# Assignment 2: Dask vs PySpark vs Koalas vs Modin *

Emanuel Tomé [†], João Ferreira [‡], Ricardo Faria [§], Vânia Guimarães [¶]

DCC-FCUP, July 2021

## 1 Introduction

Several parallel DataFrame systems alternatives to Pandas DataFrame are available nowadays. However, the differences between the available libraries is not always clear in which situations one is preferable over the others. Moreover, the cost of transition from Pandas to another library can be different not only for the need (or not) of learning a new syntax, but also in the compatibility of those libraries with others largely used by Data Scientists such as Scikit-Learn and Numpy. In this context, the main objectives of this work are:

- Identify performance bottlenecks when using a specific library.

- Identify syntactic differences among the different libraries.

- Identify operations that are best suited to one particular library.

- Get acquainted with these libraries and knowing what is supported from Pandas, Scikit-learn and Numpy.

## 2 Brief background on PySpark, Dask, Modin, JobLib, Rapids and Koalas

### 2.1 Dask

Libraries such as NumPy, Pandas and Scikit-learn, are designed to run on a single core. Therefore, all the data will be temporarily loaded onto the RAM of our local system. However, if we deal with extreme large datasets, we will inevitably problems to run it. Dask appeared mainly to solve this problem. Dask [1] is a parallel computing library that works by distributing larger computations and breaking it down into smaller computations through a task scheduler and task workers. Designed to parallelize in python ecosystems, Dask is suited to solve a wide variety of problems including structured data analysis, large-scale simulations used in scientific computing and general-purpose [2]. However, for data that fits into RAM, Pandas can often be faster and easier to use than Dask DataFrame. The popularity of the Dask is due to the union of the power of distributed computing for data science with the good integration to common Python data tools. Another advantage is that without change the interface, the users can run on clusters with multiple cores or on a common machine using a single process. Otherwise, the users don't really need to worry about the low-level internals, Dask provides several collections for wrapping low-level tasks into high-level workflows.

### 2.2 Modin

Modin [3] uses Ray or Dask to provide an effortless way to speed up pandas notebooks, scripts and libraries. The authors claim that Modin is able to scale the pandas workflow by changing a single line of code. That is, to use Modin and take advantage of its speedup, the user only has to change the line of code which imports pandas, `import pandas as pd`, to `import modin.pandas as pd` and continues using their previous pandas notebooks.

Moreover, to use Modin the user does not have to know how many cores their system has and does not need to specify how to distribute the data. Data Scientists spend their time extracting value from their data than on tools that extract data.

## 2.3   JobLib

Joblib provides solution for several activities, such as loading up large Numpy arrays, persisting python object or performance of custom python function, with the help of parallel computing, memoization and caching mechanism. Joblib [4] provides a set of tools for making the pipeline lightweight to a great extent in Python [5], without dependency on other libraries. Joblib allows to use cache which avoids recomputation of some of the steps and execute parallelization to fully utilise all the cores of CPU/GPU. So, function called with same argument will not be re-compute, instead, output loads back from cache using memmapping. It also provides a compressor during persistence for large data, to save space on disk. A fantastic library that became popular because of its optimized time-complexity feature.

## 2.4   Rapids

Rapids [6] is a suite of open source software libraries and APIs that provides the ability to execute end-to-end data science and analytics pipelines entirely on GPUs. Rapids utilises NVidia Cuda primitives for low-level compute optimisation, and exposes GPU parallelism and high-bandwidth memory speed through user-friendly Python interfaces. It also focuses on common data preparation tasks for analytics and data science, and includes support for multi-node, multi-GPU deployments, enabling vastly accelerated processing and training on much larger dataset sizes.

## 2.5   PySpark

Apache Spark is a unified computing engine and a set of libraries for parallel data processing on computer clusters [7]. Spark supports multiple widely used programming languages, such as Python, Java, Scala and R, includes libraries for diverse tasks ranging from SQL to streaming and machine learning and runs anywhere (from a laptop to a cluster of thousands of servers).

PySpark [8] is an interface for Apache Spark in Python. It allows to write Spark applications using Python APIs and also provides the PySpark shell for interactively analyse data in a distributed environment. PySpark supports most of Spark's features such as Spark SQL, DataFrame, Streaming, MLlib (Machine Learning) and Spark Core.

## 2.6   Koalas

Pandas syntax and PySpark syntax differ to a considerable degree because PySpark has been notably influenced by SQL syntax. Regular pandas users will argue that it is much less intuitive. This is where Koalas enters the picture. Koalas is a data science library that implements the Pandas APIs on top of Apache Spark so, data scientists can use their favourite APIs on datasets of all sizes [9], and be more productive. Pandas does not scale well to big data since it was designed for small data sets. Using Koalas it is possible make the transition from a single machine to a distributed environment without needing to learn a new framework.

# 3   Materials and methods

## 3.1   Machines used and their characteristics

As required, we started by creating a cluster similar to i3.4xlarge AWS and a machine similar to i3.16xlarge, on Databricks and Dataproc. However, trial accounts for both platforms do not allow work with such power machines. This information is not provided to us, therefore after we started the experiments, both counts were cancelled. To solve the problem, we asked for a machine at DCC. With advanced project, the server stopped to work. The next alternative was create a single node on Dataproc with 8 CPUs and 32 Gb of memory with 300 dollars, the coupon given by Google for each new member. Although, there was some incompatibilities with Dask and Koalas, that made it impossible to continue the execution of the assignment.

The finally alternative was to create two Virtual Machines. The experiments which required Pandas, Dask, Modin, Koalas and Joblib were executed in a Virtual Machine of Google Cloud Platform with 16 CPUs (Intel Cascade Lake) and 64 GB of memory. Selected zone was us-central1-a.

Another Virtual Machine was created to run the experiments with Rapids. This machine contains 8 CPUs, 30 GB of memory and 1 GPU Nvidia Tesla T4. Selected zone was us-central1-b.

## 3.2 Datasets description

As we mentioned before, we started with all dataset required. On each new machine created, the data was decreasing, according to the power of the respective machine. At the end, both experiments and the Machine Learning task were executed with 2009, January of Yellow Taxi Trip Records from NYC Taxi and Limousine Commission (TLC) Trip Record Data.

# 4 Experiment #1: repeat NYC taxi driver dataset study

In this section we present the results obtained for one month of the NYC taxi driver dataset study [9], namely for data from January 2009. Note that only a month of data was considered due to limitations of the available computational capacity. We present the computation times using the following libraries: Pandas (sequential), Joblib, Dask, Modin (with Dask), Koalas (PySpark) and Rapids.

In Figures 1 and 2 are presented the obtained computation times of the standard operations. The first thing to notice is that Rapids has an outstanding performance for all the operations considered. With the exception of operation "count index length" where Pandas is the faster (note that all data fit in memory), Rapids is the faster library. In the opposite side, Joblib seems to be generally the library with worst performance, having poor performance in the majority of operations and the highest total execution time (see Figure 4). However, we should remark that the parallelization with Joblib was made using Pandas dataframes, which should not be the best choice for this library since the Pandas dataframes are stored in columns and not in rows, as for example the Numpy arrays. However, this choice was made since our objective was compare the performance of different libraries when using dataframes.

It is also possible to notice that Modin, followed by Koalas, has by far the worst computation time for the operations "complex arithmetic ops" and "mean of complex artithmetic ops". Indeed, Modin has the second worst total execution time of all benchmarks (see Figure 4).

Considering the geometric mean of the execution times of the standard operations (Figure 4), one can state that Joblib had the worst geometric mean (2.92s) and Rapids and Pandas the best ones. The low computation times for Pandas is only possible since the dataframe fits in memory. We can also notice that although Dask as a lower total computation time than Koalas (Figure 3), the obtained geometric mean for both libraries is similar.

It should be remarked that care should be take when comparing the different computation times since they could not be completely fare or even comparable. For instance, Dask has the lower computation time for reading data. However, in Dask the reading is lazy, which means that the computation time for reading files will be the same regardless the size of the files.



Figure 1: Computation times of the standard operations.

Considering now the computation times obtained for the filtered dataset, that is, for a small dataset, the obtained results are presented in Figures 5 and 6. One can noticed that now Dask has the worst computation times for almost all operations. This confirms what is stated in the documentation of Dask, where the developers confirm that for small datasets does not worth using Dask. Rapids for this smaller dataset is still the faster

Figure 2: Computation times of the standard operations - logarithmic scale.



Figure 3: Total execution time of the standard operations



Figure 4: Geometric mean of the execution times of the standard operations.

library closely followed by Pandas. Modin and Koalas have similar total computation times, while Joblib has slightly lower total computation time (see Figure 7.



Figure 5: Computation times of the standard operations with filtering.



