# 01 Schneider Hackathon

Jordi Arellano Ballestero
Jordi Segura Pons

# 02

# Dataset

We had to deal with Data splitted into CSV and JSON files.

While CSV read was pretty straightforward, JSON contained 2 more columns than CSV.

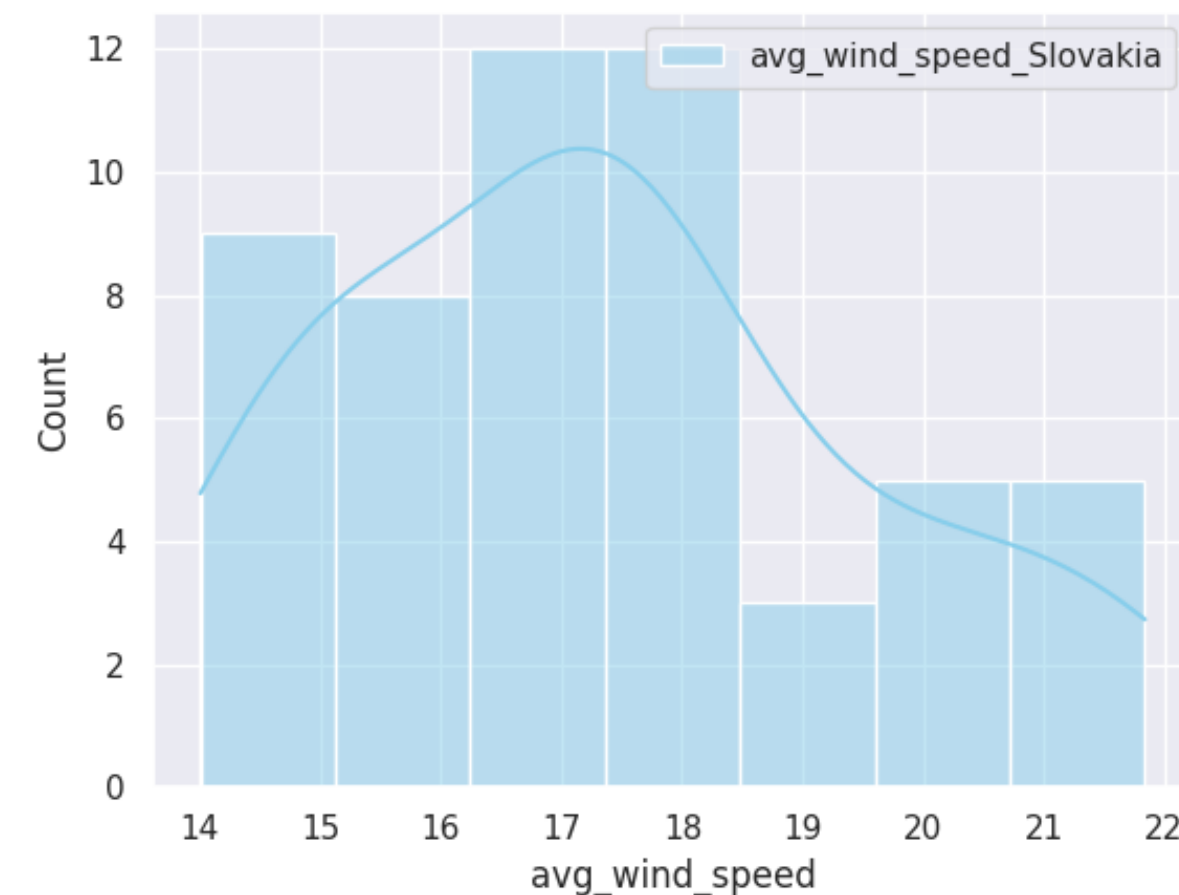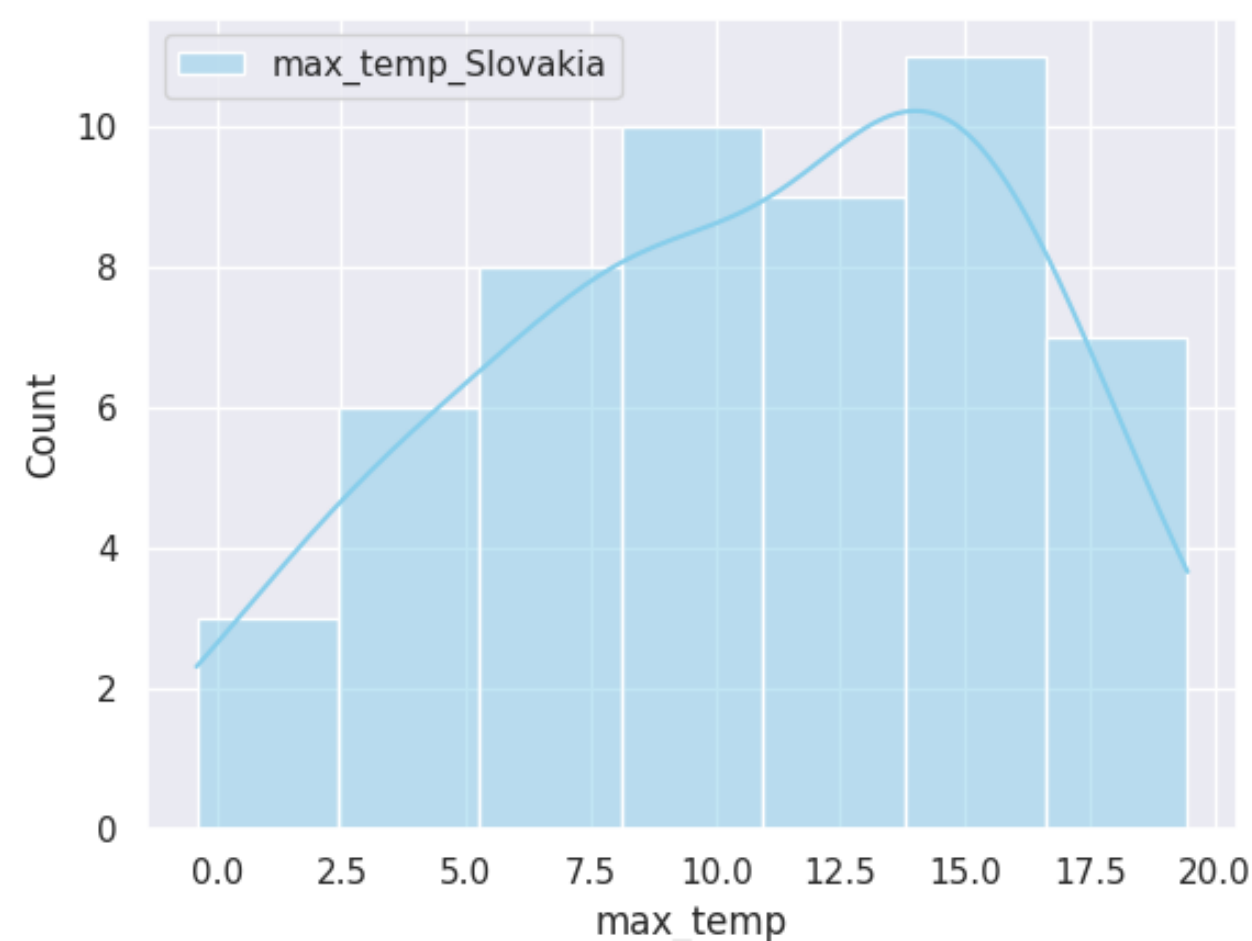