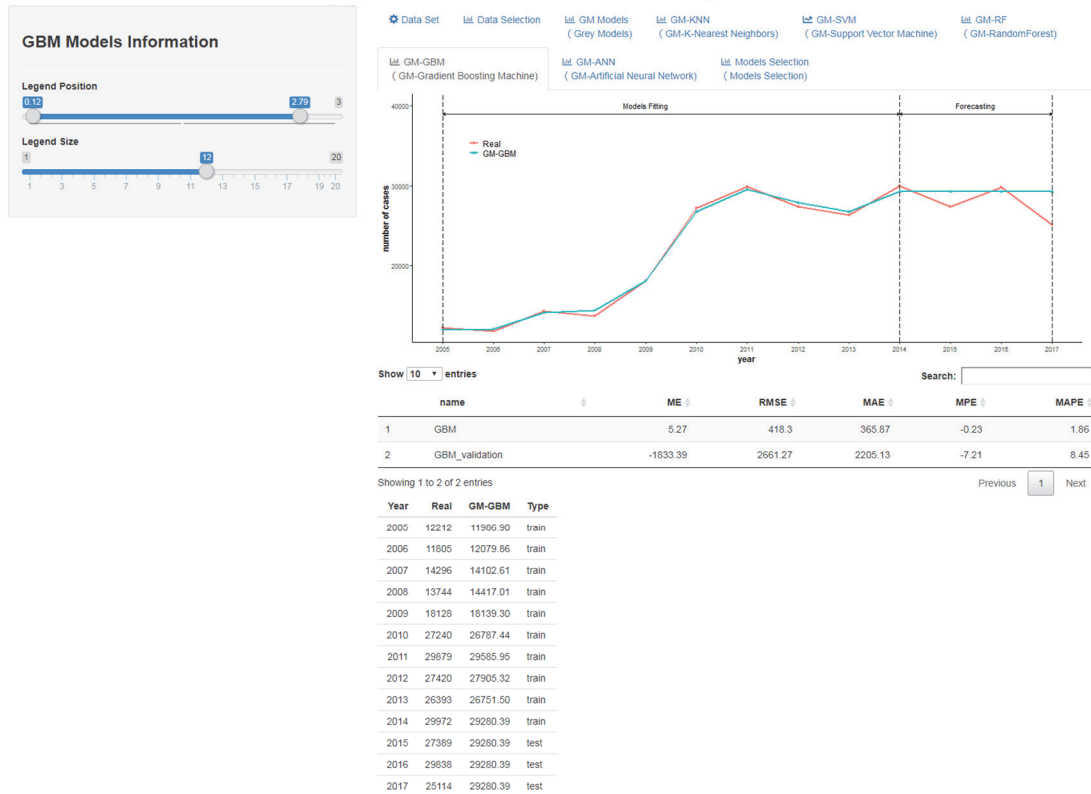
**Figure 6. GM-RF Module**

## 2.7 GM-GBM (Gradient Boosting Machine)

GM-GBM model is built with  $\alpha=0.1$  and  $\gamma=0.5$  by resampling method.

Users can see plot of GM- RF models, prediction accuracy and fitted value here. It will update the data automatically according to user's choice (Figure 7).

## Occupational Diseases Prediction Online Analysis Platform



**Figure 7. GM-GBM Module**

## 2.8 GM-ANN (Artificial Neural Network)

GM-ANN model is built from error accuracy of 0.00000001, learning times of 10,000 and neurons number of 5.

Users can see plot of GM-ANN model, prediction accuracy and fitted value here. It will update the data automatically according to user's choice (Figure 8).



