# PROJECT

## Telstra Network Disruptions Competition

PREPARED BY:   TOMÁS QUINTEROS
PREPARED FOR:      POUYA YOUSEFI

06/06/17

# TABLE OF CONTENTS

# 1. INTRODUCTION

In their first recruiting competition that ended on 29th February 2016, the company Telstra, Australia's largest telecommunications network, is challenging Kagglers to predict the severity of service disruptions on their network. Telstra is on a journey to enhance the customer experience - ensuring everyone in the company is putting customers first. In terms of its expansive network, this means continuously advancing how it predicts the scope and timing of service disruptions.

This report is based on the Telstra challenge development for the BI: Data Analytics lecture of MoTIS 8 with Pouya Yousefi. It was developed on Dataiku DSS which is a collaborative data science software platform for teams of data scientists, data analysts, and engineers. Is very intuitive, easy-to-use, friendly and allows to explore, prototype, build, and deliver your own data products efficiently.

# 2. UNDERSTANDING THE PROBLEM



**Telstra Network Disruptions**

Predict service faults on Australia's largest telecommunications network

To begin with, Telstra is the largest telecommunication network in Australia providing services such as mobile phones, broadband, landline as well as digital TV. The problem is to find out whether if a disruption occurred is a momentary glitch or total interruption of connectivity. In other words, the task was to predict the severity of service disruptions on Telstra's network.

The accurate predictions of service disruptions will help Telstra serve customers better and enhance customer experience by providing better connectivity. Telstra wants to see how we would help it drive customer advocacy by developing a more advanced predictive model for service disruptions by using a dataset of features from their service logs.

# 3. DATA DESCRIPTIONS

The goal of the problem is to predict Telstra network's **fault severity at a time at a particular location** based on the log data available. Each row in the main dataset (train.csv, test.csv) represents a **location** and a **time point**. They are identified by the "id" column, which is the key "id" used in other data files.

Fault severity ("fault_severity" in train.csv) is a measurement of actual reported faults from users of the network, is the target variable of the problem and it has 3 categories:

> 0: meaning no fault,
> 1: meaning only a few,
> 2: meaning many.

This means it was a multi-classification problem (3-class).

Different types of features are extracted from log files and other sources:

> event_type.csv - event type related to the main dataset,
> log_feature.csv - features extracted from log files,
> resource_type.csv - type of resource related to the main dataset,
> severity_type.csv - severity type of a warning message coming from the log.

Note: "severity_type" is a feature extracted from the log files (in severity_type.csv). Often this is a severity type of a warning message coming from the log. "severity_type" is categorical. It does not have an ordering.

Other datasets:

> train.csv - the training set for fault severity,
> test.csv - the test set for fault severity,
> sample_submission.csv - a sample submission file in the correct format.

# 4. APPROACH

## 4.1. Train and Test

First of all, both the train and test files were prepared by adding a column "fault_severity" into the test schema and erasing the word "location" in the location column to clean and enrich both dataset. The type of certain columns was changed e.g., ID string to integer, this required updating recipes and datasets downstream. At the same time each dataset was copied to a MySQL database connection and then stacked concatenating both prepared datasets for analysis.



Image 1. Train and test files prepared and stacked.



Image 2. Train and test data.

## 4.2. Severity_type

Afterwards, the severity_type dataset was analyzed. It contains the same number of rows as train and test files combined, and for each id there is one value for severity_type. So, first the preparation was done (same procedure as with train and test datasets) and then the merge with train and test stacked. A left join on ID, keeping all rows of the left dataset (id and severity_type) and adding information from the right dataset (location and fault_severity).



Image 3. Severity_type data.
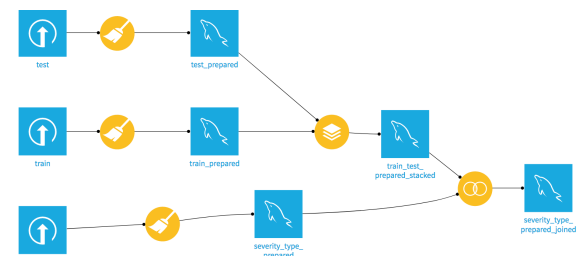


Image 4. Severity_type_joined data.



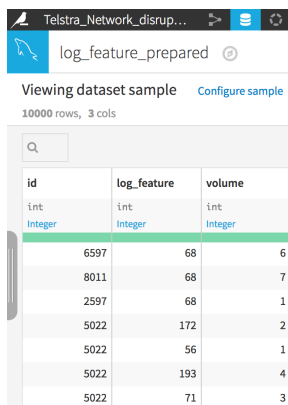Image 5. Train, test and severity_type joined.