Figure 6: Computation times of the standard operations with filtering - logarithmic scale.
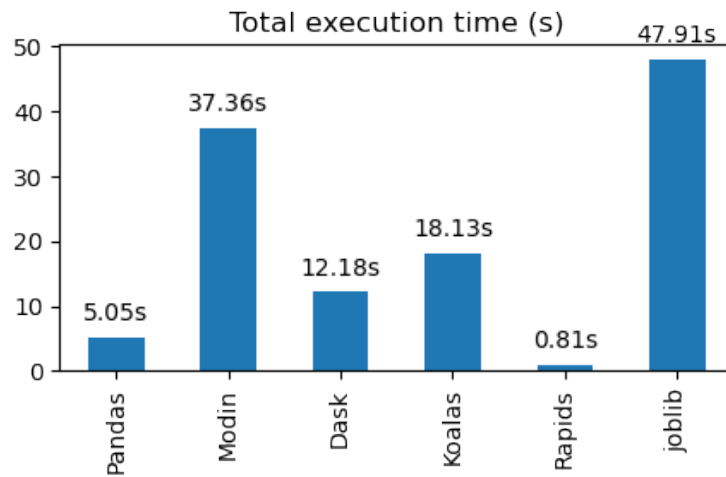
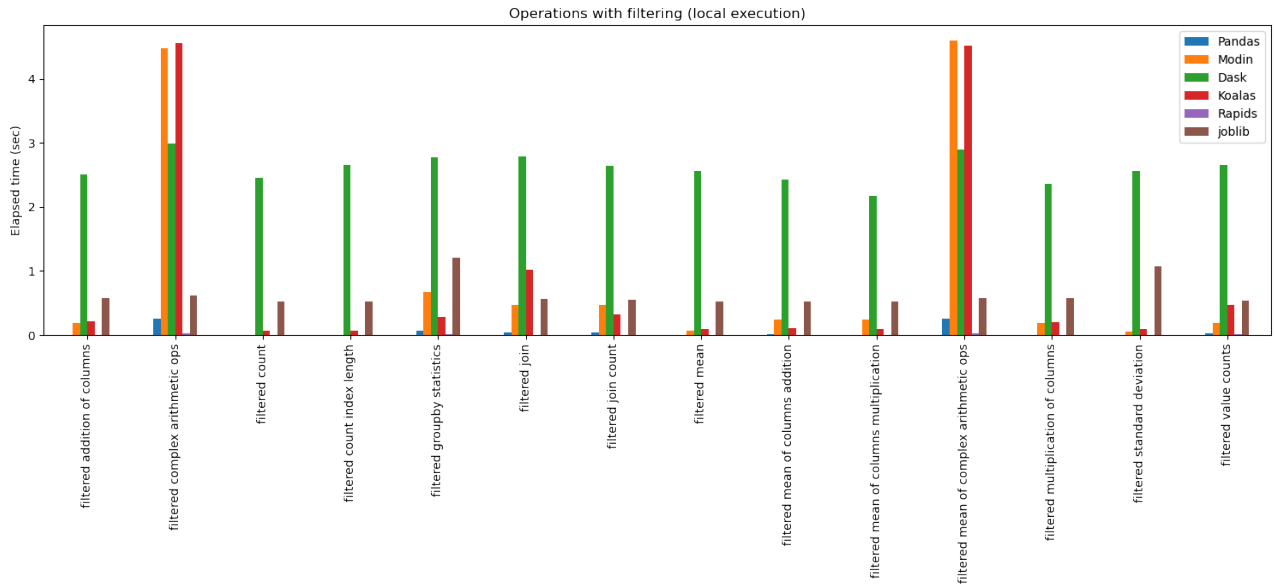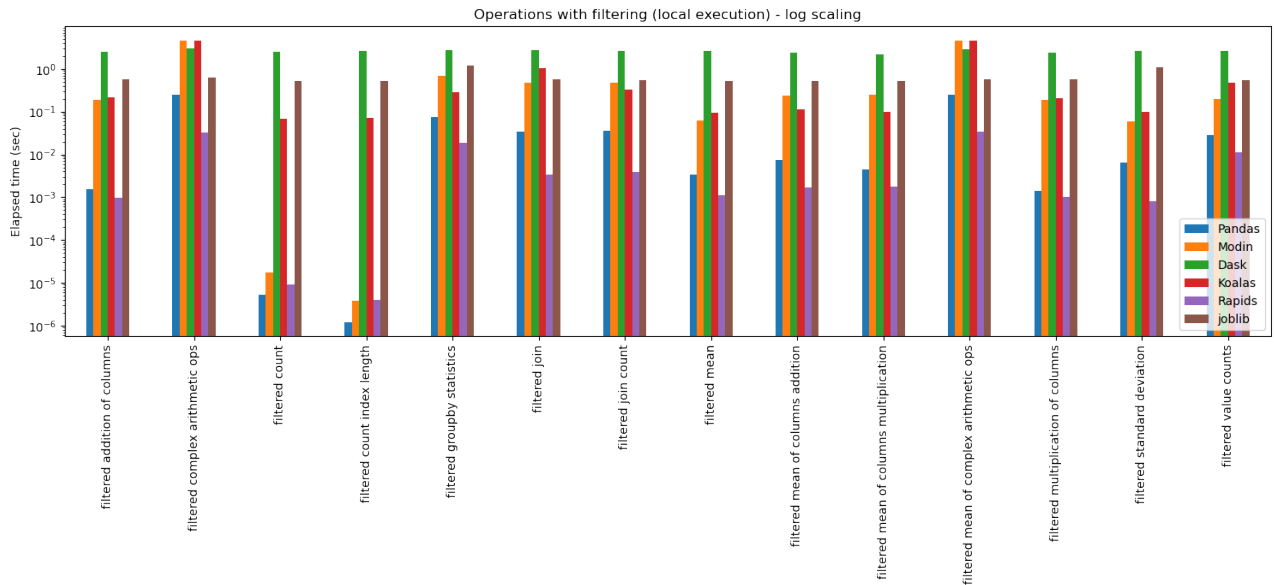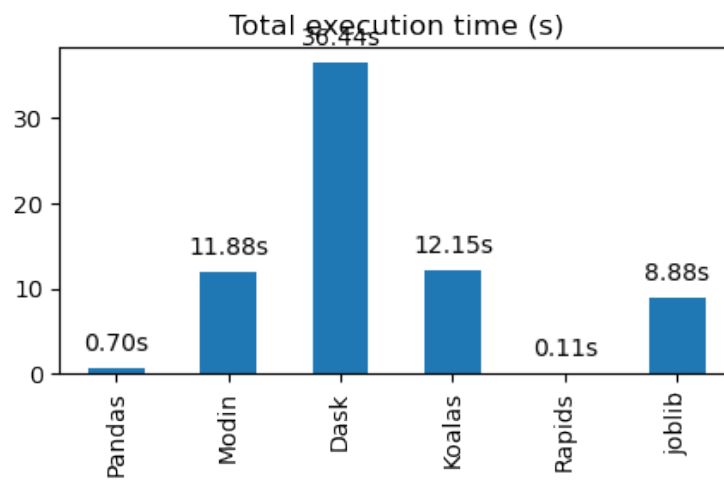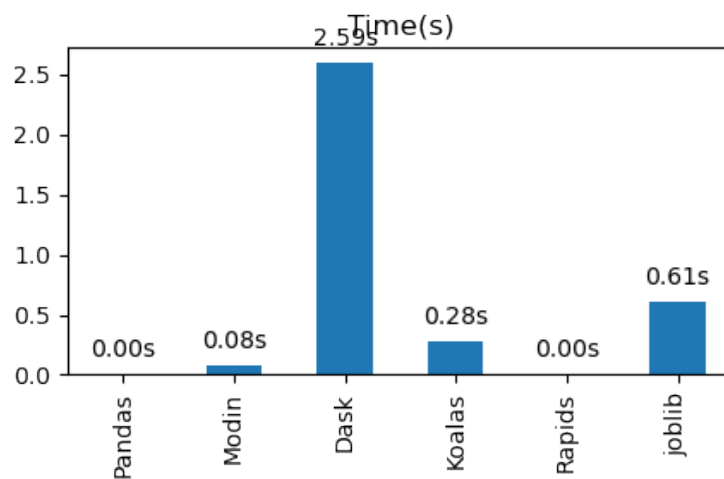Figure 7: Total execution time of the filtered standard operations.



Figure 8: Geometric mean of the execution times of the filtered standard operations.

# 5 Experiment #2

## 5.1 Description of the experiment

The aim of Experiment #2 is to build a full machine learning pipeline and predict the target variable "fare_amount" using **Dask**, **Joblib**, **Koalas** and **Rapids**.

A full machine learning pipeline has the following steps:

- Reading the data

- Pre-processing

- Training

- Testing

### 5.1.1 Read Data

Due to problems of budget and computing resources, the dataset used in the pipeline corresponds to only one month of the NYC taxi driver dataset. The chosen month was January of 2009. With exception of Joblib, which reads the data stored as a Pickle file, all the other libraries read files in the parquet format.

### 5.1.2 Pre-processing

Depending on the type of dataframe used, the pre-processing may vary but these steps are common to all libraries. Using the time in "trip_pickup_datetime" it was added two new columns: "dayofweek" (0 to 6 according to the day of the week) and "hour". A new variable, "weekday" was created with dayofweek (0 if workday, 1 if weekend). Variable "trip_duration" was created by subtrating "trip_pickup_datetime" from "trip_dropoff_datetime". In "payment_type" some categories had different values with the same mean so these cases were recoded. All values with "CASH" recoded to "Cash" and "CREDIT" recoded to "Credit". Lines with abnormal values were removed according to this selection:

- "fare_amt" $> 0$

- "trip_distance" $> 0$

- "trip_duration" $> 0$

- "tip_amt" $\geq 0$

An One Hot Enconding was applied to the categorical variables 'vendor_name' and 'payment_type'. Several variables with no interest were removed from the dataset, some of them are not related with the target, and others had no information. These variables are: 'index', 'trip_pickup_datetime', 'trip_dropoff_datetime', 'rate_code', 'mta_tax','store_and_forward', 'vendor_name' and 'payment_type'. The 'total_amt' was also removed since the 'fare_amt' can be easy computed just to subtracting in the 'total_amt' the variables 'tols_amt' and 'tip_amt'. Finally, it was applied a MinMaxScaler to all numeric variables including the target variable.

### 5.1.3 Traning

After preprocessing, randomSplit was applied to the dataset:

- 70% Train

- 30% Test

The model used by most libraries was XGBRegressor (Extreme Gradient Boosting) of Scikit-learn. We only setup parameter "objective='reg:squarederror'", which is the loss function that is minimized during the training of the model.

XGBRegression is a decision-tree-based ensemble Machine Learning algorithm that apply the principle of boosting. A decision tree is the weak learner, the resulting algorithm is called gradient boosted trees. The "eXtreme" refers to speed enhancements obtained by parallel computation.

For the experiments with Koalas, the Extreme Gradient Boosting was also tested[10] with Koalas but there were some errors that could not be solved. As described in PySpark Example, the spark-xgboost was installed and imported. However, the function XGBRegressor could not be loaded. Further, inspecting the package installed[11], several Python versions up to 3.6 are supported. Although, it was not possible work with XGBRegressor API even using Python 3.5 and 3.6.

Instead of XGB, the model used was GBTRegressor (Gradient-Boosted Trees), which is worse in terms of computational time, since it is not optimised with parallel processing. The loss function to be minimised the same as the XGB models.

## 5.2 Obtained metrics

Because this is a regression problem, in order to rank the performance of the models that were tested above, in the different pipelines, the metrics used were:

- Mean squared error - MSE.

- Root mean squared error - RMSE.

- Mean absolute error - MAE.

- Max error - ME.

- Coefficient of determination - $R^2$.

The MSE was the metric used to fit the models. The results obtained are shown in Table 1.

Table 1: Performance metrics of the fitted models on the different libraries.

|  | MSE | RMSE | MAE | Max error | $R^2$ | Pipeline execution time (s) |
|---|---|---|---|---|---|---|
| *Rapids* | 0.000244 | 0.015611 | 0.013758 | 0.954854 | 0.799853 | 97.75 |
| *Koalas* | 0.000089 | 0.009446 | 0.003511 | 0.989797 | 0.926805 | 461.52 |
| *Dask* | 0.000075 | 0.008655 | 0.002869 | 0.992073 | 0.938329 | 385.35 |
| *Joblib* | 0.000028 | 0.005292 | 0.001278 | 0.792122 | 0.977025 | 3422.39 |

As we can visualise, the library that obtained the best results (even if the differences are almost insignificant), was the *Joblib*, having the best results among the other ones. The rapids library was the one with the worse results and one possible reason for this was a limitation of the cuML library, which doesn't currently have an XGBRegressor available implemented. To circumvent this, the one used was XGBRegressor from Python XGBoost package, applied on the dataframe that was saved in GPU.

