# Configuration Manual

MSc Research Project
Data Analytics

## Tiernan Barry
Student ID: x20199121

School of Computing
National College of Ireland

Supervisor:     Dr. Catherine Mulwa

**National College of Ireland**

**MSc Project Submission Sheet**

**School of Computing**

| | |
|---|---|
| **Student Name:** | Tiernan Barry |
| **Student ID:** | X20199121 |
| **Programme:** | Master of Science - Data Analytics | **Year:** 2021 |
| **Module:** | Master of Science - Research Project |
| **Lecturer:** | Dr. Catherine Mulwa |
| **Submission Due Date:** | 16/08/2021 |
| **Project Title:** | A Comparative Approach between Batch and Online Machine Learning for Predicting Next-Minute Cryptocurrency Price Direction |
| **Word Count: 10,003** | **Page Count: 33** |

I hereby certify that the information contained in this (my submission) is information pertaining to research I conducted for this project. All information other than my own contribution will be fully referenced and listed in the relevant bibliography section at the rear of the project.

<u>ALL</u> internet material must be referenced in the bibliography section. Students are required to use the Referencing Standard specified in the report template. To use other author's written or electronic work is illegal (plagiarism) and may result in disciplinary action.

**Signature:** Tiernan Barry…………………………………………………………………

**Date:** 15/08/2021…………………………………………………………………

**PLEASE READ THE FOLLOWING INSTRUCTIONS AND CHECKLIST**

| | |
|---|---|
| Attach a completed copy of this sheet to each project (including multiple copies) | □ |
| **Attach a Moodle submission receipt of the online project submission,** to each project (including multiple copies). | □ |
| **You must ensure that you retain a HARD COPY of the project**, both for your own reference and in case a project is lost or mislaid. It is not sufficient to keep a copy on computer. | □ |

Assignments that are submitted to the Programme Coordinator Office must be placed into the assignment box located outside the office.

| **Office Use Only** | |
|---|---|
| Signature: | |
| Date: | |
| Penalty Applied (if applicable): | |

# Configuration Manual

Tiernan Barry, Student ID: x20199121

## Contents

# 1 Overview

This document provides a detailed, step-by-step manual for deploying this research project. Because this research utilises 3 different operating systems for each tier in the design (Visual Tier, Analytics Tier and Data Persistent Tier), this manual provides procedures for each, starting with the Analytics Tier. **Note: To replicate the results, only the Analytics Tier needs to be configured.**

# 2 Analytics Tier Configuration:

This is local desktop machine where the vast majority of analytics was developed:

- Feature selection
- Grid Search
- Batch Machine Learning
- Online Machine Learning

## 2.1 Hardware:

The following hardware is configured by default on current laptop (Analytics Tier). These are therefore not prerequisites:

- Laptop/Desktop Computer: HP Pavilion Power Laptop 15-cb0xx
- CPU/Processor: Intel(R) Core(TM) i5-7300HQ CPU @ 2.50GHz, 2496 Mhz, 4 Core(s), 4 Logical Processor(s)
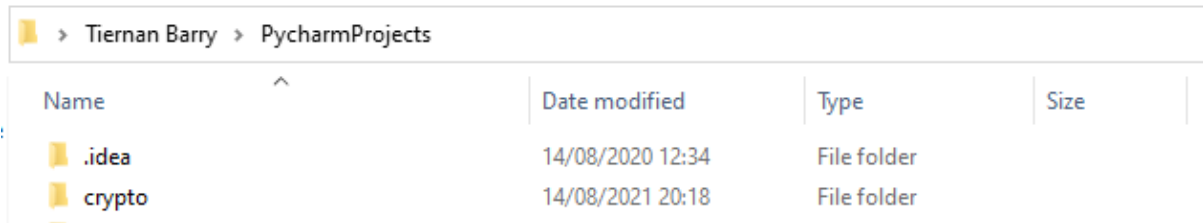- RAM: 16GB
- Graphics Card: Nividia GEFORCE GTX

## 2.2 Software:

Similarly, the following software is configured on current laptop (Analytics Tier). While the following are not prerequisites, it will make life easier for replicating results.

- Operating system: Microsoft Windows 10 Home
- Interactive Development Environment (IDE): Pycharm Community Edition 2019.2.3
  - o Note: Any other IDE will work fine too, but this document is PyCharm centric.
- Anaconda Python 3 Distribution:
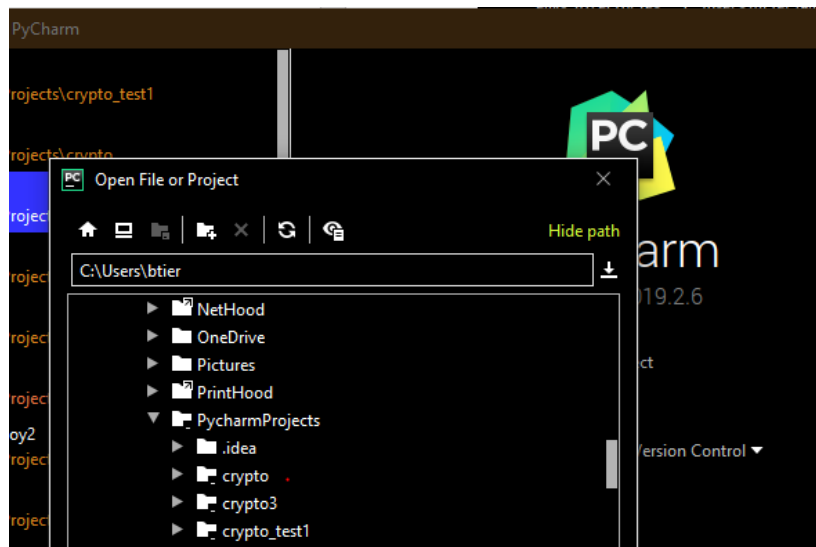  - o Version: 4.10.1