## Occupational Diseases Prediction Online Analysis Platform

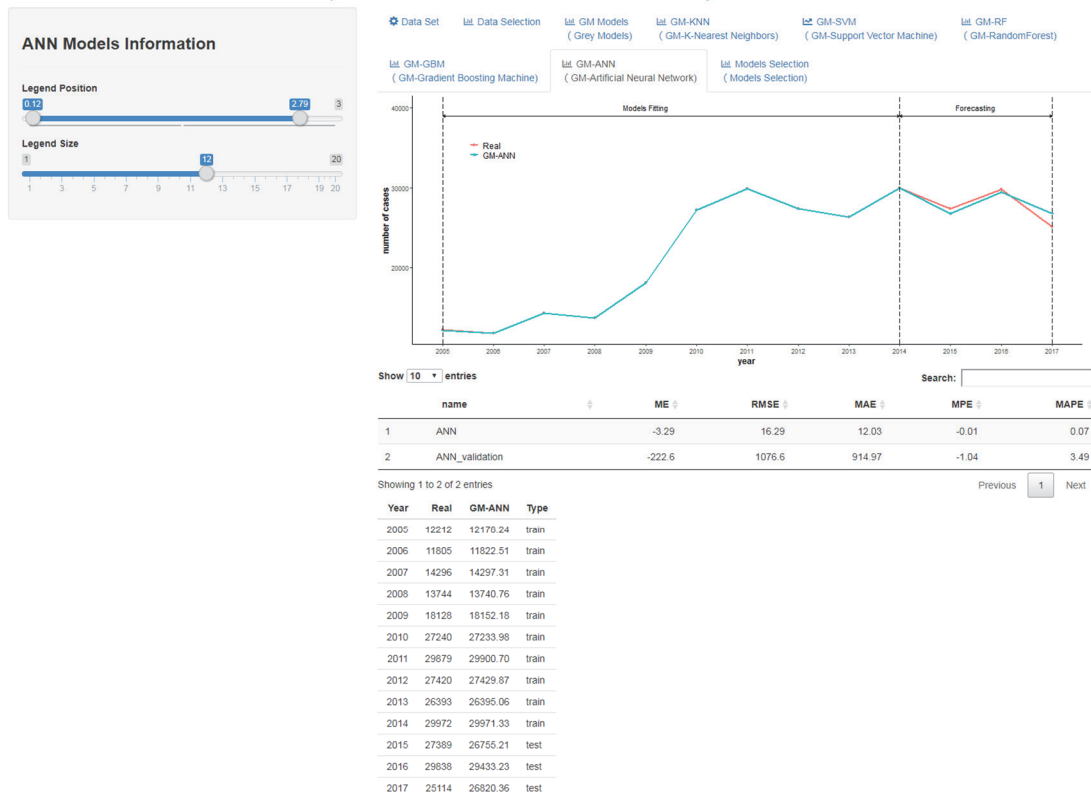


Figure 8. GM-KNN Module

## 2.9 Models Selection

Users can see plot of GM-KNN, GM-SVM, GM-RF, GM-GBM and GM-ANN models, prediction accuracy and fitted value here. It will update the data automatically according to user's choice (Figure 9).

## Occupational Diseases Prediction Online Analysis Platform

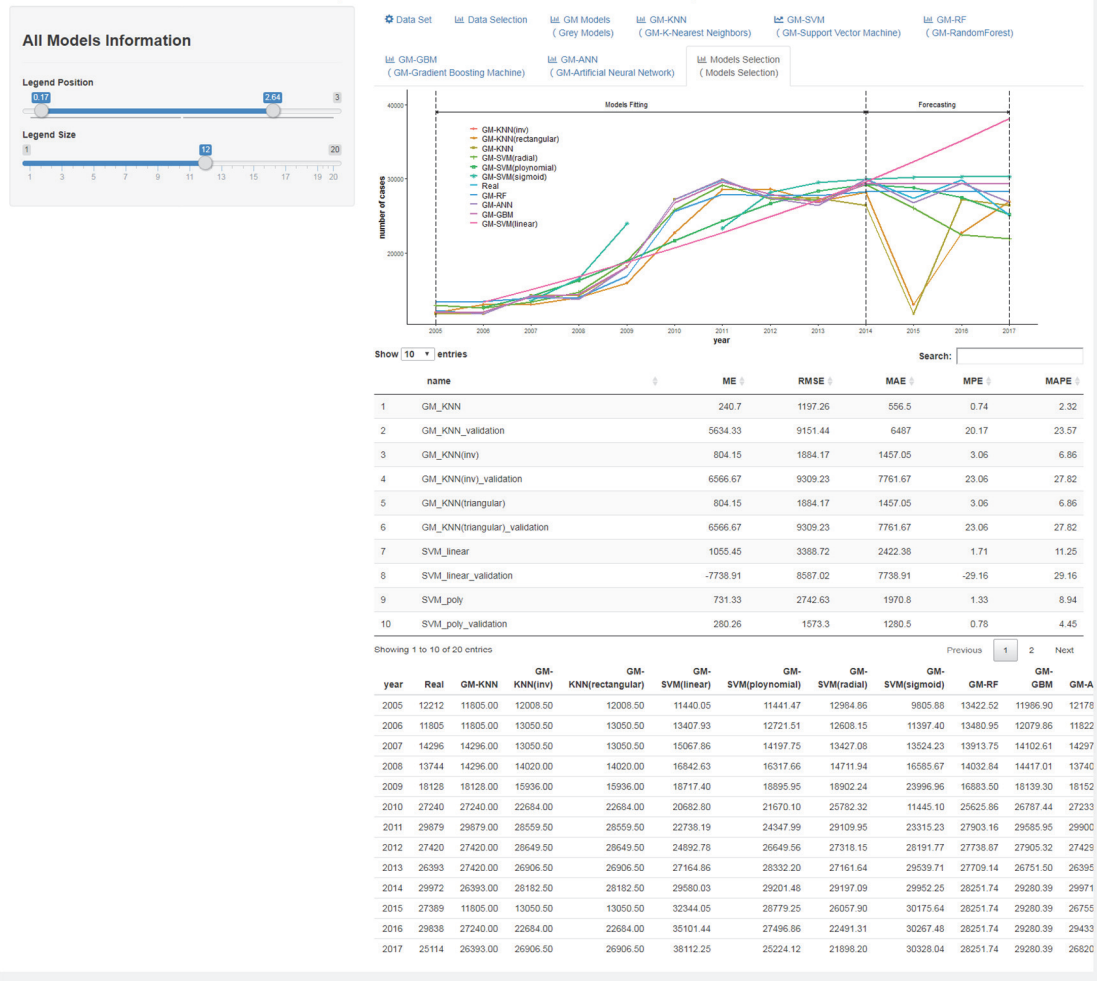


Figure 9. Models Selection Module