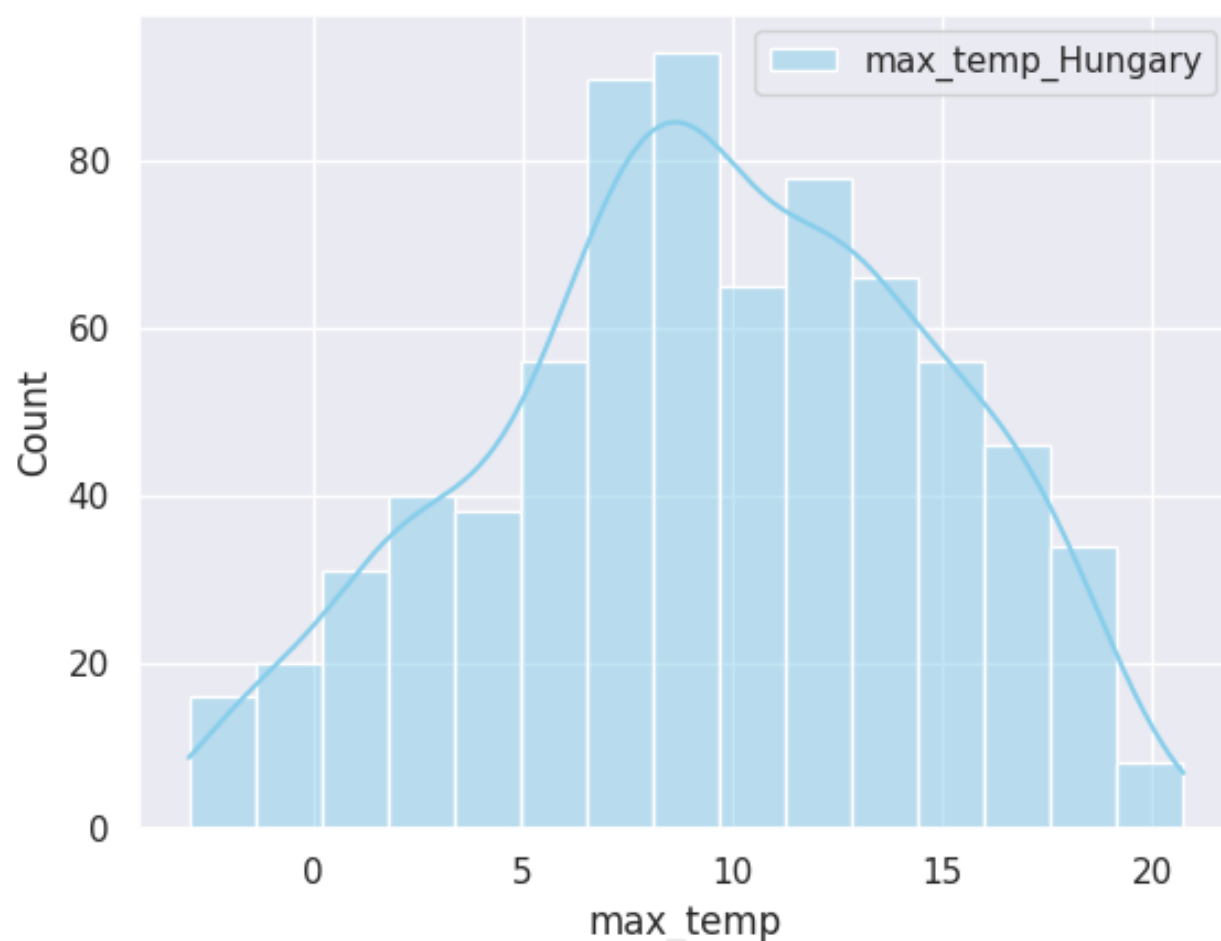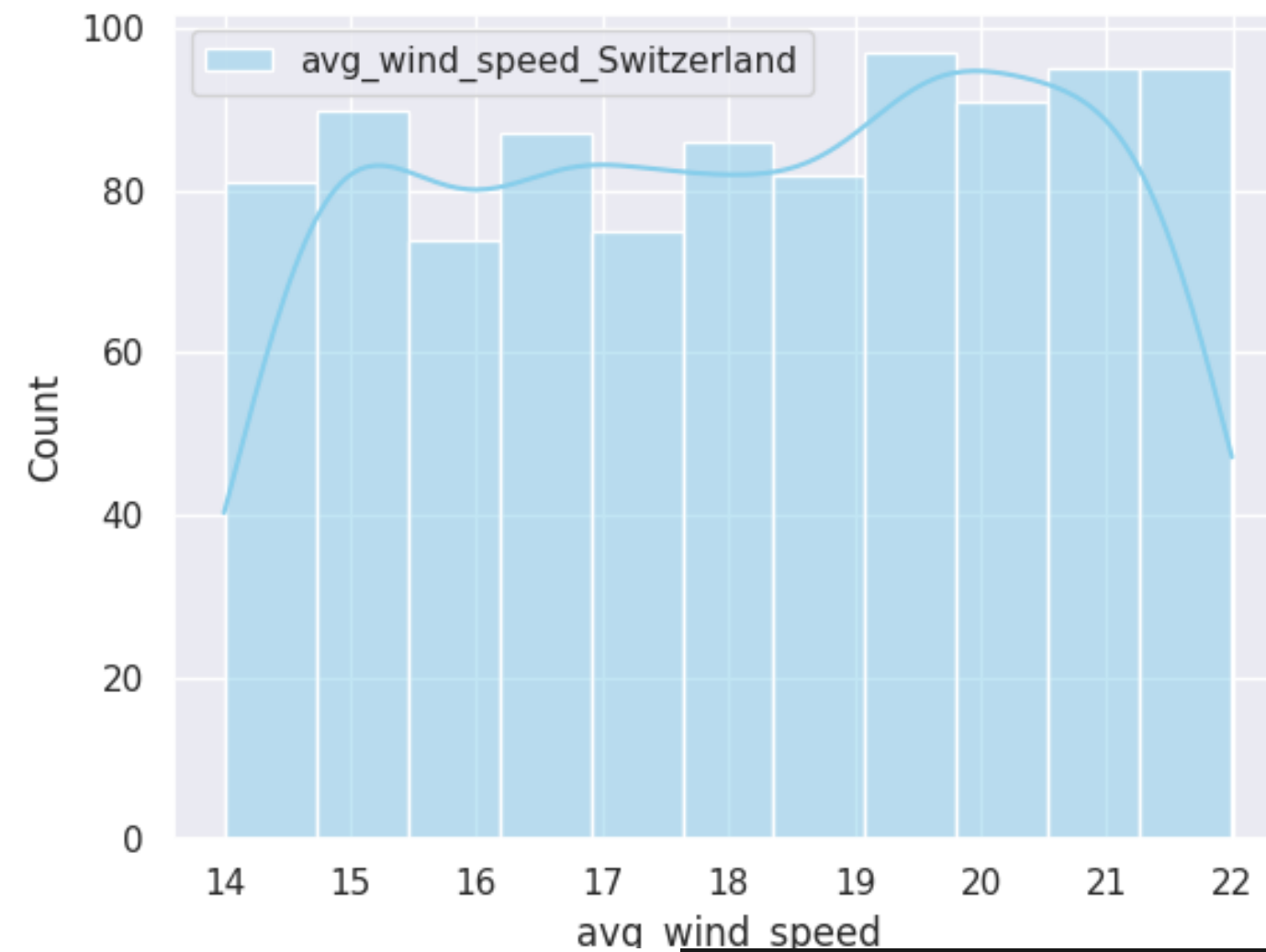After having both dataframes unified, we encoded the categorical columns and resorted them.
We also changed the date into a datetime due to facility of our models to ingest it.

```
Index(['CITY ID', 'FacilityInspireID', 'DateTime', 'CONTINENT', 'City',
