



COVID 19 REPORT

Danhong Tang

(200471292)

Introduction:

This database is developed by Toronto Public Health to report an ongoing COVID-19 outbreak. This data set contains demographic, geographic, and hospitalization information for all cases reported by Toronto Public Health since January 2020. The data are extracted from the provincial Case & Contact Management System (CCM). The data will be completely refreshed and overwritten on a weekly basis.

Background

Ontario health officials are reporting 3,301 new cases of COVID-19 on 18th Dec. The Ontario government announced new COVID-19 restrictions on 17th Dec due to the unprecedented spread of the Omicron variant. The province's hospital system is dealing with the overloaded Omicron patient for bracing for an overwhelming wave of COVID-19 patients with diminished staffing levels and limited medical resource. So a simulator that can predict number of patients that may possibly enter into a hospital and fatalities will help hospitals to prepare and allocate their resources.

Objective

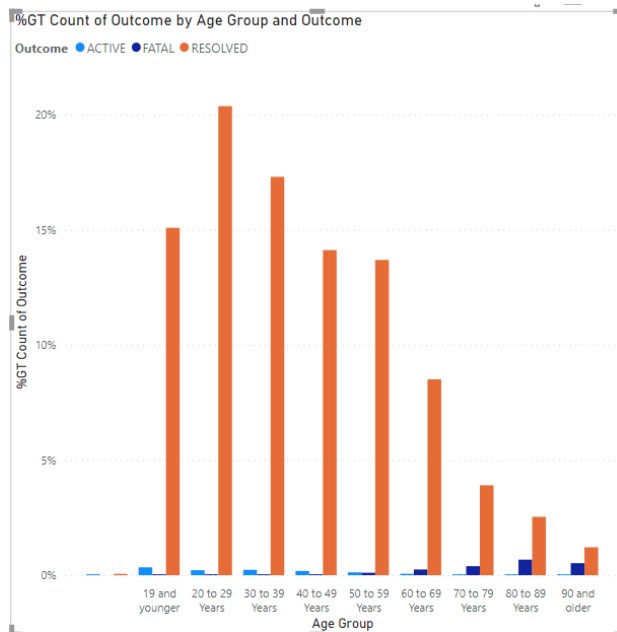
Build up an algorithm model to evaluate and predict number of patients that may possibly enter into a hospital and fatalities of them.

Tool

- Use Spark on Google Cloud Platform to build up evaluation and prediction algorithm model.
- Use PowerBI to finish visualization.

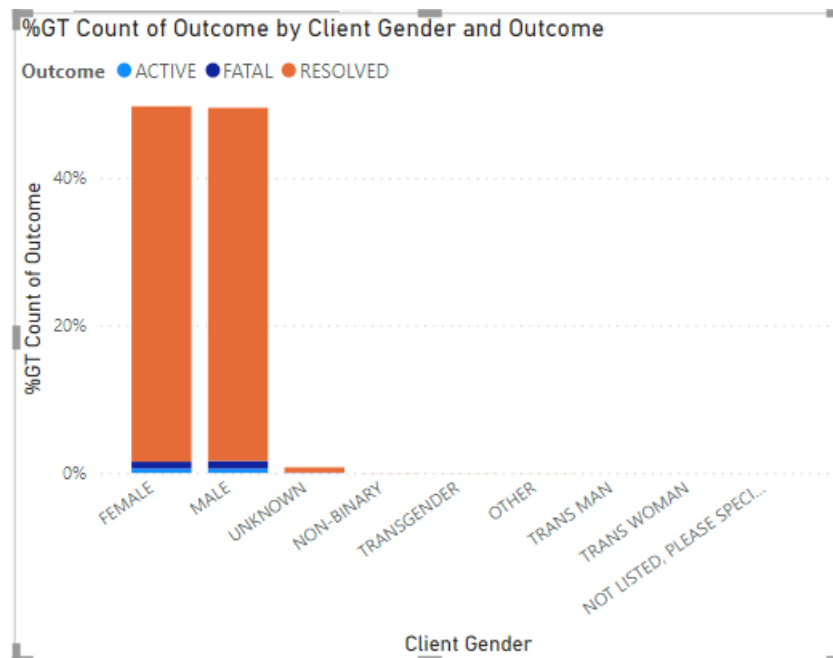
Visualization and Analysis

- Age group & Outcome Chart



- It can be seen from the chart that age from 20-29 has the highest infection rate among all age group, which might be caused by they are working and have frequent social activities.
- The fatal rate generally increases with age and group of 80-89 years has the highest fatal rate 0.68%.