One additional limitation was the holdout performed on the GPU dataframe, since there isn't currently available a test_train_split function to perform the split of the dataframe on the GPU, so the way this was surpassed was by using the indices of the dataframe corresponding to x% of the whole data, without shuffling the data. As a consequence, the method decreased the independence of the data, and increased the bias. Regarding Koalas, and as already explained before, since we couldn't use the XGBRegressor model with that library, the GBTRegressor was used. Therefore, care should be taken when comparing the obtained metrics since the models are not exactly the same over the different libraries.

Overall, all the models had a similar performance independently of the library used, however, the execution time of the pipeline is quite different along the librarys. We can see that the Joblib was by far the one that had the longest execution time of them all, counting a total of 3422 seconds (57 minutes). Then, we had the Koalas library, that counted a total of 462 seconds (approximately 8 minutes), followed by the Dask library, which had took about 368 seconds (6 minutes, approximately), and finally, rapids which took only 85 seconds (1 minute and 25 seconds). This comes to reinforce the fact that using GPU for machine learning tasks that deal with large datasets is highly efficient, lowering the execution times of almost every standard operation we could think of.

## 5.3 Profiling

Regarding the profiling of the pipelines, for the majority of the libraries the pipeline part that takes more time to compute is the training of the machine learning model. The exception was Koalas, where more computation time was spent in predicting and computing the evaluation metrics. Joblib was by far the library that took more time to conclude the computation time. However, it should be noted that this is not directly related with Joblib, but with the library/model used to build the model, even though the used model and library trains the model in parallel. If only the pre-processing and "prediction and testing" steps of the pipeline are considered, Joblib has a very good performance.

With exception of Koalas, where the second bottleneck was the training of the machine learning model, the second botleneck for all libraries was pre-processing of the data. Indeed, this is a step in the pipeline where several computations take place, such as transformation of the variables, scaling, drop of variables, computation of new variables from others, merging of dataframes, etc. For this reason, it was expected that this stage of the pipeline was one of the stages with higher computation time. The framework with higher computation time was Dask. One possible reason for that is the need to create indices in the dataframe of numeric, categorical and target variables, in order to guarantee the correct concatenation of the three. Operations like "set_index" and 'merge/join' are harder to do in a parallel or distributed setting than if they are in-memory on a single machine. Shuffling operations that rearrange data become much more communication intensive.

The operation related to import data, Dask was faster than Rapids that use a GPU. This is due to the fact that one of the features of Dask is lazy execution. That means it loads and then processes the data in chunks, so that only a subset of the data needs to be in memory at any given time. However, cuDF stores Dataframes in GPU memory and uses the GPU to perform computations. Regarding Joblib, this was the library with the worst performance. However, it should be stated that for this library the data was readed from a pickle file because the library documentation refers that this is the more efficient way to load data to joblib. If the dataframe was read using, for instance, a lower computation time for this step would be obtained.

Table 2: Line profiling on each of the pipeline's operations.

|          | Import data | Pre-processing | Training | Predicting and Testing | Pipeline |
|----------|-------------|----------------|----------|------------------------|----------|
| *Rapids* | 0.47        | 2.50           | 94.16    | 0.63                   | 97.75    |
| *Koalas* | 2.7         | 19.04          | 391.50   | 443.08                 | 856.34   |
| *Dask*   | 0.02        | 23.92          | 348.12   | 13.27                  | 385.35   |
| *Joblib* | 9.51        | 18.07          | 3388.54  | 6.26                   | 3422.39  |

# 6 Main difficulties and challenges

In this work we faced many difficulties and challenges. The main ones are listed and described bellow:

- **Rapids limitations** - The lack of some vital functions on the cuDF and cuML libraries made some tasks not as easy as they should. For example, the train_test_split function was not implemented and it was needed to implement some kind of algorithm that would split the data into different parts. The problem here was to maintain efficiency, which was not possible due to the sequential programming used to perform this task.

- **XGBRegressor with Koalas** - as already referred, although all the attempts, it was not possible to use the XGBRegressor with Koalas/PySpark. As an alternative, we used the GBTRegressor algorithm.

- **Configuration of the Clusters** - In the begging of this assignment we tried to configure a cluster in the Google Cloud Platform (GCP) using Databricks. However, our account was cancelled two times due to the limitation of GCP described in the next point.

- **Limitation in the GCP of using only a CPU with no more than 8 cores** - when using the educational vouchers, it is not possible to use virtual machines with more than 8 cores. If the user does that, the billing and voucher are cancelled.

- **Necessity of changing platform several times** - We start by trying to do our experiments with Databricks. However, in the beginning our account was cancelled and countless options were tried. We then moved to a DCC server that was made available to the authors. However, when we started to implement our experiments related with the machine learning part, the kernel was always stopping. We then moved to a single node in the Dataprocs of the GCP where we were not able to run all the libraries. Namely, we where not able to read the parquet files that was created before using koalas. We then finally configured two virtual machines in GCP and these was our final solution. Note that each platform has their own characteristics which requires configuring the system, libraries dependencies and sometimes the code.

- **Scarce resources for a so big dataset** - Besides we have to find ways to load data related to 4 years of records and convert them to parquet format, deal with a system that goes down frequently or fit the data to our current resources without lose the purpose of the project was a big challenge. Additionally, the lack of credits in GCP made it difficult to setup a proper machine to run the tasks.

- **Impossibility of doing distributed execution** - As mentioned on the guidelines of the assignment, we tried to implement a distributed execution on Databricks, tirelessly looking for information that would guide us in this direction as we couldn't associate the cluster with the Dask client. After a day of intense search, we were informed that it is not possible to do it and that it was no longer necessary to apply the distributed execution in our experiments.

- **Interpretation of the results obtained using the Yappi profiler** - The output from the Yappi profiler is not as simple as the line profiler. While the line profiler outputs the time spent in each command, the yappi profiler outputs the time spent in each operation which makes it hard to identify, without knowing the meaning of each operation, where the program is spending more time. Although yappi may be harder to analyse, it is possible to profile each thread in details.

- **Time limitation** - Despite the extension of the deadline for the assignment, the difficulties encountered along the way were mainly due to the infrastructure needed to execute the experiments. This situation took our time in the analysis of the core of the project.

# 7    Discussion and Conclusions

In this work we applied different libraries that allow the parallelization of operations in DataFrames and data science pipelines. The obtained results in the benchmarks showed that the use of Rapids/GPUs outperforms by far the results using the other libraries. Dask and Koalas had similar performances and Joblib was the one with the worst computation times. However, it should be state that we used Joblib with dataframes, while it is more suited to use Numpy arrays, which are stored in memory in rows, while the dataframes are stored in columns. Moreover, the kind of operations done in our experiments may not be the ones for which this library is more suited for.

Rapids had an outstanding performance since it used an Nvidia GPU to perform all of the dataframe's operations as well all the machine learning tasks, having some limitations as previously said. Overall, the use of Rapids had an enormous impact in the execution times, decreasing it in relation to the other libraries tested. This was possible because the used GPU made possible to process multiple computations simultaneously, as it has a large number of cores, which allows for better computation of multiple parallel processes. This can be more specifically explained by the fact that GPUs are bandwidth optimised while CPUs are latency optimised, i.e the GPU can fetch much more memory at once, while CPU's can only fetch small amounts of memory quickly, having to perform a lot of memory operations to achieve the same goal. When you add this up to thread parallelism, GPUs get a lot of advantage from the fact that now it can perform large amounts of data operations, simultaneously and multiple times.

Koalas has the disadvantage to be the library with the more distinct syntax of the ones tested, specially when it is necessary to work in PySpark dataframes. Indeed, the Koalas syntax is similar to pandas syntax. However, there was the need of working directly with PySpark since there are functions that are not implemented in Koalas. Even though, since this is a brand new technology this may improve over the time and more functionalities and documentation may be provided.

Regarding the line profiling, we concluded that, at least for the model used, for the majority of the libraries the biggest part of the pipeline computation time is spent training the model. The exception was the Koalas library, which has more computation time in prediction and testing. The prediction and testing step was computed again, but this time only for predicting and it's computation time was 0.38s. The prediction time is low but the time used to compute RegressionEvaluator for each metric is quite large.

Finally, it should be referred that a bigger dataset may have been used to due a more fair comparison between the different libraries. Indeed, a data set of about 2.5GB may not be enough to evaluate the potential of the studied libraries. Therefore, this study may be complemented in the future with other analysis, operations and a bigger datasets.

# References

[1]    Dask Development Team. *Dask: Library for dynamic task scheduling*. 2016. URL: https://dask.org.

[2]    J. Daniel. *Data Science with Python and Dask*. Manning Publications, 2019. ISBN: 9781638353546. URL: https://books.google.pt/books?id=KTkzEAAAQBAJ.

[3]    Devin Petersohn et al. "Towards Scalable Dataframe Systems". In: *CoRR* abs/2001.00888 (2020). arXiv: 2001.00888. URL: http://arxiv.org/abs/2001.00888.

[4]    Joblib Development Team. *Joblib: running Python functions as pipeline jobs*. 2020. URL: https://joblib.readthedocs.io/.

[5]    Pratik Gandhi. *Using joblib to speed up your Python pipelines*. Sept. 2020. URL: https://towardsdatascience.com/using-joblib-to-speed-up-your-python-pipelines-dd97440c653d.

[6]    Rapids Development Team. *About Rapids*. 2021. URL: https://rapids.ai/about.html.

[7]    B. Chambers and M. Zaharia. *Spark: The Definitive Guide: Big Data Processing Made Simple*. O'Reilly Media, 2018. ISBN: 9781491912300. URL: https://books.google.pt/books?id=oitLDwAAQBAJ.

[8]    *PySpark Documentation¶*. URL: https://spark.apache.org/docs/latest/api/python/index.html.

[9]    *How fast Koalas and PySpark are compared to Dask - The Databricks Blog*. Apr. 2021. URL: https://databricks.com/blog/2021/04/07/benchmark-koalas-pyspark-and-dask.html.

[10]   *Databricks XGBoost Example*. URL: https://databricks.com/notebooks/xgboost/xgboost4j-spark-example.html.

[11]  *Github Spark-xgboost.* URL: https://github.com/sllynn/spark-xgboost/blob/master/setup.py.