## 4.3. Log_feature

The next step was to analyze and prepare log_feature. This dataset includes for each log_feature a volume column which specifies a quantity associated with each variable. A particular log can contain many different types of log features. And each feature might occur multiple times in the log. For example: ID=3952, log_feature='feature 82', volume='4'.

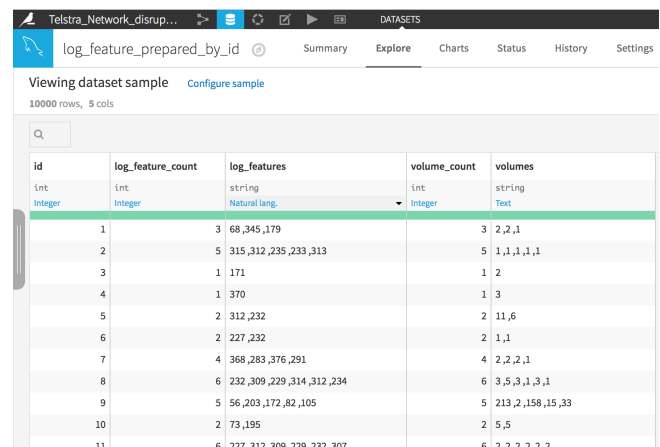The interpretation was: For the location & date associated with id 3952, feature 82 appeared in the log 4 times.

The next SQL recipe was created for count and concatenate log_features and volumes of each ID:

```
SELECT
    `id` AS `id`,
    COUNT(`log_feature`) AS `log_feature_count`,
    GROUP_CONCAT(log_feature,' ') AS `log_features`,
    COUNT(`volume`) AS `volume_count`,
    GROUP_CONCAT(volume,' ') AS `volumes`
 FROM `TELSTRA_NETWORK_DISRUPTIONS_1_log_feature_prepared`
  GROUP BY `id`
```



Image 6. Log_feature data



Image 7. Log_feature data grouped by id

In order not to mixed the analysis, the log_feature dataset was duplicated and prepared it by setting the log_feature and volume fields from string to integer to compute various aggregates of feature and volume logarithm by id (distinct, min, max, avg, sum). Count was not included due to its already in the first log_feature dataset.

Image 8. Log_feature_1 data grouped by id.

## 4.4. Resource_type and event_type

Like log_feature, this datasets contain more than one value (resource_type and event_type) for each ID. So the same procedure will be applied for both. First the classic preparation and connection to the database, then an SQL recipe to count and concatenate the values for each ID and finally the dataset duplication for changing the type, group and compute new values by id.