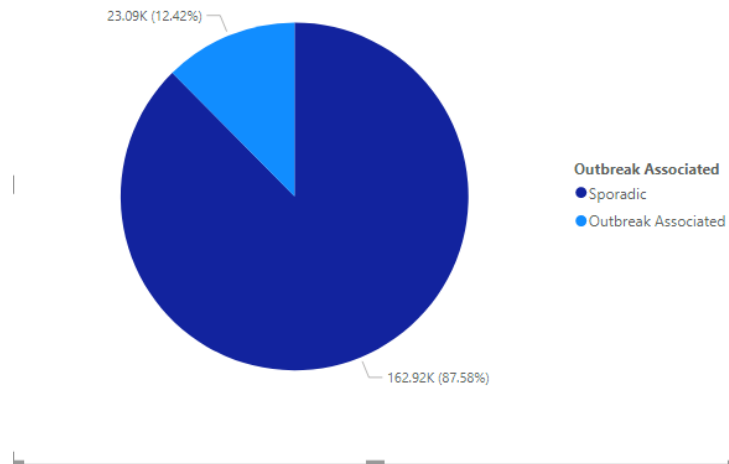
➤ Gentle & Outcome Chart



- The male and female patients are almost the same with each other, but male fatal rate 1.04% is slightly higher than female fatal rate 0.93%.

➤ Outbreak Associated Chart

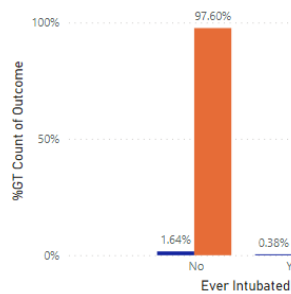
Count of _id by Outbreak Associated



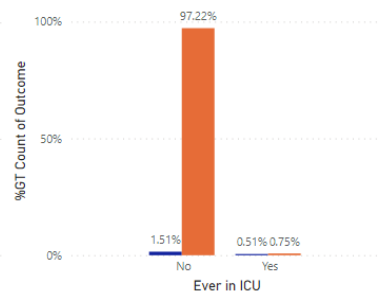
- We can see that sporadic cases (87.58%) is the most common outbreak associated type.

➤ Fatality Rate in 3 Situations Chart

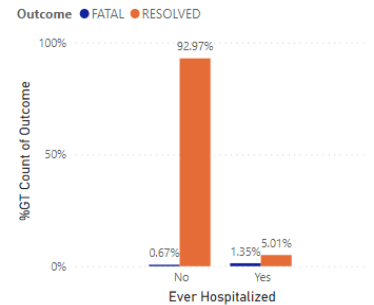
%GT Count of Outcome by Ever Intubated and Outcome
Outcome ● FATAL ● RESOLVED



%GT Count of Outcome by Ever in ICU and Outcome
Outcome ● FATAL ● RESOLVED



%GT Count of Outcome by Ever Hospitalized and Outcome
Outcome ● FATAL ● RESOLVED



	Among all cases%	Fatal among them%
Hospitalized	6.36%	21.23%
ICU	1.26%	40.48%
Intubated	0.75%	50.67%

- This meaningful comparison tells us that 6.36% cases are hospitalized, among them fatality rate 21.23%; Only 1.26% cases are sent to ICU, but the fatality rate increase to 40.48% ; Even worse, people who were intubated had almost half of them died.
- This analysis result is scary that it told people the best way to protect is to keep social distance and wear mask, and when getting serious symptoms need to be hospitalized, people have to face higher fatality rate.

