

Project: 1

Prediction of mortality of ICU patients based on random forest algorithm

Background

Because the basic causes of critical illness are different, critical patients are scattered to different medical specialties, which makes the lack of unity for critical illness. ICU (Intensive Care Unit) is a clinical base of a professional team specialized in critical care medicine and specialized in critical care treatment. It is a centralized management unit for critical patients and high-risk patients after surgery from various clinical departments. The birth of ICU has directly promoted the development of critical care medicine. It pays more attention to the characteristics of patients in critical care and the common threats and damages they face, so that many patients who previously thought that they could not be cured can survive or extend their survival time.

Develop the population mortality prediction in the intensive care unit (ICU) and the in-hospital mortality prediction method for specific patients. The information collected at the time of admission to the ICU will be used to predict which patients can survive during the hospitalization and which patients cannot.

Technology

Programming language: python

Import: pandas, imbearn, sklearn, pyecharts, matplotlib, math, numpy

Analysis: sequence classification

Introduction of Algorithm

In machine learning, random forest is a classifier containing multiple decision trees, and its output category is determined by the mode of the category output by individual trees.

Leo Breiman and Adele Cutler developed algorithms to deduce random forests. "Random Forests" is their trademark. This term is derived from the random decision forests proposed by Tin Kam Ho of Bell Laboratories in 1995. This method combines Breimans' "Bootstrap aggregating" idea and Ho's "random subspace method" to build a set of decision trees.

The advantages of random forest include:

- 1) For many kinds of data, it can produce high accuracy classifiers.
- 2) It can handle a large number of input variables;
- 3) It can evaluate the importance of variables when determining categories.
- 4) When building a forest, it can produce an unbiased estimate of the generalized error internally.
- 5) It contains a good method to estimate the lost data, and if a large part of the data is lost, the accuracy can still be maintained.
- 6) It provides an experimental method to detect variable interactions.

- 7) For unbalanced classification data sets, it can balance errors.
- 8) It calculates the proximity in each example, which is very useful for data mining, outlier detection and data visualization.
- 9) Use the above. It can be extended to unmarked data, which usually uses unsupervised clustering. It can also detect deviators and view data.
- 10) The learning process is very fast
- 11) Significant effect on high latitude data

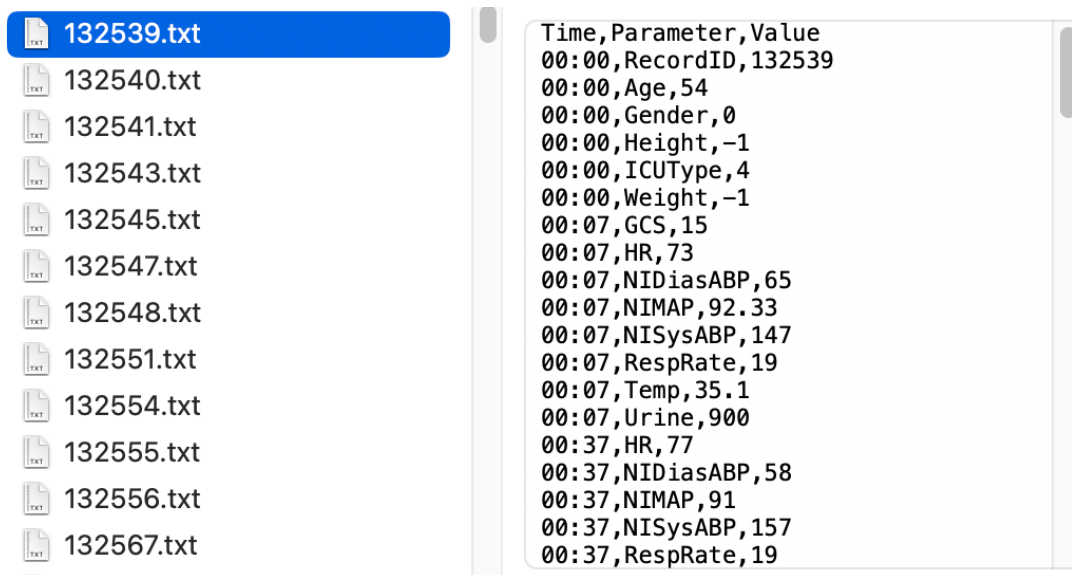
Data source

<https://www.physionet.org/content/challenge-2012/1.0.0/>

Project process

1: Data acquisition:

The original data is a .txt file. Each file is used for the various data collected by a patient since entering the hospital.



The screenshot shows a file explorer interface. On the left, a list of files is displayed, with '132539.txt' selected and highlighted in blue. The other files in the list are 132540.txt, 132541.txt, 132543.txt, 132545.txt, 132547.txt, 132548.txt, 132551.txt, 132554.txt, 132555.txt, 132556.txt, and 132567.txt. On the right, the content of the selected file '132539.txt' is displayed. The content is a text file with a header 'Time,Parameter,Value' followed by several lines of data. The data is organized into groups by time, with each group starting with a timestamp followed by a comma and then a parameter and its value. The parameters include RecordID, Age, Gender, Height, ICUType, Weight, GCS, HR, NIDiasABP, NIMAP, NISysABP, RespRate, Temp, and Urine.

```
Time,Parameter,Value
00:00,RecordID,132539
00:00,Age,54
00:00,Gender,0
00:00,Height,-1
00:00,ICUType,4
00:00,Weight,-1
00:07,GCS,15
00:07,HR,73
00:07,NIDiasABP,65
00:07,NIMAP,92.33
00:07,NISysABP,147
00:07,RespRate,19
00:07,Temp,35.1
00:07,Urine,900
00:37,HR,77
00:37,NIDiasABP,58
00:37,NIMAP,91
00:37,NISysABP,157
00:37,RespRate,19
```

- *RecordID* (defined as above)
- *SAPS-I score* (Le Gall et al., 1984)
- *SOFA score* (Ferreira et al., 2001)
- *Length of stay* (days)
- *Survival* (days)
- *In-hospital death* (0: survivor, or 1: died in-hospital)

2: Preliminary arrangement of data:

Read the .txt file through the open function loop and format the data after slicing, segmentation and statistics to form a set_a_data.csv、set_b_data.csv、set_c_Data.csv

Three tables with the same column name, 12000 data in total.

