

Analysis of COVID-19 survival rate in Toronto

Abstract

Coronaviruses are a large group of viruses known to cause the common cold and even more serious illnesses, such as Middle East Respiratory syndrome (MERS) and Severe acute Respiratory syndrome (SARS).mNovel Coronavirus (CoVID-19) is a kind of novel coronavirus that has not been found in humans before 2019.To date, the total number of confirmed cases worldwide has exceeded 18 million. Seven hundred thousand people die of the disease, and the number is rising. How to control COVID-19 and reduce mortality has become a top priority for countries around the world. We now have data sets from Toronto on COVID-19 cases, and we want to analyze these data sets through the KDD process to develop models that actually mitigate the impact of COVID-19 globally and protect specific populations.

Key words: COVID-19 cases, Toronto, KDD process

目录

| | |
|--|----|
| 1.1 Identify the objectives of the business..... | 5 |
| 1.1.1 Background..... | 5 |
| 1.1.2 Business objectives..... | 6 |
| 1.1.3 Business success criteria..... | 7 |
| 1.2 Assess the situation..... | 7 |
| 1.2.1 Assumption..... | 7 |
| 1.2.2 Constraints..... | 8 |
| 1.3 Determine data mining objectives..... | 8 |
| 1.3.1 Data mining goals..... | 8 |
| 1.3.2 Data mining success criteria..... | 8 |
| 1.4 Produce a project plan..... | 9 |
| 1.4.1 Project Plan Overview..... | 9 |
| 1.4.2 Resources..... | 10 |
| 1.4.3 Risks..... | 10 |
| 2.1 Collect initial data..... | 11 |
| 2.1.1 Source: Kaggle website..... | 11 |
| 2.1.2 Original source:..... | 11 |

| | |
|---|----|
| 2.2 Describe the Data..... | 12 |
| 2.2.1 Data quantity..... | 12 |
| 2.2.2 Data Quality..... | 12 |
| 2.3 Explore data | 17 |
| 2.3.1 Data Overview | 17 |
| 2.3.2 Data Hypotheses | 18 |
| 2.3.3 Data analyses | 18 |
| 2.4 Verify the data quality | 24 |
| 2.4.1 Data missing | 24 |
| 2.4.2 Data errors and metric errors | 25 |
| 3.1 Select the Data | 25 |
| 3.2 Clean the Data..... | 29 |
| 3.3 Construct the data | 30 |
| 3.4 Integrate various data sources..... | 31 |
| 3.5 Format the data as required..... | 31 |
| 3.5.1 Format the data to fit decision tree model..... | 31 |
| 3.5.2 Check again for useless fields | 33 |
| 4.1 Reduce the Data | 34 |
| 4.1.1 Feature selection | 34 |

| | |
|---|----|
| 4.1.2 Reduce unimportant attribute | 35 |
| 4.2 Project the Data..... | 36 |
| 4.2.1 Balance the Data | 36 |
| 4.2.2 Distribution of target attribute | 37 |
| 5.1 Match and discuss the objectives of data mining to data mining methods | 41 |
| 5.1.1 Supervised and Unsupervised Learning | 41 |
| 5.1.2 Classification, Association and Segmentation | 42 |
| 5.2 Select the appropriate data-mining method (s) based on discussion .. | 43 |
| 5.2.1 Choose supervised learning | 43 |
| 5.2.2 Choose classification | 43 |
| 6.1 Conduct exploratory analysis and discuss | 44 |
| 6.1.1 Algorithm discussion | 44 |
| 6.2 Select data-mining algorithms based on discussion | 46 |
| 6.2.1 Algorithm requirements :..... | 46 |
| 6.2.2 Select data-mining algorithms | 48 |
| 6.3 Select appropriate model(s) and choose relevant parameter(s) | 49 |
| 7.1 Logical test designs..... | 51 |
| 7.2 Data mining must be conducted (the model must run)..... | 52 |
| 1 Run Bayesian Network model..... | 52 |

| | |
|---|----|
| 7.3 Search for patterns and document the model's output. | 53 |
| 8.1 Study and discuss the mined patterns. | 55 |
| (1) data and result | 55 |
| (2) models and patterns..... | 55 |
| 8.2 - Visualize the data, results, models and patterns in a clear and effective manner..... | 56 |
| 8.3 Interpret the results, models and patterns showing a clear understanding of the results. | 60 |
| 8.4 - Assess and evaluate the results, models and patterns using the appropriate methods/processes. | 61 |
| 8.4.1 assess the results, models and patterns | 61 |
| 8.4.2 evaluation the results, models and patterns | 63 |
| 8.5 Iterate prior steps (1 – 7) as required | 63 |
| 8.5.1 Business understanding | 63 |
| Reference | 66 |

1.1 Identify the objectives of the business

1.1.1 Background

As a known large group of viruses, coronavirus can cause the common cold and even more serious diseases, such as the Middle East respiratory syndrome (MERS) and severe acute respiratory syndrome (SARS). The novel coronavirus (COVID-19) is a new coronavirus which has not been found in humans before 2019.

The World Health Organization (who) defines coronavirus disease as a pandemic. So far, the total number of confirmed cases in the world has exceeded 18 million. 700000 people have died of the disease, and the number is still on the rise. It not only has a profound impact on human health, but also has an indelible impact on the global economic and social environment. How to curb the novel coronavirus pneumonia and reduce mortality has become the primary problem to be solved.

1.1.2 Business objectives

(1) Analyze the impact of age, gender, source of infection, outbreak associated, hospitalization status. Identify key factors and characteristics that contribute to survive in patients with COVID-19.

(2) Control these characteristics and factors that lead to death in patients with COVID-19, and guide the country, institution or family to locate groups that are more likely to survive from COVID-19.

(3) Protect those at risk for COVID-19 by pre-positioning and taking preventive measures.

(4) Reduce national or regional mortality from COVID-19 through prevention.

1.1.3 Business success criteria

(1) Significant increasing in survival rate of novel coronavirus pneumonia. Obtain characteristics of COVID-19 cases, including sex, age, source of infection, hospital status, etc. Implementing targeted prevention, and successfully improve COVID-19 survival rate.

(2) The project can be completed on time without exceeding the budget.

1.2 Assess the situation

1.2.1 Assumption

(1) Survival rate for COVID-19 patients was associated with age. Generally speaking, young children are at the greatest risk of infection. For example, approximately 57% of malaria occurs in children under 5 years of age. However, in the face of the new coronavirus, the elderly are the most at risk. This may be because older people have potential health problems, especially cardiovascular diseases and respiratory diseases. The elderly are more likely to have these health problems than the young, which may be one of the important reasons why the elderly are not likely to survive the risk of COVID-19.

(2) The gendered impact on health outcomes. In many cases, since most of the world's health workers are women, women seem to be more likely to be diagnosed with covid-19. At the same time, compared with women, the male mortality rate in each country has maintained a higher growth trend. This may be due to the fact that men have a higher smoking rate and are more likely to suffer from cardiopulmonary diseases.

(3) Whether or not a person is treated, the level of treatment and the patient's ability to recover from the disease all affect whether or not he will die of a disease.

1.2.2 Constraints

(1) The source of data is single. The data of COVID-19 cases for this study are from Toronto Public Health in January 2020. The loss of some data will still cause certain errors in the overall statistics of outcome data.

(2) In order to protect personal privacy, some more detailed data cannot be obtained, and the lack of critical information may skew the overall objectivity of the results.

1.3 Determine data mining objectives

1.3.1 Data mining goals

(1) Get a set of decision rules that determine the survival rate of COVID-19.

(2) Use decision rules to predict or classify the COVID-19 survival rate of a person or group of people in the future

1.3.2 Data mining success criteria

(1) Model quality

More than 85% accuracy; Faster response speed; The output is easy to understand

(2) Engineering dimension

Flexible model; easy to use; tight layout, embeddable and extensible.

(3) Logistical constraints

Simple calculation; less development time.

1.4 Produce a project plan

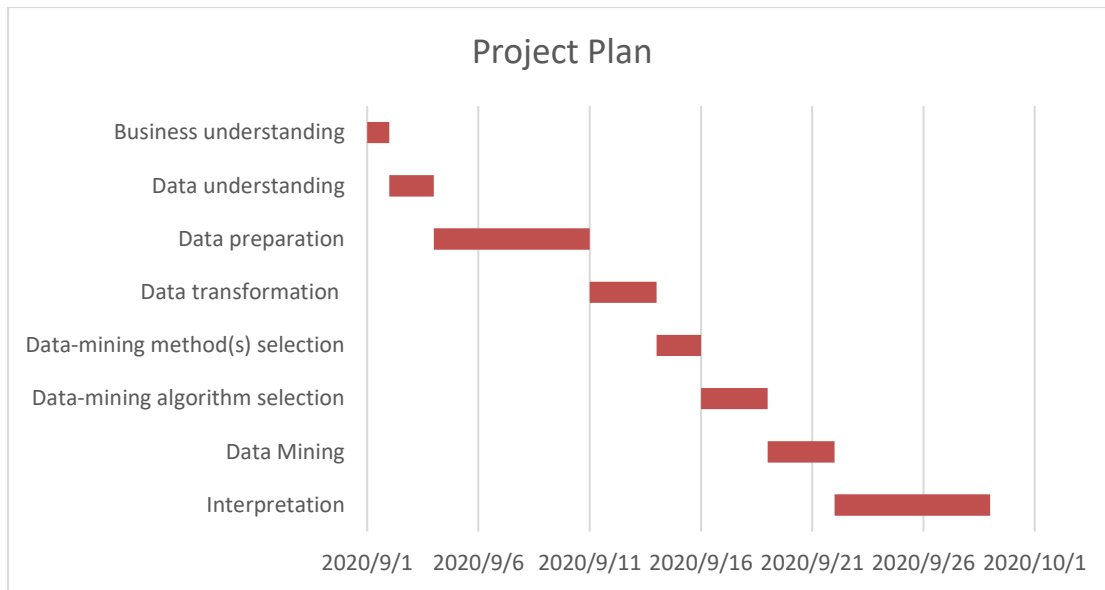
1.4.1 Project Plan Overview

The whole project will take about a month, which is divided into eight steps. Data preparation and interpraction accounted for the largest proportion. This is because data preparation is the basis of the following steps. Detailed data preparation can help us speed up the data transformation, the selection of methods and algorithms, and data mining. The interpraction step requires us to evaluate and review the previous seven steps, which will affect the accuracy and relevance of the modeling.

Note :Modeling is a phase in which multiple iterations usually occur.

The overview plan for the study is as shown in the table below.

Table 1.1 project plan overview



1.4.2 Resources

Because the data set is simple and the project needs less resources, I only need my time and myself as the resources for data analysis.

1.4.3 Risks

There will be some risks in some stages of the project, which will affect the progress, quality and results of the project.

Table 1.2 project plan risk

| Task | Days | Start Date | Risk |
|------------------------|------|------------|--|
| Business understanding | 1 | 2020/9/1 | Economic change |
| Data understanding | 2 | 2020/9/2 | Data problems, technology problems |
| Data preparation | 7 | 2020/9/4 | Data problems, technology problems |

| | | | |
|------------------------------------|---|-----------|--|
| Data transformation | 3 | 2020/9/11 | Data problems, technology problems |
| Data-mining method(s) selection | 2 | 2020/9/14 | Technology problems, inability to find adequate model |
| Data-mining algorithm selection | 3 | 2020/9/16 | Technology problems, inability to find adequate model |
| Data Mining | 3 | 2020/9/19 | Technology problems, inability to find adequate model |
| Interpretation | 7 | 2020/9/22 | Economic change, inability to implement results |

2.1 Collect initial data

2.1.1 Source: Kaggle website

<https://www.kaggle.com/divyansh22/toronto-covid19-cases>

2.1.2 Original source:

Toronto Public Health

Collection methodology: The data was collected published by Toronto Public Health under Open Government License - Toronto

<https://open.toronto.ca/dataset/covid-19-cases-in-toronto/>

License
Open Government License - Toronto

Publisher
Published by
Toronto Public Health
Contact
cdsu@toronto.ca

| DATA PREVIEW | | | | | | | |
|--------------|-------------|---------------------|----------------|--------------------|-----|---------------------|----------------|
| _id | Assigned_ID | Outbreak Associated | Age Group | Neighbourhood Name | FSA | Source of Infection | Classification |
| 126705 | 1 | Sporadic | 50 to 59 Years | Willowdale East | M2N | Travel | CONFIRMED |
| 126706 | 2 | Sporadic | 50 to 59 Years | Willowdale East | M2N | Travel | CONFIRMED |
| 126707 | 3 | Sporadic | 20 to 29 Years | Parkwoods-Donalda | M3A | Travel | CONFIRMED |

figure 1. About COVID-19 Cases in Toronto

2.2 Describe the Data

2.2.1 Data quantity

Amount of data: 14910 records and 16 attributes, the data used in this analysis is to record the COVID- 19 survival rate of different populations.

```

      _id  Outbreak Associated ... Ever in ICU Ever Intubated
14907  59201 Outbreak Associated ...           No           No
14908  59202 Outbreak Associated ...           No           No
14909  59203 Outbreak Associated ...           No           No
14910  59204 Outbreak Associated ...           No           No

[4 rows x 16 columns]
```

figure 2. Amount of data

2.2.2 Data Quality

(1) Relevant Attributes: _id, Outbreak Associated, Age Group, Neighbourhood Name, Client Gender, Classification, Source of Infection, Episode Date, Reported Date, Outcome, Currently Hospitalized, Currently in ICU, Currently Intubated, Ever Hospitalized, Ever in ICU, Ever Intubated

| Column | Description |
|---------------------|--|
| _id | Unique row identifier for Open Data database |
| Outbreak Associated | Outbreak associated cases are associated with outbreaks of COVID-19 in Toronto healthcare institutions and healthcare settings (e.g. long-term care homes, retirement homes, hospitals, etc.) and other Toronto congregate settings (such as homeless shelters). |
| Age Group | Age groups (in years): ≤19, 20-29, 30-39, 40-49, 50-59, 60-69, 70-79, 80-89, 90+, unknown (blank) |
| Neighbourhood Name | To help government and community agencies with local planning, Toronto is divided into 140 geographically diverse communities. This is about the specific location of the case |

| | |
|---------------------|--|
| Source of Infection | The most likely routes of covid-19 infection may include: travel, close contact with a case, institutional setting, healthcare setting, community, unknown / missing, Pending (Information on source of infection pending) , N/A (Outbreak-associated cases) |
| Classification | The provincial case definition is applied to classify the cases as confirmed or possible cases according to the standards. |
| Episode Date | The date of onset is a derivative variable, which refers to the onset of symptoms (the first day of covid-19 symptoms), the date of collection of laboratory specimens or the date of reporting. |
| Reported Date | The date on which the case was reported to Toronto Public Health. |
| Client Gender | People are classified according to the specified physiological sex. |

| | |
|------------------------|--|
| Outcome | <p>Fatal: Cases with a fatal outcome reported.</p> <p>Resolved: Cases with no reported deaths, and those reported as "recovered" or reported more than 14 days after the onset of symptoms, and the case is not currently hospitalized.</p> <p>Active: All other cases</p> |
| Currently Hospitalized | Cases that are currently admitted to hospital |
| Currently in ICU | Cases that are currently admitted to the intensive care unit (ICU) |
| Currently Intubated | Cases that were intubated related to their COVID-19 infection |
| Ever Hospitalized | Cases that were hospitalized related to their COVID-19 infection |
| Ever in ICU | Cases that were admitted to the intensive care unit (ICU) related to their COVID-19 infection. |
| Ever Intubated | Cases that were intubated related to |

| | |
|--|--------------------------|
| | their COVID-19 infection |
|--|--------------------------|