## 2.3 Open code in PyCharm/IDE:

- Open Windows Explorer, and navigate to PyCharm folder (if using another IDE, go to wherever the desired/default location is)
- Unzip the source code into the folder, until you can see the project root folder (crypto) directly under PyCharmProjects like so:



- Open the PyCharm IDE, and then open the 'crypto' project:



- Now, you will see the folder structure on the LHS of the screen. Next, we need to install packages in the following section using conda.
- Note: All code is also version controlled using a Github private repository. Please reach out if access to this is needed.

## 2.4 Create conda environment:

Once Anaconda is installed, a conda environment can now be created using the 'env.yml' file provided in the source code repository (crypto\env.yml). This will install all required packages for the Analytics Tier, and will avoid having to manually install packages individually (Anaconda, 2021):

- Launch Anaconda Prompt as follows
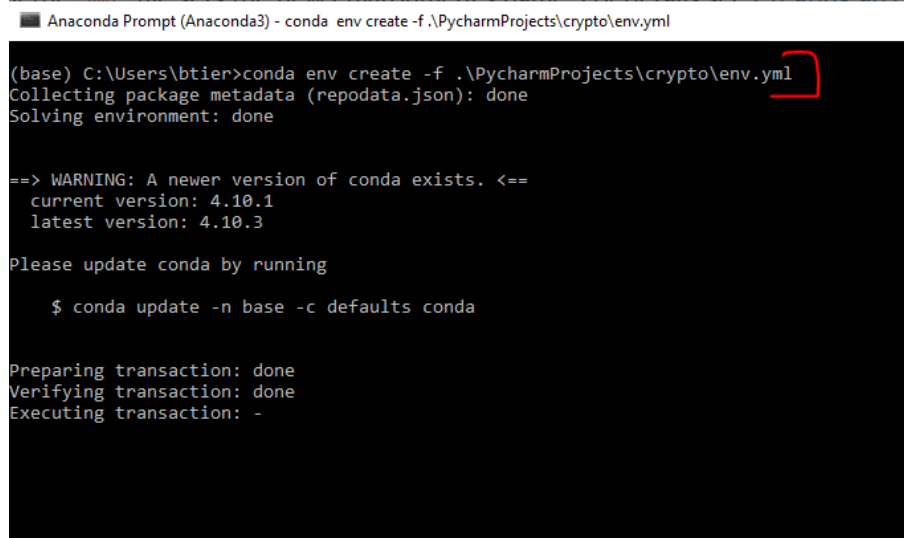- Open 'Anaconda Prompt' from the Start Menu. A terminal will pop up.

- Depending on your folder structure, run the following terminal command by providing the path to env.yml as follows:

  conda env create -f .\path\to\env.yml

- If using PyCharm it should look something like this:

  conda env create -f .\PycharmProjects\crypto\Scripts\env.yml
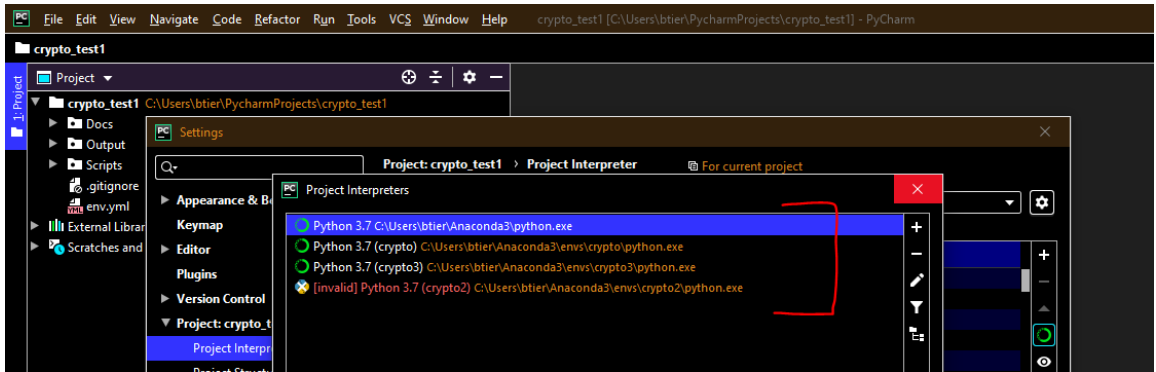
- The packages will begin installing as follows:



- Once complete (few minutes), run the following command to activate the conda environment:
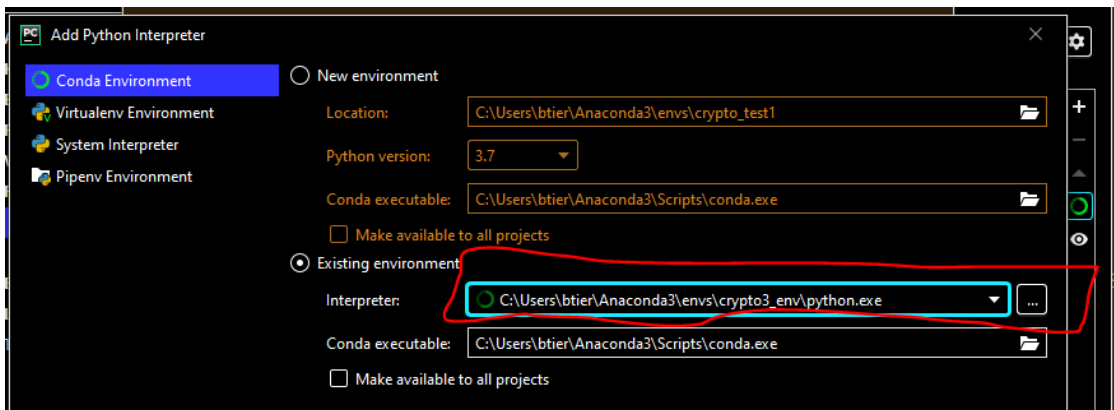
  conda activate crypto3

- crypto3 is the name of the conda environment (as per the env.yml file), which now needs to be applied as the interpreter in PyCharm as follows.

## 2.5 Set PyCharm project interpreter:

- In Pycharm, go to File > Settings > Project Interpreter
- Show all interpreters (drop-down):

- If crypto3 is not present as shown above, click the + sign to the RHS as follows. This will allow the interpreted to be added PyCharm. Find the required environment from the drop down, and apply this by clicking 'Okay'.



- Now, all packages will be available within PyCharm, and we can now deploy the analytics.

# 3   Analytics Tier Deployment:

## 3.1   Please note:

- The code is ran from the PyCharm Console
- **Because of this, the working directory automatically defaults to the root project folder of 'crypto'.**
- **Working directory being set to 'crypto' is a prerequisite for running the code. In my case, the working directory is:**
  - **C:\Users\btier\PycharmProjects\crypto**
  - **Check your working directory is 'crypto':**

## 3.2 Update config.py file with output location:

To make the results easier to reproduce, a config python script is defined in the below location. Please update this with a local output folder location (just the folder, not a file name) for writing results to. This is needed to upload files to S3:

- crypto\Scripts\Config\config.py
- Note that there is a double '\\' at the end. Please ensure this is applied.

## 3.3 Feature Selection:

Feature selection is run for all 3 alt-coin as follows:

- crypto\Scripts\FeatureSelection\prod_feature_selection.py



- As you can see at the end, the selected features are sent to AWS S3, and are then queried back during the following analyses. Please see below AWS S3 screenshot:

## 3.4 Batch Grid Search:

As discussed in the technical report, grid search is used to find the final parameters which are to be used. Grid search scripts are located in the following location for each batch model (runs for all 3 alt-coins):

- crypto\Scripts\Prediction\Batch\Classification\Production:
    - classf_grid_search_decision_tree.py
    - classf_grid_search_log_reg.py
    - classf_grid_search_random_for.py

- For each script, once they have completed (can take significant time, 10+ hours in the case of classf_grid_search_log_reg.py), the results will be loaded into the AWS S3 buckets, which are then visualised through the Visual Tier. Below is a snippet of the code running, with time series splitting and print outs:

- Results output to S3:

## 3.5   Batch Machine Learning:

For each alt-coin, 1 batch machine learning script is developed for all 3 models under the following files. Within each file, a batch decision tree, logistic regression and random forest was developed.

- crypto\Scripts\Prediction\Batch\Classification\Production:
  - \BNB\bnb_batch_pipeline.py
  - \DOGE\ doge_batch_pipeline.py
  - \ETH\eth_batch_pipeline.py
- As each of these scripts run, some results are printed out as follows, while the final results are also sent to AWS S3. Final results also feed into the Visual Tier from AWS S3.

## 3.6 Online Machine Learning:

Likewise, a script for running online learning models is developed for each alt-coin under the following paths, and as shown in following screenshot:

- crypto\Scripts\Prediction\Online\Classification\Production:
    - \BNB\bnb_online_pipeline.py
    - \DOGE\ doge_online_pipeline.py
    - \ETH\eth_online_pipeline.py

- As can be seen in above example of BinanceCoin, the online learning is running and predicting at each time step (1-minute) and printing to console. Once it is done, the results export to AWS S3:

# 4   Visual Tier Configuration:

As discussed in the technical report, this is hosted on an Amazon EC2 cloud instance. Please note the following:

- To access the Visual Tier dashboard, you only need to open the following URL link to see the results (this is now kept running until grading is complete):

  http://ec2-3-236-15-25.compute-1.amazonaws.com:3838/crypto/

- The source code is in the below R script, and looks as follows:
  crypto/Scripts/Results/Shiny/crypto/app.r

```
←  →  C   ⚠ Not secure | ec2-3-236-15-25.compute-1.amazonaws.com:8787

  R      File  Edit  Code  View  Plots  Session  Build  Debug  Profile  Tools  Help

  📋 ▾ | 📄 | 🎁 ▾ | 🔲 🔲 | 🖨 | ▶ Go to file/function  | 🔳 ▾ Addins ▾

  app.r ×   Untitled1* ×   concept_drift_aggregate.r* ×   Untitled3* ×   Untitled4* ×   Untitled6* ×

  ←→ | 🔳 | 🔲 | 🔍 ⚡ ▾ | 🔲

    1 ▾ ##########################################################
    2   # DES: Shiny dashboard for Msc
    3   #       restart shiny: sudo systemctl restart shiny-server
    4   # BY: Tiernan Barry
    5 ▾ ##########################################################
    6
    7   library("aws.s3")
    8   library("shiny")
    9   library("zoo")
   10   library("shinydashboard")
   11   library("dashboardthemes")
   12   library("tidyverse")
   13   library("utils")
   14   library("plotly")
   15   library("DT")
   16   library("corrplot")
   17   library("ggcorrplot")
   18   library("grDevices")
   19   library("stats")
   20   source("/home/rstudio/Analytics/crypto/Scripts/HelperFunctions/GetAWSS3Bucket.r")
   21
   22 ▾ ##########################################################
   23   # Define UI
   24 ▾ ##########################################################
   25
   26 ▸ DefineUserInterface <- function(){⬅⮕}
  531
  532 ▾ ##########################################################
  533   # Define Server
  534 ▾ ##########################################################
  535
  536 ▸ DefineServer <- function(input, output){⬅⮕}
 2375   |
 2376 ▾ ##########################################################
 2377   # Launch app
 2378 ▾ ##########################################################
 2379
```

- To access the RStudio login, go to below link (password is needed, please reach out for details):

  http://ec2-3-236-15-25.compute-1.amazonaws.com:8787/

- However, for completeness, this section will also outline the required steps taken to set up an AWS EC2 instance, and also to install R and RStudio.
- **Note: if required and/or easier, please reach out for gaining user access to the exact EC2 instances used in this project, and to avoid setting everything up from scratch. These currently need my AWS credentials to get the SSH address and keys.**

## 4.1 Setting up and SSH into an AWS EC2 instance:

- Log into AWS management Console (or sign up first)



- Click into EC2, and click Launch Instance
- Choose AMI: Ubuntu Server 20.04 LTS (HVM), SSD Volume Type
- Instance type: t.2 large
- Next: Configure Instance (skip)
- Add storage: 16GB of Elastic Block Storage (EBS).
- Configure Security Group:
  - Create and apply new security group as follows:
  - Apply 3838 port to be available to all IP addresses (this allows anyone to open the Visual Tier dashboard)

Inbound rules for sg-067c830f5eff605a7 (Selected security groups: sg-067c830f5eff605a7)

| Type ⓘ | Protocol ⓘ | Port Range ⓘ | Source ⓘ |
|---|---|---|---|
| HTTP | TCP | 80 | 86.40.55.206/32 |
| SSH | TCP | 22 | 86.40.55.206/32 |
| Custom TCP Rule | TCP | 8787 | 86.40.55.206/32 |
| Custom TCP Rule | TCP | 3838 | 0.0.0.0/0 |
| Custom TCP Rule | TCP | 3838 | ::/0 |
| RDP | TCP | 3389 | 86.40.55.206/32 |

  - I suggest allows custom IP addresses (ie. only your own) for other ports (ie RStudio, SSH, etc) as shown above.
- Review and launch
- Create new or use existing key pair:

16

o   Ensure you save the path as we need this to SSH into EC2



- Now, we need to SSH into the server.
- Navigate to EC2 instance console to get the IP address



- By clicking into the instance ID of the visual tier, we can now see the following, and can copy the IP address which has just been generated:

- With the IP address and the key pair saved, we can now SSH into the EC2 instance using a command as below:

ssh -i .\Scripts\Keys\tbarry_msc_key.pem ubuntu@ec2-3-236-15-25.compute-1.amazonaws.com

- Click yes when prompted, and you will be logged in as follows:



## 4.2   Installing R and RStudio onto Ubuntu AWS EC2

Following this very useful guide, R and RStudio can be installed (Zwitch, 2013). A few slight adjustments are made to ensure latest versions are available (crypto\Docs\rserver_install.txt):

```
# notes for R studio install
# source: https://www.r-bloggers.com/2013/04/instructions-for-installing-using-r-
on-amazon-ec2/

sudo useradd rstudio
sudo mkdir /home/rstudio
sudo passwd rstudio
sudo chmod -R 0777 /home/rstudio

#Update all files from the default state
sudo apt-get update
sudo apt-get upgrade

#Add CRAN mirror to custom sources.list file using vi
sudo vi /etc/apt/sources.list.d/sources.list

#Add following line (or your favorite CRAN mirror)
deb http://lib.stat.cmu.edu/R/CRAN/bin/linux/ubuntu precise/

#Update files to use CRAN mirror
#Don't worry about error message
sudo apt-get update

#Install latest version of R
#Install without verification
sudo apt-get install r-base

#Install in order to use RCurl & XML
sudo apt-get install libcurl4-openssl-dev
sudo apt-get install libxml2-dev

#Install a few background files
sudo apt-get install gdebi-core
sudo apt-get install libapparmor1

#Change to a writeable directory
#Download & Install RStudio Server
cd /tmp



wget https://download2.rstudio.org/server/bionic/amd64/rstudio-server-1.4.1106-
amd64.deb
sudo gdebi rstudio-server-1.4.1106-amd64.deb

# run in terminal for packages
# https://stackoverflow.com/questions/55855898/installing-aws-s3-r-package
sudo apt-get install -y build-essential libssl-dev libxml2-dev libcurl4-openssl-
dev
```

```
# restart shiny service
sudo systemctl restart shiny-server
```