- On the other side, China had much lower fatality rate among patients who were hospitalized among the same period of time, the reason of that might be there were much lower number of total cases so patients could get more medical resources.

Machine learning algorithm

➤ Preparation dataset(Appendix 1,2)

1. Download datasets
2. Index the string variables
3. Upload file to Google cloud
4. Launch Spark shell and import libraries
5. Load dataset
6. Select and cast relative columns into Int, also filter off "AVTIVE" cases in "outcome" that are not meaningful.

➤ Random Forest machine learning (Appendix 3,4,5)

1. Split our dataset into training and test data typical 80 20
2. Assemble features using vector assembler
3. Random forest: Create a new random forest object, give feature as a vector, give the label as "outcome"
4. Set up pipeline
5. Evaluate the model: using MulticlassClassificationEvaluator to compare "outcome" to the prediction column
6. hyper parameters:
 - MaxDepth is an array with two values, 5, 10 which is the limit on how deep we can construct the tree
 - Impurity which we give as entropy
7. Cross validate model
 - Cross validator will divide the training data set into 3
 - Each fold is coupled with the paramters for each type: Fold 1 is tried with max depth 3 and entropy and then fold 1 is again tried but this time with max depth 5 and entropy
 - the best model is picked
8. Training dataset, give the best model
9. Evaluate the model and print accuracy 0.98

Reference:

<https://toronto.ctvnews.ca/ontario-reports-3-301-new-covid-19-cases-four-additional-deaths-1.5712753>

Appendix

1

```
darhontang1@cluster-data3-m: ~ - Google Chrome
ssh.cloud.google.com/projects/engaged-reducer-329903/zones/us-central-1-a/instances/cluster-data3-m?authuser=0&hl=en_GB&projectNumber=790215681936&useAdminProxy=true&troubleshoot4005Ena...
jshoot255En...

scala> :paste
// Entering paste mode (ctrl-D to finish)

import org.apache.spark.sql.functions._
import org.apache.spark.sql.expressions.Window
import org.apache.spark.ml.feature.VectorAssembler, StringIndexer
import org.apache.spark.ml.Pipeline
import org.apache.spark.ml.classification.RandomForestClassificationModel, RandomForestClassifier
import org.apache.spark.ml.tuning.(CrossValidator, CrossValidatorModel, ParamGridBuilder)
import org.apache.spark.ml.evaluation.MulticlassClassificationEvaluator
import org.apache.spark.ml.param.ParamMap
import org.apache.spark.sql.types.{IntegerType, DoubleType}

import org.apache.spark.ml.regression.LinearRegression
import org.apache.spark.ml.evaluation.RegressionEvaluator
import org.apache.spark.ml.param.ParamMap
import org.apache.spark.ml.evaluation.BinaryClassificationEvaluator

val cov = spark.read
  .format("csv")
  .option("header", "true")
  .load("hdfs://10.128.0.6/BigData/Case2.csv")

val rawdata = cov.select(
  col("Gender_index").cast(IntegerType),
  col("Outbreak_index").cast(IntegerType),
  col("C_Hospitalized_index").cast(IntegerType),
  col("C_ICU_index").cast(IntegerType),
  col("C_Intubated_index").cast(IntegerType),
  col("Outcome").alias("outcome"))
  .filter(! (col("Outcome") === lit("ACTIVE")))

rawdata.show()

// Exiting paste mode, now interpreting.

+-----+-----+-----+-----+-----+
|Gender_index|Outbreak_index|C_Hospitalized_index|C_ICU_index|C_Intubated_index|outcome|
+-----+-----+-----+-----+-----+
|2|1|2|2|2|2 (RESOLVED)|
|1|1|2|2|2|2 (RESOLVED)|
|2|1|2|2|2|2 (RESOLVED)|
|2|1|2|2|2|2 (RESOLVED)|
|1|1|2|2|2|2 (RESOLVED)|
|1|1|2|2|2|2 (RESOLVED)|
|1|1|2|2|2|2 (RESOLVED)|
|1|1|2|2|2|2 (RESOLVED)|
|1|1|2|2|2|2 (RESOLVED)|
|1|1|2|2|2|2 (RESOLVED)|
|1|1|2|2|2|2 (RESOLVED)|
|1|1|2|2|2|2 (RESOLVED)|
+-----+-----+-----+-----+-----+
only showing top 20 rows
```