# Appendixes

# A    Detailed description of the used Machines

**Machine 1 - No GPU**
CPU:

- Architecture: x86_64

- CPU op-mode(s): 32-bit, 64-bit

- Byte Order: Little Endian

- CPU(s): 16

- On-line CPU(s) list: 0-15

- Thread(s) per core: 2

- Core(s) per socket: 8

- Socket(s): 1

- NUMA node(s): 1

- Vendor ID: GenuineIntel

- CPU family: 6

- Model: 85

- Model name: Intel(R) Xeon(R) CPU

- Stepping: 7

- CPU MHz: 3100.240

- BogoMIPS: 6200.48

- Hypervisor vendor: KVM

- Virtualization type: full

- L1d cache: 32K

- L1i cache: 32K

- L2 cache: 1024K

- L3 cache: 25344K

- NUMA node0 CPU(s): 0-15

RAM: 64GB

**Machine 2 - with GPU**

- Architecture: x86_64

- CPU op-mode(s): 32-bit, 64-bit

- Byte Order: Little Endian

- CPU(s): 8

- On-line CPU(s) list: 0-7

- Thread(s) per core: 2

- Core(s) per socket: 4

- Socket(s): 1

- NUMA node(s): 1

- Vendor ID: GenuineIntel

- CPU family: 6

- Model: 63

- Model name: Intel(R) Xeon(R) CPU @ 2.30GHz

- Stepping: 0

- CPU MHz: 2299.998

- BogoMIPS: 4599.99

- Hypervisor vendor: KVM

- Virtualization type: full

- L1d cache: 32K

- L1i cache: 32K

- L2 cache: 256K

- L3 cache: 46080K

- NUMA node0 CPU(s): 0-7

RAM: 30GB
GPU: NVIDIA Tesla T4 (16GB)

# B    Benchmark results

# C    Source code of the setup functions for JobLib

```python
import os
import joblib
from joblib import dump, load
from joblib import Parallel, delayed
import pandas as pd
import numpy as np
from tqdm import tqdm
import functools
# import databricks.koalas as ks
import time
import math

def get_results(benchmarks):
    """Return a pandas DataFrame containing benchmark results."""
    return pd.DataFrame.from_dict(benchmarks)

def benchmark(f, n_jobs=-1, df=None, benchmarks=None, name=None, **kwargs):
    """Benchmark the given function against the given DataFrame.

    Parameters
    ----------
    f: function to benchmark
    df: data frame
    benchmarks: container for benchmark results
    name: task name

    Returns
    -------
    Duration (in seconds) of the given operation
    """
    start_time = time.time()
```

Table 3: Computation times of the benchmark operations (in seconds).

| | Pandas | Modin | Dask | Koalas | Rapids | Joblib |
|---|---|---|---|---|---|---|
| read file | 1.097639 | 1.933903 | 0.021552 | 0.116975 | 0.440296 | 8.386324 |
| count | 0.000021 | 0.000017 | 0.161452 | 0.048637 | 0.000008 | 2.341892 |
| count index length | 0.000002 | 0.000004 | 1.84509 | 0.053336 | 0.000003 | 2.322545 |
| mean | 0.017941 | 0.090053 | 0.163021 | 0.144577 | 0.00205 | 2.328575 |
| standard deviation | 0.035824 | 0.136144 | 0.173782 | 0.102439 | 0.001807 | 4.737141 |
| mean of columns addition | 0.029497 | 0.395677 | 0.19121 | 0.172799 | 0.004861 | 2.321197 |
| addition of columns | 0.011179 | 0.315962 | 0.493689 | 1.09793 | 0.001438 | 2.600978 |
| mean of columns multiplication | 0.029217 | 0.384626 | 0.202049 | 0.141332 | 0.004349 | 2.319647 |
| multiplication of columns | 0.011427 | 0.315781 | 0.508529 | 1.015454 | 0.001162 | 2.606427 |
| value counts | 0.156386 | 0.576157 | 0.168916 | 0.835458 | 0.016973 | 2.336677 |
| mean of complex arithmetic ops | 1.504419 | 14.310768 | 0.784181 | 5.617018 | 0.112237 | 2.557495 |
| complex arithmetic ops | 1.431074 | 13.929066 | 1.032044 | 6.947288 | 0.114976 | 2.783355 |
| groupby statistics | 0.372646 | 2.109774 | 2.039868 | 0.481217 | 0.082091 | 5.423327 |
| join count | 0.177912 | 1.429141 | 2.177256 | 0.420335 | 0.015542 | 2.403003 |
| join | 0.177711 | 1.428062 | 2.215343 | 0.932832 | 0.011545 | 2.440177 |
| filtered count | 0.000005 | 0.000017 | 2.459231 | 0.069556 | 0.000009 | 0.520608 |
| filtered count index length | 0.000001 | 0.000004 | 2.660978 | 0.069706 | 0.000004 | 0.519908 |
| filtered mean | 0.003317 | 0.061901 | 2.558253 | 0.095373 | 0.001124 | 0.523736 |
| filtered standard deviation | 0.006458 | 0.058423 | 2.55637 | 0.09805 | 0.000816 | 1.067333 |
| filtered mean of columns addition | 0.00735 | 0.241343 | 2.422405 | 0.113718 | 0.00173 | 0.525848 |
| filtered addition of columns | 0.001515 | 0.185699 | 2.506612 | 0.218986 | 0.000979 | 0.577525 |
| filtered mean of columns multiplication | 0.004525 | 0.245141 | 2.175726 | 0.100009 | 0.001773 | 0.525475 |
| filtered multiplication of columns | 0.00138 | 0.191738 | 2.362341 | 0.208008 | 0.001014 | 0.576129 |
| filtered mean of complex arithmetic ops | 0.252189 | 4.605596 | 2.900001 | 4.521427 | 0.033652 | 0.578335 |
| filtered complex arithmetic ops | 0.250223 | 4.477522 | 2.987706 | 4.55545 | 0.032899 | 0.619596 |
| filtered value counts | 0.027867 | 0.192825 | 2.652701 | 0.468087 | 0.011148 | 0.532077 |
| filtered groupby statistics | 0.073863 | 0.677037 | 2.770664 | 0.283749 | 0.018422 | 1.203191 |
| filtered join count | 0.03491 | 0.474413 | 2.639927 | 0.328058 | 0.00389 | 0.544891 |
| filtered join | 0.033922 | 0.470783 | 2.7875 | 1.020233 | 0.003434 | 0.561686 |

```python
32        ret = func_parallel(f, n_jobs, df, **kwargs)
33        benchmarks['duration'].append(time.time() - start_time)
34        benchmarks['task'].append(name)
35        print(f"{name} took: {benchmarks['duration'][-1]} seconds")
36
37        return benchmarks['duration'][-1]
38
39
40  def func_parallel(f, n_jobs=-1, df=None, **kwargs):
41      n_jobs = n_jobs if n_jobs>0 else os.cpu_count()
42
43      if df is not None:
44          dim0 = df.shape[0] # numero de linhas
45          chunk = math.ceil(dim0/n_jobs) #if n_jobs else math.ceil(dim0/os.cpu_count()) #
    tamanho do chunk
46          nb_chunks = math.ceil(dim0/chunk)
47
48          index_list = [(i*chunk,min(chunk*(i+1),dim0)) for i in range(nb_chunks)]
49
50
51      if f.__name__ in ['standard_deviation']:
52          average = func_parallel(mean, n_jobs=n_jobs, df=df)
53
54          res_list = Parallel(n_jobs=n_jobs)(delayed(f)(df.iloc[index[0]:index[1],:],average, **
    kwargs) for index in index_list)
55
56          res = (np.sum([res1[0] for res1 in res_list])/(np.sum([res1[1] for res1 in res_list])
    -1))**(1/2)
57
58      elif f.__name__ in ['write_file_joblib']:
59          Parallel(n_jobs=n_jobs)(delayed(f)(df.iloc[index[0]:index[1],:], i, pickle_file=kwargs
    ['pickle_file']) for i,index in enumerate(index_list))
60          res = 1
61
62      elif f.__name__ in ['read_file_parquet']:
63          res_list=Parallel(n_jobs=n_jobs)(delayed(f)(pickle_file) for pickle_file in kwargs['
    files_list'])
64          res = pd.concat(res_list)
65
66      else:
67          res_list = Parallel(n_jobs=n_jobs)(delayed(f)(df.iloc[index[0]:index[1],:], **kwargs)
    for index in index_list)
68
69          if f.__name__ in ['count', 'count_index_length']:
70              res = sum(res_list)
71
72          elif f.__name__ in ['mean','mean_of_sum','mean_of_product','
    mean_of_complicated_arithmetic_operation']:
73              res = np.sum([res1[0] for res1 in res_list])/np.sum([res1[1] for res1 in res_list
    ])
74
75          elif f.__name__ in ['sum_columns','product_columns','complicated_arithmetic_operation'
    ]:
76              res = pd.concat(res_list)
77
78          elif f.__name__ in ['value_counts']:
79              res = functools.reduce(lambda a, b: a.add(b, fill_value=0), res_list)
80
81          elif f.__name__ in ['groupby_statistics']:
82              # mean computation
83              res = functools.reduce(lambda a, b: a.add(b, fill_value=0), res_list)
84              mean_fare = pd.DataFrame(res.fare_amt['sum']/res.fare_amt['count'], columns=['
    fare_amt_mean'])
85              mean_tip = pd.DataFrame(res.tip_amt['sum']/res.tip_amt['count'], columns=['
    tip_amt_mean'])
86
87              # std computation
88              def groupby_statistics_2(df, mean_fare, mean_tip):
89                  gb = df.groupby(by='passenger_count')
90                  df_2 = pd.DataFrame()
91                  for i, data in gb:
92                      df_2.loc[i,'fare_amt_std'] = sum((data['fare_amt'] - mean_fare.loc[i,'
    fare_amt_mean'])**2)
93                      df_2.loc[i,'tip_amt_std'] = sum((data['tip_amt'] - mean_tip.loc[i,'
    tip_amt_mean'])**2)
94                  return df_2
```

```
95
96          res_list_std = Parallel(n_jobs=n_jobs)(delayed(groupby_statistics_2)(df.iloc[index
     [0]:index[1],:], mean_fare, mean_tip) for index in index_list)
97          res_std = functools.reduce(lambda a, b: a.add(b, fill_value=0), res_list_std)
98          res_std['fare_amt_std']=(res_std['fare_amt_std']/(res['fare_amt']['count']-1))
     **(1/2)
99          res_std['tip_amt_std']=(res_std['tip_amt_std']/(res['tip_amt']['count']-1))**(1/2)
100
101          res = pd.concat([mean_fare, res_std['fare_amt_std'], mean_tip, res_std['
     tip_amt_std']],axis=1)
102          res.index.name= 'passenger_count'
103
104      elif f.__name__ in ['join_data', 'join_count']:
105          res = res_list[0]
106
107      return res
```