       'countryName', 'EPRTRAnnexIMainActivityCode',
       'EPRTRAnnexIMainActivityLabel', 'EPRTRSectorCode', 'eprtrSectorName',
       'avg_temp', 'avg_wind_speed', 'max_temp', 'max_wind_speed', 'min_temp',
       'min_wind_speed', 'targetRelease', 'facilityName', 'reportingYear',
       'MONTH', 'DAY', 'REPORTER NAME', 'DAY WITH FOGS', 'pollutant'],
      dtype='object')
```

```
Diff columns are: Index(['', 'EPRTRAnnexIMainActivityCode', 'EPRTRSectorCode'], dtype='object')
```
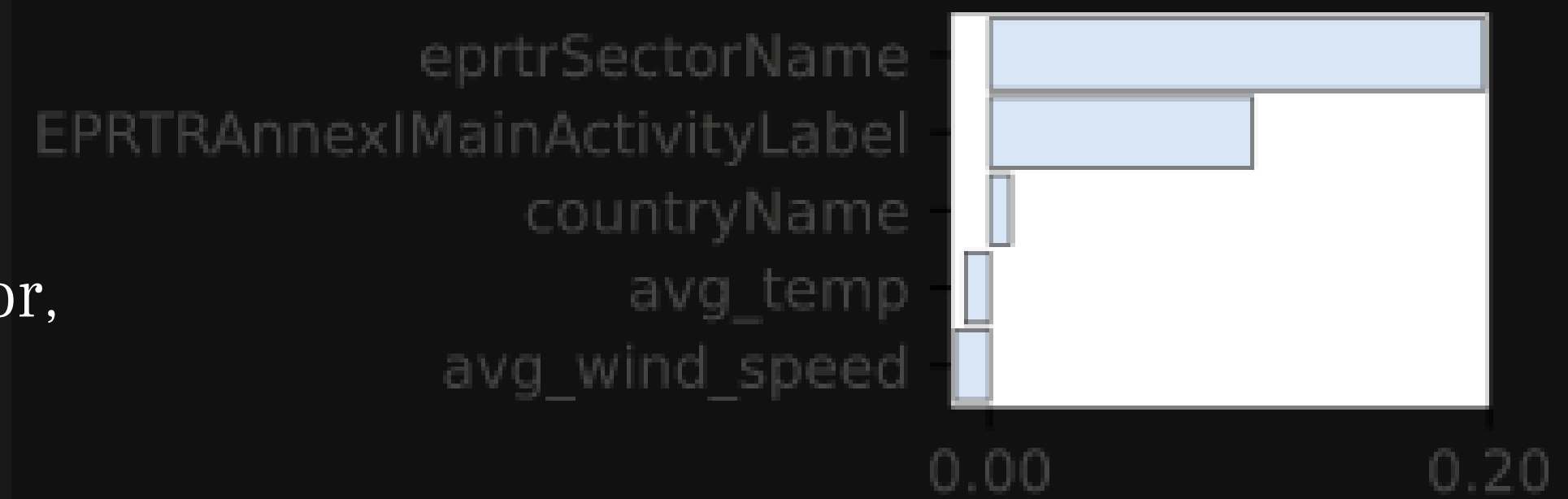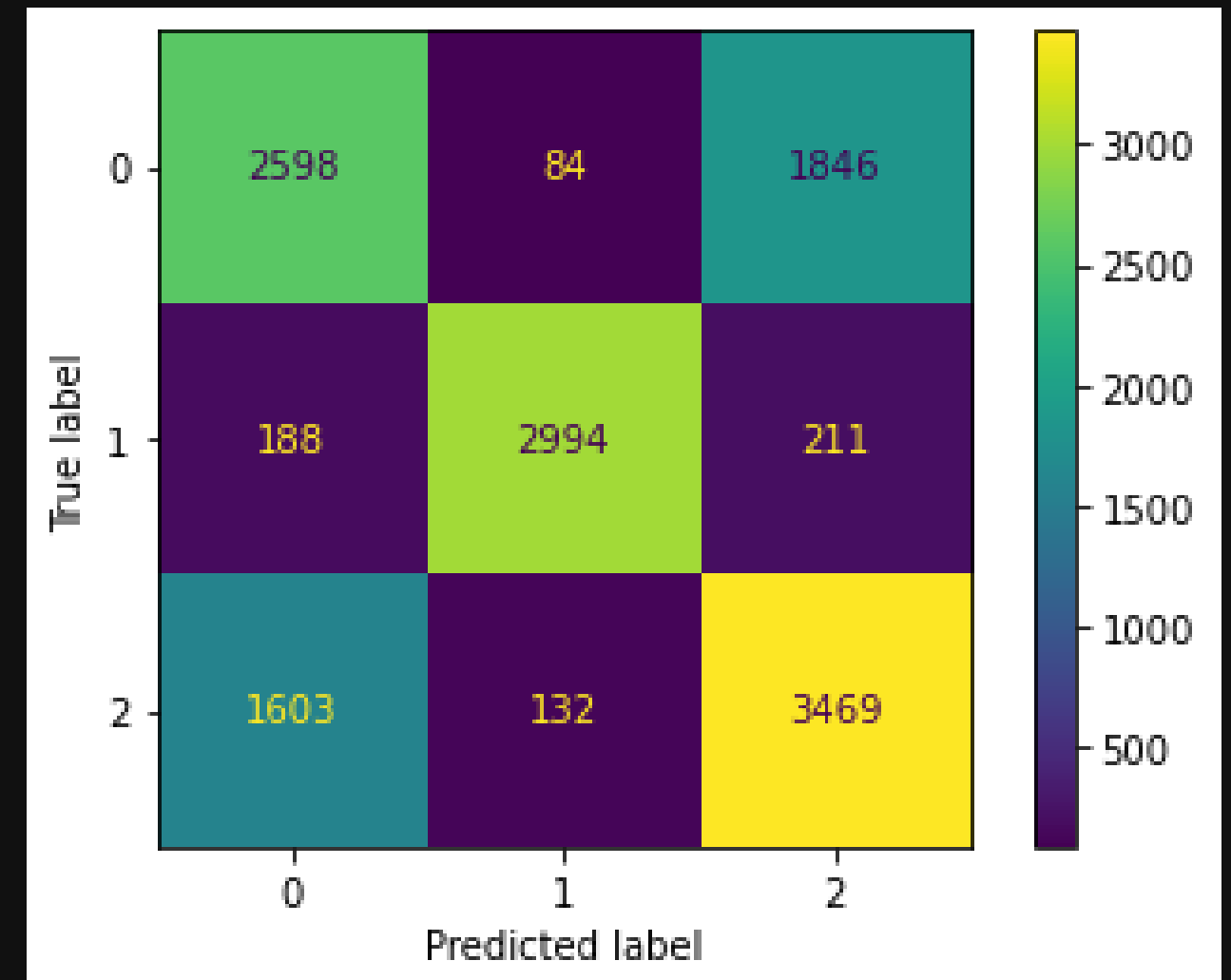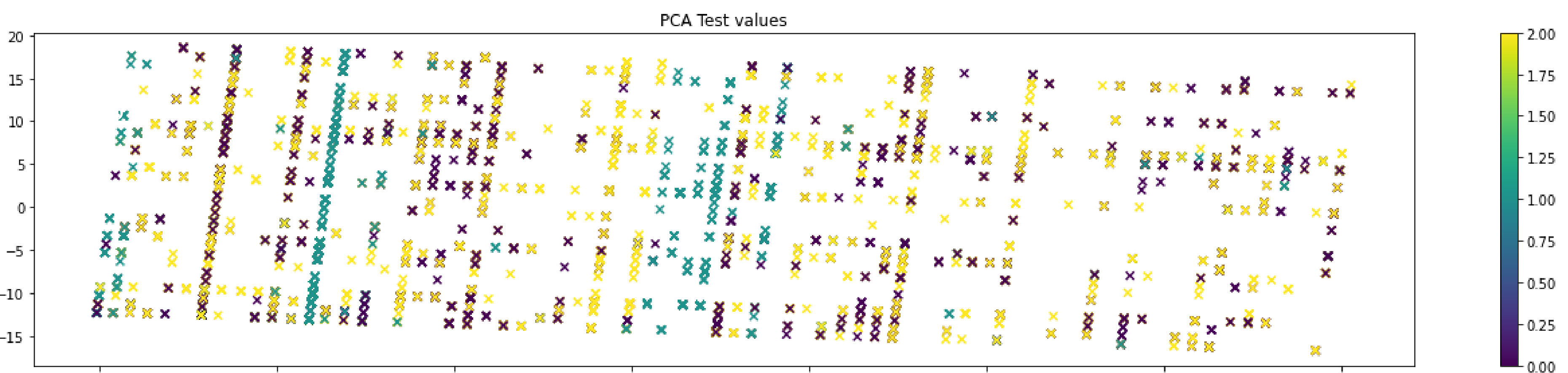
# Models

Feature Selection

We end up using only the country, activity, sector, average temperature and wind speed.

We wanted to try classic ML Multiclassification algorithms, such as: RF, Decision Trees, SVC, kNN and XGBoost.

The best results were from a Decision Tree with LightGBM and kNN with PCA.

PCA Test values

We can appreciate a structure in the data, where pollunt 1 is really well defined but 0 and 2 are mixed each other. We believe it should be useful to iterate over this idea in order to get more thorough information.

In the future, might be worth trying a resampling and DL algorithms.