2

```
darhontang1@cluster-data3-m: ~ - Google Chrome
ssh.cloud.google.com/projects/engaged-reducer-329903/zones/us-central-1-a/instances/cluster-data3-m?authuser=0&hl=en_GB&projectNumber=790215681936&useAdminProxy=true&troubleshoot4005Ena...
jshoot255En...

import org.apache.spark.sql.expressions.Window
import org.apache.spark.ml.feature.VectorAssembler, StringIndexer
import org.apache.spark.ml.classification.RandomForestClassificationModel, RandomForestClassifier
import org.apache.spark.ml.tuning.(CrossValidator, CrossValidatorModel, ParamGridBuilder)
import org.apache.spark.ml.evaluation.MulticlassClassificationEvaluator
import org.apache.spark.ml.param.ParamMap
import org.apache.spark.sql.types.{IntegerType, DoubleType}
import org.apache.spark.ml.regression.LinearRegression
import org.apache.spark.ml.evaluation.RegressionEvaluator
import org.apache.spark.ml.param.ParamMap
import org.apache.spark.ml.evaluation.BinaryClassificationEvaluator

var data=rawdata.withColumn("outcome",when(col("outcome") === "RESOLVED",lit("1").cast(IntegerType)).otherwise(lit("2").cast(IntegerType)))
data: org.apache.spark.sql.Dataset = [Gender_index: int, Outbreak_index: int ... 4 more fields]

scala> data.show()

+-----+-----+-----+-----+-----+
|Gender_index|Outbreak_index|C_Hospitalized_index|C_ICU_index|C_Intubated_index|outcome|
+-----+-----+-----+-----+-----+
|1|1|2|2|2|1|
|2|1|2|2|2|1|
|2|1|2|2|2|1|
|1|1|2|2|2|1|
|1|1|2|2|2|1|
|1|1|2|2|2|1|
|1|1|2|2|2|1|
|1|1|2|2|2|1|
|1|1|2|2|2|1|
|1|1|2|2|2|1|
|1|1|2|2|2|1|
|1|1|2|2|2|1|
|1|1|2|2|2|1|
|1|1|2|2|2|1|
|1|1|2|2|2|1|
|1|1|2|2|2|1|
|1|1|2|2|2|1|
|1|1|2|2|2|1|
|1|1|2|2|2|1|
|1|1|2|2|2|1|
+-----+-----+-----+-----+-----+
only showing top 20 rows

scala> :paste
// Entering paste mode (ctrl-D to finish)
```

