# Data Science Hackathon

Lucas Sulzbach Rilho

# Introduction

- The main goal of this work was to create a model to predict what kind of pollutant will be produced based on the training data.

# Data Descritpion

- Data was collected from 2 CSV files, 3 JSON RESTapi and from 83 PDF reports.
- All sources contained the same kind of information with a little variation on names and field shown.

# Methodology

- Firstly the training data was cleaned and joined in a single training dataframe;
- From this dataframe was made a train/test group to test the fields used on the Decision Tree classification;
- It was decided to keep only the fields 'EPRTRSectorCode' and 'EPRTRAnnexIMainActivityCode' as a main indicator of the company production type, that was directly related with the pollutant emited.
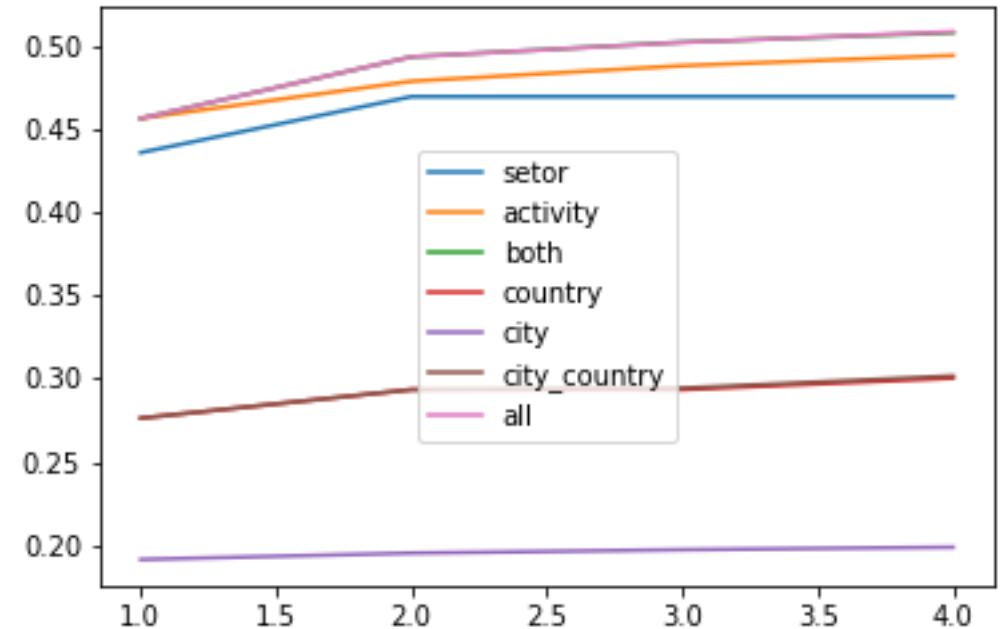


Fig 1. – Training with different fields

# Results

- This model have a F1 score of 0.5077352830457011.