This data is sequence data, and features are extracted based on six basic statistical methods: mean value, variance, extreme value (maximum, minimum), initial value and current value,

Feature extraction: For RecordID, Age, Gender, Height, ICUType and Weight, these six fields are read directly,

The remaining 37 fields:

- *Albumin* (g/dL)
- *ALP* [Alkaline phosphatase (IU/L)]
- *ALT* [Alanine transaminase (IU/L)]
- *AST* [Aspartate transaminase (IU/L)]
- *Bilirubin* (mg/dL)
- *BUN* [Blood urea nitrogen (mg/dL)]
- *Cholesterol* (mg/dL)
- *Creatinine* [Serum creatinine (mg/dL)]
- *DiasABP* [Invasive diastolic arterial blood pressure (mmHg)]
- *FiO2* [Fractional inspired O₂ (0-1)]
- *GCS* [Glasgow Coma Score (3-15)]
- *Glucose* [Serum glucose (mg/dL)]
- *HCO3* [Serum bicarbonate (mmol/L)]
- *HCT* [Hematocrit (%)]
- *HR* [Heart rate (bpm)]
- *K* [Serum potassium (mEq/L)]
- *Lactate* (mmol/L)
- *Mg* [Serum magnesium (mmol/L)]
- *MAP* [Invasive mean arterial blood pressure (mmHg)]
- *MechVent* [Mechanical ventilation respiration (0:false, or 1:true)]
- *Na* [Serum sodium (mEq/L)]
- *NIDiasABP* [Non-invasive diastolic arterial blood pressure (mmHg)]
- *NIMAP* [Non-invasive mean arterial blood pressure (mmHg)]
- *NISysABP* [Non-invasive systolic arterial blood pressure (mmHg)]
- *PaCO2* [partial pressure of arterial CO₂ (mmHg)]
- *PaO2* [Partial pressure of arterial O₂ (mmHg)]
- *pH* [Arterial pH (0-14)]
- *Platelets* (cells/nL)
- *RespRate* [Respiration rate (bpm)]
- *SaO2* [O₂ saturation in hemoglobin (%)]
- *SysABP* [Invasive systolic arterial blood pressure (mmHg)]
- *Temp* [Temperature (°C)]
- *TropI* [Troponin-I (μg/L)]
- *TropT* [Troponin-T (μg/L)]
- *Urine* [Urine output (mL)]
- *WBC* [White blood cell count (cells/nL)]
- *Weight* (kg)⁺

In addition to the weight, each field is subdivided into 6 features, and the total number of features is 222.

(The data processing and feature extraction code is in the open_file.py file).

