

# Data Intake Report

Name: APS Failure from Scania trucks

Report date: 2021-10-25

Internship Batch: NLP01

Version: 1.0

Data intake by: Juan Carlos Gutiérrez

Data intake reviewer: NA

Data storage location: <https://archive.ics.uci.edu/ml/machine-learning-databases/00421/>

## Tabular data details:

*File name: aps\_failure\_training\_set.csv*

<b>Total number of observations</b>	60000
<b>Total number of files</b>	1
<b>Total number of features</b>	171
<b>Base format of the file</b>	.csv
<b>Size of the data</b>	44.7 MB

*File name: aps\_failure\_test\_set.csv*

<b>Total number of observations</b>	16000
<b>Total number of files</b>	1
<b>Total number of features</b>	171
<b>Base format of the file</b>	.csv
<b>Size of the data</b>	11.9 MB

## Proposed Approach:

The goal is to have a definitive dataset, ready to analyse and/or use for modelling purposes. These are the steps I followed:

1. Review the data files: values, variables, dimension, etc.
2. Look for null values: in this case, those are represented by the string “na” in the data frames.
3. Encode the target from “neg” and “pos”, to 0 and 1.
4. Have the correct data types: integer for the target variable and float for the rest.
5. Standardize the values, since there are some outliers, I preferred to standardize instead of normalize, expecting to have a better scaling before modelling.
6. Work on the standardized data sets, given the variability of the different variables.