(3) Data type

1. Basic data of the people:

_id: numeric

Age Group: categorical (string)

Client Gender: categorical (string)

Neighbourhood Name: categorical (string)

2. Basic data on virus infection

Outbreak Associated: categorical (string)

Classification: categorical (string)

Source of Infection: categorical (string)

Episode Date: numeric

Reported Date: numeric

3. Basic data of severity information

Outcome: categorical (string)

Currently Hospitalized: categorical (string)

Currently in ICU: categorical (string)

Currently Intubated: categorical (string)

Ever Hospitalized: categorical (string)

Ever in ICU: categorical (string)

Ever Intubated: categorical (string)

```
_id                int64
Outbreak Associated object
Age Group          object
Neighbourhood Name object
Source of Infection object
Classification     object
Episode Date       object
Reported Date      object
Client Gender      object
Outcome            object
Currently Hospitalized object
Currently in ICU   object
Currently Intubated object
Ever Hospitalized  object
Ever in ICU        object
Ever Intubated     object
dtype: object
```

figure 3. type of data

2.3 Explore data

2.3.1 Data Overview

This data set contains demographic, geographic, and severity information for all confirmed and probable cases reported to and managed by Toronto Public Health since the first case was reported in January 2020.

| data - DataFrame | | | | | | | | | | | | | |
|------------------|-------|---------------------|----------------|--------------------|---------------------------|----------------|--------------|---------------|----------------|----------|------------------------|------------------------|---------------------|
| Index | Id | Outbreak Associated | Age Group | Neighbourhood Name | Source of Infection | Classification | Episode Date | Reported Date | Patient Gender | Outcome | Currently Hospitalized | Currently in Isolation | Currently Intubated |
| 0 | 44294 | Sporadic | 50-59 | Malvern | Institutional | CONFIRMED | 2020/3/25 | 2020/3/27 | MALE | RESOLVED | No | No | No |
| 1 | 44295 | Sporadic | 20-29 | Malvern | Community | CONFIRMED | 2020/3/20 | 2020/3/28 | MALE | RESOLVED | No | No | No |
| 2 | 44296 | Sporadic | 60-69 | Malvern | Travel | CONFIRMED | 2020/3/4 | 2020/3/8 | FEMALE | RESOLVED | No | No | No |
| 3 | 44297 | Outbreak Associated | 50-59 | Rouge | N/A - Outbreak associated | CONFIRMED | 2020/5/2 | 2020/5/4 | FEMALE | RESOLVED | No | No | No |
| 4 | 44298 | Sporadic | 30-39 | Rouge | Close contact | CONFIRMED | 2020/5/31 | 2020/6/6 | FEMALE | RESOLVED | No | No | No |
| 5 | 44299 | Sporadic | 20-29 | Rouge | Close contact | CONFIRMED | 2020/6/1 | 2020/6/6 | MALE | RESOLVED | No | No | No |
| 6 | 44300 | Sporadic | 60-69 | Rouge | Community | CONFIRMED | 2020/5/22 | 2020/6/1 | MALE | RESOLVED | No | No | No |
| 7 | 44301 | Sporadic | 30-39 | Rouge | Close contact | PROBABLE | 2020/5/26 | 2020/6/2 | MALE | RESOLVED | No | No | No |
| 8 | 44302 | Sporadic | 30-39 | Malvern | Close contact | CONFIRMED | 2020/5/11 | 2020/5/16 | MALE | RESOLVED | No | No | No |
| 9 | 44303 | Sporadic | 19 and younger | Malvern | Close contact | PROBABLE | 2020/6/6 | 2020/6/9 | MALE | RESOLVED | No | No | No |
| 10 | 44304 | Sporadic | 30-39 | Malvern | Close contact | CONFIRMED | 2020/5/17 | 2020/5/21 | MALE | RESOLVED | No | No | No |

figure 4. data set of COVID-19 cases

2.3.2 Data Hypotheses

A total of 16 attributes can help us explore factors that can influence the outcome of recovery in COVID-19 cases. It is assumed that gender, age, source of infection and medical degree are all factors that will influence the condition of the case, and geographical location and source of infection are useless factors for further exploration.

2.3.3 Data analyses

This project uses Spyder to analyze the basic data and conduct preliminary statistical analysis of the data set.

```
import pandas as pd
import seaborn as sns
import matplotlib.pyplot as plt
# from sklearn import preprocessing

cases_data = pd.read_csv("G:/我的云端硬盘/semester2/722/assignment/2/COVID19 cases Toronto.csv")
print(cases_data)
```

figure 5. Code for reading data

(1) Target analyses

In order to measure the recovery of a case, the outcome attribute can be used as the target value. The outcome property contains three values, FATAL, ACTIVE and RESOLVED. It can be seen from the figure that

the number of recovered cases is far more than that of death.

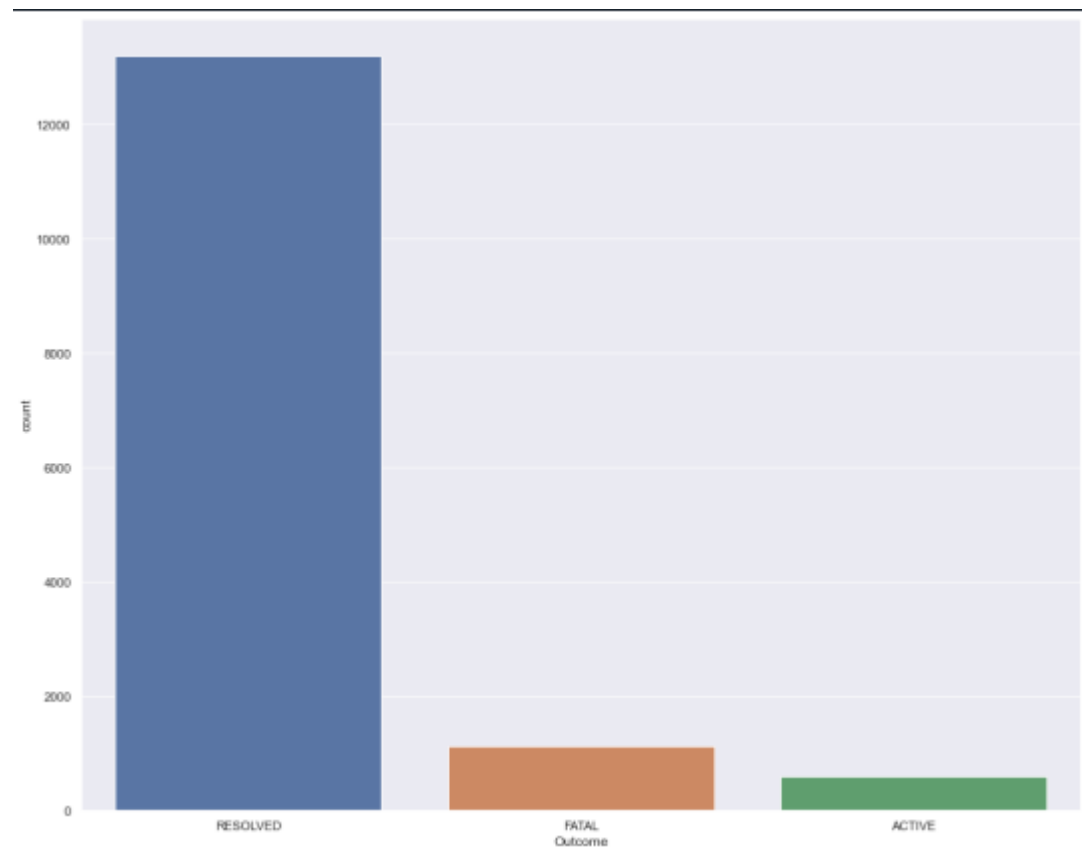


figure 6. Target value distribution

(2) Taking age as an example, it can be seen from the figure that among the confirmed or possibly confirmed cases, there are more cases of rehabilitation between 20 and 60 years old, almost all of which are more than 1500, while the cases of other age groups are less.

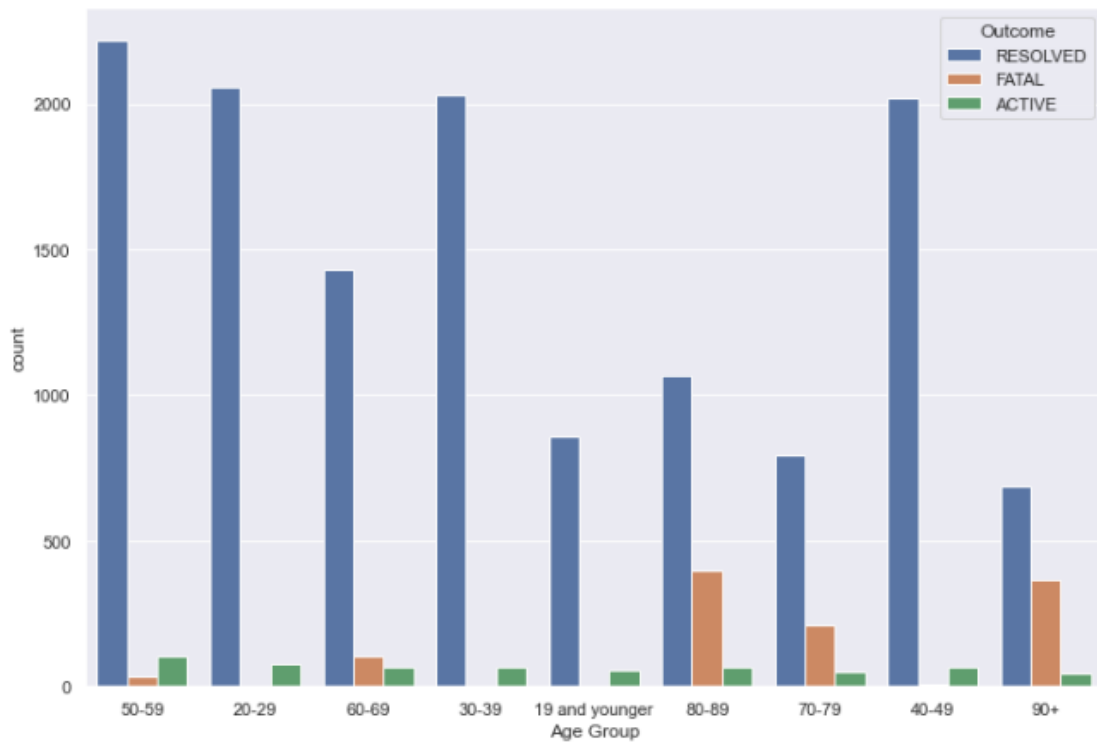


figure 7 confirmed or suspected cases over different age group

(3) It can be seen from the figure that women have more rehabilitation cases than men, which shows that gender may be one of the factors affecting the level of rehabilitation. It is also possible that there are more cases in women than men.

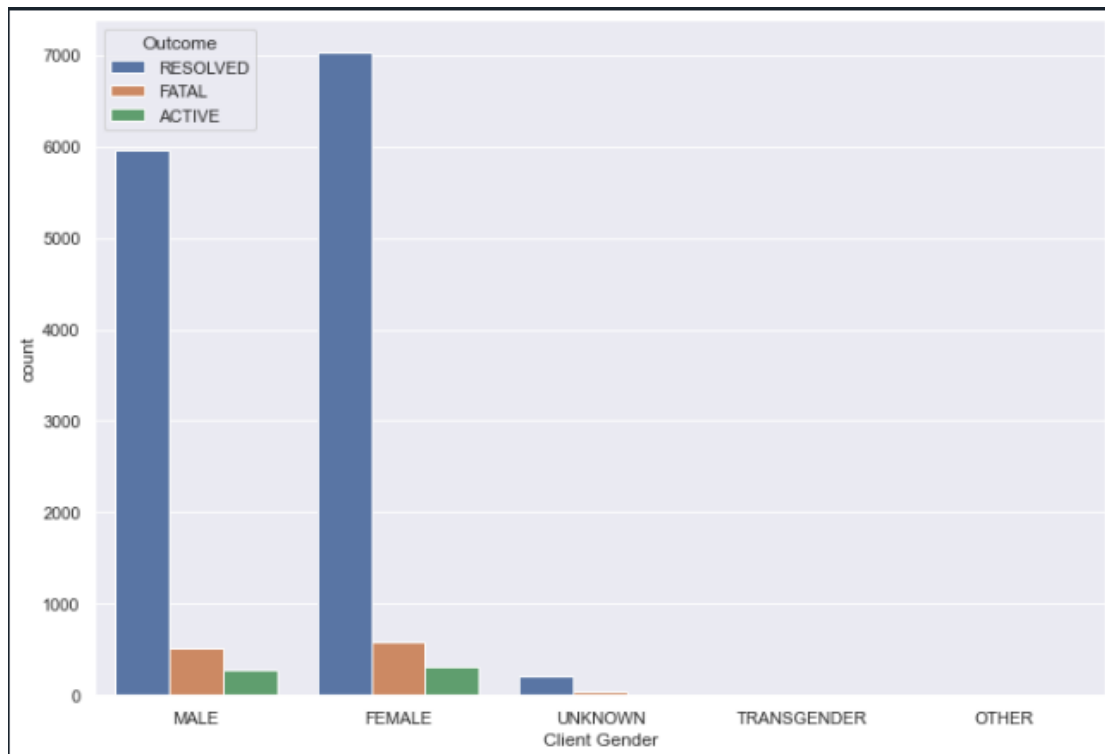


figure 8 confirmed or suspected cases over different gender group

(3) Before I did novel coronavirus pneumonia, I think that the medical conditions of the cases will affect their rehabilitation level of the new crown pneumonia cases. However, through visual analysis, it seems that most of the cases without treatment or hospitalization still have a high recovery rate. Currently, novel coronavirus pneumonia related treatments are rare, and only a few cases have been hospitalized.