MS_MLData.csv																									使用Microsoft Excel打开	
RecordID	Age	Gender	Height	ICUType	Weight	GCS_jun	GCS_fang	GCS_da	GCS_xiao	GCS_zhu	GCS_dang	HR_jun	HR_fang	HR_da	HR_xiao	HR_zhu	HR_dang	NIDiasABP_jun	NIDiasABP_fang	NIDiasABP_da	NIDiasABP_xiao	NIDiasABP_zhu	NIDiasABP_dang	NIMAP_jun	NE	
132539.0	54.0	0.0	-1.0	-1.0	4.0	-1.0	14.923079823079823	0.971059171397633	15.0	14.0	15.0	70.810108108108	72.0452883177501	86.0	58.0	73.0	86.0	50.14705882352941	61.77249134948096	67.0	39.0	65.0	55.0	71.55911764705883	83	
132540.0	0	54.0	1.0	175.5	2.0	78.0	13.333333333333334	0.9555555555555556	11.0	3.0	15.0	80.76415428571429	81.34888888888889	56.742857142857143	67.0	58.0	73.0	51.34888888888889	61.77249134948096	67.0	39.0	65.0	55.0	71.55911764705883	11	
132541.0	0	64.0	0.0	-1.0	3.0	56.7	5.937676307676308	1.3017751476389943	8.0	5.0	7.0	83.7594525925926	136.827206869328	113.0	57.0	100.0	71.0	79.0	65.05172413793103	44.38387637693214	81.0	66.0	84.0	84.0	96.7513137894737	53
132542.0	0	68.0	1.0	180.3	3.0	84.6	14.944444444444445	0.9246913380246914	15.0	14.0	15.0	70.98333333333333	58.38035555555555	88.0	57.0	79.0	79.0	65.05172413793103	44.38387637693214	81.0	45.0	63.0	81.0	83.88531162179332	38	
132543.0	0	68.0	0.0	-1.0	3.0	-1.0	15.0	0.0	15.0	15.0	15.0	74.95833333333333	54.41495555555555	94.0	65.0	93.0	68.0	45.7209932325814	152.8058409951325	96.0	26.0	41.0	42.0	74.94651162179332	85	
132547.0	0	64.0	1.0	180.3	1.0	134.0	8.666666666666666	0.8888888888888887	11.0	7.0	7.0	8.0	88.5131489161703	63.100393932364	101.1	71.0	78.0	92.0	70.5	342.25	89.0	52.0	89.0	52.0	81.893	41
132548.0	0	68.0	0.0	162.6	3.0	87.0	15.0	0.0	15.0	15.0	15.0	68.26868686868687	41.27661256196657	80.0	50.0	73.0	60.0	72.0	298.85714285714286	98.0	31.0	68.0	77.0	102.3474285714286	54	
132551.0	1	78.0	0.0	162.6	3.0	46.4	11.84613964537847	8.899482643643658	10.0	8.0	15.0	9.0	70.94526547946206	132.0249462297658	111.1	55.0	111.0	58.0	30.6797441860465	135.327238939677	56.0	14.0	51.0	19.0	55.17796897644715	11
132554.0	0	64.0	0.0	-1.0	3.0	60.7	15.0	0.0	15.0	15.0	15.0	127.23913104347826	22.52977315889961	137.0	115.0	127.0	122.0	64.47826868656522	83.3364393154707	92.0	47.0	71.0	53.0	84.47739130434783	77	
132555.0	0	74.0	1.0	175.5	2.0	66.1	14.983333333333334	3.4097222222222228	15.0	10.0	10.0	15.0	85.189051174138	64.0847267307491	99.0	67.0	67.0	78.0	53.0	16.0	57.0	49.0	57.0	49.0	75.67	46
132556.0	0	64.0	0.0	-1.0	3.0	65.0	15.0	0.0	15.0	15.0	15.0	110.5625	465.371787575	161.1	57.0	101.0	91.0	48.166666666666664	11.348888888888889	60.0	52.0	46.0	45.0	57.38615555555556	22	
132557.0	0	71.0	0.0	137.5	2.0	50.0	14.181818181818182	2.50781329660422	15.0	10.0	15.0	15.0	95.22272727272727	53.02410464519359	112.0	82.0	84.0	95.0	43.43333333333333	38.945555555555555	57.0	32.0	45.0	46.0	64.68733333333335	46
132560.0	0	66.0	0.0	157.5	3.0	84.5	15.0	0.0	15.0	15.0	15.0	88.5	34.06818181818182	100.7	77.0	88.0	93.0	56.8974388974389	109.8356345191906	79.0	36.0	62.0	46.0	71.21621621621621	81	
132570.0	0	84.0	1.0	170.2	1.0	102.6	14.6	0.24	15.0	14.0	15.0	70.60975609756098	143.567218913711	106.0	55.0	70.0	73.0	48.675	125.06637300000001	66.0	0.0	61.0	47.0	72.38461538461539	7	
132573.0	0	77.0	1.0	162.6	1.0	50.1	15.0	0.0	15.0	15.0	15.0	73.66071428571429	19.04509849759322	84.0	64.0	80.0	68.0	48.53357142857143	60.8185867346038	82.0	37.0	49.0	44.0	82.83614285714287	81	
132575.0	0	78.0	1.0	167.6	2.0	63.0	10.87743857142857	27.9112040616328	15.0	3.0	3.0	15.0	84.4325373132383	146.2494642630918	119.0	65.0	73.0	106.0								
132577.0	0	65.0	1.0	-1.0	3.0	63.3	12.125	6.109375	15.0	9.0	9.0	15.0	84.67916666666667	25.66623688888888	96.0	71.0	94.0	88.0	74.90909090909091	201.3373908264463	98.0	48.0	62.0	65.0	95.65181818181819	24
132582.0	0	84.0	1.0	182.9	3.0	82.5	14.84613964157847	0.13017751476389943	15.0	14.0	15.0	84.07142857142857	45.53061224468795	108.0	77.0	101.0	83.0	43.870370370370374	48.58432773203191	60.0	24.0	57.0	51.0	66.0740274074028	46	
132584.0	0	78.0	0.0	-1.0	3.0	72.8	9.506666666666666	0.4355555555555556	11.0	3.0	11.0	84.125	299.339375	124.0	53.0	124.0	73.0	64.29411764705883	158.2076124587475	129.0	44.0	75.0	52.0	84.74233294117646	57	
132585.0	0	40.0	0.0	165.1	2.0	84.7	13.78714285714286	10.336938775103207	15.0	3.0	3.0	15.0	80.5614350587719	87.09011594133671	101.1	66.0	79.0	82.0	41.57142857142857	38.81632835931225	50.0	36.0	36.0	50.0	57.265714285714285	23
132586.0	0	68.0	0.0	154.9	3.0	42.3	15.0	0.0	15.0	15.0	15.0	91.85	166.875	115.0	76.0	115.0	78.0	50.981333333333336	77.85416666666666	54.0	31.0	47.0	46.0	64.54861111111112	62	
132590.0	0	58.0	1.0	188.0	2.0	96.0	12.0	14.181818181818182	15.0	3.0	3.0	15.0	97.54089360557378	90.51061542540675	131.1	84.0	119.0	88.0	55.0	55.0	55.0	55.0	55.0	55.0	75.67	71
132591.0	0	81.0	1.0	-1.0	3.0	63.7	15.0	0.0	15.0	15.0	15.0	70.13725400196079	10.78812027272991	109.0	59.0	62.0	61.0	35.740888321584	54.033642606746	70.0	28.0	62.0	54.0	74.143264117647	36	
132592.0	0	35.0	0.0	-1.0	3.0	71.8	15.0	0.0	15.0	15.0	15.0	97.07142857142857	80.76861228571429	115.1	77.0	112.0	82.0	60.69756087560975	105.11600237935937	85.0	43.0	43.0	55.0	80.9256973609756	16	
132595.0	0	24.0	0.0	-1.0	3.0	1.0																				
132597.0	0	66.0	0.0	137.2	3.0	82.0	15.0	0.0	15.0	15.0	15.0	68.7942304301886	179.7871128154859	143.0	56.0	76.0	65.0	39.3613584613585	48.6710059711598	60.0	24.0	37.0	25.0	60.4930762307692	57	
132598.0	1	80.0	0.0	-1.0	4.0	60.0	7.719913043678261	0.3178813042646504	8.0	6.0	6.0	8.0	76.7086075949367	116.95587088428984	158.0	49.0	51.0	72.0	66.7814386404511	247.6566469442924	115.0	43.0	69.0	61.0	80.0376744804045	18
132599.0	0	53.0	0.0	177.8	4.0	73.5	8.416666666666666	2.9097222222222222	14.0	7.0	14.0	85.39285714285714	97.667801878247	104.0	67.0	98.0	94.0	63.0	18.5	69.0	57.0	69.0	64.0	79.5	18	
132601.0	0	74.0	1.0	177.8	2.0	79.9	12.58461359461538	16.55613011775148	15.0	3.0	3.0	15.0	99.24193483970796	37.89794633962315	109.1	87.0	103.0	95.0	86.0	66.0	66.0	66.0	66.0	66.0	93.0	83
132602.0	1	80.0	1.0	180.3	3.0	70.0	15.0	0.0	15.0	15.0	15.0	74.6328018867624	286.3967154651263	144.0	59.0	67.0	85.0	66.38	60.915600000000001	82.0	47.0	72.0	61.0	80.79166666666667	12	
132603.0	1	80.0	0.0	-1.0	3.0	75.0	7.416666666666667	0.2430555555555555	8.0	7.0	7.0	72.12880776743486	65.9677467769285	98.0	58.0	72.0	98.0	36.57142857142857	174.3520481813054	75.0	22.0	66.0	72.0	64.095	24	
132610.0	0	72.0	1.0	172.9	3.0	102.6	14.9	0.09	15.0	14.0	15.0	74.11627906976744	12.73931038991765	104.0	70.0	75.0	78.0	39.0	12.666666666666666	42.0	34.0	34.0	42.0	34.0	57.33333333333336	1
132612.0	0	50.0	1.0	-1.0	4.0	109.0	9.222222222222222	4.99061738930617	11.0	3.0	11.0	8.0	74.61111111111111	257.42289595617286	126.1	57.0	114.0	73.0	74.0	64.0	43.0	43.0	52.0	73.35666666666667	11	
132614.0	0	77.0	1.0	162.6	1.0	59.0	15.0	0.0	15.0	15.0	15.0	70.78479484181952	1.146209861495465	74.0	70.0	73.0	78.0	39.42857142857143	29.123464181635702	49.0	32.0	41.0	36.0	57.42374285714285	28	
132615.0	0	46.0	0.0	152.4	3.0	88.4	8.111111111111111	10.76543296476543	11.0	3.0	11.0	70.95797104462975	101.387947030475	107.0	55.0	80.0	55.0	43.13333333333333	69.44888888888887	61.0	30.0	61.0	36.0	61.0	63.52000000000001	86
132617.0	1	77.0	1.0	170.2	1.0	75.0	15.0	0.0	15.0	15.0	15.0	82.74285714285715	83.9035362244889	100.0	66.0	66.0	79.0	40.0	46.5411254610256	61.0	19.0	56.0	37.0	58.59846135846134	5	
132618.0	0	72.0	0.0	152.4	4.0	58.38	11.71																			