# D  Source code of standard operations for JobLib

```python
1  def write_file_joblib(df, i, pickle_file):
2      with open(pickle_file+'_'+str(i), 'wb') as f:
3          dump(df, f)#, compress='zlib')
4
5  def read_file_parquet(pickle_file):
6      with open (pickle_file, 'rb') as f:
7          df = load(f)
8      return df
9
10 def count(df=None): # DONE
11     return len(df)
12
13 def count_index_length(df=None): # DONE
14     return len(df.index)
15
16 def mean(df): # DONE
17     return (df.fare_amt.sum(), df.shape[0])
18
19 def standard_deviation(df, average): # DONE
20     res = (df.fare_amt-average)**2
21     return (res.sum(), res.shape[0])
22
23 def mean_of_sum(df): # DONE
24     return ((df.fare_amt + df.tip_amt).sum(), df.shape[0])
25
26 def sum_columns(df): # DONE
27     x = df.fare_amt + df.tip_amt
28     return x
29
30 def mean_of_product(df): # DONE
31     return ((df.fare_amt * df.tip_amt).sum(), df.shape[0])
32
33 def product_columns(df): # DONE
34     x = df.fare_amt * df.tip_amt
35     return x
36
37 def value_counts(df): # DONE
38     val_counts = df.fare_amt.copy().value_counts()
39     return val_counts
40
41 def complicated_arithmetic_operation(df): # DONE
42     theta_1 = df.start_lon
43     phi_1 = df.start_lat
44     theta_2 = df.end_lon
45     phi_2 = df.end_lat
46     temp = (np.sin((theta_2 - theta_1) / 2 * np.pi / 180) ** 2
47             + np.cos(theta_1 * np.pi / 180) * np.cos(theta_2 * np.pi / 180) * np.sin((phi_2 -
     phi_1) / 2 * np.pi / 180) ** 2)
48     ret = np.multiply(np.arctan2(np.sqrt(temp), np.sqrt(1-temp)),2)
49     return ret
50
51 def mean_of_complicated_arithmetic_operation(df): # DONE
52     theta_1 = df.start_lon
53     phi_1 = df.start_lat
54     theta_2 = df.end_lon
```

```
55     phi_2 = df.end_lat
56     temp = (np.sin((theta_2 - theta_1) / 2 * np.pi / 180) ** 2
57            + np.cos(theta_1 * np.pi / 180) * np.cos(theta_2 * np.pi / 180) * np.sin((phi_2 -
       phi_1) / 2 * np.pi / 180) ** 2)
58     ret = np.multiply(np.arctan2(np.sqrt(temp), np.sqrt(1-temp)),2)
59     return (ret.sum(), ret.shape[0])
60
61 def groupby_statistics(df): # DONE
62     gb = df.groupby(by='passenger_count').agg(
63       {
64          'fare_amt': ['sum', 'count'],
65          'tip_amt': ['sum', 'count']
66       }
67     )
68     return gb
69
70 def join_count(df, **kwargs): # DONE
71     return len(pd.merge(df, kwargs['other'], left_index=True, right_index=True))
72
73 def join_data(df, **kwargs): # DONE
74     return pd.merge(df, kwargs['other'], left_index=True, right_index=True)
```

# E   Source code of Experiment #2

## E.1   Dask

```
1 # Read data
2 import databricks.koalas as ks
3 def read_data():
4
5     t0 = time.time()
6     dask_data = dd.read_parquet('../../data/2009-01/', index_col='index')
7     t1 = time.time()
8     print(f'[INFO] Computation time: {t1-t0}s')
9
10     return dask_data
11
12 def pre_processing(df):
13
14     # Create variables dayofweek and hour
15     df=df.assign(dayofweek = df.trip_pickup_datetime.dt.dayofweek,
16               hour = df.trip_pickup_datetime.dt.hour)
17
18     # Creation of the variable weekend (0-weekend, 1-work day)
19     df['weekend'] = (df['dayofweek'] < 5).astype(int)
20
21     # Creation of the variables trip_duration
22     df['trip_duration'] = (df['trip_dropoff_datetime'] - df['trip_pickup_datetime'])/timedelta
       (minutes=1)
23     #df['trip_duration'] = df['trip_duration']/timedelta(seconds=1)
24
25     # Drop variables with no interest
26     df = df.drop(['trip_pickup_datetime',
27              'trip_dropoff_datetime',
28              'rate_code', 'mta_tax',
29              'store_and_forward','dayofweek', 'total_amt'], axis = 1)
30
31
32     # Corretion of the classes of the variable payment_type
33     df = df.categorize(columns=['vendor_name', 'payment_type'])
34     df['payment_type'] = df['payment_type'].replace('CASH', 'Cash')
35     df['payment_type'] = df['payment_type'].replace('CREDIT', 'Credit')
36
37
38     # Remove lines with abnormal values
39     df = df[df['fare_amt']>0]
40     df = df[df['trip_distance']>0]
41     df = df[df['trip_duration']>0]
42     df = df[df['tip_amt']>=0]
43
44     # Min-max scaling of numerical variables
45     numeric_variables = df[['passenger_count', 'trip_distance', 'start_lon', 'start_lat', '
       end_lon',
46                    'end_lat', 'surcharge', 'tip_amt', 'tolls_amt',
```

```python
                            'trip_duration', 'fare_amt']]

    scaler = MinMaxScaler()
    scaler.fit(numeric_variables)
    num_data = scaler.transform(numeric_variables)
    num_data = num_data.assign(idx=1)
    num_data = num_data.set_index(num_data.idx.cumsum()-1)
    num_data = num_data.drop('idx', axis=1)

    # One hot encondig of variables vendor_name and payment_type
    categorical_variables = df[['vendor_name', 'payment_type']]
    cat_data = dd.get_dummies(categorical_variables)
    cat_data = cat_data.assign(idx=1)
    cat_data = cat_data.set_index(cat_data.idx.cumsum()-1)
    cat_data = cat_data.drop('idx', axis=1)

    weekend = df[['weekend']]
    weekend = weekend.assign(idx=1)
    weekend = weekend.set_index(weekend.idx.cumsum()-1)
    weekend = weekend.drop('idx', axis=1)

    # Concatenate numerical variables and categorical variables
    num_cat_data = dd.merge(num_data, cat_data, left_index=True, right_index=True)
    dask_data = dd.merge(num_cat_data, weekend,  on = 'idx')#left_index=True, right_index=True
    ) #

    return dask_data

def training(dask_data):
    # Split into training and testing data
    train, test = dask_data.random_split([0.7, 0.3], random_state=0)

    # Separate labels from data
    y_train = train.fare_amt
    y_test = test.fare_amt

    del train['fare_amt']  # remove informative column from data
    del test['fare_amt']  # remove informative column from data

    model = XGBRegressor(objective='reg:squarederror')

    t0=time.time()

    model.fit(train, y_train)

    t1 = time.time()
    print(f'[INFO] Computation time: {t1-t0}')

    return model, train, y_train, test, y_test

folder_out = '2009-01/'

def save_model(model, model_name):
    io.save_obj(model,model_name)

#save_model(model, folder_out + 'Dask_model'+ datetime.now().strftime("%H%M%S"))

# Prediction and test
from dask_ml.metrics import mean_absolute_error

def pred_test(model, test, y_test):
    # Predictions
    prediction = model.predict(test)

    # Evaluation metrics computation
    y_test = np.asarray(y_test)
    prediction = np.asarray(prediction)

    # Metrics
    df_metrics = pd.DataFrame()
    df_metrics.loc[0,'max_error'] = metrics.max_error(y_test, prediction)
    df_metrics.loc[0,'R^2'] = metrics.explained_variance_score(y_test,prediction)
    df_metrics.loc[0,'MAE'] = metrics.mean_absolute_error(y_test,prediction)
    df_metrics.loc[0,'MSE'] = metrics.mean_squared_error(y_test,prediction)
    df_metrics.loc[0,'RMSE'] = np.sqrt(metrics.mean_squared_error(y_test,prediction))
```

```python
122        return df_metrics
123
124
125 def pipeline_dask(n_jobs=-1):
126        t0 = time.time()
127        # Read data
128        print('[INFO] Reading data...')
129        df = read_data()
130        print('[INFO] Data loaded.')
131
132        # Pre-processing
133        print('[INFO] Pre-processing')
134        dask_data = pre_processing(df)
135
136        #Train the model
137        print('[INFO] Training the model...')
138        t2 = time.time()
139        model, train, y_train, test, y_test = training(dask_data)
140        t3 = time.time()
141        print(f'[INFO] Model trained in {t3-t2} s.')
142
143        # Save the model
144        print('[INFO] Saving the model...')
145        save_model(model, folder_out + 'Dask_model'+ datetime.now().strftime("%H%M%S"))
146        print('[INFO] Model saved.')
147
148        # Predict and model metrics
149        print('[INFO] Predicting and computing model metrics...')
150        df_metrics = pred_test(model, test, y_test)
151
152        # Save metrics
153        io.save_obj(df_metrics,folder_out + 'Dask_model_metrics'+ datetime.now().strftime("%H%M%S"))
154        t1 = time.time()
155        print('[INFO] All metrics computed.\n[INFO] End of the pipeline.')
156        print(f'[INFO] Computation time for training the model: {t1-t0}s')
```