From severity_type_prepared_joined (joined with train and test) a left join with all the new datasets and the flow looks like this:
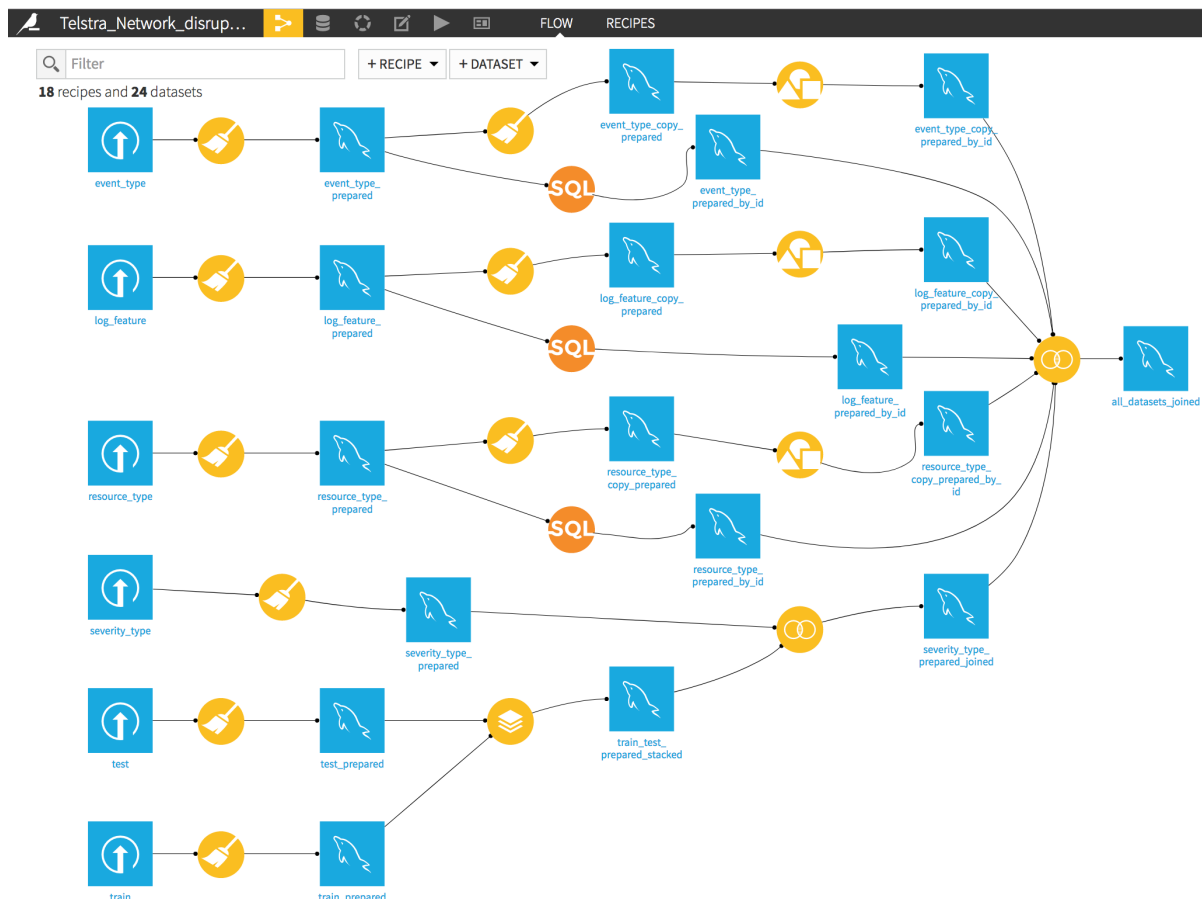


Image 9. All_datasets_joined.

At this point we are able to split the joined dataset into train and test according to the fault_severity column, and then prepare the predictive model to obtain the test dataset scored. However, we can still improve the dataset with a couple of more features.

# 5. IMPROVEMENTS AND PREDICTIVE MODELS

In the first place we will use the string concatenated columns of log_feature, resource_type and event_type created: "log_features", "volumes", "resource_types" and "event_types". For each column we will use extract numbers processor that finds and extracts numerical values from a text column. In this way we are adding the exact value of an event, resource, volume or log feature into a column that can represent for example a date/time or a category. After this the string columns makes no sense so they can be erased.



Image 10. Extracted numbers

Another nice feature is the location sort. In the first part of the analysis we can realize that each of the four datasets come in the same order of ID.

| id | severity_type |
|---|---|
| bigint | string |
| Integer | Text |
| 6597 | severity_type 2 |
| 8011 | severity_type 2 |
| 2597 | severity_type 2 |
| 5022 | severity_type 1 |
| 6852 | severity_type 1 |
| 5611 | severity_type 2 |
| 14838 | severity_type 1 |

Image 11. Severity_type

| id | event_type |
|---|---|
| int | string |
| Integer | Text |
| 6597 | event_type 11 |
| 8011 | event_type 15 |
| 2597 | event_type 15 |
| 5022 | event_type 15 |
| 5022 | event_type 11 |
| 6852 | event_type 11 |
| 6852 | event_type 15 |
| 5611 | event_type 15 |
| 14838 | event_type 15 |

Image 12. Event_type

| id | log_feature | volume |
|---|---|---|
| bigint | string | bigint |
| Integer | Text | Integer |
| 6597 | feature 68 | 6 |
| 8011 | feature 68 | 7 |
| 2597 | feature 68 | 1 |
| 5022 | feature 172 | 2 |
| 5022 | feature 56 | 1 |
| 5022 | feature 193 | 4 |
| 5022 | feature 71 | 3 |
| 6852 | feature 201 | 2 |
| 6852 | feature 56 | 1 |
| 6852 | feature 80 | 2 |
| 5611 | feature 80 | 2 |
| 14838 | feature 203 | 5 |
| 14838 | feature 82 | 8 |
| 14838 | feature 80 | 9 |

| id | resource_type |
|---|---|
| bigint | string |
| Integer | Text |
| 6597 | resource_type 8 |
| 8011 | resource_type 8 |
| 2597 | resource_type 8 |
| 5022 | resource_type 8 |
| 6852 | resource_type 8 |
| 5611 | resource_type 8 |
| 14838 | resource_type 8 |