Missing value processing: remove the feature data with missing ratio greater than 20%, and finally get 109 features, as shown below:

```
data1 = pd.read_csv('set_a_data.csv')
data2 = pd.read_csv('set_b_data.csv')
data3 = pd.read_csv('set_b_data.csv')
data = pd.concat([data1, data2, data3])
data_list = data.dropna(thresh=9600, axis=1)
print(data_list.isnull().sum())
print(data_list.head())
print(data_list)
```

```
RecordID    0
lab         0
Age         0
Gender      0
Height      0
...
WBC_fang    205
WBC_da      205
WBC_xiao    205
WBC_chu     205
WBC_dang    205
Length: 109, dtype: int64
```

```
Length: 109, dtype: int64
RecordID  lab  Age  Gender  ...  WBC_da  WBC_xiao  WBC_chu  WBC_dang
0  132539.0   0  54.0    0.0  ...   11.2    9.4    11.2    9.4
1  132540.0   0  76.0    1.0  ...   13.3    7.4    7.4   13.3
2  132541.0   0  44.0    0.0  ...    6.2    3.7    4.2    6.2
3  132543.0   0  68.0    1.0  ...   11.5    7.9   11.5    7.9
4  132545.0   0  88.0    0.0  ...    4.8    3.8    3.8    4.8
[5 rows x 109 columns]
```

```
RecordID  lab  Age  Gender  ...  WBC_da  WBC_xiao  WBC_chu  WBC_dang
0  132539.0   0  54.0    0.0  ...   11.2    9.4    11.2    9.4
1  132540.0   0  76.0    1.0  ...   13.3    7.4    7.4   13.3
2  132541.0   0  44.0    0.0  ...    6.2    3.7    4.2    6.2
3  132543.0   0  68.0    1.0  ...   11.5    7.9   11.5    7.9
4  132545.0   0  88.0    0.0  ...    4.8    3.8    3.8    4.8
...      ...  ...  ...    ...  ...    ...    ...    ...    ...
3995  152849.0   0  78.0    1.0  ...   20.0   14.5   20.0   15.8
3996  152851.0   0  90.0    1.0  ...   41.8   18.0   27.4   21.4
3997  152858.0   0  70.0    0.0  ...   15.1   13.1   14.8   15.1
3998  152862.0   0  49.0    0.0  ...   16.6   13.6   13.6   16.6
3999  152864.0   0  82.0    0.0  ...   17.3   10.3   17.3   10.3
[12000 rows x 109 columns]
```

Missing value padding:

```
for i in data_list.isnull().any().items():
    name = str(i[0])
    val = str(i[1])
    if 'True'==val:
        data_list[name] = data_list[name].fillna(data_list[name].median())
print(data_list.isnull().any())
```