- Once RStudio is installed, we can log into RStudio through a web browser using the IP address, followed by the RStudio port of 8787.

http://ec2-3-236-15-25.compute-1.amazonaws.com:8787/

**Sign in to RStudio**

Username:

rsudio

Password:

☐ Stay signed in when browser closes

You will automatically be signed out after 60 minutes of inactivity.

**Sign In**

## 4.3 Get code onto server (clone from Github)

- Code was managed through Github in a private repo. Please request access to repository if needed to clone the code into the Visual Tier.
- Once logged into RStudio, you can use the terminal plugin on the bottom left to run linux commands.
- Create an Analytics folder:
  cd ~
  mkdir Analytics

```
rstudio@ip-172-31-72-239:~$ ls
Analytics  R
```

- Clone Github repository (access needs to be granted)

```
rstudio@ip-172-31-72-239:~/Analytics/crypto$ ls -l
total 20
drwxr-xr-x  2 rstudio rstudio 4096 Jun 22 20:34 Docs
drwxr-xr-x  4 rstudio rstudio 4096 Jun 22 14:05 Input
drwxr-xr-x  2 rstudio rstudio 4096 Jul 17 11:52 Output
drwxr-xr-x 11 rstudio rstudio 4096 Aug 15 00:16 Scripts
-rw-r--r--  1 rstudio rstudio 2612 Aug 15 00:16 env.yml
rstudio@ip-172-31-72-239:~/Analytics/crypto$ 
```

## 4.4 Installing and Configuring Shiny service

The Shiny service needs to be configured to enable Shiny apps to be published to the web. This is done by configuring the following file as follows:

- /etc/shiny-server/shiny-server.conf
- Change the default settings to the ones provided below (run as, site dir).

```
# Instruct Shiny Server to run applications as the user "shiny"
run_as :HOME_USER: rstudio;

# Define a server that listens on port 3838
server {
  listen 3838;

  # Define a location at the base URL
  location / {

    # Host the directory of Shiny Apps stored in this directory
    #site_dir /srv/shiny-server;
    site_dir /home/rstudio/Analytics/crypto/Scripts/Results/Shiny;

    # Log all Shiny output to files in this directory
    log_dir /var/log/shiny-server;

    # When a user visits the base URL rather than a particular application,
    # an index of the applications available in this directory will be shown.
    directory_index on;
  }
}
```

- Save and ensure it is correct in server:

```
rstudio@ip-172-31-72-239:~/Analytics/crypto$ cat /etc/shiny-server/shiny-server.conf
# Instruct Shiny Server to run applications as the user "shiny"
run_as :HOME_USER: rstudio;

# Define a server that listens on port 3838
server {
  listen 3838;

  # Define a location at the base URL
  location / {

    # Host the directory of Shiny Apps stored in this directory
    #site_dir /srv/shiny-server;
    site_dir /home/rstudio/Analytics/crypto/Scripts/Results/Shiny;

    # Log all Shiny output to files in this directory
    log_dir /var/log/shiny-server;

    # When a user visits the base URL rather than a particular application,
    # an index of the applications available in this directory will be shown.
    directory_index on;
  }
}
rstudio@ip-172-31-72-239:~/Analytics/crypto$ []
```

## 4.5   Install R Packages:

Once logged into an R session, the following packages must be installed first in the R console:

```
install.packages('curl')
install.packages('httr')
install.packages('xml2')
```

Then, the following packages are needed to be installed (using the above command in R) launch the web-app:

install.packages("aws.s3")
install.packages("shiny") # install shiny R package now
install.packages("zoo")
install.packages("shinydashboard")
install.packages("dashboardthemes")
install.packages("tidyverse")
library("utils") # should be installed already – can skip this
install.packages("plotly")
install.packages("DT")
install.packages("corrplot")
install.packages("ggcorrplot")
install.packages("grDevices")

library("stats") # should be installed already – can skip this

## 4.6 Launch Visual Tier Dashboard (Shiny App):

- First, from the command line, restart the Shiny Service:

  sudo systemctl restart shiny-server

- Now, we can navigate to the Web App from a browser by changing the port to 3838 as follows:

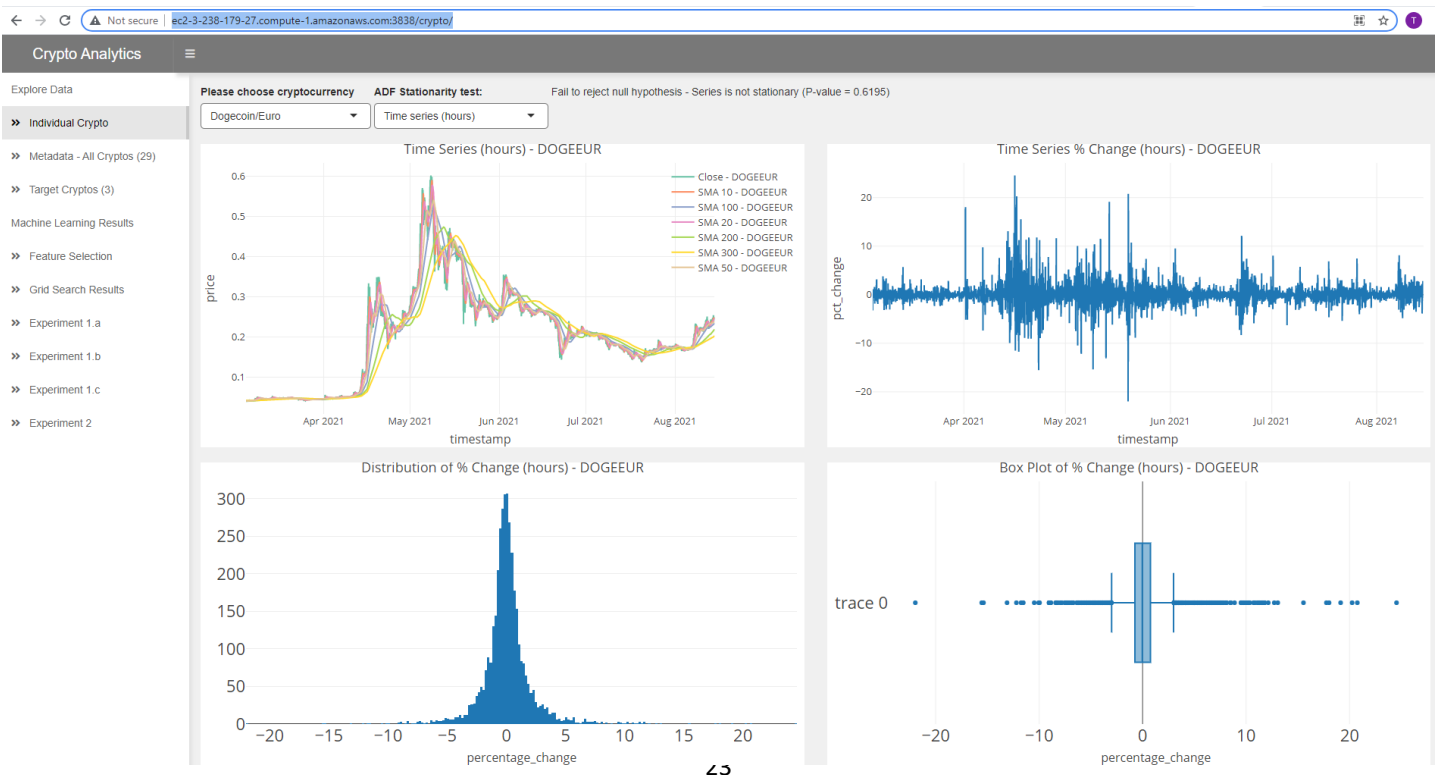  http://ec2-3-238-179-27.compute-1.amazonaws.com:3838/



- Now, we can click on 'crypto' to open the app and complete the URL link:

  http://ec2-3-238-179-27.compute-1.amazonaws.com:3838/crypto/

# 5   Data Persistent Tier Configuration:

Likewise, this is also hosted on an Amazon EC2 cloud instance, and combines a number of AWS Services as discussed here. Similar to the Visual Tier, please reach out if access credentials are required for logging into the exact instance, or AWS account (due to the large number of steps, etc). The following AWS Services need to be configured (useful Youtube video here: (TotalCloud, 2020)):

- AWS EC2
- AWS CloudFormation
- AWS Dynamo DB
- Apply Tag to AWS EC2 instance

## 5.1   Create a new EC2 Ubuntu Instance:
- Follow same steps as 4.1.
- Log into the EC2 instance

| Instances (1/5) Info | | | | |
|---|---|---|---|---|
| **Name** ▽ | **Instance ID** | **Instance state** ▽ | **Instance type** ▽ | **Status check** |
| ☑ data_persistent_tier_python_server | i-02f7031eaa4700ea4 | ⊖ Stopped ⊕⊖ | t2.large | – |
| ☐ rstudio_server | i-0dd66cf5f2c594beb | ⊖ Stopped ⊕⊖ | t2.medium | – |
| ☐ python_server | i-06f4e3c80cf997aba | ⊖ Stopped ⊕⊖ | t2.medium | – |
| ☐ r-shiny-server | i-001565ee2ec75ef39 | ⊖ Stopped ⊕⊖ | t2.large | – |
| ☐ visual_tier_rstudio_server | i-0a304250b1d5c44a4 | ⊘ Running ⊕⊖ | t2.large | ⊘ 2/2 checks passed |

## 5.2   Create CloudFormation template:
- Create Stack
- Using the below S3 URL, load the instance scheduler template:

https://s3.amazonaws.com/solutions-reference/aws-instance-scheduler/latest/instance-scheduler.template

