## Deliverable C

### Steps 1 and 2

```matlab
tab = readtable('fisheriris.csv');
tab = tab(randperm(size(tab,1)), :);
```
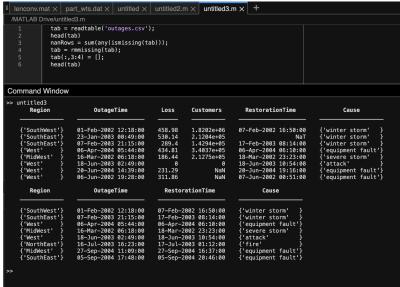




We decided to remove the Loss and Customer columns from the table because we thought it was best to use this data to focus on only the cause/location of the outages and the time it took to fix them. Loss and number of customers were irrelevant to that goal so we removed them.

Step 3

Our initial table had 200 columns of data, and many were unnecessary. We cut this down to 20 columns. Additionally, all missing values were removed. At the end of the cleaning, the dataset size went from 641491 to 300254 rows, with 87% of the reduction due to missing values.

Data import

```
data = readtable("./Dataset/combined.csv")
```

Warning: Column headers from the file were modified to make them valid MATLAB identifiers
The original column headers are saved in the VariableDescriptions property.
Set 'VariableNamingRule' to 'preserve' to use the original column headers as table variabl

data = 631491×200 table

|    | ReportSubmissionSource | MultipleRowsPerIncident |                              |
|----|------------------------|-------------------------|------------------------------|
| 1  | 'Paper'                | 'No'                    | '<a href = https://portal.phm |
| 2  | 'Paper'                | 'No'                    | '<a href = https://portal.phm |
| 3  | 'Paper'                | 'No'                    | '<a href = https://portal.phm |
| 4  | 'Paper'                | 'No'                    | '<a href = https://portal.phm |
| 5  | 'Paper'                | 'No'                    | '<a href = https://portal.phm |
| 6  | 'Paper'                | 'No'                    | '<a href = https://portal.phm |
| 7  | 'Paper'                | 'No'                    | '<a href = https://portal.phm |
| 8  | 'Paper'                | 'No'                    | '<a href = https://portal.phm |
| 9  | 'Paper'                | 'No'                    | '<a href = https://portal.phm |
| 10 | 'Paper'                | 'No'                    | '<a href = https://portal.phm |
| 11 | 'Paper'                | 'No'                    | '<a href = https://portal.phm |
| 12 | 'Paper'                | 'No'                    | '<a href = https://portal.phm |
| 13 | 'Paper'                | 'No'                    | '<a href = https://portal.phm |
| 14 | 'Paper'                | 'No'                    | '<a href = https://portal.phm |

Keeping only select columns from the 200 column dataset:

- First 3 columns used for removing rows, then are removed
- Next 7 are inputs
- Last 13 are potential outputs

```
data = data(:,{'MultipleRowsPerIncident', ...
    'ReportType', ...
    'IncidentCountry', ...
    'DateOfIncident', ...
    'ModeOfTransportation', ...
    'TransportationPhase', ...
    'CommodityLongName', ...
    'IdentificationNumber', ...
    'PackagingType', ...
    'GeneralPackageType', ...
    'WhatFailedCode', ...
    'HowFailedCode', ...
    'FailureCauseCode', ...
    'SpillageResultInd', ...
```

```
    'FireResultInd', ...
    'ExplosionResultInd', ...
    'WaterSewerResultInd', ...
    'GasDispersionResultInd', ...
    'EnvironmentalDamageResult', ...
    'OtherCleanupInd', ...
    'MaterialLoss', ...
    'RemediationCleanupCost', ...
    'SeriousIncidentInd'});
```

Removal of select rows due to row content:

- MultipleRowsPerIncident == 'Yes' (Removal of 3% of dataset, n=19390) [For simplicity and clarity of counting rows and not double-counting]
- IncidentCountry ~= 'US' (Removal of 0.179%, n=1130) [Negligible reduction in rows, allows for removal of IncidentCountry column for file size]
- ReportType ~= 'A hazardous material incident' (Removal of 2.936%, n=18538) [For uniformity, the removed 2.936% were nonhazardous reports]

Overall reduction from 631491 to 596921 rows, n=34570, 5.5%

```
data(ismember(data.MultipleRowsPerIncident,'Yes'),:)=[];
data(~ismember(data.IncidentCountry,'US'),:)=[];
data(~ismember(data.ReportType,'A hazardous material incident'),:)=[];
data(:,1:3) = [];
```

Removal of rows with missing values:

Reduction from 596921 to 300254 rows, n=296667, 50%

```
data = rmmissing(data);
data = data(randperm(size(data,1)),:)
```

data = 300254×20 table

| | DateOfIncident | ModeOfTransportation | TransportationPhase | |
|---|---|---|---|---|
| 1 | 17-Nov-2017 | 'Highway' | 'In Transit' | 'SODIL |
| 2 | 11-Jul-2020 | 'Highway' | 'Unloading' | 'DISINI |
| 3 | 10-Sep-1995 | 'Highway' | 'Unloading' | 'INK PI |
| 4 | 17-Sep-2012 | 'Highway' | 'Unloading' | 'SODIL |
| 5 | 26-May-2021 | 'Highway' | 'Unloading' | 'CORR |
| 6 | 01-Apr-2022 | 'Highway' | 'Loading' | 'ISOPF |
| 7 | 19-Jan-2007 | 'Highway' | 'Unloading' | 'POTA: |
| 8 | 10-Aug-2017 | 'Highway' | 'Unloading' | 'HYDR |
| 9 | 12-Jul-2019 | 'Highway' | 'Loading' | 'ISOPF |
| 10 | 10-Nov-2005 | 'Highway' | 'Loading' | 'PAINT |
| 11 | 27-May-2023 | 'Highway' | 'Loading' | 'SODIL |
| 12 | 25-Aug-2023 | 'Highway' | 'Unloading' | 'HYDR |
| 13 | 12-Oct-2023 | 'Highway' | 'Unloading' | 'CORR |
| 14 | 15-Aug-2012 | 'Highway' | 'Unloading' | 'PAINT |