Fill in the missing values of each column in the way of median, and the results show.

```
RecordID      False
lab           False
Age           False
Gender        False
Height        False
...
WBC_fang      False
WBC_da        False
WBC_xiao      False
WBC_chu       False
WBC_dang      False
Length: 109, dtype: bool
```

3. Analysis of correlation coefficient for each feature:

```
def xiangguanxishu(X, Y):
    """计算相关系数"""
    XY = X * Y
    X2 = X ** 2
    Y2 = Y ** 2
    n = len(XY)
    numerator = n * XY.sum() - X.sum() * Y.sum() # 分子
    denominator = math.sqrt(n * X2.sum() - X.sum() ** 2) * math.sqrt(n * Y2.sum() - Y.sum()
** 2) # 分母
    if denominator == 0:
        return 'NaN'
    rhoXY = numerator / denominator
    return rhoXY
```

Print results:

```
['Age', 0.13048980922783202]
['Gender', -0.015418852639407147]
['Height', -0.028738795253975824]
['ICUType', 0.0702577782808679]
.....
['WBC_xiao', 0.05710239848080716]
['WBC_chu', 0.03335311353058349]
['WBC_dang', 0.09195582784538718]
```

4. Data processing:

According to statistics, in the data sample, the tag value survivors (lab=0) account for about 85%, and the death (lab=1) account for about 15%. According to the binary classification algorithm, it belongs to the unbalanced data, and the final prediction result

will deviate greatly.

Solution: Reduce sampling method。

```
cc = ClusterCentroids()
X_resampled, y_resampled = cc.fit_resample(data_txt, data_lab.values.reshape([-1]))
data_train_X, data_train_Y = X_resampled, y_resampled
```

The final data is divided into training set (iris_x_train, iris_y_train), test set (iris_x_test, iris_y_test) and verification set (iris_x_yan, iris_y_yan) according to the two-eighth rule. The following is the division of data code block and display.

```
x_train_all, iris_x_test, y_train_all, iris_y_test = train_test_split(data_train_X, data_train_Y,
random_state=7, test_size=0.2)
iris_x_train, iris_x_yan, iris_y_train, iris_y_yan = train_test_split(x_train_all, y_train_all,
random_state=11, test_size=0.2)
print(iris_x_test)
print(iris_x_train)
print(iris_x_yan)
```

Test set:

```
      Age  Gender  Height  ...  WBC_xiao  WBC_chu  WBC_dang
1395  52.20000  0.60  34.260000  ...  31.220000  34.860000  31.860000
1062  42.00000  0.00  165.100000  ...   9.700000   9.700000  10.500000
440   70.75000  0.50  62.875000  ...   5.856250   8.187500   6.643750
1417  72.09375  0.50  47.540625  ...  11.053125  13.496875  12.053125
3333  64.00000  0.00  157.500000  ...   8.100000   8.100000  10.700000
...      ...      ...      ...      ...      ...      ...
962   44.00000  1.00  -1.000000  ...   3.533333   6.100000   3.666667
3102  88.00000  1.00  182.900000  ...   6.400000   6.400000  13.400000
2478  78.00000  1.00  182.900000  ...  14.000000  14.000000  33.500000
1601  75.50000  1.00  -1.000000  ...   5.400000   7.100000   5.550000
429   65.00000  0.25  107.700000  ...  10.825000  15.400000  13.675000

[676 rows x 107 columns]
```

Training set:

```
      Age  Gender  Height  ...  WBC_xiao  WBC_chu  WBC_dang
1460  76.116279  0.395349  59.862791  ...  10.102326  12.653488  11.241860
1857  89.000000  1.000000  -1.000000  ...   9.900000  10.000000   9.900000
3016  80.000000  1.000000  -1.000000  ...   8.500000  11.700000   9.000000
18    64.000000  1.000000  -1.000000  ...   8.800000  10.100000  12.700000
2206  87.000000  1.000000  -1.000000  ...   7.300000  11.300000   7.300000
...      ...      ...      ...      ...      ...      ...
1215  67.000000  0.000000  -1.000000  ...   2.600000   2.600000   5.900000
760   58.857143  0.857143  151.885714  ...   7.785714   9.700000   9.685714
3295  63.000000  0.000000  -1.000000  ...   8.800000   9.700000   8.800000
1595  62.444444  0.666667  171.588889  ...   7.944444   9.933333   9.622222
466   69.000000  1.000000  -1.000000  ...  11.300000  17.000000  11.300000

[2163 rows x 107 columns]
```

Validation set:

	Age	Gender	Height	...	WBC_xiao	WBC_chu	WBC_dang
1044	65.666667	0.555556	37.044444	...	13.333333	17.40	13.833333
3052	80.000000	1.000000	175.300000	...	3.900000	15.00	13.700000
1565	37.500000	1.000000	-1.000000	...	10.900000	14.45	10.900000
2989	66.000000	1.000000	170.200000	...	3.900000	3.90	9.800000
1169	70.000000	0.666667	141.216667	...	8.733333	9.30	10.766667
...
2739	76.000000	1.000000	182.900000	...	15.200000	15.20	16.100000
2063	55.000000	1.000000	-1.000000	...	28.200000	33.60	28.200000
2808	73.000000	1.000000	160.000000	...	31.300000	31.30	60.200000
2380	76.000000	0.000000	-1.000000	...	12.500000	12.50	16.300000
2656	50.000000	0.000000	167.600000	...	9.400000	16.10	15.400000
[541 rows x 107 columns]							