3

```
danhongtang1@cluster-data3-m: ~ - Google Chrome
ssh.cloud.google.com/projects/engaged-reducer-329903/zones/us-central1-a/instances/cluster-data3-m?authuser=0&hl=en_GB&projectNumber=790215681936&useAdminProxy=true&troubleshoot4005Enabled=true&troubleshoot255En...

at org.apache.spark.ml.evaluation.Evaluator.set(Evaluator.scala:28)
at org.apache.spark.ml.evaluation.MulticlassClassificationEvaluator.setMetricName(MulticlassClassificationEvaluator.scala:65)
... 58 elided

scala> val Array(trainingData, testData) = data.randomSplit(Array(0.8, 0.2), 754)
trainingData: org.apache.spark.sql.Dataset[org.apache.spark.sql.Row] = [Gender_index: int, Outbreak_index: int ...
4 more fields]
testData: org.apache.spark.sql.Dataset[org.apache.spark.sql.Row] = [Gender_index: int, Outbreak_index: int ... 4 mo
re fields]

scala> :paste
// Entering paste mode (ctrl-D to finish)

val assembler = new VectorAssembler()
.setInputCols(Array("Gender_index", "Outbreak_index", "C_Hospitalized_index", "C_ICU_index", "C_Intubated_index"))
.setOutputCol("assembled-features")

// Exiting paste mode, now interpreting.

assembler: org.apache.spark.ml.feature.VectorAssembler = VectorAssembler: uid=vecAssembler_e934611e383f, handleInva
lid-error, numInputCols=5

scala> :paste
// Entering paste mode (ctrl-D to finish)

val rf = new RandomForestClassifier()
.setFeaturesCol("assembled-features")
.setLabelCol("outcome")
.setSeed(1234)

// Exiting paste mode, now interpreting.

rf: org.apache.spark.ml.classification.RandomForestClassifier = rfc_5aad73e7d82c

scala> :paste
// Entering paste mode (ctrl-D to finish)

val pipeline = new Pipeline()
.setStages(Array(assembler, rf))

// Exiting paste mode, now interpreting.

pipeline: org.apache.spark.ml.Pipeline = pipeline_c6cf962f47ba

scala> :paste
// Entering paste mode (ctrl-D to finish)
```

4

```
danhongtang1@cluster-data3-m: ~ - Google Chrome
ssh.cloud.google.com/projects/engaged-reducer-329903/zones/us-central1-a/instances/cluster-data3-m?authuser=0&hl=en_GB&projectNumber=790215681936&useAdminProxy=true&troubleshoot4005Enabled=true&troubleshoot255En...

// Entering paste mode (ctrl-D to finish)

val evaluator = new MulticlassClassificationEvaluator()
.setLabelCol("outcome")
.setPredictionCol("prediction")
.setMetricName("accuracy")

// Exiting paste mode, now interpreting.

evaluator: org.apache.spark.ml.evaluation.MulticlassClassificationEvaluator = MulticlassClassificationEvaluator: ui
d=mcEval_7a6e4b09a7e, metricName=accuracy, metricLabel=0.0, beta=1.0, eps=1.0E-15

scala> :paste
// Entering paste mode (ctrl-D to finish)

val paramGrid = new ParamGridBuilder()
.addGrid(rf.maxDepth, Array(3, 5))
.addGrid(rf.impurity, Array("entropy", "gini")).build()

// Exiting paste mode, now interpreting.

paramGrid: Array[org.apache.spark.ml.param.ParamMap] =
Array((
  rfc_5aad73e7d82c-impurity: entropy,
  rfc_5aad73e7d82c-maxDepth: 3
), (
  rfc_5aad73e7d82c-impurity: gini,
  rfc_5aad73e7d82c-maxDepth: 3
), (
  rfc_5aad73e7d82c-impurity: entropy,
  rfc_5aad73e7d82c-maxDepth: 5
), (
  rfc_5aad73e7d82c-impurity: gini,
  rfc_5aad73e7d82c-maxDepth: 5
))

scala> :paste
// Entering paste mode (ctrl-D to finish)

val cross_validator = new CrossValidator()
.setEstimator(pipeline)
.setEvaluator(evaluator)
.setEstimatorParamMaps(paramGrid)
.setNumFolds(3)

// Exiting paste mode, now interpreting.

cross_validator: org.apache.spark.ml.tuning.CrossValidator = cv_62730ed22023

scala> :paste
```

5

```
scala> val Model = cross_validator.fit(trainingData)
Model: org.apache.spark.ml.tuning.CrossValidatorModel = CrossValidatorModel: uid=cv_62730ed22023, bestModel=pipeline_c6cf962f47ba, numFolds=3

scala> val predict = Model.transform(testData)
predict: org.apache.spark.sql.DataFrame = [Gender_index: int, Outbreak_index: int ... 8 more fields]

scala> val accuracy = evaluator.evaluate(predict)
accuracy: Double = 0.9802446933158941

scala> println("Accuracy based on test data = " + accuracy)
Accuracy based on test data = 0.9802446933158941

scala> |
```