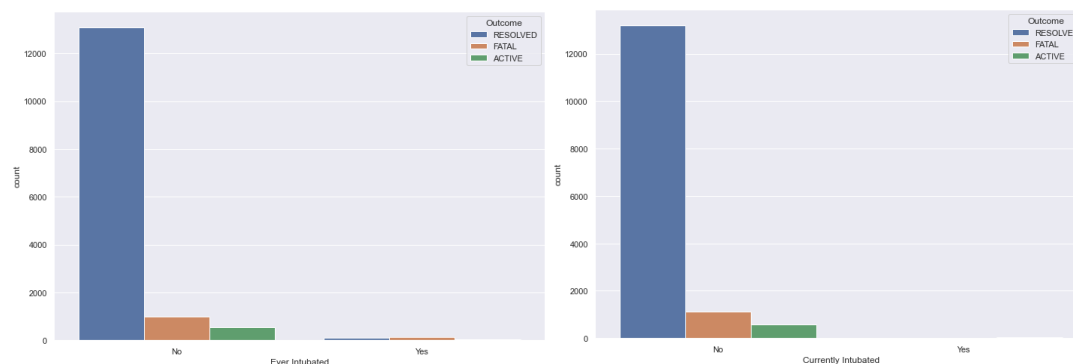


figure 9 Ever Intubated and currently Intubated

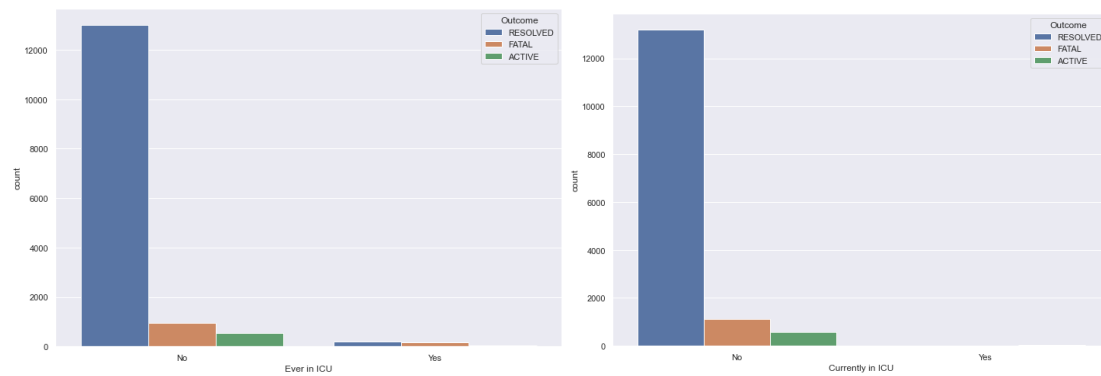


figure 10 Ever in ICU and currently in ICU

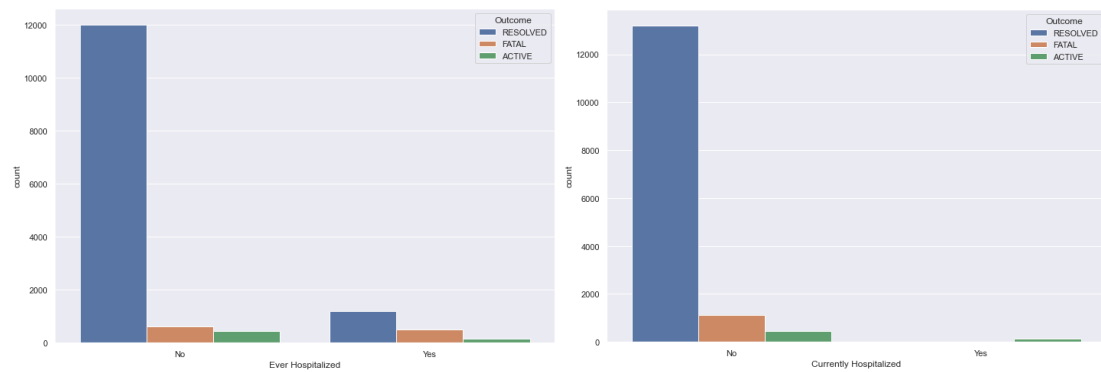


figure 11 Ever Hospitalized and Hospitalized

(4) The source of infection can be roughly divided into two categories: outbreak related cases and sporadic cases. The number of patients recovered from sporadic infection was more than that of outbreak related cases. No matter from Figure 11 or Figure 12, it can be seen that the number of death cases related to outbreak was more than that of sporadic infection. This suggests that the source of infection may be an important factor in the death of the case.

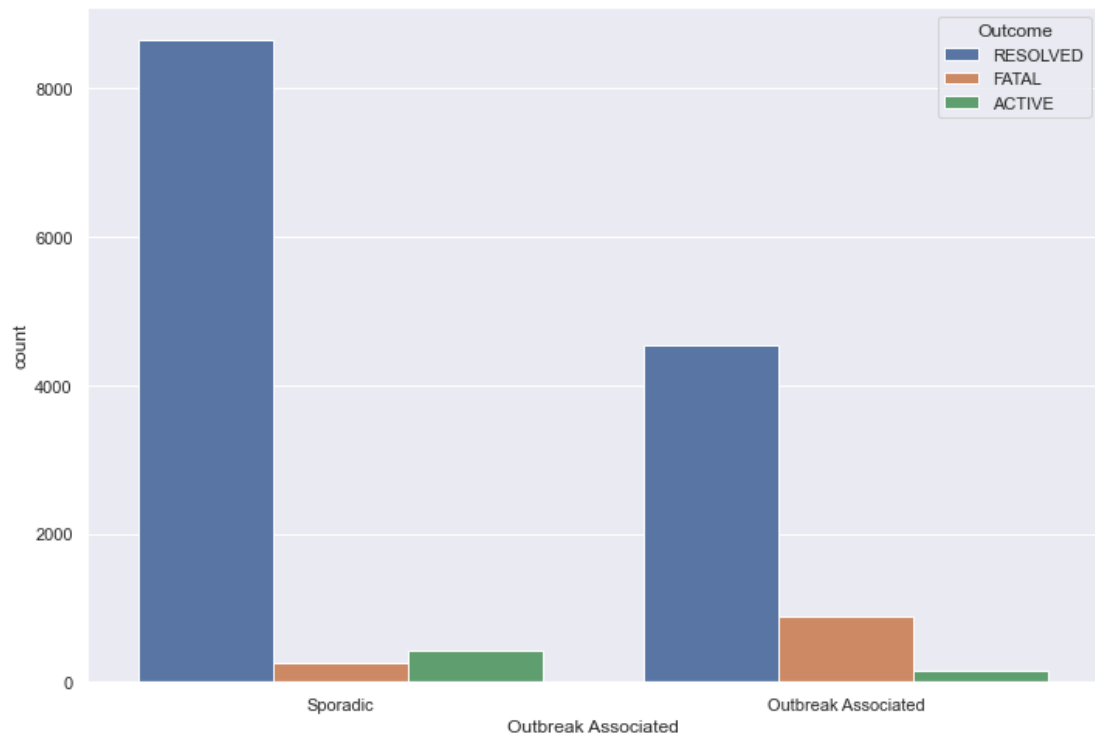


figure 11. Sporadic and Outbreak Associated cases

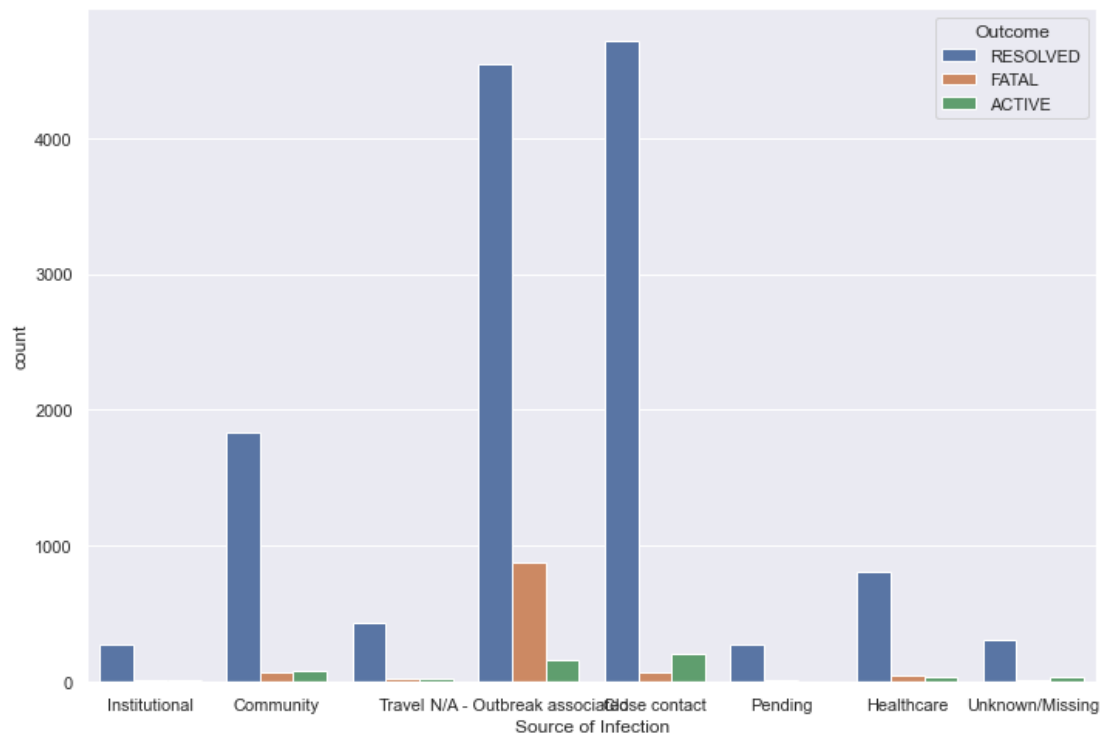


figure 12. Source of Infection

2.4 Verify the data quality

2.4.1 Data missing

First of all, all magnetic field measurements are correct. As can be seen from Figure 14, the data is relatively complete, only the age group and Neighbourhood Name has a small number of null value. This may be because it is difficult to collect complete information during the COVID-19 outbreak. There are some values in the dataset that require further processing.

```
nulldata=cases_data.isnull().sum()  
print(nulldata)
```

figure 13. Code for data missing

```
id 0  
Outbreak Associated 0  
Age Group 32  
Neighbourhood Name 613  
Source of Infection 0  
Classification 0  
Episode Date 0  
Reported Date 0  
Client Gender 0  
Outcome 0  
Currently Hospitalized 0  
Currently in ICU 0  
Currently Intubated 0  
Ever Hospitalized 0  
Ever in ICU 0  
Ever Intubated 0  
dtype: int64  
<class 'pandas.core.frame.DataFrame'>
```

figure 14. Number of missing values

2.4.2 Data errors and metric errors

Since most of the data is a string type and cannot be calculated, there are basically no extreme values and outliers in the data. Only `_id` the only data of type int, `_id` can detect outliers and extreme values. However, due to `_id` itself is just a case number, and will not have any impact on the recovery or death of cases, so it has no practical significance. In addition to these defects and limitations about the dataset, no other data errors were found.

```
detection=cases_data.describe()
print(detection)
```

figure 15. Code for outliers and extremes detection

```
count    14911.000000
mean      51749.000000
std       4304.579267
min       44294.000000
25%       48021.500000
50%       51749.000000
75%       55476.500000
max       59204.000000
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 14911 entries, 0 to 14910
Data columns (total 16 columns):
 #   Column                Non-Null Count  Dtype
---  -
```

figure 16. outliers and extremes

3.1 Select the Data

By observing complete data and the distribution of missing data, I

decided to select ten attributes in the dataset.

Select attribute: Age Group, Client Gender, Source of Infection, Outcome, Currently Hospitalized, Currently in ICU, Currently Intubated, Ever Hospitalized, Ever in ICU and Ever Intubated. Because these attributes may have different degrees of influence and response in analyzing COVID- 19 survival rate of different populations.

```
# select data
select_data=cases_data[['Outbreak Associated',
                        'Age Group','Client Gender','Source of Infection',
                        'Outcome','Currently Hospitalized',
                        'Currently in ICU','Currently Intubated',
                        'Ever Hospitalized','Ever in ICU',
                        'Ever Intubated']]
print(select_data.head)
```

figure 16. code for selecting attributes

(1) Source of Infection

The condition of confirmed cases may be related to different sources of infection, so different sources of infection may increase or decrease the survival rate of COVID-19.

(2) Client Gender

Gender may lead to different COVID- 19 survival rate because of men and women differences in personality and physical conditions.

(3) Age Group

People in different age groups have large differences in their physical, mental, or economic conditions, so this attribute should be retained, and its effect on COVID- 19 survival rate.

(4) Outcome

Data on the status of all confirmed and probable cases at present, which is the direct data of COVID- 19 survival rate.

(5) Currently Hospitalized and Ever Hospitalized

Whether or not a patient is admitted to hospital represents to some extent, the level of treatment for a confirmed case, which may affect the patient's level of recovery.

(6) Currently in ICU and Ever in ICU

The ICU represents the level of treatment for severe cases and affects the survival rate of COVID-19 cases.

(7) Currently Intubated and Ever Intubated

Intubated as treatment may play a decisive role in the rehabilitation of COVID-19 cases and increase the survival rate of COVID-19.

| Outbreak Associated | Age Group | Source of Infection | Outcome | ntly Hospit. | rrently in IC | ently Intub. | r Hospitali. | Ever in ICU | er Intubate |
|---------------------|----------------|---------------------------|----------|--------------|---------------|--------------|--------------|-------------|-------------|
| Sporadic | 50-59 | Institutional | RESOLVED | No | No | No | No | No | No |
| Sporadic | 20-29 | Community | RESOLVED | No | No | No | Yes | No | No |
| Sporadic | 60-69 | Travel | RESOLVED | No | No | No | Yes | Yes | Yes |
| Outbreak Associated | 50-59 | N/A - Outbreak associated | RESOLVED | No | No | No | No | No | No |
| Sporadic | 30-39 | Close contact | RESOLVED | No | No | No | No | No | No |
| Sporadic | 20-29 | Close contact | RESOLVED | No | No | No | No | No | No |
| Sporadic | 60-69 | Community | RESOLVED | No | No | No | No | No | No |
| Sporadic | 30-39 | Close contact | RESOLVED | No | No | No | No | No | No |
| Sporadic | 30-39 | Close contact | RESOLVED | No | No | No | No | No | No |
| Sporadic | 19 and younger | Close contact | RESOLVED | No | No | No | No | No | No |
| Sporadic | 30-39 | Close contact | RESOLVED | No | No | No | No | No | No |
| Outbreak Associated | 20-29 | N/A - Outbreak associated | RESOLVED | No | No | No | No | No | No |

figure 17. Data after select

Table 3.1. Useless data and reason

| Useless Data | reason |
|----------------|---|
| _id | no research value |
| Classification | It had nothing to do with survival after covid-19. |
| Episode Date | Means the time of coVID-19 diagnosis is independent of survival rate. |
| Reported Date | Means the time of coVID-19 diagnosis is independent of survival rate. |