5. For data after down sampling:

feature standardization, maximum and minimum value standardization, converted value range (0,1)

```
from sklearn.preprocessing import MinMaxScaler
min_max_scaler = MinMaxScaler(copy=True, feature_range=(0, 1))
new_X_train = iris_x_train
new_X_test = iris_x_test
new_X_yan = iris_x_yan
from sklearn.preprocessing import Normalizer
normalizer = Normalizer(copy=True, norm='l2').fit(new_X_train)
new_X_train = normalizer.transform(new_X_train)
new_X_test = normalizer.transform(new_X_test)
new_X_yan = normalizer.transform(new_X_yan)
```

6. Model Building:

Model 1:

MLPClassifier (multilayer perceptron classifier) and try to use this algorithm to build the model. The parameters are set to three hidden layers, the first layer is 100 neurons, the second layer is 50 neurons, the third layer is 50 neurons, clf=MLPClassifier (solver='lbfgs', alpha=1e-5, hidden_layer_sizes=(100, 50, 50), random_State=1) # 3 hidden layers.

```
# 拟合 (模型训练)
print(iris_x_train)
clf.fit(iris_x_train, iris_y_train)
iris_y_predict = clf.predict(iris_x_test)
score = clf.score(iris_x_test, iris_y_test, sample_weight=None)
print('iris_y_predict=')
print(iris_y_predict)
```

```

ddd = []
for i in iris_y_predict:
    ddd.append(i)
print(ddd)
print('iris_y_test=')
print(iris_y_test)
print('Accuracy:', score)

```

Finally, the accuracy rate is obtained through the test set test:
Accuracy: 0.5754437869822485.

Model 2 (final):

Random forest algorithm: It is very suitable for multi-dimensional data and has a good effect on dealing with binary classification problems.

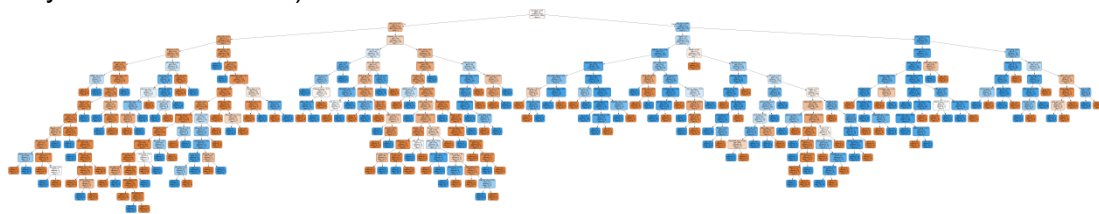
```

RandomForest Classifier model under sklearn
clf1 = RandomForestClassifier(n_estimators=200)
clf1.fit(new_X_train, iris_y_train)

```

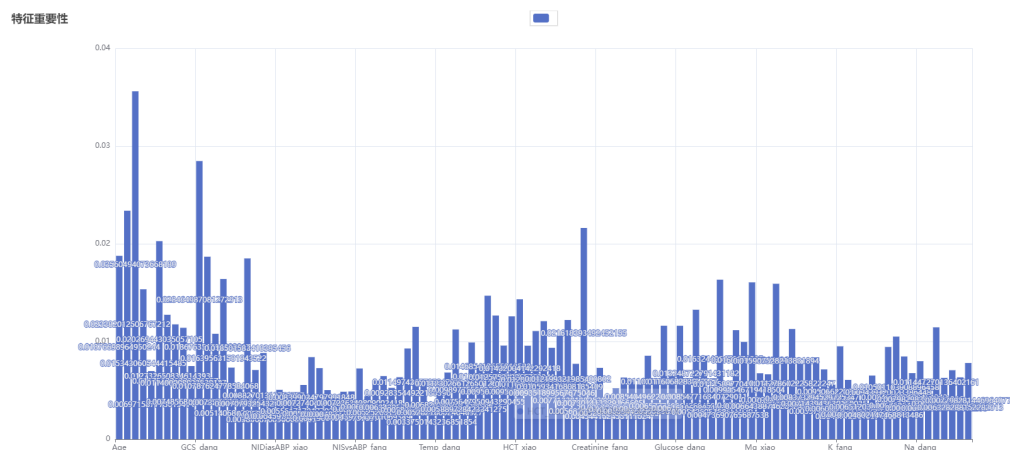
Set the number of decision trees to 200 (the best effect)

The following is a visual diagram of each decision tree (200 complete under the project, only one is shown here).



7. Model results

Visualization of feature importance:



The impact of features on the whole model can be clearly analyzed from the figure (file all.html), and some features with low impact can be eliminated.

Predicted value: (1: death, 0: survivor)

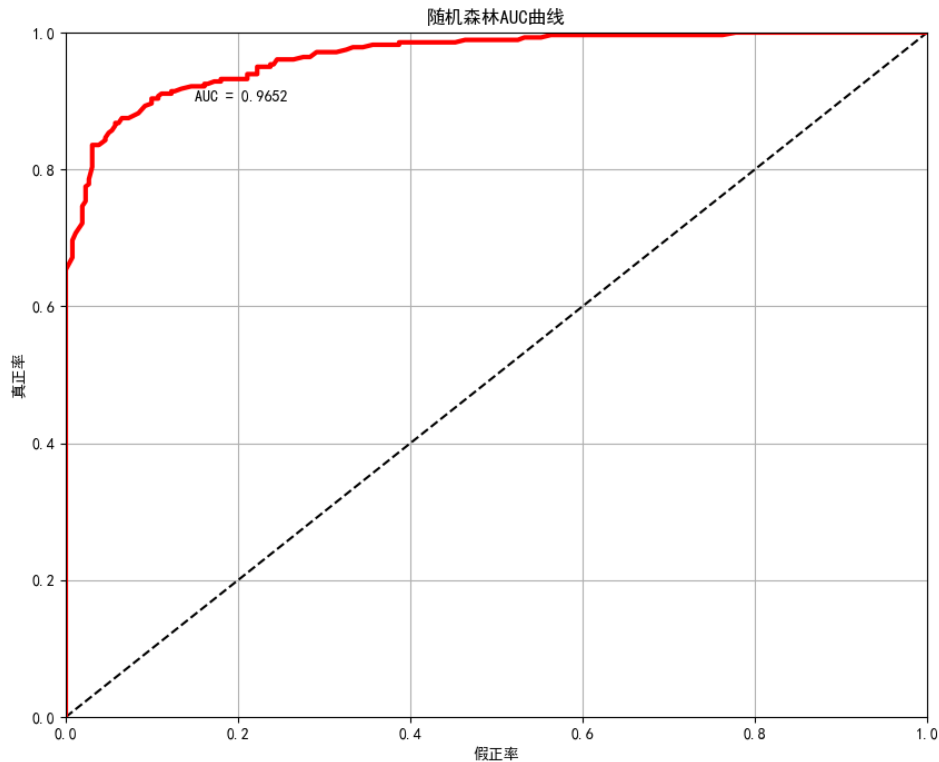
[0, 1, 1, 0, 1, 0, 0, 0, 1, 1, 1, 1, 0, 1, 0, 1, 1, 1, 0, 0, 1, 0, 0, 1, 1, 1, 1, 0, 0, 0, 0, 1, 1,