Image 13. Resource_type

Image 12. Log_feature

Once you merge one of this files with the train and test files stacked, it turns out that the records are sorted by location, but not location as a number, but location as an ASCII string, this means 1000 comes before 999 and the order is like 1, 10, 100…

So, for keeping this particular order of IDs in one of the ordered datasets, a row_number was created in severity_type_prepared, with the next SQL recipe:

```
SELECT      *,
            @curRank := @curRank + 1 AS `row_number`
FROM        `TELSTRA_NETWORK_DISRUPTIONS_1_severity_type_prepared`, (SELECT @curRank := 0) r
```

And then joined to train_test_prepared_stacked:

| id | location | fault_severity | row_number |
|---|---|---|---|
| int | int | int | int |
| Integer | Integer | Integer | Integer |
| 4787 | 1 | | 62 |
| 16586 | 1 | 0 | 63 |
| 8796 | 1 | | 64 |
| 6954 | 1 | | 65 |
| 17648 | 1 | | 66 |
| 1849 | 1 | 0 | 67 |
| 16949 | 1 | 0 | 68 |
| 14593 | 1 | | 69 |
| 12587 | 10 | | 70 |
| 2240 | 10 | 0 | 71 |
| 5356 | 10 | | 72 |
| 1285 | 100 | | 73 |
| 17155 | 100 | | 74 |
| 1308 | 100 | 0 | 75 |
| 3251 | 100 | | 76 |

Image 14. Prepared, stacked, ordered, joined

Since its order by location, it could be also order by time (maybe represented by the ID). In this case, a rank order was added to each location individually using the Window Visual Recipe of Dataiku. First, partitioning is made on location and order by row_number to maintain the original order of IDs.

| id | location | fault_severity | row_number | rank |
|---|---|---|---|---|
| int | int | int | int | bigint |
| Integer | Integer | Integer | Integer | Integer |
| 4787 | 1 | | 62 | 62 |
| 16586 | 1 | 0 | 63 | 63 |
| 8796 | 1 | | 64 | 64 |
| 6954 | 1 | | 65 | 65 |
| 17648 | 1 | | 66 | 66 |
| 1849 | 1 | 0 | 67 | 67 |
| 16949 | 1 | 0 | 68 | 68 |
| 14593 | 1 | | 69 | 69 |
| 12587 | 10 | | 70 | 1 |
| 2240 | 10 | 0 | 71 | 2 |
| 5356 | 10 | | 72 | 3 |
| 1285 | 100 | | 73 | 1 |
| 17155 | 100 | | 74 | 2 |
| 1308 | 100 | 0 | 75 | 3 |
| 3251 | 100 | | 76 | 4 |
| 9313 | 100 | 0 | 77 | 5 |
| 14926 | 100 | | 78 | 6 |

Image 15. Prepared, stacked, ordered, ranked by location

Prepared and joined to the rest of the datasets, keeping the location order, the flow looks like this:



Image 16. Last preparation

## 5.1. Predictive Models

Now we are ready to split the data and train the predictive model. First, a split was made in the last prepared dataset with the next settings:



Image 17. Split settings

Next step is to create a prediction model in "to_train dataset" with the fault_severity variable to predict. A **Performance** model was chosen to train a highly cross-validated model to attempt to obtain the best possible predictions.



Image 18. Performance models

The XGBoost is an advanced gradient boosted tree algorithm. It has support for parallel processing, regularization, early stopping which makes it a very fast, scalable and accurate algorithm. The settings where used by default:

Image 19. XGBoost

The other model was the Random Forest for decision tree classification and is a simple algorithm which builds a decision tree. Each node of the decision tree includes a condition on one of the input features. A Random Forest classifier is made of many decision trees. When predicting a new record, it is predicted by each tree, and each tree "votes" for the final answer of the forest.

The forest chooses the class having the most votes. When "growing" (ie, training) the forest:
    for each tree, a random sample of the training set is used;
    for each decision point in the tree, a random subset of the input features is considered.
Random Forests generally provide good results, at the expense of "explainability" of the model.

Once the models are trained, the next step is to apply them on the test dataset to predict the fault severities. Then, the result is prepared one last time before the submission in the website of the competition. The final preparation consist on leaving just the ID and the predict_0, predict_1 and predict_2 columns, then export the .csv file and upload it to the platform to obtain both private and public score.