3.2 Clean the Data

After selecting the data, only the age group attribute has missing data. As the number of missing values is very small, only 32, which has little impact on the sample population, so I choose to fill the missing data with the previous data. After filling, the missing quantity is zero.

```
[14911 rows x 11 columns]>
Outbreak Associated      0
Age Group                32
Client Gender            0
Source of Infection      0
Outcome                  0
Currently Hospitalized    0
Currently in ICU          0
Currently Intubated       0
Ever Hospitalized         0
Ever in ICU               0
Ever Intubated            0
dtype: int64
```

figure 18. check missing values

```
select_data['Age Group'].fillna(method = 'pad',inplace = True)
nulldata=select_data.isnull().sum()
print(nulldata)
```

figure 19. fill null value.

```
Outbreak Associated      0
Age Group                0
Client Gender            0
Source of Infection      0
Outcome                  0
Currently Hospitalized    0
Currently in ICU          0
Currently Intubated       0
Ever Hospitalized         0
Ever in ICU               0
Ever Intubated            0
dtype: int64
```

figure 20. check missing values after filling

3.3 Construct the data

Create four new attributes and generate the field Outcome2, ICU, Intubated and Hospitalized. Both the former hospitalization and the current hospitalization belong to the hospitalization treatment, which can be combined into one field. Similarly, Intubated experiences can also compose a field. There are three values in Outcome , ACTIVE, RESOLVED and FATAL. ACTIVE and RESOLVED both represent patients still alive and can be combined into a single value. In the "client gender" attribute, "unknown" and "other" values have similar meanings and can be classified into one class.

```
# Construct the data
frame = pd.DataFrame(select_data, columns=['Outbreak Associated',
                                           'Age Group', 'Client Gender', 'Source of Infection',
                                           'Outcome', 'Currently Hospitalized',
                                           'Currently in ICU', 'Currently Intubated',
                                           'Ever Hospitalized', 'Ever in ICU',
                                           'Ever Intubated'])

def function(a, b):
    if a == 'Yes' or b == 'Yes':
        return 'Yes'
    else:
        return 'No'

frame['Hospitalized'] = frame.apply(lambda x: function(x["Currently Hospitalized"], x["Ever Hospitalized"]))

frame['ICU'] = frame.apply(lambda x: function(x["Currently in ICU"], x["Ever in ICU"]), axis=1)
frame['Intubated'] = frame.apply(lambda x: function(x["Currently Intubated"], x["Ever Intubated"]), axis=1)
frame['Outcome2'] = frame['Outcome']

frame.loc[frame['Outcome'] == 'ACTIVE', 'Outcome2'] = 'NONFATAL'
frame.loc[frame['Outcome'] == 'RESOLVED', 'Outcome2'] = 'NONFATAL'

frame.loc[frame['Client Gender'] == 'UNKNOWN', 'Client Gender'] = 'OTHER'
```

Figure22. generate new attributes.

| Index | break Associated | Age Group | Patient Gender | Source of Infection | Outcome | Initially Hospitalized | Currently in ICU | Currently Intubated | Previously Hospitalized | Ever in ICU | Ever Intubated | Hospitalized | ICU | Intubated | Outcome2 |
|-------|---------------------|----------------|----------------|---------------------------|----------|------------------------|------------------|---------------------|-------------------------|-------------|----------------|--------------|-----|-----------|----------|
| 0 | Not Associated | 50-59 | MALE | Institutional | RESOLVED | No | No | No | No | No | No | No | No | No | NONFATAL |
| 1 | Not Associated | 20-29 | MALE | Community | RESOLVED | No | No | No | Yes | No | No | Yes | No | No | NONFATAL |
| 2 | Not Associated | 60-69 | FEMALE | Travel | RESOLVED | No | No | No | Yes | Yes | Yes | Yes | Yes | Yes | NONFATAL |
| 3 | Outbreak Associated | 50-59 | FEMALE | N/A - Outbreak associated | RESOLVED | No | No | No | No | No | No | No | No | No | NONFATAL |
| 4 | Not Associated | 30-39 | FEMALE | Close contact | RESOLVED | No | No | No | No | No | No | No | No | No | NONFATAL |
| 5 | Not Associated | 20-29 | MALE | Close contact | RESOLVED | No | No | No | No | No | No | No | No | No | NONFATAL |
| 6 | Not Associated | 60-69 | MALE | Community | RESOLVED | No | No | No | No | No | No | No | No | No | NONFATAL |
| 7 | Not Associated | 30-39 | MALE | Close contact | RESOLVED | No | No | No | No | No | No | No | No | No | NONFATAL |
| 8 | Not Associated | 30-39 | MALE | Close contact | RESOLVED | No | No | No | No | No | No | No | No | No | NONFATAL |
| 9 | Not Associated | 19 and younger | MALE | Close contact | RESOLVED | No | No | No | No | No | No | No | No | No | NONFATAL |
| 10 | Not Associated | 30-39 | MALE | Close contact | RESOLVED | No | No | No | No | No | No | No | No | No | NONFATAL |
| 11 | Outbreak Associated | 20-29 | MALE | N/A - Outbreak associated | RESOLVED | No | No | No | No | No | No | No | No | No | NONFATAL |

Figure22. New attributes.

3.4 Integrate various data sources

There is no need to integrate data because the project has only one data source, since the current data source already includes medical factors (hospitalization, source of infection) and personal factors (gender, age). There is already a comprehensive set of factors affecting cure rates, so we can use current data sources for data mining.

3.5 Format the data as required

3.5.1 Format the data to fit decision tree model

Since the machine learning algorithm in scikit-learn only understands the number input, we convert string data into numerical data.

(1) First, we need to layout the different categorical columns from the frame dataframe and the unique classes under each of these columns.


```
#import the necessary module
from sklearn import preprocessing

# create the Labelencoder object
le = preprocessing.LabelEncoder()

print("Outbreak Associated: ",frame['Outbreak Associated'].unique())
print("Source of Infection: ",frame['Source of Infection'].unique())
print("Age Group : ",frame['Age Group'].unique())
print("Outcome2 : ",frame['Outcome2'].unique())
print("ICU: ",frame['ICU'].unique())
print("Intubated: ",frame['Intubated'].unique())
print("Hospitalized: ",frame['Hospitalized'].unique())
print("Client Gender": ",frame['Client Gender'].unique())
```

Figure23. Layout the different categorical columns.

```
Outbreak Associated: ['Sporadic' 'Outbreak Associated']
Source of Infection: ['Institutional' 'Community' 'Travel' 'N/A - Outbreak associated'
'Close contact' 'Pending' 'Healthcare' 'Unknown/Missing']
Age Group : ['50-59' '20-29' '60-69' '30-39' '19 and younger' '80-89' '70-79' '40-49'
'90+']
Outcome2 : ['NONFATAL' 'FATAL']
ICU: ['No' 'Yes']
Intubated: ['No' 'Yes']
Hospitalized: ['No' 'Yes']
Client Gender: ['MALE' 'FEMALE' 'OTHER' 'TRANSGENDER']
```

Figure24. unique classes under each of these columns.

(2) Encode these strings into numeric labels.

```
# convert the categorical columns into numeric
frame['Outbreak Associated'] = le.fit_transform(frame['Outbreak Associated'])
frame['Source of Infection'] = le.fit_transform(frame['Source of Infection'])
frame['Age Group'] = le.fit_transform(frame['Age Group'].astype(str))
frame['Outcome2'] = le.fit_transform(frame['Outcome2'])
frame['ICU'] = le.fit_transform(frame['ICU'])
frame['Intubated'] = le.fit_transform(frame['Intubated'])
frame['Hospitalized'] = le.fit_transform(frame['Hospitalized'])
frame['Client Gender'] = le.fit_transform(frame['Client Gender'])

#display the initial records
print(frame.head())
```

Figure25. Encode these strings into numeric labels

| | Outbreak Associated | Age Group | Client Gender | ... | ICU | Intubated |
|----------|---------------------|-----------|---------------|-----|-----|-----------|
| Outcome2 | | | | | | |
| 0 | 1 | 4 | 1 | ... | 0 | 0 |
| 1 | 1 | 1 | 1 | ... | 0 | 0 |
| 1 | 1 | 5 | 0 | ... | 1 | 1 |
| 2 | 0 | 4 | 0 | ... | 0 | 0 |
| 3 | 1 | 2 | 0 | ... | 0 | 0 |
| 4 | | | | | | |
| 1 | | | | | | |

[5 rows x 15 columns]

Figure26. Result of the encode

3.5.2 Check again for useless fields

The Hospitalized attribute is generated based on the currently hospitalized and the 'ever hospitalized' attributes using the derive node. The Intubated attribute is generated based on the currently intubated and the ever the currently hospitalized and the 'ever hospitalized' attributes using the derive node. Hospitalized, currently in ICU and, 'ever in ICU' are similar and repetitive attributes. 'Outbreak associated' has been included in the 'source of infection' attribute. To avoid duplication, it can be removed. So we need to filter out duplicate properties again to reduce interference.

```
frame=frame.drop(['Currently Hospitalized','Outbreak Associated',  
                  'Currently in ICU','Currently Intubated',  
                  'Ever Hospitalized','Ever in ICU',  
                  'Ever Intubated','ICU','Outcome'],axis=1)  
print("column_name:",frame.columns)
```

Figure27.Filter useless fields

```
column_name: Index(['Age Group', 'Client Gender', 'Source of  
Infection', 'Hospitalized',  
                  'Intubated', 'Outcome2'],  
                  dtype='object')
```

Figure28.Result of filter useless fields

4.1 Reduce the Data

4.1.1 Feature selection

(1) Select features related to the predictor variable Outcome2. Use this to reduce data vertically

(2) The result is shown below, the attributes including Age Group, Hospitalized, Source of Infection and Client Gender are shown as important

(3) Intubated is shown as a coefficient of variation below the threshold, so reducing attribute Intubated.

```
# 4.1 Reduce the Data
X = frame.iloc[:,0:4]
y = frame.iloc[:, -1]
from sklearn.ensemble import ExtraTreesClassifier
import matplotlib.pyplot as plt
model = ExtraTreesClassifier()
model.fit(X,y)
#use inbuilt class feature importances of tree based classifiers
print(model.feature_importances_)
#plot graph of feature importances for better visualization
feat_importances = pd.Series(model.feature_importances_, index=X.columns)
feat_importances.nlargest(4).plot(kind='barh')
plt.show()

frame=frame.drop(['Intubated'],axis=1)
print("column_name:",frame.columns)
```

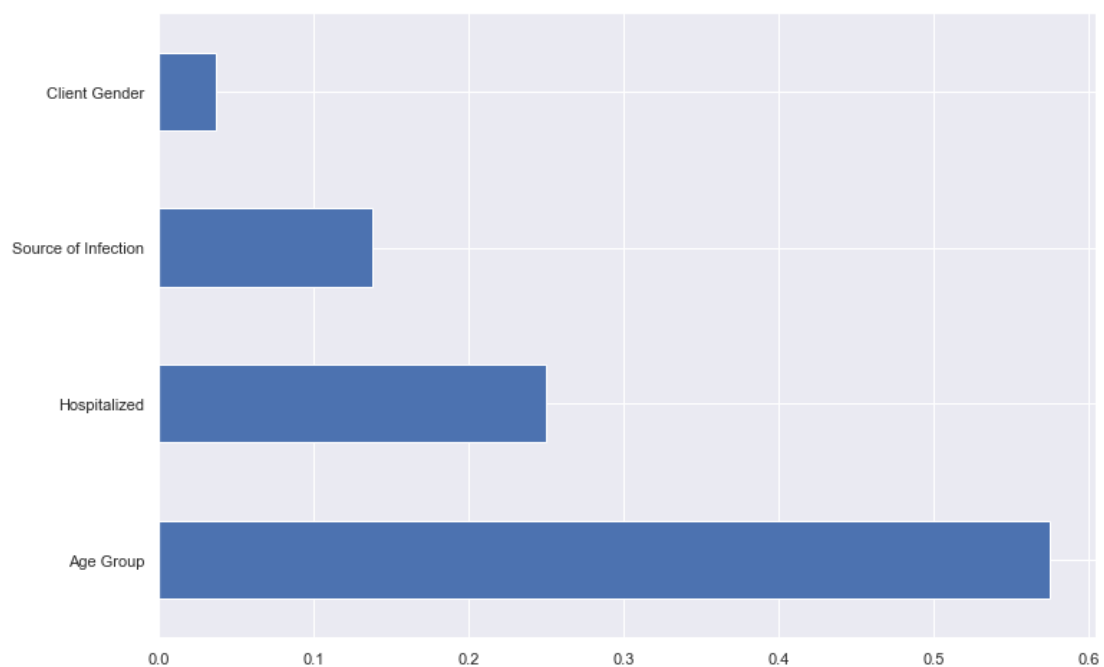


Figure29. Feature selection

4.1.2 Reduce unimportant attribute

Use the Outcome2 model to generate a filter which help to delete unimportant variation 'Intubated' selected during feature selection

```
frame=frame.drop(['Intubated'],axis=1)
print("column_name:",frame.columns)
```

```
column_name: Index(['Age Group', 'Client Gender',
'Source of Infection', 'Hospitalized',
'Outcome2'],
dtype='object')
```

Figure30. reduce unimportant attribute

4.2 Project the Data

4.2.1 Balance the Data

- (1) Select the target attribute and display its distribution. We want to predict the 'Outcome2', so it should be our 'target' rather than part of 'data'. The results are shown in the following figure. In the Outcome attribute, value 'FATAL' accounts for 7.52%, value 'NONFATAL' accounts for 92.48 %.

```
print(target.value_counts())
print("-----")
print(target.value_counts(normalize=True))
```

```
1    13790
0     1121
Name: Outcome2, dtype: int64
-----
1    0.924821
0    0.075179
Name: Outcome2, dtype: float64
```

Figure31. Distribution of target attribute

(2) Observe the data distribution of the target attribute;
balanced data is shown below. value 'FATAL' accounts for 50.4%,
value 'NONFATAL' accounts for 49.6 %. Data balance is
achieved.

```
df1=frame[frame['Outcome2']==1]
df0=frame[frame['Outcome2']==0]

df2=df1.sample(frac=0.08)

df_new=pd.concat([df0,df2])

x=df_new.iloc[:,1:-1]
y=df_new["Outcome2"]

print(y.value_counts())
print("-----")
print(y.value_counts(normalize=True))
```

```
0    1121
1    1103
Name: Outcome2, dtype: int64
-----
0    0.504047
1    0.495953
Name: Outcome2, dtype: float64
```