- Apply a stack name, RServerScheduler is used in this case
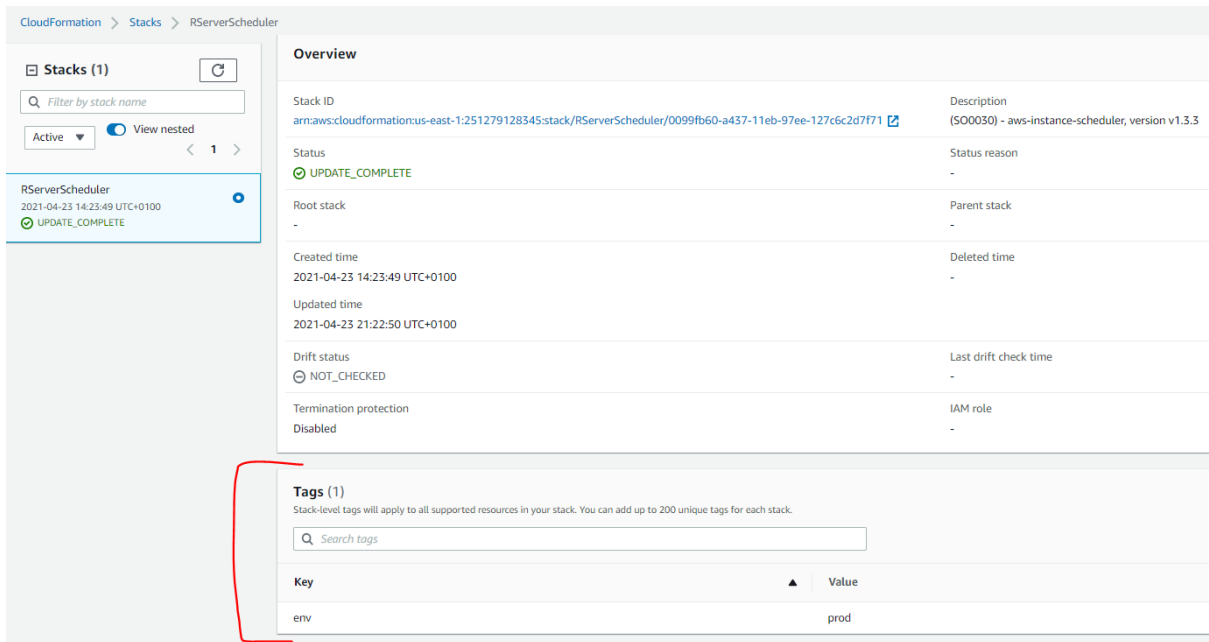- The following parameters are applied when setting up:

# RServerScheduler

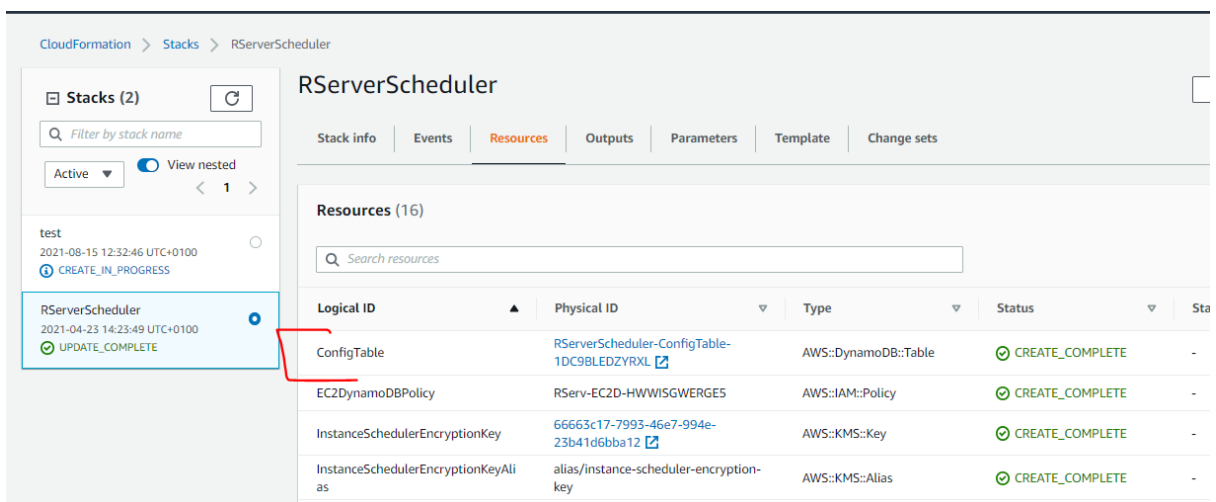Stack info | Events | Resources | Outputs | **Parameters** | Template | Change se

## Parameters (17)

Q Search parameters

| Key ▲ | Value |
|---|---|
| CreateRdsSnapshot | No |
| CrossAccountRoles | - |
| DefaultTimezone | UTC |
| LogRetentionDays | 30 |
| MemorySize | 128 |
| Regions | - |
| ScheduleLambdaAccount | Yes |
| ScheduleRdsClusters | No |
| ScheduledServices | EC2 |
| SchedulerFrequency | 2 |
| SchedulingActive | Yes |
| SendAnonymousData | Yes |
| StartedTags | started_schedule=true |
| StoppedTags | stopped_schedule=true |
| TagName | Schedule |
| Trace | No |
| UseCloudWatchMetrics | No |

- Once Stack is created, open the following config table:



## 5.3   Create Schedule in Dynamo DB:

- Open the following page:

- Here, we need to create the above files (use/edit example templates provided):
  - Period:



  - Schedule:

## 5.4 Apply Instance Scheduling tag to EC2 instance:

- Go to EC2 dashboard and open the EC2 tags as follows:



- Apply the Key/Value of the scheduler which was created in previous steps:

- This is now scheduled to launch at the defined times.
- More information on Aws instance scheduling can be found here (Amazon, 2020)

## 5.5 Cron-Job configuration:

- Similar to the Visual tier, the code is managed using Github to replicate the code.
- Once the code is available in the Data Persistent Tier under the Analytics directory, we then schedule a shell script to run using cron.