## E.2   Joblib

```python
1 def func_parallel(f, n_jobs=-1, df=None, **kwargs):
2        n_jobs = n_jobs if n_jobs>0 else os.cpu_count()
3
4        if df is not None:
5            dim0 = df.shape[0] # numero de linhas
6            chunk = math.ceil(dim0/n_jobs) #if n_jobs else math.ceil(dim0/os.cpu_count()) #
tamanho do chunk
7            nb_chunks = math.ceil(dim0/chunk)
8
9            index_list = [(i*chunk,min(chunk*(i+1),dim0)) for i in range(nb_chunks)]
10
11
12        if f.__name__ in ['standard_deviation']:
13            average = func_parallel(mean, n_jobs=n_jobs, df=df)
14
15            res_list = Parallel(n_jobs=n_jobs)(delayed(f)(df.iloc[index[0]:index[1],:],average, **
kwargs) for index in index_list)
16
17            res = (np.sum([res1[0] for res1 in res_list])/(np.sum([res1[1] for res1 in res_list])
-1))**(1/2)
18
19        elif f.__name__ in ['write_file_joblib']:
20            Parallel(n_jobs=n_jobs)(delayed(f)(df.iloc[index[0]:index[1],:], i, pickle_file=kwargs
['pickle_file']) for i,index in enumerate(index_list))
21            res = 1
22
23        elif f.__name__ in ['read_file_parquet']:
24            res_list=Parallel(n_jobs=n_jobs)(delayed(f)(pickle_file) for pickle_file in kwargs['
files_list'])
25            res = pd.concat(res_list)
26
27        else:
28            res_list = Parallel(n_jobs=n_jobs)(delayed(f)(df.iloc[index[0]:index[1]], **kwargs)
for index in index_list)
29
30            if f.__name__ in ['count', 'count_index_length']:
```

```python
            res = sum(res_list)

        elif f.__name__ in ['mean','mean_of_sum','mean_of_product','
    mean_of_complicated_arithmetic_operation']:
            res = np.sum([res1[0] for res1 in res_list])/np.sum([res1[1] for res1 in res_list
    ])

        elif f.__name__ in ['sum_columns','product_columns','complicated_arithmetic_operation'
    ]:
            res = pd.concat(res_list)

        elif f.__name__ in ['value_counts']:
            res = functools.reduce(lambda a, b: a.add(b, fill_value=0), res_list)

        elif f.__name__ in ['groupby_statistics']:
            # mean computation
            res = functools.reduce(lambda a, b: a.add(b, fill_value=0), res_list)
            mean_fare = pd.DataFrame(res.fare_amt['sum']/res.fare_amt['count'], columns=['
    fare_amt_mean'])
            mean_tip = pd.DataFrame(res.tip_amt['sum']/res.tip_amt['count'], columns=['
    tip_amt_mean'])

            # std computation
            def groupby_statistics_2(df, mean_fare, mean_tip):
                gb = df.groupby(by='passenger_count')
                df_2 = pd.DataFrame()
                for i, data in gb:
                    df_2.loc[i,'fare_amt_std'] = sum((data['fare_amt'] - mean_fare.loc[i,'
    fare_amt_mean'])**2)
                    df_2.loc[i,'tip_amt_std'] = sum((data['tip_amt'] - mean_tip.loc[i,'
    tip_amt_mean'])**2)
                return df_2

            res_list_std = Parallel(n_jobs=n_jobs)(delayed(groupby_statistics_2)(df.iloc[index
    [0]:index[1],:], mean_fare, mean_tip) for index in index_list)
            res_std = functools.reduce(lambda a, b: a.add(b, fill_value=0), res_list_std)
            res_std['fare_amt_std']=(res_std['fare_amt_std']/(res['fare_amt']['count']-1))
    **(1/2)
            res_std['tip_amt_std']=(res_std['tip_amt_std']/(res['tip_amt']['count']-1))**(1/2)

            res = pd.concat([mean_fare, res_std['fare_amt_std'], mean_tip, res_std['
    tip_amt_std']],axis=1)
            res.index.name= 'passenger_count'

        elif f.__name__ in ['join_data', 'join_count']:
            res = res_list[0]

    return res

def read_file_parquet(pickle_file):
    with open (pickle_file, 'rb') as f:
        df = load(f)
    return df

def mean(df): # DONE
    return (df.sum(), df.shape[0])

# Read data
def read_data():
    folder_in_out = '2009-01/'
    folder_data = '../../data/joblib_taxi_pickles/' + folder_in_out
    files_list = ['part__0', 'part__1', 'part__2', 'part__3', 'part__4', 'part__5', 'part__6',
     'part__7']
    files_list =[folder_data+file for file in files_list]

    t0 = time.time()

    joblib_data = func_parallel(read_file_parquet, n_jobs=-1, df=None, name='read file',
    files_list=files_list)

    t1 = time.time()
    print(f'[INFO] Computation time: {t1-t0}s')

    return joblib_data

# Pre-processing
```

```python
95  from datetime import timedelta
96  from sklearn.preprocessing import MinMaxScaler
97
98  def pre_processing_1(df):
99      # Create variables dayofweek and hour
100     df=df.assign(dayofweek = df.trip_pickup_datetime.dt.dayofweek,
101                  hour = df.trip_pickup_datetime.dt.hour)
102
103     # Creation of the variable weekend (0-weekend, 1-work day)
104     df['weekend'] = (df['dayofweek'] < 5).astype(int)
105
106     # Drop of variable weekend
107     df = df.drop(['dayofweek'], axis=1)
108
109     # Creation of the variables trip_duration
110     df['trip_duration'] = (df['trip_dropoff_datetime'] - df['trip_pickup_datetime'])/timedelta
        (seconds=1)
111
112     # Remove lines with abnormal values
113     df = df[df['fare_amt']>0]
114     df = df[df['trip_distance']>0]
115     df = df[df['trip_duration']>0]
116     df = df[df['tip_amt']>=0]
117
118     # One hot encondig of variables vendor_name and payment_type
119     cat_data = pd.get_dummies(df[['vendor_name', 'payment_type']])
120
121     # Target
122 #     y = df['fare_amt']
123
124     # Drop variables with no interest
125     df = df.drop(['index','trip_pickup_datetime',
126                   'trip_dropoff_datetime',
127                   'rate_code', 'mta_tax',
128                   'store_and_forward',
129                   'vendor_name', 'payment_type'], axis = 1)
130
131     # Concatenate dataframes
132     df = pd.concat([df, cat_data], axis=1)
133
134     return df
135
136
137 def pre_processing_2(df, scaler, numerical_variables):
138     # Min-max scaling of numerical variables
139 #     numerical_variables=['passenger_count', 'trip_distance', 'start_lon', 'start_lat', '
        end_lon',
140 #                          'end_lat', 'surcharge', 'tip_amt', 'tolls_amt', 'total_amt',
141 #                          'trip_duration','hour']
142
143 #     scaler = MinMaxScaler()
144 #     scaler.fit(numeric_variables)
145
146     # Concatenate numerical variables and categorical variables
147     df = pd.concat([pd.DataFrame(scaler.transform(df[numerical_variables]), columns=
        numerical_variables),
148                     df.drop(numerical_variables, axis=1).reset_index(drop=True)], axis=1)
149
150     return df
151
152
153 # Training
154 from sklearn.model_selection import train_test_split
155 from xgboost import XGBRegressor
156
157 def training(X, y):
158     # Hold out of the data
159     X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.30, random_state=0)
160
161     # Model training
162     model = XGBRegressor(objective='reg:squarederror')
163
164     # Model training
165     model.fit(X_train, y_train)
166
167     return model, X_train, y_train, X_test, y_test
```

```python
168
169
170  def save_model(model, model_name):
171      io.save_obj(model, model_name)
172
173  # Prediction and test
174
175  def pred_test(model, X_test, y_test, mean_test):
176      # Predictions
177      prediction = model.predict(X_test)
178
179      # Evaluation metrics computation
180      df_error = y_test - prediction
181
182      # Metrics
183      df_metrics = pd.DataFrame()
184      df_metrics.loc[0, 'max_error'] = max(df_error)
185      df_metrics.loc[0, 'R^2'] = sum((y_test - mean_test)**2)
186      df_metrics.loc[0, 'MAE'] = sum(abs(df_error))
187      df_metrics.loc[0, 'MSE'] = sum(df_error**2)
188      df_metrics.loc[0, 'size'] = df_error.shape[0]
189
190      return df_metrics
191
192
193  def pipeline_joblib(n_jobs=-1):
194      t0 = time.time()
195      # Read data
196      print('[INFO] Reading data...')
197      df = read_data()
198      print('[INFO] Data loaded.')
199
200      # Pre-processing 1 (parallel)
201      tp0 = time.time()
202      print('[INFO] Pre-processing 1')
203      n_jobs = n_jobs if n_jobs>0 else os.cpu_count()
204      dim0 = df.shape[0] # numero de linhas
205      chunk = math.ceil(dim0/n_jobs) #if n_jobs else math.ceil(dim0/os.cpu_count()) # tamanho do
            chunk
206      nb_chunks = math.ceil(dim0/chunk)
207      index_list = [(i*chunk,min(chunk*(i+1),dim0)) for i in range(nb_chunks)]
208
209      res_list = Parallel(n_jobs=n_jobs)(delayed(pre_processing_1)(df.iloc[index[0]:index[1],:])
             for index in index_list)
210      df = pd.concat(res_list)
211  #     X = pd.concat([r[0] for r in res_list])
212  #     y = pd.concat([r[1] for r in res_list])
213
214      tp1 = time.time()
215      print(f'[INFO] End of pre-processing 1 in {tp1-tp0}s.')
216
217      tp0 = time.time()
218      print(f'[INFO] Fit of the scaller.')
219      # Reset index (serial)
220      df = df.reset_index(drop=True)
221  #     y = y.reset_index(drop=True)
222
223      # Fit of the Min-Max scaller (serial)
224      numerical_variables=['passenger_count', 'trip_distance', 'start_lon', 'start_lat', '
            end_lon',
225                           'end_lat', 'surcharge', 'tip_amt', 'tolls_amt', #'total_amt',
226                           'trip_duration','hour','fare_amt']
227      scaler = MinMaxScaler()
228      scaler.fit(df[numerical_variables])
229
230      tp1 = time.time()
231      print(f'[INFO] End of the fit of the scaller in {tp1-tp0}s.')
232
233      # Pre-processing 2 (parallel)
234      tp0=time.time()
235      print('[INFO] Pre-processing 2')
236      dim0 = df.shape[0] # numero de linhas
237      chunk = math.ceil(dim0/n_jobs) #if n_jobs else math.ceil(dim0/os.cpu_count()) # tamanho do
            chunk
238      nb_chunks = math.ceil(dim0/chunk)
239      index_list = [(i*chunk,min(chunk*(i+1),dim0)) for i in range(nb_chunks)]
```

```
240    res_list = Parallel(n_jobs=n_jobs)(delayed(pre_processing_2)(df.iloc[index[0]:index[1],:],
        scaler=scaler, numerical_variables=numerical_variables) for index in index_list)
241    df = pd.concat(res_list)
242    tp1 = time.time()
243    print(f'[INFO] End of pre-processing 2 in {tp1-tp0}s.')
244
245    # Training (serial function, but XGBR training is parallel)
246    print('[INFO] Training the model...')
247    tp0 = time.time()
248    model, X_train, y_train, X_test, y_test = training(df.drop(['fare_amt'],axis=1), df['
        fare_amt'])
249    tp1 = time.time()
250    print(f'[INFO] Model trained in {tp1-tp0} s.')
251
252    # Save the model (serial)
253    print('[INFO] Saving the model...')
254    tp0 = time.time()
255    save_model(model, folder_out + 'Joblib_model'+ datetime.now().strftime("%H%M%S"))
256    tp1 = time.time()
257    print('[INFO] Model saved.')
258
259    # Predict and model metrics (parallel)
260    tp0 = time.time()
261    print('[INFO] Predicting and computing model metrics...')
262    mean_test = func_parallel(mean, n_jobs=-1, df=y_test)
263
264    dim0 = X_test.shape[0] # numero de linhas
265    chunk = math.ceil(dim0/n_jobs) #if n_jobs else math.ceil(dim0/os.cpu_count()) # tamanho do
        chunk
266    nb_chunks = math.ceil(dim0/chunk)
267    index_list = [(i*chunk,min(chunk*(i+1),dim0)) for i in range(nb_chunks)]
268
269
270    res_list = Parallel(n_jobs=n_jobs)(delayed(pred_test)(model,
271                                                          X_test.iloc[index[0]:index[1],:],
272                                                          y_test.iloc[index[0]:index[1]],
        mean_test) for index in index_list)
273    res = functools.reduce(lambda a, b: a.add(b, fill_value=0), res_list)
274    res = res/res['size'].values[0]
275    res['max_error'] = max([x['max_error'].values[0] for x in res_list])
276    res['R^2'] = 1 - sum([x['MSE'].values[0] for x in res_list])/sum([x['R^2'].values[0] for x
        in res_list])
277    res.loc[0,'RMSE'] = (res.loc[0,'MSE'])**(1/2)
278
279    res = res.drop(['size'],axis=1)
280
281    # Save metrics
282    io.save_obj(res,folder_out + 'Joblib_model_metrics'+ datetime.now().strftime("%H%M%S"))
283    tp1 = time.time()
284    print(f'[INFO] Metrics computed and saved in {tp1-tp0} s.')
285
286    t1 = time.time()
287    print(f'[INFO] All metrics computed.\n[INFO] End of the pipeline in {t1-t0} seconds.')
```