*Figure 32. Distribution of target attribute after the
balance*

4.2.2 Distribution of target attribute

- (1) first, separate our features and target variables. So, in the first line of the code above, we selected only the columns which do not match 'Outcome2' and assigned them to a variable 'cols'. Next, we created a new data frame with the columns in the list cols. This

will serve as our feature set. Then we took the 'Outcome2' column from the data frame and created a new data frame target.

```
# 4.2 Project the Data
#select columns other than 'Outcome2'
cols = [col for col in frame.columns if col not in ['Outcome2']]

#dropping the 'Outcome2'column
data = frame[cols]

#assigning the Oppurtunity Result column as target
target = frame['Outcome2']

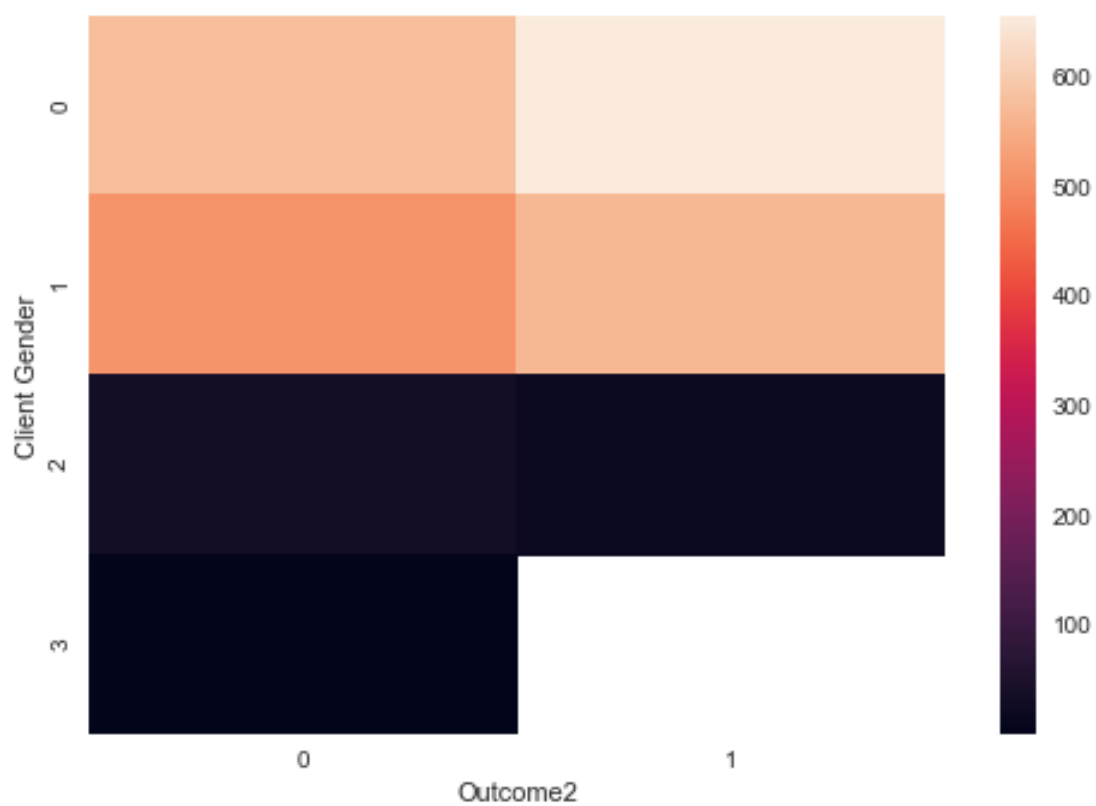
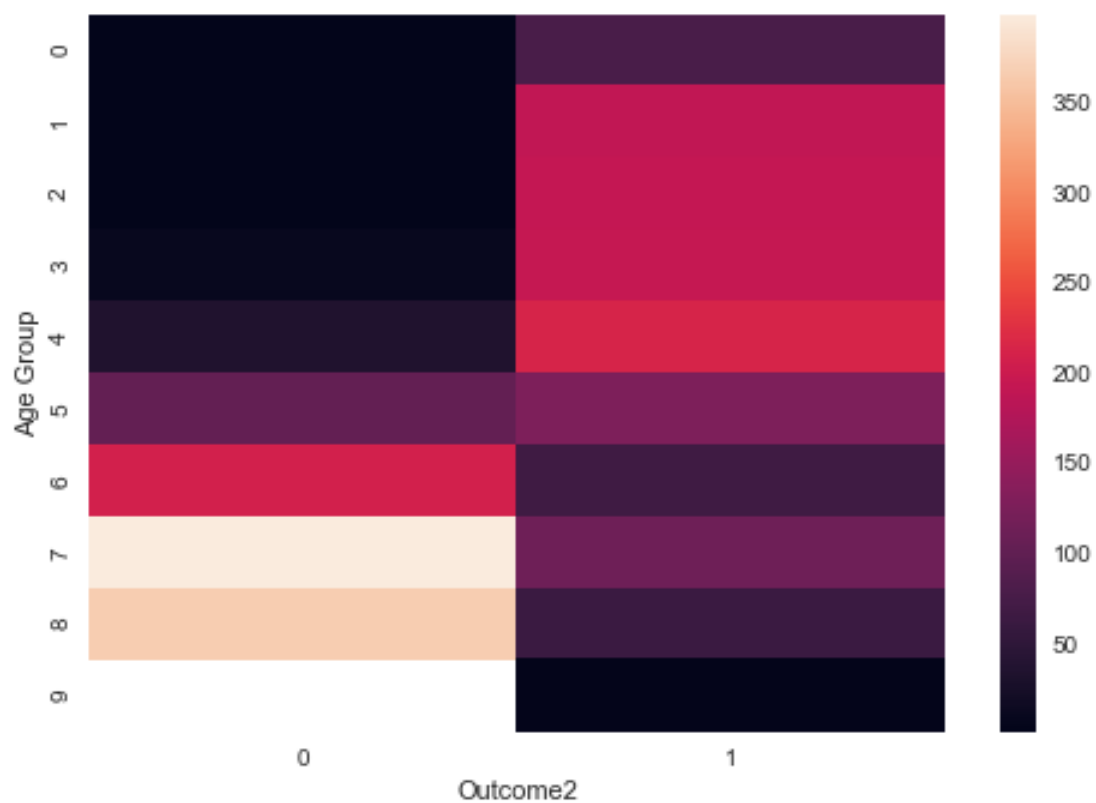
print(data.head(n=2))
```

Figure33. separate our features and target variables

| | Age Group | Client Gender | Source of Infection | Hospitalized |
|---|-----------|---------------|---------------------|--------------|
| 0 | 4 | 1 | 3 | 0 |
| 1 | 1 | 1 | 1 | 1 |

Figure34. new data frame

(2)The data correlation is obtained by creating a heat map to see the correlation between the attributes. The results show that Age Group , Client Gender, Source of Infection and Hospitalized have a correlation with Outcome2. However, the correlations of Source of Infection and Hospitalized are not as strong as age group, client gender.



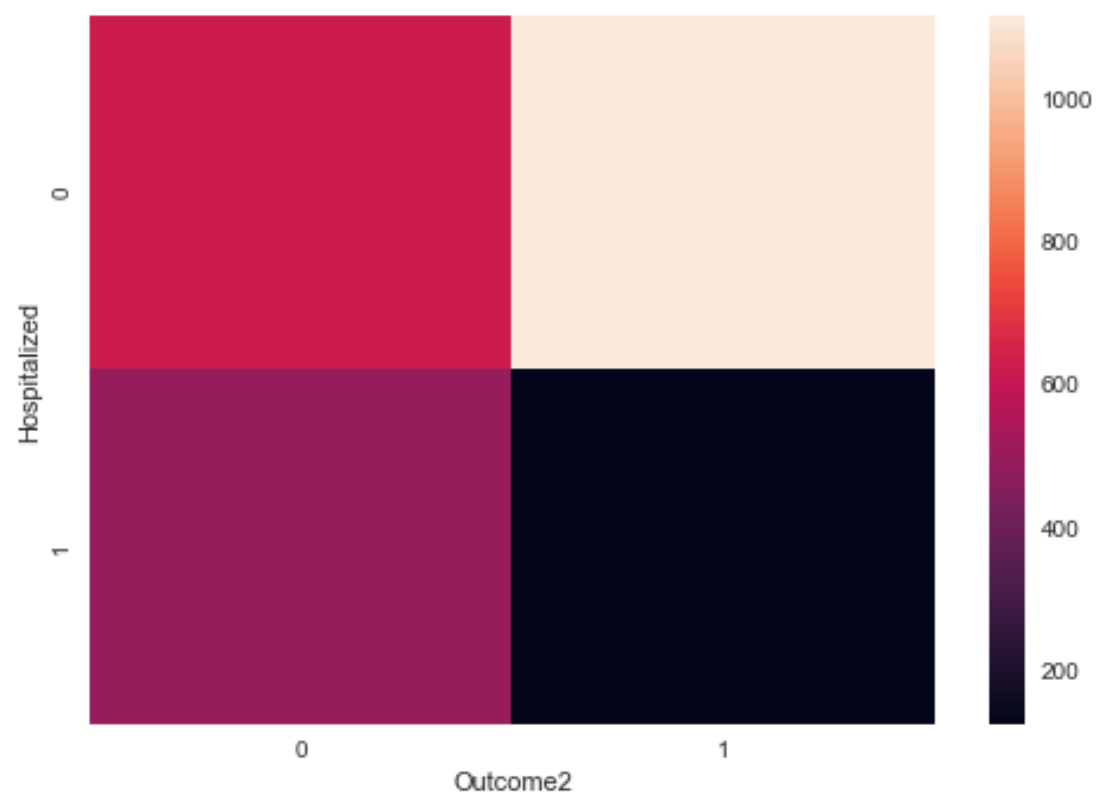
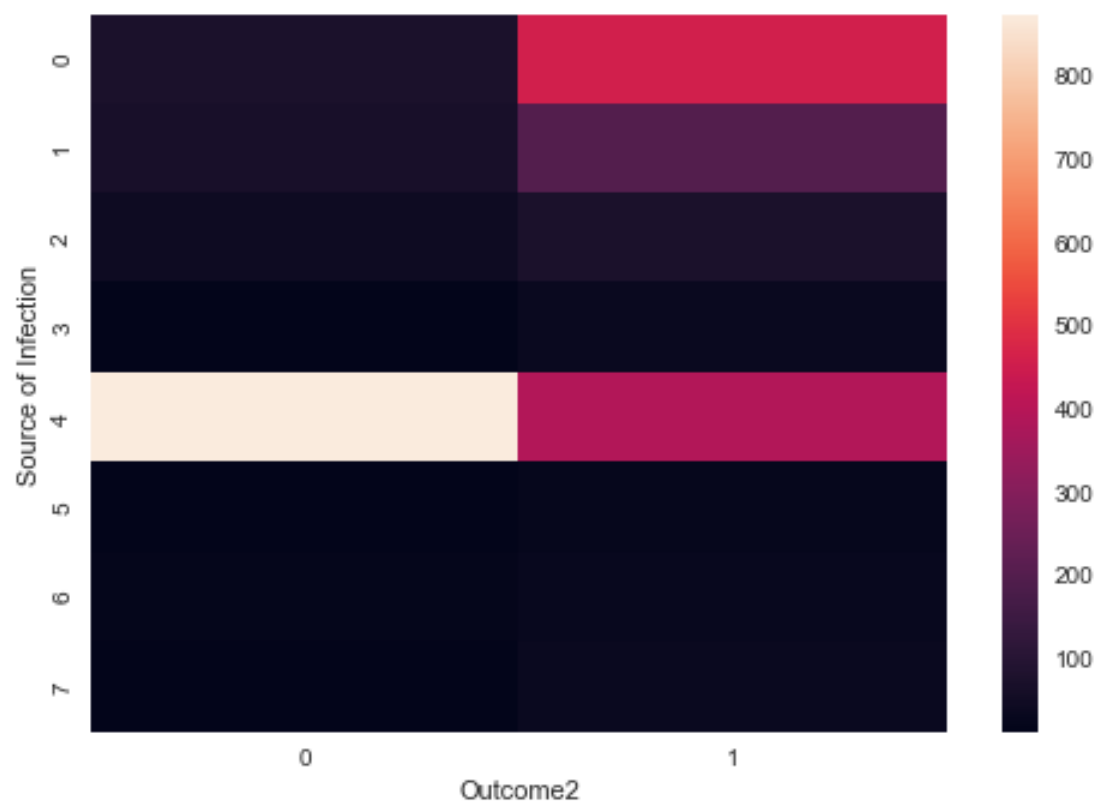


Figure34.data correlation

5.1 Match and discuss the objectives of data mining to data mining methods

5.1.1 Supervised and Unsupervised Learning

(1) How to select learning Method

Whether the data set use labeled data

(2) Supervised learning

The supervisory model uses the values of one or more input fields to predict the values of one or more output or target fields. Some examples are:

decision trees (C &R trees, CHAID, QUEST and C5.0 algorithms), regression (linear, logical, generalized linear and Cox regression algorithms), neural networks, support vector machines, and Bayesian networks.

Oversight models help organizations predict known outcomes, such as whether customers will buy or leave, or whether transactions conform to known fraud patterns. Model technology includes machine learning, rule induction, subgroup identification, statistical method and multi model generation.

(3) Unsupervised learning

Unsupervised learning groups data into unlabeled data sets according to potential hidden features. Because there is no label, the results cannot be evaluated (this is the key difference between supervised learning algorithms). By grouping the data by unsupervised learning, you can learn about some raw data that may not be visible in other situations. In high-dimensional data sets or large data sets, this problem is more obvious (Jones,2017).

5.1.2 Classification, Association and Segmentation

(1) Classification: Find the common characteristics of a group of data objects in the database and divide them into different categories according to the classification mode. The purpose is to map the data items in the database to a given category through the classification model (Data for Data Mining, n.d.).

(2) Association: The association model finds a pattern in the data in which one or more entities (such as events, purchases, or properties) are associated with one or more other entities. The model builds the set of rules that define these relationships. Here, the fields in the data can be either input or target. You can find these associations manually, but the association rules algorithm finds them faster and explores more complex patterns. (IBM Knowledge Center, n.d.).

(3) Segmentation: The segmentation model divides the data into segments or clusters with similar input field patterns. Because they are only interested in input fields, the segmentation model has no concept of output fields or target fields. Examples of segmentation models include Kohonen network, K-means clustering, two-step clustering and anomaly detection (IBM Knowledge Center, n.d.) Subdivision models (also known

as "clustering models") are useful when specific results are unknown (for example. The clustering model focuses on identifying groups of similar records and labeling them according to the group they belong to. This is done without prior knowledge of the group and its characteristics, and it distinguishes the clustering model from other modeling techniques because the model does not have predefined output or target fields for prediction (IBM Knowledge Center, n.d.).

5.2 Select the appropriate data-mining method (s) based on discussion

5.2.1 Choose supervised learning

Because there are both input variables and output variables (Outcome), so choose to use supervised learning methods. Because the prediction target is continuous, the regression method is appropriate.

5.2.2 Choose classification

Classification can be used for forecasting. The purpose of classification is to automatically derive the general description of given data from historical data records, so as to predict future data. Different from regression, the output of classification is discrete, while the output of regression is continuous.

Because the target attribute Outcome is categorical value, so we choose the classification model (IBM Knowledge Center, n.d.).

6.1 Conduct exploratory analysis and discuss

6.1.1 Algorithm discussion

(1) C & R tree

Requirements: The C & R tree model requires one or more input fields and one target field. The target field and input field can be continuous (numeric range) or classified. Fields set to all or none are ignored. Fields used in the model must fully instantiate their type, and any ordinal (ordered set) fields used in the model must have a numeric store (not a string) (IBM Knowledge Center, n.d.).

Strengths: C & R tree model is very powerful when there are problems such as data loss and a large number of fields. They usually don't need a long training time to estimate. In addition, the C & R tree model is easier to understand than some other model types - rules derived from the model have very simple explanations. Unlike C5.0, C & R trees can accommodate continuous and classified output fields(IBM Knowledge Center, n.d.).

(2) CHAID

Requirements: The target field and input field can be continuous or classified. Nodes can be divided into two or more subgroups at each level. Any ordinal field used in the model must have a numeric store (not a string). If necessary, you can use the reclassification node to transform it (IBM Knowledge Center, n.d.).

Strengths: Unlike C & R trees and quest nodes, CHAID can generate non

binary trees, which means that some splits have more than two branches. Therefore, it tends to create larger trees than the binary growth method. CHAID applies to all types of input and accepts case weight and frequency variables (IBM Knowledge Center, n.d.).

(3) QUEST

Requirements: The input fields can be contiguous (numeric range), but the target fields must be classified. All splits are binary. The weight field cannot be used. Any ordinal (ordered set) fields used in the model must have a numeric store (not a string). If necessary, you can use the Reclassification node to transform it(IBM Knowledge Center, n.d.).

Strengths: Quest uses a series of rules based on the important test to evaluate the input fields on the node. For selection purposes, you may only need to perform a test once on each input on the node. Unlike C & R trees, all splits are not checked; unlike C & R trees and CHAID, category combinations are not tested when evaluating input fields for selection. This can speed up the analysis (IBM Knowledge Center, n.d.).

(4) C 5.0

Requirements: To practice the C5.0 model, there must be a classified (i.e., nominal or ordered) target field and one or more input fields of any type. Fields set to all or none are ignored. Fields used in models must fully instantiate their types. You can also specify a weight field (IBM Knowledge Center, n.d.).