- Cron job config:

```
ubuntu@ip-172-31-13-246:~/Analytics/crypto$ sudo crontab -l
# Edit this file to introduce tasks to be run by cron.
#
# Each task to run has to be defined through a single line
# indicating with different fields when the task will be run
# and what command to run for the task
#
# To define the time you can provide concrete values for
# minute (m), hour (h), day of month (dom), month (mon),
# and day of week (dow) or use '*' in these fields (for 'any').
#
# Notice that tasks will be started based on the cron's system
# daemon's notion of time and timezones.
#
# Output of the crontab jobs (including errors) is sent through
# email to the user the crontab file belongs to (unless redirected).
#
# For example, you can run a backup of all your user accounts
# at 5 a.m every week with:
# 0 5 * * 1 tar -zcf /var/backups/home.tgz /home/
#
# For more information see the manual pages of crontab(5) and cron(8)
#
# m h  dom mon dow   command
50 18 * * * /home/ubuntu/Analytics/crypto/Scripts/main.sh
```

- As seen above, there is 1 shell script which runs at 18:50 system time.
- Therefore, once the instance is launched automatically using AWS instnace scheduling, 5 minutes later the above cron job kicks off and collects and loads latest data to S3.
- The shell script is as follows, and is provided in the code:

```
ubuntu@ip-172-31-13-246:~/Analytics/crypto/Scripts$ cat main.sh
#!/bin/sh


######################################################################
# DES: Shell script to be ran via cronjob for automating project
# BY: Tiernan Barry
######################################################################


# set WD as Output in crypto folder
cd /home/ubuntu/Analytics/crypto


##########################################
# Get data
##########################################


# get API data
PYTHONPATH=$(pwd) python3 /home/ubuntu/Analytics/crypto/Scripts/GetData/GetBinanceData/get_binance_data.py
PYTHONPATH=$(pwd) python3 /home/ubuntu/Analytics/crypto/Scripts/GetData/GetBinanceData/put_binance_data_s3.py
PYTHONPATH=$(pwd) python3 /home/ubuntu/Analytics/crypto/Scripts/GetData/GetBinanceData/get_daily_prices.py
PYTHONPATH=$(pwd) python3 /home/ubuntu/Analytics/crypto/Scripts/GetData/GetBinanceData/get_combined_dataset.py
PYTHONPATH=$(pwd) python3 /home/ubuntu/Analytics/crypto/Scripts/ExploreData/Classification/get_multiclass_data.py


##########################################
# Explore data
##########################################


# metadata
PYTHONPATH=$(pwd) python3 /home/ubuntu/Analytics/crypto/Scripts/ExploreData/Metadata/get_metadata_table.py
PYTHONPATH=$(pwd) python3 /home/ubuntu/Analytics/crypto/Scripts/ExploreData/Metadata/get_metadata_table2.py
PYTHONPATH=$(pwd) python3 /home/ubuntu/Analytics/crypto/Scripts/ExploreData/Metadata/analyse_metadata_table.py


# correlation matrices
PYTHONPATH=$(pwd) python3 /home/ubuntu/Analytics/crypto/Scripts/ExploreData/correlation_matrix.py


##########################################
# Deploy ML
##########################################


# Batch:
# PYTHONPATH=$(pwd) python3 /home/ubuntu/Analytics/crypto/Scripts/ExploreData/correlation_matrix.py
```

- This runs every day, ensuring a live, voluminous dataset.
- More information on using the linux Cron service found here (Mehra, 2017)

# 6  References

Amazon. (2020). *AWS Instance Scheduler*. Retrieved from https://docs.aws.amazon.com/:
      https://docs.aws.amazon.com/solutions/latest/instance-scheduler/security.html

Anaconda. (2021). *Managing environments - Creating an environment from an
      environment.yml file*. Retrieved from https://conda.io/projects/conda/en/latest/user-
      guide/tasks/manage-environments.html#creating-an-environment-from-an-
      environment-yml-file

Mehra, A. (2017). *The Ultimate Crontab Cheatsheet*. Retrieved from https://devhints.io/:
      https://www.codementor.io/@akul08/the-ultimate-crontab-cheatsheet-5op0f7o4r

TotalCloud. (2020). *AWS Instance Scheduler | Step by Step tutorial to Start and Stop EC2
      and RDS Instances*. Retrieved from https://www.youtube.com/watch?v=oooxZsz0hS4

Zwitch, R. (2013). *Instructions for Installing & Using R on Amazon EC2*. Retrieved from
      https://www.r-bloggers.com/2013/04/instructions-for-installing-using-r-on-amazon-
      ec2/