## E.3   Koalas

```
1
2  import pandas as pd
3  import numpy as np
4  import databricks.koalas as ks
5  from datetime import timedelta, datetime
6  from utils import io
7  import time
8
9  import matplotlib.pyplot as plt
10
11 from pyspark.ml.regression import LinearRegression
12 from pyspark.ml.feature import VectorAssembler
13 from pyspark.ml.feature import MinMaxScaler
14 from pyspark.ml.evaluation import RegressionEvaluator
15 from pyspark.ml.regression import GBTRegressor
16 from pyspark.ml.linalg import Vectors
17 from pyspark.sql.functions import udf
18 from pyspark.sql.types import DoubleType
19 from pyspark.sql.functions import monotonically_increasing_id
```

```python
20
#########################################################################
######################### Read data ####################################
#########################################################################
def read_data():
    koalas_data = ks.read_parquet('../../data/2009-01', index_col='index')
    return koalas_data


#########################################################################
######################## Pre-processing ################################
#########################################################################
def pre_processing(koalas_data):

    koalas_data=koalas_data.assign(pdayofweek = koalas_data.trip_pickup_datetime.dt.dayofweek,
                                   phour = koalas_data.trip_pickup_datetime.dt.hour)

    # Creation of the variable weekend (0-weekend, 1-work day)
    koalas_data['WEEKDAY'] =  koalas_data['pdayofweek'] //5

    # Drop of variable weekend
    koalas_data = koalas_data.drop(['pdayofweek'], axis = 1)

    # Creation of the variables trip_duration
    koalas_data['trip_duration'] = koalas_data['trip_dropoff_datetime'] - koalas_data['trip_pickup_datetime']

    # Recode same variables with different names
    koalas_data['payment_type'] = koalas_data['payment_type'].mask(koalas_data['payment_type']=='CASH', 'Cash')
    koalas_data['payment_type'] = koalas_data['payment_type'].mask(koalas_data['payment_type']=='CREDIT', 'Credit')

    # Remove lines with abnormal values
    koalas_data = koalas_data[koalas_data['fare_amt']>0]
    koalas_data = koalas_data[koalas_data['trip_distance']>0]
    koalas_data = koalas_data[koalas_data['trip_duration']>0]
    koalas_data = koalas_data[koalas_data['tip_amt']>=0]

    # One hot encondig of variables vendor_name and payment_type
    cat_data = ks.get_dummies(koalas_data[['vendor_name', 'payment_type']])

    # Transform Categorical Variables in Koalas Dataframe to Pyspark Dataframe
    pycat = cat_data.to_spark()

    # Drop variables with no interest
    num_data = koalas_data.drop(['index', 'trip_pickup_datetime', 'trip_dropoff_datetime'
                                , 'rate_code', 'mta_tax','store_and_forward', 'total_amt',
                                'vendor_name', 'payment_type'], axis = 1)

    # Transform Numerical Variables Koalas Dataframe to Pyspark Dataframe
    pynum = num_data.to_spark()

    # Adding all values exept "fare_amt" to an array called "features"
    vectorAssembler_num = VectorAssembler(inputCols =['passenger_count', 'trip_distance', 'start_lon', 'start_lat',
                                                       'end_lon', 'end_lat', 'surcharge', 'tip_amt', 'tolls_amt',
                                                       'phour', 'trip_duration', 'WEEKDAY'],
                                          outputCol = "num_features")
    pynum2 = vectorAssembler_num.transform(pynum)

    # Apply MinMaxScaler to num_features and return scaledFeatures column
    scaler = MinMaxScaler(inputCol="num_features", outputCol="scaledFeatures")
    scalerModel = scaler.fit(pynum2)
    scaledData = scalerModel.transform(pynum2)

    # Transform fare_amt to a vector to apply MinMaxScaler on target variable
    vectorAssembler_fare = VectorAssembler(inputCols =['fare_amt'], outputCol = "fare_amt_f")
    scaledData2 = vectorAssembler_fare.transform(scaledData)

    # Apply MinMaxScaler to fare_amt_f and return s_fare_amt column
    scaler2 = MinMaxScaler(inputCol="fare_amt_f", outputCol="s_fare_amt")
    scalerModel2 = scaler2.fit(scaledData2)
    scaledData3 = scalerModel2.transform(scaledData2)
```

```python
 91     # Create id column on both categorical and numerical dataframes and merge by id
 92     df1 = scaledData3
 93     df2 = pycat
 94
 95     df1 = df1.withColumn("id", monotonically_increasing_id())
 96     df2 = df2.withColumn("id", monotonically_increasing_id())
 97
 98     df3 = df2.join(df1, "id").drop("id")
 99
100     # Add categorical values and scaledFeatures into new column features
101     vectorAssembler = VectorAssembler(inputCols =['vendor_name_CMT', 'vendor_name_DDS', '
        vendor_name_VTS',
102                                                   'payment_type_Cash', 'payment_type_Credit',
103                                                   'payment_type_Dispute', 'payment_type_No Charge'
        ,
104                                                   'scaledFeatures'],        outputCol = "
        features")
105     df4 = vectorAssembler.transform(df3)
106
107     # Convert vector s_fare_amt to Double scaled_fare_amt
108     unlist = udf(lambda x: float(list(x)[0]), DoubleType())
109     pyko2 = df4.withColumn("scaled_fare_amt", unlist("s_fare_amt"))
110
111     return koalas_data, pyko2
112
113
114 ######################################################################
115 ############################ Training ################################
116 ######################################################################
117
118 def train(pyko2):
119     # Hold out of the data
120     splits = pyko2.randomSplit([0.7,0.3])
121     train_df = splits[0]
122     test_df = splits[1]
123
124     # Model training
125     gbt = GBTRegressor(featuresCol="features", labelCol="scaled_fare_amt", lossType = "squared
        ")
126
127     # Model training
128     gbt_model = gbt.fit(train_df)
129
130     return gbt_model, train_df, test_df
131
132
133 ######################################################################
134 ##################### Prediction and test ###########################
135 ######################################################################
136
137 def pred_test(gbt_model, test_df):
138     # Predictions
139     gbt_predictions = gbt_model.transform(test_df)
140
141     # Root Mean Square Error
142     gbt_evaluator_rmse = RegressionEvaluator(labelCol="scaled_fare_amt", predictionCol="
        prediction", metricName="rmse")
143     gbt_rmse = gbt_evaluator_rmse.evaluate(gbt_predictions)
144
145     # Mean Square Error
146     gbt_evaluator_mse = RegressionEvaluator(labelCol="scaled_fare_amt", predictionCol="
        prediction", metricName="mse")
147     gbt_mse = gbt_evaluator_mse.evaluate(gbt_predictions)
148
149     # Mean Absolute Error
150     gbt_evaluator_mae = RegressionEvaluator(labelCol="scaled_fare_amt", predictionCol="
        prediction", metricName="mae")
151     gbt_mae = gbt_evaluator_mae.evaluate(gbt_predictions)
152
153     # R2
154     gbt_evaluator_r2 = RegressionEvaluator(labelCol="scaled_fare_amt", predictionCol="
        prediction", metricName="r2")
155     gbt_r2 = gbt_evaluator_r2.evaluate(gbt_predictions)
156
157     # Kolas Dataframe with predictions
158     koalas_predict = ks.DataFrame(gbt_predictions)
```

```python
    # Metrics
    df_metrics = pd.DataFrame()
    df_metrics.loc[0,'max_error'] = ((koalas_predict['prediction']-koalas_predict['
        scaled_fare_amt']).abs()).max()
    df_metrics.loc[0,'R^2'] = gbt_r2
    df_metrics.loc[0,'MAE'] = gbt_mae
    df_metrics.loc[0,'MSE'] = gbt_mse
    df_metrics.loc[0,'RMSE'] = gbt_rmse

    return koalas_predict, df_metrics

####################################################################
############################# Pipeline #############################
####################################################################

def pipeline_koalas():
    t0 = time.time()

    # Read data
    tr0 = time.time()
    print('[INFO] Reading data...')
    koalas_data = read_data()
    tr1 = time.time()
    print(f'[INFO] Data loaded in {tr1-tr0} s.\n')

    # Pre processing
    tp0 = time.time()
    print('[INFO] Pre-processing')
    koalas_data, pyko2 = pre_processing(koalas_data)
    tp1 = time.time()
    print(f'[INFO] Pre-processing done in {tp1-tp0} s.\n')

    # Traning
    tt0 = time.time()
    print('[INFO] Training the model...')
    gbt_model, train_df, test_df = train(pyko2)
    tt1 = time.time()
    print(f'[INFO] Model trained in {tt1-tt0} s.\n')

    # Predict and model metrics
    tm0 = time.time()
    print('[INFO] Predicting and computing model metrics...')
    koalas_predict, df_metrics = pred_test(gbt_model, test_df)
    tm1 = time.time()
    print(f'[INFO] Predictions and all metrics computed in {tm1-tm0} s.\n')

    # Save metrics
    io.save_obj(df_metrics, 'Koalas_model_metrics'+ datetime.now().strftime("%H%M%S"))

    t1 = time.time()
    print(f'[INFO] Pipeline computed in {t1-t0} s.\n[INFO] End of the pipeline.')
    return koalas_predict, df_metrics

####################################################################
############################# Yappi ################################
####################################################################

import yappi

# Yappi Read
yappi.clear_stats()
yappi.start()
koalas_data = read_data()
yappi.stop()
stats = yappi.get_func_stats(filter_callback=lambda x: x.ttot>0.5)#.print_all()
stats.save('../outputs/callgrind.kread.prof', type = "ystat")

# Yappi Pre Processing
yappi.clear_stats()
yappi.start()
koalas_data, pyko2 = pre_processing(koalas_data)
yappi.stop()
stats = yappi.get_func_stats(filter_callback=lambda x: x.ttot>0.5)#.print_all()
stats.save('../outputs/callgrind.kprepros.prof', type = "ystat")
```

```
234  # Yappi Train
235  yappi.clear_stats()
236  yappi.start()
237  gbt_model, train_df, test_df = train(pyko2)
238  yappi.stop()
239  stats = yappi.get_func_stats(filter_callback=lambda x: x.ttot>0.5)#.print_all()
240  stats.save('../outputs/callgrind.ktrain.prof', type = "ystat")
241
242  # Yappi Pred
243  yappi.clear_stats()
244  yappi.start()
245  koalas_predict, df_metrics = pred_test(gbt_model, test_df)
246  yappi.stop()
247  stats = yappi.get_func_stats(filter_callback=lambda x: x.ttot>0.5)#.print_all()
248  stats.save('../outputs/callgrind.kpred.prof', type = "ystat")
```