Strengths: The C5.0 model is very powerful when there are problems such as data loss and a large number of input fields. They usually don't

need a long training time to estimate. In addition, because the rules derived from this model have very direct interpretations, the C5.0 model is easier to understand than some other model types. C5.0 also provides a powerful enhancement method to improve the accuracy of classification (IBM Knowledge Center, n.d.).

(5) Bayesian Network

Requirements: The target field must be classified and can have nominal, *Ordinal* or tagged measurement levels. The input can be any type of field. Continuous (numeric range) input fields are automatically discarded; however, if the distribution is skewed, you can get better results by manually binding fields with a binding node before the Bayesian network node.

Strengths: It helps you understand causality. As a result, it enables you to understand the problem area and predict the consequences of any intervention. The network provides an effective method to avoid data over fitting. It is easy to observe a clear visualization of the relationships involved (IBM Knowledge Center, n.d.).

6.2 Select data-mining algorithms based on discussion

6.2.1 Algorithm requirements :

(1) Objective: Find out the relationship between the target attribute and the predictors, and use it to predict the probability of the target variable

occurring.

(2) Target attribute requirement: flag.

(3) Data type requirements: nominal and flag(string)

(4) Result requirements: prefer model results which are easy to demonstrate

| requirements model | Objective | Target attribute | Input data type | Result |
|-----------------------|--|------------------|------------------|---------------------|
| | relationship between the target and predictors | flag | nominal and flag | easy to demonstrate |
| CHAID | √ | × | √ | √ |

| | | | | |
|---------------------|---|---|---|---|
| QUEST | √ | × | √ | √ |
| C 5.0 | √ | √ | √ | √ |
| Bayesian Network | √ | √ | √ | √ |
| C & R tree | √ | × | √ | √ |

Table 3

6.2.2 Select data-mining algorithms

Through the above analysis, considering the data characteristics of this data set, which is the some attributes contain more than two variables. So we need to use the Multiple Classification method.

In addition, in order to use scikit learn in combination, the scheme sheet helps select the model that can be used by checking the following

requirements. Naive Bayes and Linear SVC.

- More than 50 samples – Check
- predicting a category – Check
- labeled data – Check
- Less than 100k samples – Check

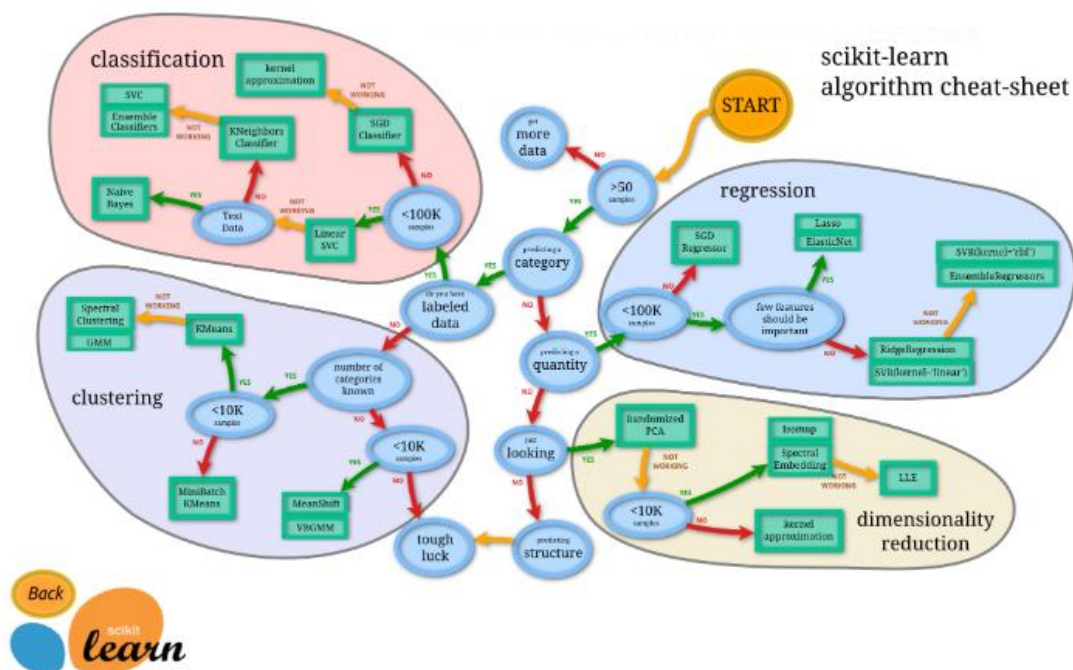


Figure34. Choosing the right estimator

(Pedregosa et al., 2011)

6.3 Select appropriate model(s) and choose relevant parameter(s)

(1) Build Bayesian Network model

Scikitlearn provides a set of classification algorithms, assuming that each pair of features in the dataset is independent. This

assumption is the basic principle of Bayes theorem. Naive Bayes algorithm calculates the probability of the feature connecting with the target variable, and then selects the feature with the highest probability. Now let's novel coronavirus pneumonia algorithm in scikit, and create our prediction model: we have a new set of data for new crown pneumonia cases, which will be our feature set. The probability of "Outcome2" will be our goal.

```
#import the necessary module
from sklearn.naive_bayes import GaussianNB
from sklearn.metrics import accuracy_score

#create an object of the type GaussianNB
gnb = GaussianNB()

#train the algorithm on training data and predict using the testing data
pred = gnb.fit(data_train, target_train).predict(data_test)

# print(pred.tolist())

#print the accuracy score of the model
print("Naive-Bayes accuracy : ",accuracy_score(target_test, pred, normalize = True))
```

Figure 35. Bayesian Network model

(2) Build LinearSVC model

First, the required modules are imported. Then we create an object SVC of type linearsvc_ model, where random_ State is 0. Random_ State is placed in the instructions of the built-in random number generator to scramble the data in a specific order. Next, we train the linear SVC on the training data, and then use the test data to predict the target. Finally, we use accuracy_ score () method was used to check the accuracy score.

```
#import the necessary modules
from sklearn.svm import LinearSVC
from sklearn.metrics import accuracy_score

#create an object of type LinearSVC
svc_model = LinearSVC(random_state=0)

#train the algorithm on training data and predict using the testing data
pred = svc_model.fit(data_train, target_train).predict(data_test)

#print the accuracy score of the model
print("LinearSVC accuracy : ",accuracy_score(target_test, pred, normalize = True))
```

Figure 36. LinearSVC model

7.1 Logical test designs

7.1.1 create logical test design

(1) When an algorithm use a limited data set to search for the best parameters for a particular model, it may model not only the general pattern in the data, but also any noise specific to that data set, resulting in poor performance of the model on the test data, that is, overfitting (Usama, Gregory & Padhraic,1996). So we chose cross-validation to solve this problem

(2) Due to the small number of data samples and the need for sufficient data to test results, we set 70% training set and 30% testing set. Use partition node to divide the data into 70% training data and 30% test data

```
#import the necessary module
from sklearn.model_selection import train_test_split

#split data set into train and test sets
data_train, data_test, target_train, target_test = train_test_split(data,target,
                                                                    test_size = 0.30, random_state = 10)
```

Figure 37. training data and test data

7.2 Data mining must be conducted (the model must run).

1 Run Bayesian Network model

```
#import the necessary module
from sklearn.naive_bayes import GaussianNB
from sklearn.metrics import accuracy_score

#create an object of the type GaussianNB
gnb = GaussianNB()

#train the algorithm on training data and predict using the testing data
pred = gnb.fit(data_train, target_train).predict(data_test)

# print(pred.tolist())

#print the accuracy score of the model
print("Naive-Bayes accuracy : ",accuracy_score(target_test, pred, normalize = True))
```

Figure 38. Bayesian Network model

2. Build LinearSVC model

```
#import the necessary modules
from sklearn.svm import LinearSVC
from sklearn.metrics import accuracy_score

#create an object of type LinearSVC
svc_model = LinearSVC(random_state=0)

#train the algorithm on training data and predict using the testing data
pred = svc_model.fit(data_train, target_train).predict(data_test)

#print the accuracy score of the model
print("LinearSVC accuracy : ",accuracy_score(target_test, pred, normalize = True))
```

Figure 39. LinearSVC model

7.3 Search for patterns and document the model's output.

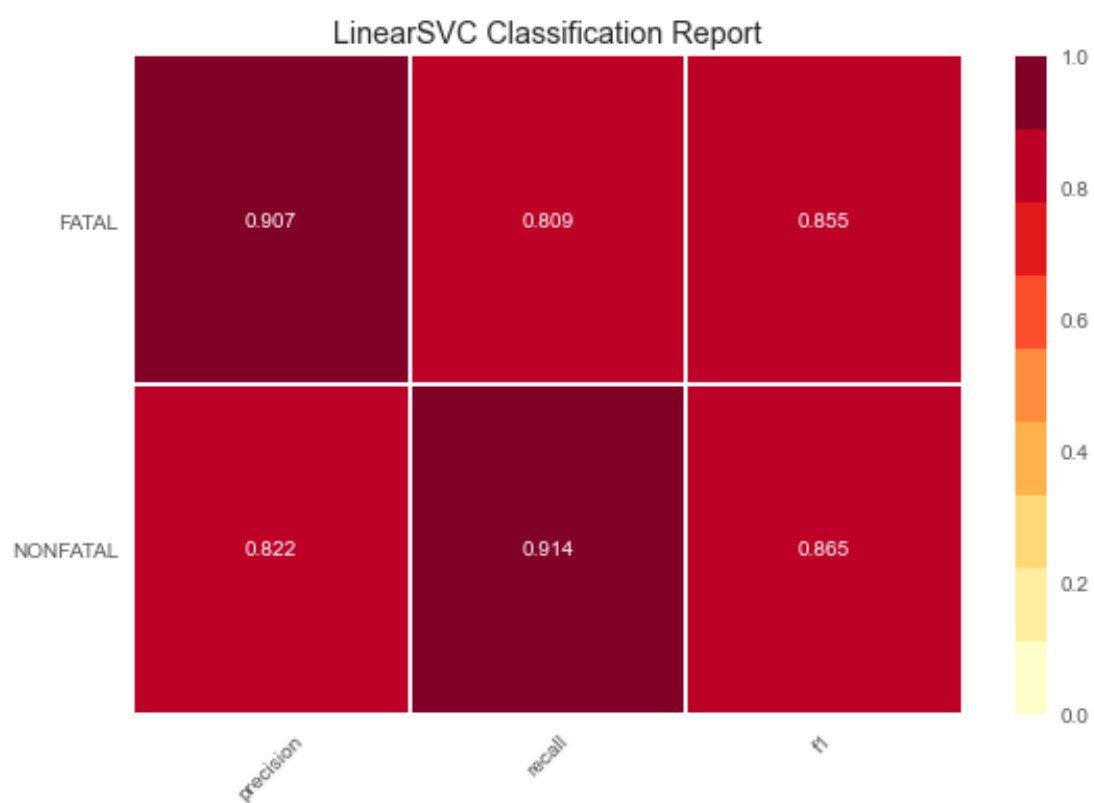
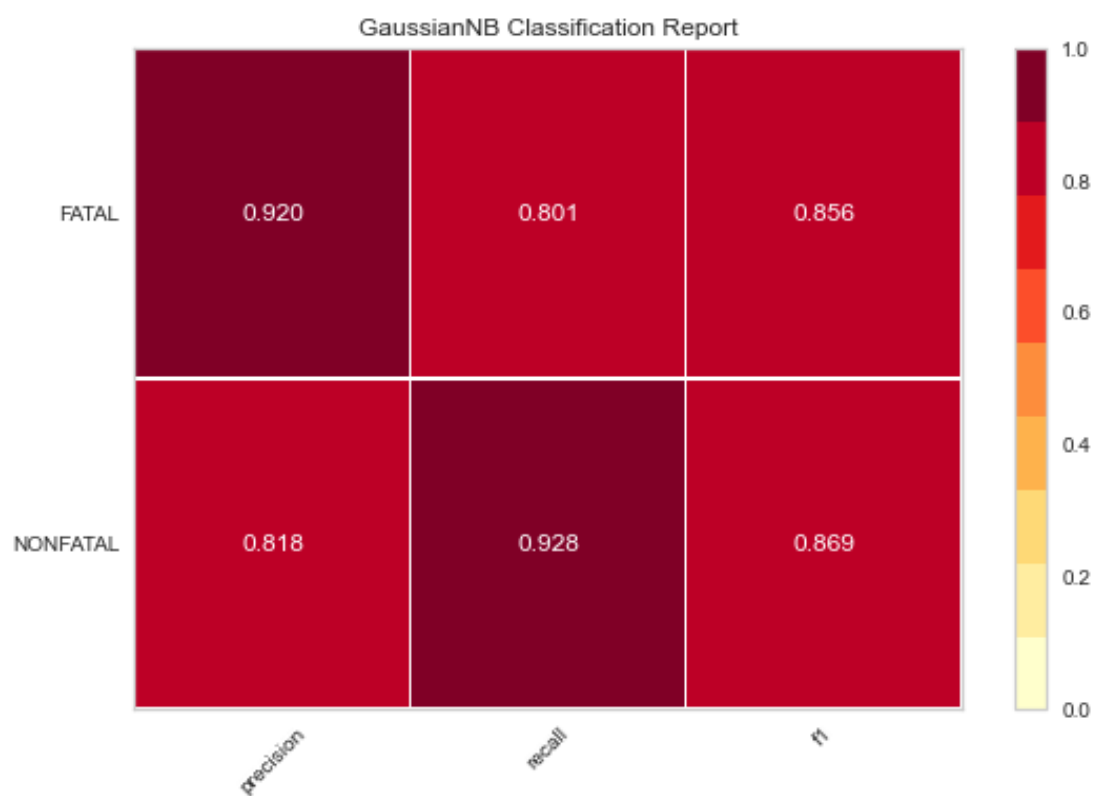
There are many patterns to choose from. In order to make the model as simple and understandable as possible, I only chose two models LinearSVC and Bayesian Network Model, both of which meet the requirements of the data set, cooperate so that I have many discoveries.

As can be seen from the figure below, the accuracy of this model is relatively high. For he prediction of the results is more accurate, parameter setting is also relatively simple, easy to understand. Novel coronavirus pneumonia is a very useful model for predicting the survival and mortality of new crown pneumonia cases.

model's

output

```
Naive-Bayes accuracy : 0.8617771509167842  
LinearSVC accuracy : 0.8603667136812412
```



8.1 Study and discuss the mined patterns.

Carry out an in-depth discussion about the data, results, models and patterns.

(1) data and result

The data range is only from the Toronto area, and the geographical scope is relatively small. In addition, due to the different cultural habits, living habits and protection measures for COVID-19, the infected population may be different. The outcome of the forecast may change depending on the location. The selection of data increases the spatial limitation of the prediction.

In addition, this data set mainly includes data from January 2020. In the months when the epidemic is very severe, data statistics may be affected by the epidemic, resulting in information loss or incorrect data input, which may increase the error.

(2) models and patterns

Object data dominate my data set, and the target value is a value in a flag format. Identify the factors that influence the mortality or survival of COVID-19 by distinguishing between the high mortality and the low mortality groups. Predict which populations are vulnerable, and help people increase their survival rates. So I chose the Bayesian Network Model and LinearSVC Model. These two models help us to identify which types of characteristics people are more likely to die, with an accuracy of more than 85%, which shows the performance of the models.

