

$$\text{Let } N = \bar{n}_i = \sum_{i=2}^{d=90} \frac{n_i}{d}$$

be the mean number of rows for a single Day. We expect $N \approx 150 \times 10^6$.

Our dataset presents $F = 11$ columns.

Suppose the proposed procedure manages to cut down the datasize to $p = 0.01$ of its former size.

We can confidently process $\approx 25 \times 11 \times 10^6$ datapoints, so that is $\approx \frac{11N}{6} = T$ points.

Applying the procedure, we get:

$$N \times F \times p \times (d-1) \approx \frac{11 \times 89N}{100} = 11N \times 0.89 = 5.34T \quad (1)$$

This is 53.4 times greater than what we can feasibly process.

A second strategy is to remove some of the features. We can at most hope to remove 4 of them: 'Src_Bytes', 'Dst_Bytes', 'Src_Packets' and 'Dst_Packets'. We'd then get:

$$N \times (F-4) \times p \times (d-1) \approx \frac{7N \times 89}{100} \approx 0.566 \times 11N = 3.396T \quad (2)$$

This is 33.96 times greater than desired still. Moreover, the assumption $p = 1\%$ is rather generous, as early observations presented the following:

$$p_1 = 0.2, p_2 = 3.4, p_3 = 0.4, p_4 = 7.8, p_5 = 1.2$$

We then get a mean preliminary $\bar{p} = \sum_{j=1}^5 \frac{p_j}{5} = 2.6 = 2.6p$, 2.6 times our desired reduction.

With the current observations, we're expecting an oversize of a factor between $2.6 \times 3.396 = 8.8296$ and $2.6 \times 5.34 = 13.884$ times the target datasize.