## E.4   Rapids

```
1   def import_data(path = None):#"../../data/ks_taxi_parquet"
2       t0 = time.time()
3
4       try:
5           rapids_data = cudf.io.parquet.read_parquet(path)
6           global train_ind
7           train_ind = round(len(rapids_data)*0.7)
8       except:
9           print("No path provided.")
10      t1 = time.time()
11
12      print(f'[INFO] Computation time for loading the file: {t1-t0}s')
13      return rapids_data
14
15
16  def preprocessing(rapids_data):
17      t0 = time.time()
18
19      rapids_data = rapids_data.assign(pdayofweek = rapids_data.trip_pickup_datetime.dt.
    dayofweek,
20                              phour = rapids_data.trip_pickup_datetime.dt.hour)
21
22      # computation of the variable weekend
23      #zero weekend, 1 week
24      rapids_data['weekend'] = (rapids_data['pdayofweek'] < 5).astype(int)
25
26      # computation of the variable trip_duration
27      rapids_data['trip_duration'] = rapids_data['trip_dropoff_datetime'] - rapids_data['
    trip_pickup_datetime']
28      rapids_data['trip_duration'] = rapids_data['trip_duration']/timedelta(seconds=1)
29
30      # Corretion of the classes of the variable payment_type
31      rapids_data['payment_type'] = rapids_data['payment_type'].mask(rapids_data['payment_type'
    ]=='CASH', 'Cash')
32      rapids_data['payment_type'] = rapids_data['payment_type'].mask(rapids_data['payment_type'
    ]=='CREDIT', 'Credit')
33
34      # Drop concatentated timestamp columns
35      rapids_data = rapids_data.drop(['trip_pickup_datetime', 'trip_dropoff_datetime', '
    rate_code', 'mta_tax','store_and_forward','pdayofweek'], axis = 1)
36
37      # Remove lines with abnormal values
38      rapids_data = rapids_data[rapids_data['fare_amt']>0]
39      rapids_data = rapids_data[rapids_data['trip_distance']>0]
40      rapids_data = rapids_data[rapids_data['trip_duration']>0]
41      rapids_data = rapids_data[rapids_data['tip_amt']>=0]
42      rapids_data = rapids_data.drop('total_amt', axis = 1)
43
44      rapids_data = rapids_data.iloc[:,1:].reset_index()
45
46      # One hot enconding of variables vendor_name and payment_type
47      rapids_data["vendor_name"] = rapids_data["vendor_name"].astype("category")
48      rapids_data["payment_type"] = rapids_data["payment_type"].astype("category")
49
50      categorical_variables = rapids_data[['vendor_name', 'payment_type']]
51
52      # encoder = OneHotEncoder(sparse=False)
```

```python
53      # encoder.fit(categorical_variables)
54      # cat_data = encoder.transform(categorical_variables).compute()
55
56      cat_data = cudf.get_dummies(categorical_variables)
57
58      # Min-max scaler of numerical variables
59      numeric_variables = rapids_data[['passenger_count', 'trip_distance', 'start_lon', '
        start_lat', 'end_lon',
60                                      'end_lat', 'surcharge', 'tip_amt', 'tolls_amt', '
        trip_duration', 'fare_amt']]
61      scaler = MinMaxScaler()
62      scaler.fit(numeric_variables)
63      num_data = scaler.transform(numeric_variables)
64
65      num_data = num_data.rename(columns={0: 'passenger_count', 1: 'trip_distance', 2: '
        start_lon', 3: 'start_lat', 4: 'end_lon', 5: 'end_lat',
66                                      6: 'surcharge', 7: 'tip_amt', 8: 'tolls_amt', 9: '
        trip_duration', 10: 'fare_amt'})
67
68      rapids_data = cudf.concat([num_data, cat_data], axis=1) #rapids_data['fare_amt']]
69
70      t1 = time.time()
71      print(f'[INFO] Computation time for preprocessing the data: {t1-t0}s')
72      return rapids_data
73
74
75  def training(X, y):#X = rapids_data.drop('fare_amt', axis = 1), y = rapids_data['fare_amt']
76      # Hold out of the data
77      X_train = X[0:train_ind]
78      X_test = X[train_ind:]
79      y_train = y[0:train_ind]
80      y_test = y[train_ind:]
81
82      #Convertion to a DMatrix
83      X_train_DM = xgb.DMatrix(X_train.values, label=y_train.values)
84
85      # Parameters and Model training
86      param = {'objective': 'reg:squarederror'}
87
88      t0 = time.time()
89
90      bst = xgb.train(param, X_train_DM)
91
92      t1 = time.time()
93
94      print(f'[INFO] Computation time for training the model: {t1-t0}s')
95      return bst, X_train, y_train, X_test, y_test
96
97
98  from datetime import datetime
99  def save_model(model, model_name):
100     io.save_obj(model, model_name + datetime.now().strftime("%H%M%S"))
101
102
103 import cuml.metrics.regression as cuml_metrics
104 import sklearn.metrics as skl_metrics
105 # Prediction and test
106
107 def pred_test(model, X_test, y_test):
108     #Convertion to a DMatrix
109     X_test_DM = xgb.DMatrix(X_test.values, label = y_test.values)
110
111     # Predictions
112     prediction = model.predict(X_test_DM)
113
114     # Evaluation metrics computation
115     max_error = skl_metrics.max_error(y_test.to_array(), prediction)
116     variance = cuml_metrics.r2_score(y_test.to_array(), prediction)
117     mean_absolute_error = cuml_metrics.mean_absolute_error(np.float64(y_test.to_array()), np.
        float64(prediction))
118     mean_squared_error = cuml_metrics.mean_squared_error(np.float64(y_test.to_array()), np.
        float64(prediction))
119     root_mean_squared = np.sqrt(cuml_metrics.mean_squared_error(np.float64(y_test.to_array()),
         np.float64(prediction)))
120
121     print('Max Error:', max_error)
```

```
122      print('R^2:', variance)
123      print('MAE:', mean_absolute_error)
124      print('MSE:', mean_squared_error)
125      print('RMSE:', root_mean_squared)
126
127      # Metrics
128      gdf_metrics = cudf.DataFrame({'Max Error' : max_error,
129                                    'Explained Variance': variance,
130                                    'Mean absolute error': mean_absolute_error,
131                                    'Mean squared error': mean_squared_error,
132                                    'Root mean Squared error': root_mean_squared})
133
134      io.save_obj(gdf_metrics, '../outputs/Metrics/Rapids_model_metrics'+ datetime.now().
         strftime("%H%M%S"))
135      return gdf_metrics
136
137
138 def pipeline_rapids():
139      t0 = time.time()
140      # Read data
141      df = import_data("../../data/2009-01")
142
143      df = preprocessing(df)
144      # Training (serial function, but XGBR training is parallel)
145      model, X_train, y_train, X_test, y_test = training(df.drop('fare_amt', axis = 1), df['
         fare_amt'])
146
147      # Save the model (serial)
148      save_model(model, '../outputs/Models/Rapids_model')
149
150      # Predict and model metrics
151      df_metrics = pred_test(model, X_test, y_test)
152
153      t1 = time.time()
154      print(f'[INFO] Computation time to perform the pipeline: {t1-t0}s')
155      return df_metrics
156
157 pipeline_rapids()
158
159 %load_ext line_profiler
160
161 %lprun -f pipeline_rapids pipeline_rapids()
162
163
164 import yappi
165
166 yappi.clear_stats()
167
168 yappi.start()
169 import_data("../../data/2009-01")
170 yappi.stop()
171
172 stats_import_data = yappi.get_func_stats(filter_callback=lambda x: x.ttot>2.5)#.print_all()
173 stats_import_data.save('../outputs/Yappi/ystat.import_data' + datetime.now().strftime("%H%M%S"
     ) + '.prof', type = "ystat")
174
175 #stats.save('../outputs/Dask/callgrind.preproc.' + datetime.now().isoformat(), 'callgrind')
176
177
178 import yappi
179
180 df = import_data("../../data/2009-01")
181
182 yappi.clear_stats()
183
184 yappi.start()
185 df = preprocessing(df)
186 yappi.stop()
187
188 stats_preprocessing = yappi.get_func_stats(filter_callback=lambda x: x.ttot>2.5)#.print_all()
189 stats_preprocessing.save('../outputs/Yappi/ystat.preprocessing' + datetime.now().strftime("%H%
     M%S") + '.prof', type = "ystat")
190
191 #stats.save('../outputs/Dask/callgrind.preproc.' + datetime.now().isoformat(), 'callgrind')
192
193
```

```python
194  import yappi
195
196  yappi.clear_stats()
197
198  yappi.start()
199  model, X_train, y_train, X_test, y_test = training(df.drop('fare_amt', axis = 1), df['fare_amt
         '])
200  yappi.stop()
201
202  stats_training = yappi.get_func_stats(filter_callback=lambda x: x.ttot>2.5)#.print_all()
203  stats_training.save('../outputs/Yappi/ystat.training'+ datetime.now().strftime("%H%M%S") + '.
         prof', type = "ystat")
204
205  #stats.save('../outputs/Dask/callgrind.preproc.' + datetime.now().isoformat(), 'callgrind')
206
207
208  import yappi
209
210  yappi.clear_stats()
211
212  yappi.start()
213  save_model(model, 'Rapids_model_yappi')
214  yappi.stop()
215
216  stats_save_model = yappi.get_func_stats(filter_callback=lambda x: x.ttot>0.1)#.print_all()
217  stats_save_model.save('../outputs/Yappi/ystat.save_model' + datetime.now().strftime("%H%M%S")
         + '.prof', type = "ystat")
218
219  #stats.save('../outputs/Dask/callgrind.preproc.' + datetime.now().isoformat(), 'callgrind')
220
221
222  import yappi
223
224  yappi.clear_stats()
225
226  yappi.start()
227  pred_test(model, X_test, y_test)
228  yappi.stop()
229
230  stats_prediction = yappi.get_func_stats(filter_callback=lambda x: x.ttot>0.5)#.print_all()
231  stats_prediction.save('../outputs/Yappi/ystat.prediction' + datetime.now().strftime("%H%M%S")
         + '.prof', type = "ystat")
232
233  #stats.save('../outputs/Dask/callgrind.preproc.' + datetime.now().isoformat(), 'callgrind')
234
235
236  import yappi
237
238  yappi.clear_stats()
239
240  yappi.start()
241  pipeline_rapids()
242  yappi.stop()
243
244  stats_pipeline = yappi.get_func_stats(filter_callback=lambda x: x.ttot>10)#.print_all()
245  stats_pipeline.save('../outputs/Yappi/ystat.pipeline' + datetime.now().strftime("%H%M%S") + '.
         prof', type = "ystat")
246
247  #stats.save('../outputs/Dask/callgrind.preproc.' + datetime.now().isoformat(), 'callgrind')
```