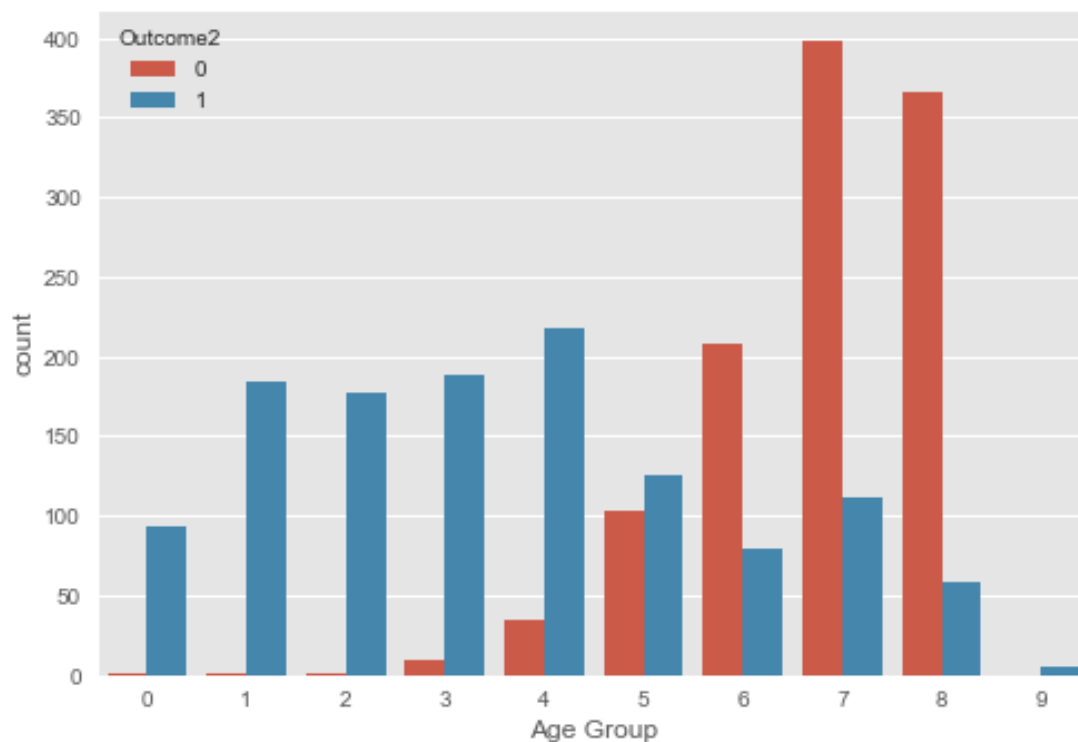
8.2 - Visualize the data, results, models and patterns in a clear and effective manner.

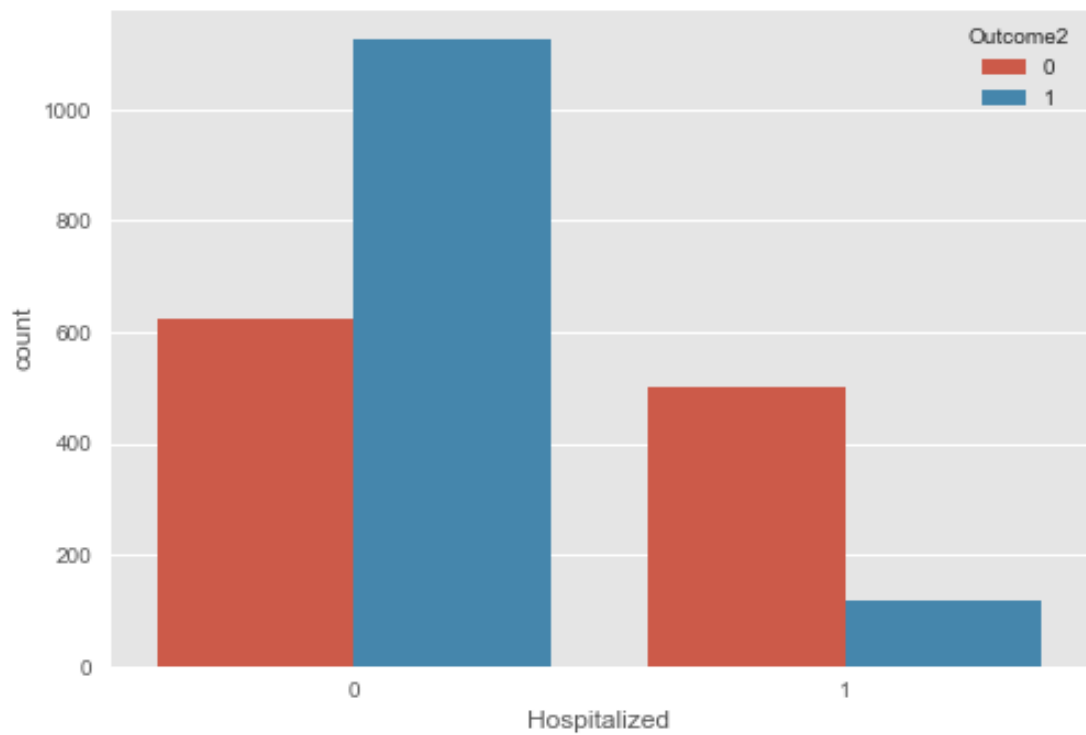
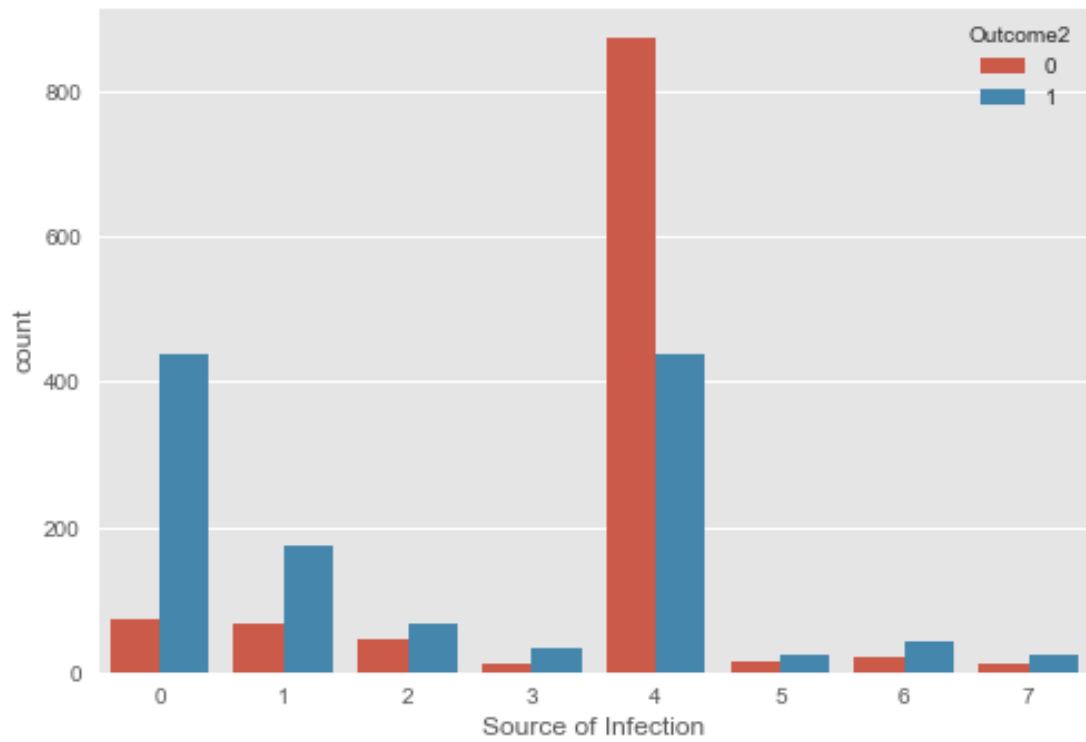
(1) Data distribution

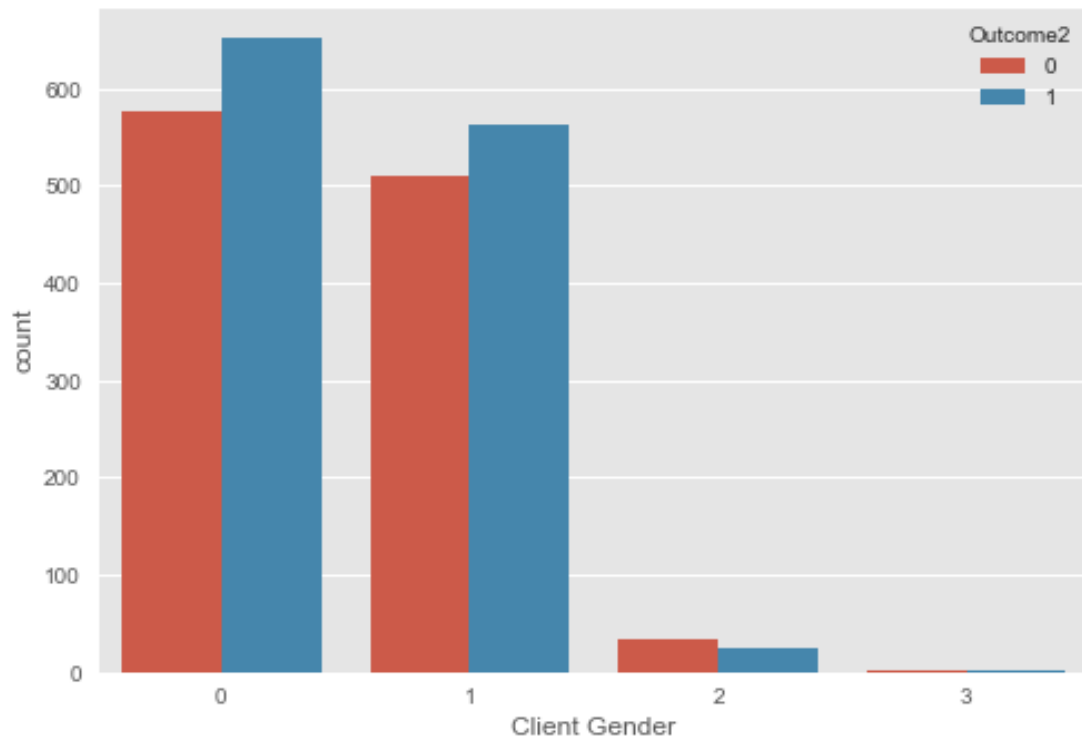
1= “ NONFATAL”

0= “ FATAL”

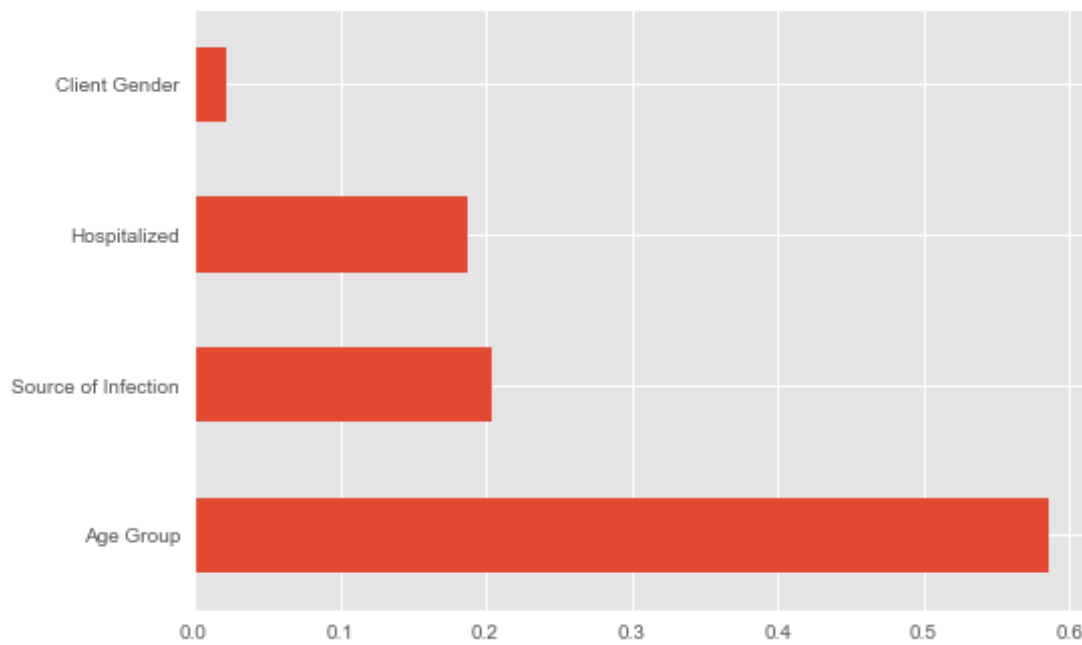
```
[14511 rows x 13 columns]
Age Group: ['50-59' '20-29' '60-69' '30-39' '19 and younger' '80-89' '70-79' '40-49'
nan '90+']
Age Group: [4 1 5 2 0 7 6 3 9 8]
Source of Infection: ['Institutional' 'Community' 'Travel' 'N/A - Outbreak associated'
'Close contact' 'Pending' 'Healthcare' 'Unknown/Missing']
Source of Infection: [3 1 6 4 0 5 2 7]
Hospitalized : ['No' 'Yes']
Hospitalized : [0 1]
Client Gender : ['MALE' 'FEMALE' 'OTHER' 'TRANSGENDER']
Client Gender : [1 0 2 3]
```