1, 0, 1, 0, 0, 1, 1, 0, 0, 1, 1, 1, 0, 0, 0, 1, 1, 1, 1, 0, 1, 0, 0, 0, 1, 1, 1, 1, 1, 1, 0, 0, 0,
0, 0, 0, 1, 0, 0, 0, 1, 1, 0, 1, 0, 0, 0, 0, 0, 0, 0, 0, 1, 0, 1, 1, 1, 0, 0, 0, 1, 0, 1, 1, 1, 0,
1, 1, 1, 0, 0, 0, 0, 1, 0, 1, 0, 1, 0, 0, 0, 0, 1, 1, 0, 0, 0, 0, 0, 0, 0, 1, 0, 0, 1, 0, 1, 0, 1,
0, 1, 1, 0, 0, 0, 0, 1, 0, 0, 0, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 0, 0, 1, 1, 1, 0, 1, 1, 1,
1, 0, 1, 0, 0, 0, 1, 0, 1, 0, 0, 0, 0, 1, 0, 0, 1, 1, 0, 1, 1, 0, 0, 0, 1, 1, 1, 1, 1, 0, 0, 1, 1,
1, 1, 1, 1, 1, 0, 1, 1, 0, 0, 1, 0, 1, 0, 1, 0, 1, 0, 0, 1, 0, 1, 0, 1, 1, 1, 1, 1, 1, 0, 1, 0, 1,
0, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 0, 0, 1, 0, 0, 0, 1, 1, 1, 1, 1, 0, 1, 0, 0, 0, 0, 1, 0, 0, 1,
0, 1, 0, 1, 0, 0, 0, 0, 1, 0, 0, 0, 0, 1, 0, 0, 0, 1, 1, 1, 0, 1, 0, 0, 0, 0, 0, 0, 0, 1, 0, 0, 1,
1, 0, 0, 1, 0, 0, 0, 1, 0, 0, 0, 0, 1, 1, 1, 0, 1, 0, 1, 1, 0, 1, 1, 1, 1, 1, 0, 0, 0, 1, 1, 0,
1, 1, 1, 1, 0, 1, 0, 1, 1, 0, 1, 0, 1, 0, 0, 1, 0, 0, 0, 1, 0, 1, 0, 0, 1, 0, 1, 0, 0, 1, 0, 1, 1,
0, 1, 1, 1, 0, 0, 0, 0, 0, 0, 1, 0, 1, 1, 1, 0, 0, 1, 0, 1, 1, 1, 0, 1, 0, 1, 0, 0, 1, 0, 0, 0,
0, 0, 1, 0, 0, 0, 1, 0, 0, 0, 1, 0, 1, 1, 1, 1, 0, 0, 1, 0, 0, 1, 0, 0, 1, 1, 1, 0, 1, 1, 1, 1,
0, 1, 0, 1, 0, 1, 1, 1, 0, 1, 0, 1, 0, 1, 0, 0, 1, 1, 1, 0, 1, 1, 1, 0, 1, 0, 0, 0, 1, 0, 0, 0, 1,
1, 1, 0, 0, 0, 1, 1, 1, 1, 0, 0, 0, 1, 1, 1, 1, 1, 1, 0, 0, 1, 1, 0, 1, 1, 1, 0, 0, 1, 1, 1, 0,
1, 0, 1, 1, 1, 0, 1, 0, 1, 1, 0, 1, 0, 0, 0, 1, 0, 1, 1, 0, 0, 1, 0, 1, 0, 1, 1, 0, 0, 0, 0, 0, 0,
1, 0, 0, 0, 0, 0, 0, 1, 1, 0, 0, 0, 0, 1, 1, 0, 0, 1, 0, 0, 1, 0, 1, 1, 0, 1, 0, 1, 0, 0, 1,
1, 0, 0, 0, 0, 1, 0, 1, 1, 0, 0, 0, 0, 1, 0, 1, 0, 0, 1, 1, 0, 0, 1, 1, 0, 0, 0, 1, 0, 1, 1, 1,
1, 0, 1, 1, 0, 0, 1, 0, 0, 1, 0, 0, 1, 0, 0, 0, 0, 1, 1, 0, 0, 1, 0, 0, 1, 1, 0, 0, 1, 1, 1, 0, 0,
0, 0, 0, 0, 0, 1, 0, 0, 0, 1, 1, 1, 1, 0, 0, 1, 1, 0, 0, 1, 1, 0, 0, 0, 0, 1, 1, 1, 1, 0, 1, 1, 1,
1, 1, 0, 1, 1, 0, 1, 1, 1, 0, 0, 0, 1, 1, 0, 1]