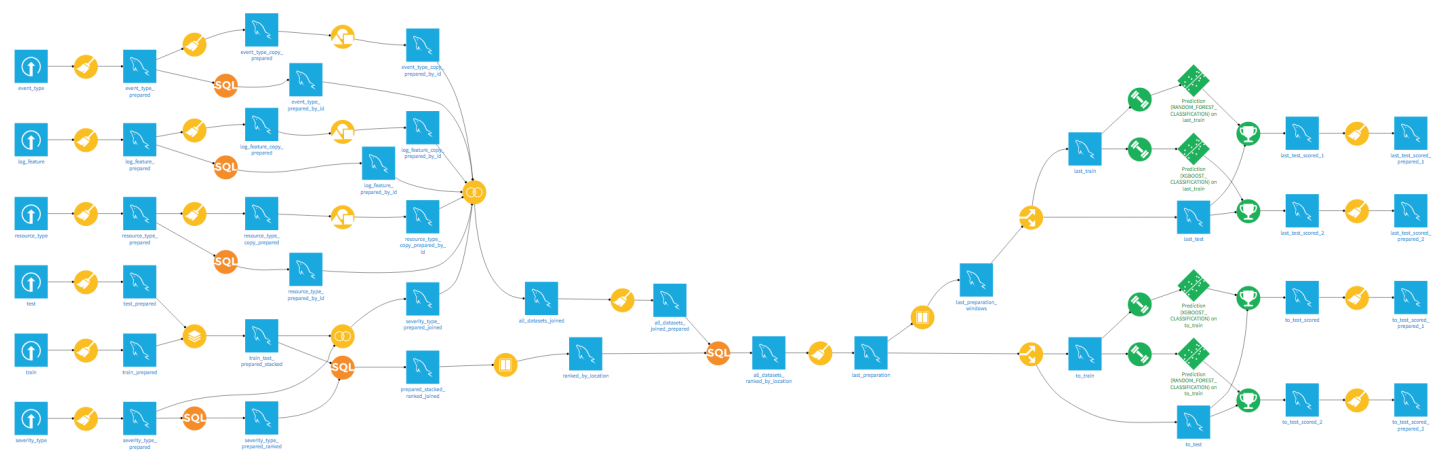This is how the flow looks at the end and with some other predictive model trials on the same datasets:



Image 20. Final flow

# 6. RESULTS AND CONCLUSION

At first, the results of the submissions were above the 500th positions in the leaderboard. There were just a few features and some basic preparation in the datasets. However, being curious with the DataScienceStudio and in general with what Data Science is about, there was a considerable improvement in the model and features besides a very exiting and fun learning process.

| Submission and Description | Private Score | Public Score | Use for Final Score |
|---|---|---|---|
| **test_all_info_scored_XGB2_prepared.csv** <br> a few seconds ago by TAQP <br> *add submission details* | 0.53711 | 0.54779 | ☐ |
| **test_all_info_scored_XGB_prepared.csv** <br> 2 hours ago by TAQP <br> *add submission details* | 1.01797 | 1.01958 | ☐ |
| **test_scored_prepared.csv** <br> 2 hours ago by TAQP <br> *add submission details* | 0.55882 | 0.55744 | ☐ |
| **test_all_info_scored_prepared.csv** <br> 2 hours ago by TAQP <br> *add submission details* | 0.55849 | 0.56671 | ☐ |
| **test_all_info_scored.csv** <br> 2 hours ago by TAQP <br> *add submission details* | NULL | Error ⓘ | ☐ |
| **TAQP_test_dataset_scored_prepared.csv** <br> 7 hours ago by TAQP <br> *add submission details* | 0.61736 | 0.63562 | ☐ |
| **test_dataset_XGB_scored_prepared.csv** <br> 8 hours ago by TAQP <br> *add submission details* | 0.67018 | 0.68375 | ☐ |
| **test_dataset_scored_prepared.csv** | 0.63818 | 0.65213 | ☐ |

Image 21. Competition submissions

The datasets for this competition seemed really simple at first but turned out to have some hidden depths and tricky but significant steps. It was more about feature engineering to get a good result, rather than tuning a powerful predictive model.

If there is one take away from this exercise, is that you have to be creative and the most important thing to do is to first capture all the possible information presented in the data. Data Science is about working out how to solve not so straight forward problems, with exploration, feature engineering and finally Predictive Modeling using advanced machine learning algorithms.

Personally, the fact that competitors who stand out were considered for data science roles in Telstra's Big Data team is really amazing and means that Data Scientists are being highly requested by the industry.

The final and best result on the submissions was a private score of ~0.52 (position 414). However, there are more features in development to keep improving the entire model.