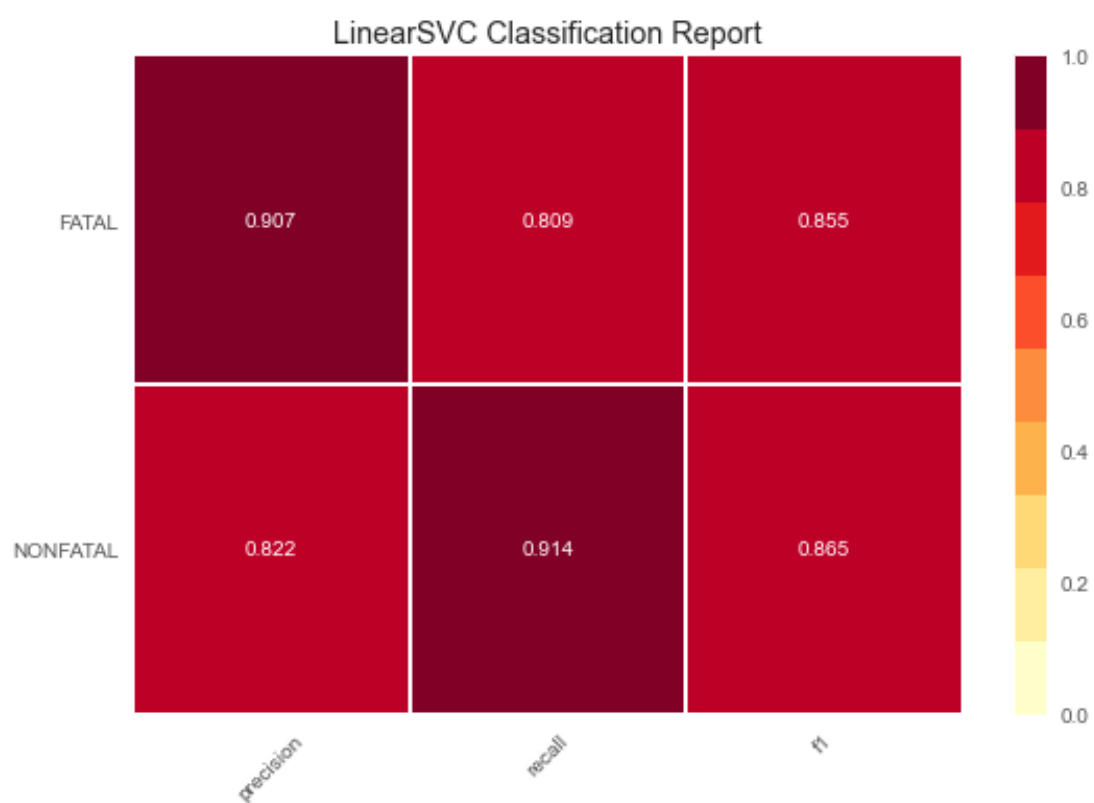
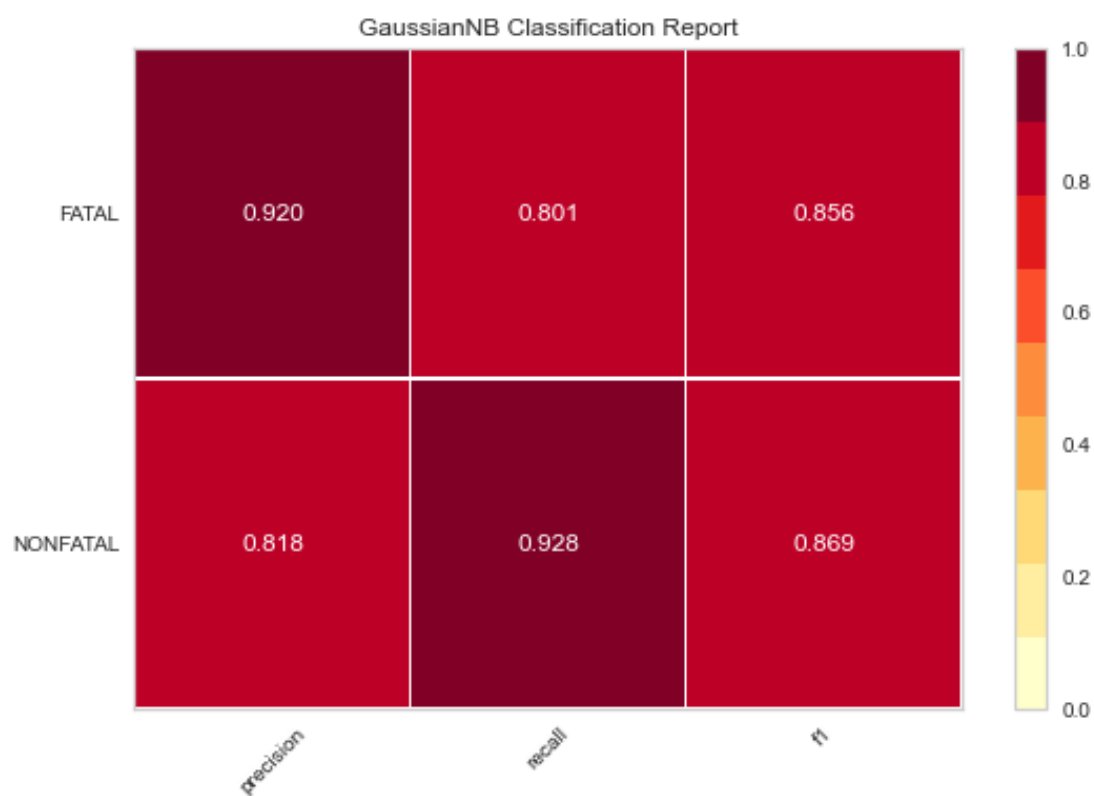




(2).Attribute importance ranking



(2) model performance



8.3 Interpret the results, models and patterns showing a clear understanding of the results.

Through the above steps, I have a clear understanding of the project's results, models, and patterns.

(1) Novel coronavirus pneumonia is the most important factor affecting the survival rate of patients with new crown pneumonia, and the number of deaths over 50 years old has even exceeded the survival rate from the results of 8.2. This suggests that we must pay attention to the key protection of the middle-aged and the elderly, which is very important.

Among the factors of infection source, N / a - outbreak associated had the greatest influence. This novel coronavirus pneumonia case is mainly from concentrated infection, or that the cases of concentrated infection are more likely to die. That's what we need to control the concentration of people to control the death toll.

(2) In the construction of the model and understanding of the model, I encountered many difficulties, especially in the application of Python statements, often can not achieve the effect I want. Therefore, I chose two relatively easy to operate models, Bayesian network model and linear SVC model. Through the evaluation results of the two models, we can find that the precision and recall of the two models are relatively high, so the F1 score is not low. This gives us more confidence in using this model to predict the results.

(3) Achieve business goals

According to the bar chart of predictor importance, the characteristics that affect suicide rate are Age Group, Source of Infection, Hospitalized.

The rule set can be used to guide the state, institutions or families to focus on protecting groups of older persons. These groups with lower survival rates are regularly checked and equipped with more effective protective measures.

8.4 - Assess and evaluate the results, models and patterns using the appropriate methods/processes.

8.4.1 assess the results, models and patterns

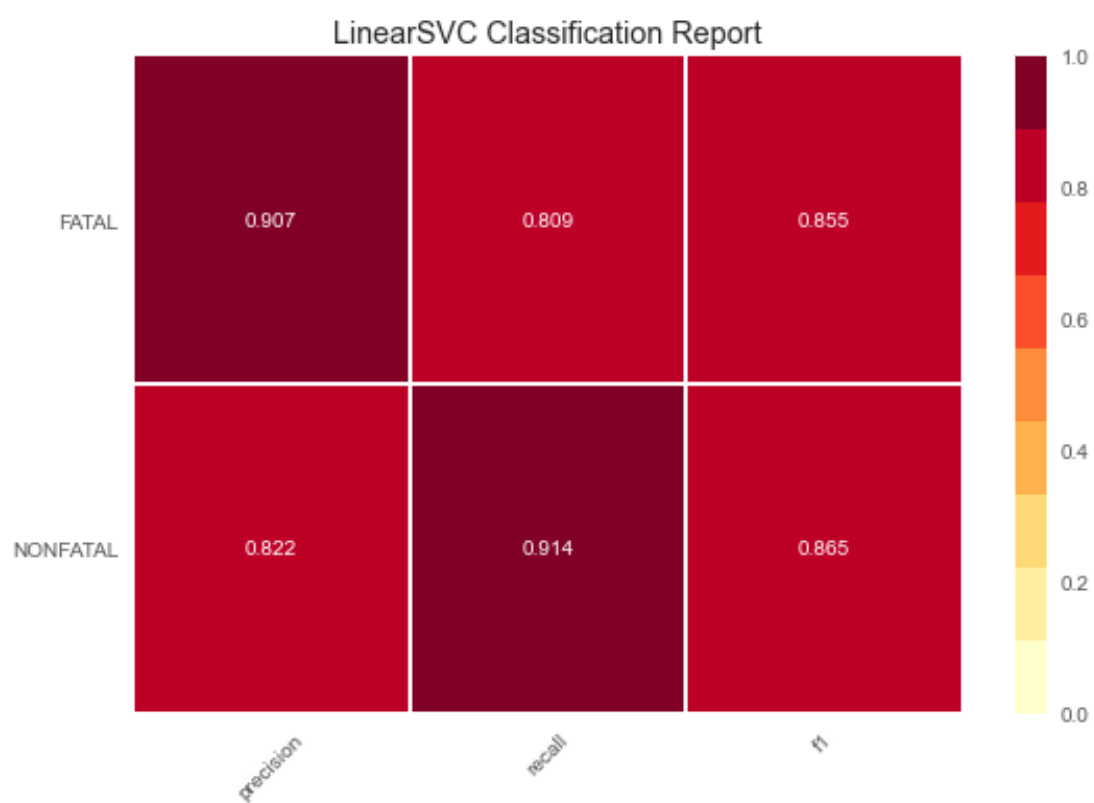
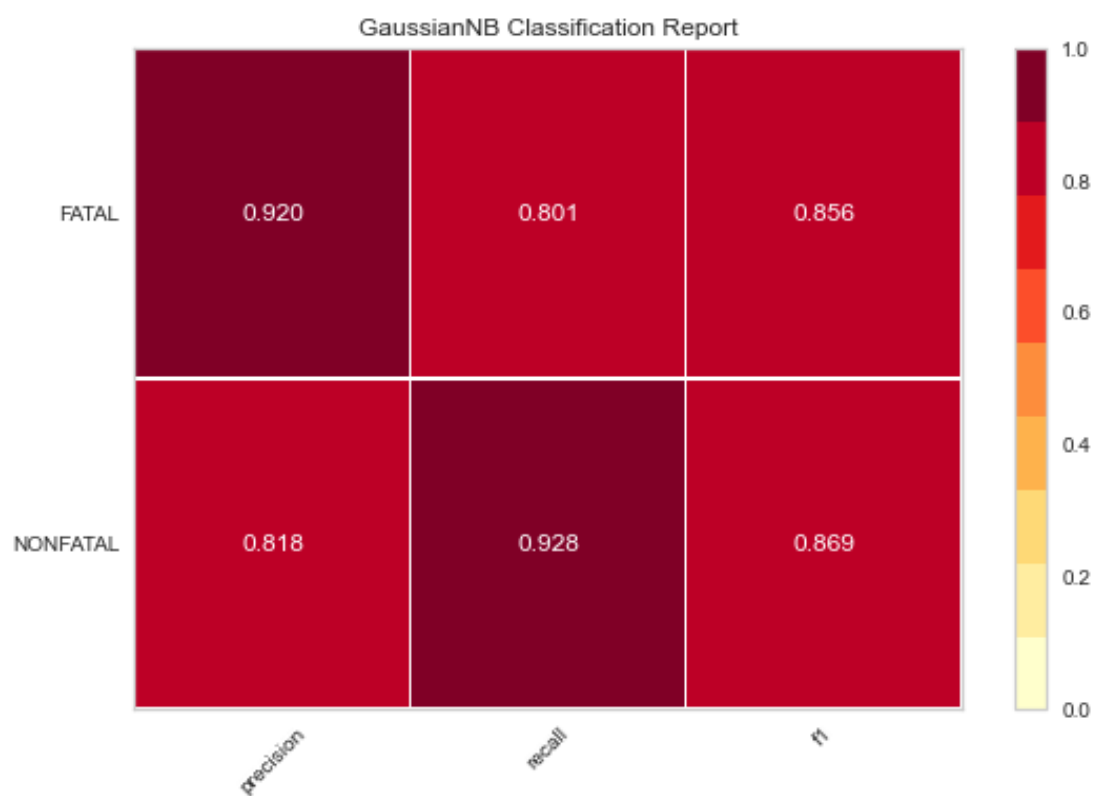
(1) Test Accuracy

Use the analysis node to evaluate the model, and the results are shown below. Both of the model have a high accuracy rate. The accuracy rate of LinearSVC model is 86%. The accuracy rate of Bayesian Network is 86.2%.

```
Naive-Bayes accuracy : 0.8617771509167842
LinearSVC accuracy : 0.8603667136812412
```

(2) Evaluation Graph

From this figure, we can see that the accuracy of the model is relatively high, which is more than 80%, and some of them are still up to 90%. It has certain credibility.



8.4.2 evaluation the results, models and patterns

(1) Business goals

People at high risk of COVID-19 can be protected by pre-locating and taking preventive measures.

(2) Data mining goals

Decision rules can be used to predict or classify a group of people with high COVID-19 survival rates, achieving the required model accuracy, but the data quality is poor and may affect the use of the model. The result of the model is logical. The results are easy to understand and easy to deploy

8.5 Iterate prior steps (1 – 7) as required

8.5.1 Business understanding

Covid-19 is one of the most difficult diseases in the world, and there are still tens of thousands of living cases every day in countries that are at risk of death every day. However, the vaccine development process still needs a lot of time. How to reduce the increase of cases in the meantime, who should we care about most. To find out who is more deserving of additional protection and special care, we used a dataset of COVID-19 cases in the Toronto area as a data source to analyze which factors contribute to the survival of infected persons and which factors are associated with case fatality.

8.5.2 data understanding

Data collection, preliminary cognition and processing of data, and certain cognition of data type, data size and processing method. There are assumptions and analysis.

8.5.3 Data preparation

Select data, clean data, merge data, deal with missing value, and format data. Prepare the data required for subsequent modelling.

8.5.4 Data transformation

To further simplify the data set and remove the unimportant and disturbing data, balance the data.

8.5.5 Data-mining method selection

The output variable and the input variable are both String values. In consideration of the target values and methods to predict, I choose supervisory learning and classification.

8.5.6 Data-mining algorithms selection

SPSS Modeler provides many reliable algorithms. After investigation, it was found that these algorithms have their own advantages and disadvantages, so I finally chose To use Land the Bayesian network model and LinearSVC model.

8.5.7 Data-mining

Through the construction of the model, let me have a feeling of suddenly enlightened. The model clearly shows the importance of the attributes and their relationship to the target value, shaping the prediction.

8.5.8 Interpretation

This step can also involve visualizing the extracted patterns and models, and we evaluate and evaluate the models, results, and their reliability.

Reference

Data for Data Mining. (n.d.). *Principles of Data Mining*, 11–21.

doi:10.1007/978-1-84628-766-4_1.

IBM Knowledge Center (n.d.) Overview of modeling nodes. Retrieved from

https://www.ibm.com/support/knowledgecenter/en/SS3RA7_18.1.0/modeler_mainhelp_client_ddita/clementine/modeling_nodes.html

Jones, M.(2017). Unsupervised learning for data classification. Retrieved from <https://developer.ibm.com/articles/ccunsupervised-learning-data-classification/>

Pal, S. (2018). Scikit-learn Tutorial: Machine Learning in Python. Retrieved from <https://www.dataquest.io/blog/sci-kit-learn-tutorial/>

Pedregosa *et al.*, (2011) . Scikit-learn: Machine Learning in Python. *JMLR* 12, pp. 2825-2830. Retrieved from https://scikit-learn.org/stable/tutorial/machine_learning_map/index.html

Usama, F., Gregory, P.-S., & Padhraic, S. (1996) Knowledge discovery and data mining: toward a unifying framework. In *Proceeding of The Second Int. Conference on Knowledge Discovery and Data Mining*, pages 82-88, 1996.

"I acknowledge that the submitted work is my own original work in accordance with the University of Auckland guidelines and policies on academic integrity and copyright.

(See: <https://www.auckland.ac.nz/en/students/forms-policies-and-guidelines/student-policies-and-guidelines/academic-integrity-copyright.html>Links to an external site.).

I also acknowledge that I have appropriate permission to use the data that I have utilized in this project. (For example, if the data belongs to an organization and the data has not been published in the public domain then the data must be approved by the rights holder.) This includes permission to upload the data file to Canvas. The University of Auckland bears no responsibility for the student's misuse of data."