True value: (1: death, 0: survivor)

[0, 0, 0, 0, 1, 0, 0, 0, 1, 1, 1, 1, 0, 1, 0, 0, 1, 1, 0, 0, 1, 0, 0, 1, 1, 1, 1, 0, 1, 0, 0, 0, 1,
1, 0, 1, 0, 0, 1, 1, 0, 0, 1, 0, 0, 0, 0, 0, 1, 1, 1, 1, 0, 1, 0, 0, 0, 1, 1, 1, 1, 1, 1, 0, 0, 0,
0, 0, 0, 1, 0, 0, 0, 1, 0, 0, 1, 0, 0, 1, 0, 0, 0, 0, 0, 1, 0, 1, 1, 1, 0, 0, 0, 1, 0, 1, 1, 1, 0,
1, 1, 1, 0, 0, 0, 0, 1, 0, 1, 0, 1, 0, 0, 1, 0, 0, 1, 0, 0, 0, 0, 1, 0, 0, 1, 0, 1, 1, 0, 1, 0, 1,
0, 1, 1, 0, 0, 0, 0, 1, 1, 1, 0, 1, 1, 0, 1, 1, 1, 1, 0, 1, 1, 1, 1, 0, 0, 1, 1, 1, 0, 1, 1, 1,
0, 0, 1, 0, 0, 0, 1, 0, 1, 0, 0, 0, 0, 1, 0, 0, 1, 1, 1, 0, 1, 0, 0, 0, 1, 1, 1, 1, 1, 0, 0, 1, 1,
1, 1, 1, 1, 1, 0, 1, 1, 0, 0, 1, 0, 1, 0, 1, 1, 1, 0, 1, 1, 0, 1, 0, 1, 1, 0, 1, 1, 1, 0, 1, 0, 1,
0, 0, 1, 1, 1, 1, 1, 1, 1, 0, 1, 0, 0, 1, 0, 0, 0, 1, 1, 1, 1, 0, 0, 0, 1, 0, 1, 0, 0, 1, 0, 0, 1,
0, 1, 0, 1, 0, 0, 0, 0, 1, 0, 1, 0, 0, 1, 0, 0, 0, 1, 1, 0, 0, 1, 0, 0, 0, 0, 0, 0, 0, 1, 0, 0, 1,
1, 0, 0, 1, 1, 0, 0, 1, 0, 0, 0, 0, 1, 0, 1, 0, 1, 1, 1, 1, 0, 1, 1, 1, 1, 0, 0, 0, 1, 1, 0,
1, 1, 1, 1, 0, 1, 0, 1, 1, 0, 1, 0, 1, 0, 0, 1, 1, 0, 0, 1, 0, 1, 0, 0, 1, 0, 1, 0, 0, 1, 0, 1, 1,
0, 1, 1, 1, 0, 0, 0, 0, 0, 0, 1, 0, 0, 1, 1, 1, 0, 1, 0, 1, 1, 1, 0, 1, 1, 1, 0, 0, 1, 1, 0, 0,
0, 1, 1, 1, 0, 0, 0, 0, 0, 0, 1, 0, 1, 0, 1, 1, 0, 0, 0, 0, 0, 1, 0, 0, 1, 1, 1, 0, 1, 0, 1, 1, 1,
1, 1, 0, 1, 0, 1, 1, 1, 0, 1, 0, 1, 1, 0, 1, 1, 1, 0, 1, 1, 1, 0, 1, 0, 0, 1, 1, 1, 0, 0, 1,
1, 1, 0, 0, 0, 1, 0, 1, 1, 0, 0, 1, 1, 1, 1, 1, 1, 0, 0, 1, 1, 0, 1, 1, 1, 0, 0, 1, 1, 1, 0,
1, 0, 0, 1, 1, 0, 1, 0, 1, 1, 0, 1, 0, 0, 0, 1, 0, 1, 1, 0, 0, 1, 0, 1, 0, 1, 1, 0, 0, 0, 0, 0,
1, 0, 0, 0, 0, 0, 0, 1, 0, 0, 0, 0, 0, 1, 0, 0, 0, 0, 0, 0, 1, 0, 1, 1, 0, 1, 0, 1, 0, 1, 0, 0, 1,
1, 0, 0, 0, 0, 0, 0, 1, 1, 0, 0, 0, 0, 0, 1, 1, 0, 1, 1, 1, 0, 0, 1, 1, 0, 0, 1, 1, 0, 1, 0, 1, 1,
1, 0, 1, 1, 0, 0, 1, 0, 0, 0, 0, 1, 1, 0, 0, 0, 1, 1, 1, 0, 0, 1, 0, 0, 1, 1, 1, 0, 1, 1, 1, 0, 0,
1, 0, 0, 0, 0, 0, 0, 0, 1, 1, 1, 1, 1, 0, 0, 1, 1, 1, 0, 0, 0, 1, 1, 1, 1, 0, 1, 1, 1, 1,
1, 1, 0, 1, 1, 0, 1, 1, 1, 0, 0, 0, 1, 1, 0, 0]

Use the validation set to generate the AUC visualization:



Evaluation index F1: 0.89583333333334

Accuracy: 0.89583333333334

Response rate: 0.89583333333334

Conclusion

The research and analysis of this medical data, for predicting the mortality rate of patients in ICU wards, adopts the random forest algorithm to build the model, and the accuracy rate of the results in the test set reaches 89.5%, which has certain medical reference value. The performance of the model in the validation set, draw the AUC curve, and get AUC=0.9652 (the closer the value is, the better).

Improvements:

For data processing and feature engineering, the extraction and segmentation of features are not detailed enough, this data can also mine deeper features, which need to be optimized.