

# **WATER QUALITY PREDICTION**

## **ABSTRACT:**

We all know water is one of the most essential resource for our living. But as the development is increasing, we are exploiting water by wasting it and treating it with harmful materials which makes water impure and unfit for use. This is the reason it is very important to know the quality of water. This project is based on water quality prediction. In this project, water quality index (WQI) and quality status of water is predicted through some parameters that affects water quality. In this notebook we have performed Data Cleaning steps and did Exploratory Data Analysis. Then we have did some calculations as the data does not contain the column which can be used for prediction. Then we have created 2 models for prediction. The first model is Linear Regression model and The second model is Logistic Regression model. We have done this project with pyspark and other spark packages.

## **ABOUT THE DATA:**

The data contains water quality parameters of different rivers of India. There are 8 parameters and each parameter is the average values measured over a period of time.

## **AIM:**

The main aim of this project is to predict the water quality index , the quality status of the water and various parameter influencing it, using spark machine learning packages.

## SPARK & PYSPARK INSTALLATION:

- Install packages required for Spark

```
sudo apt install default-jdk  
Scala git -y
```

- Verify the installed packages

```
java -version;  
javac -version;  
scala -version;  
git --version;
```

- Use the wget command and the direct link to download spark

```
wget https://downloads.apache.org/spark/spark-3.0.1/spark-3.0.1-bin-hadoop2.7.tgz
```

- Now extract using tar

```
tar xvf spark-*
```

- Use this command to move the unpacked directory

```
sudo mv spark-3.0.1-bin-hadoop2.7 /opt/spark
```

- Now configure spark environment using these commands

```
echo "export SPARK_HOME=/opt/spark" && ~/.profile  
echo "export  
PATH=$PATH:$SPARK_HOME/bin:$SPARK_HOME/sbin"  
&& ~/.profile
```

```
echo "export PYSPARK_PYTHON=/usr/bin/python3" &&>
~/.profile
```

- Open profile using nano and paste these commands

Nano .profile

```
export SPARK_HOME=/opt/spark
export
PATH=$PATH:$SPARK_HOME/bin:$SPARK_HOME/sbin
export PYSPARK_PYTHON=/usr/bin/python3
```

- Start standalone spark master server

Start-master.sh

- To view the Spark Web user interface use

<http://127.0.0.1:8080/>

- Test spark shell

spark-shell

- Install pyspark using

Python3 pip install pyspark

- Test pyspark

pyspark

```
Activities Terminal Nov 9 21:52 hdoop@srn-pc ~/sparkpro

hdoop@srn-pc:~$ wget https://downloads.apache.org/spark/spark-3.3.1/spark-3.3.1-bin-hadoop3.tgz
--2022-11-05 12:56:57-- https://downloads.apache.org/spark/spark-3.3.1/spark-3.3.1-bin-hadoop3.tgz
Resolving downloads.apache.org (downloads.apache.org)... failed: Temporary failure in name resolution.
wget: unable to resolve host address 'downloads.apache.org'

hdoop@srn-pc:~$ wget https://downloads.apache.org/spark/spark-3.3.1/spark-3.3.1-bin-hadoop3.tgz
--2022-11-05 12:57:01-- https://downloads.apache.org/spark/spark-3.3.1/spark-3.3.1-bin-hadoop3.tgz
Resolving downloads.apache.org (downloads.apache.org)... 135.181.214.104, 88.99.95.219, 2a01:4f8:10a:201a::2, ...
Connecting to downloads.apache.org (downloads.apache.org)|135.181.214.104|:443... connected.
HTTP request sent, awaiting response... 200 OK
Length: 299350810 (285M) [application/x-gzip]
Saving to: 'spark-3.3.1-bin-hadoop3.tgz'

spark-3.3.1-bin-had 100%[=====] 285.48M 117KB/s in 34m 33s

2022-11-05 13:31:35 (141 KB/s) - 'spark-3.3.1-bin-hadoop3.tgz' saved [299350810/299350810]

hdoop@srn-pc:~$ tar xvf spark-*.tar.gz
spark-3.3.1-bin-hadoop3/
spark-3.3.1-bin-hadoop3/LICENSE
spark-3.3.1-bin-hadoop3/NOTICE
spark-3.3.1-bin-hadoop3/R/
spark-3.3.1-bin-hadoop3/R/lib/
spark-3.3.1-bin-hadoop3/R/lib/sparkR/
spark-3.3.1-bin-hadoop3/R/lib/sparkR/DESCRIPTION
spark-3.3.1-bin-hadoop3/R/lib/sparkR/INDEX
spark-3.3.1-bin-hadoop3/R/lib/sparkR/Meta/
spark-3.3.1-bin-hadoop3/R/lib/sparkR/Meta/Rd.rds
spark-3.3.1-bin-hadoop3/R/lib/sparkR/Meta/Features.rds
spark-3.3.1-bin-hadoop3/R/lib/sparkR/Meta/Research.rds
spark-3.3.1-bin-hadoop3/R/lib/sparkR/Meta/Links.rds
spark-3.3.1-bin-hadoop3/R/lib/sparkR/Meta/NSInfo.rds
spark-3.3.1-bin-hadoop3/R/lib/sparkR/Meta/package.rds
spark-3.3.1-bin-hadoop3/R/lib/sparkR/Meta/vignette.rds
spark-3.3.1-bin-hadoop3/R/lib/sparkR/NAMESPACE
spark-3.3.1-bin-hadoop3/R/lib/sparkR/R/
spark-3.3.1-bin-hadoop3/R/lib/sparkR/R/sparkR
spark-3.3.1-bin-hadoop3/R/lib/sparkR/R/sparkR.rdb
spark-3.3.1-bin-hadoop3/R/lib/sparkR/R/sparkR.rdx
spark-3.3.1-bin-hadoop3/R/lib/sparkR/doc/
spark-3.3.1-bin-hadoop3/R/lib/sparkR/doc/index.html
spark-3.3.1-bin-hadoop3/R/lib/sparkR/doc/sparkR-vignettes.R
spark-3.3.1-bin-hadoop3/R/lib/sparkR/doc/sparkR-vignettes.Rmd
spark-3.3.1-bin-hadoop3/R/lib/sparkR/doc/sparkR-vignettes.html
spark-3.3.1-bin-hadoop3/R/lib/sparkR/help/
spark-3.3.1-bin-hadoop3/R/lib/sparkR/help/antIndex
spark-3.3.1-bin-hadoop3/R/lib/sparkR/help/sparkR.rdb
spark-3.3.1-bin-hadoop3/R/lib/sparkR/help/sparkR.rdx
spark-3.3.1-bin-hadoop3/R/lib/sparkR/help/aliases.rds
```

```
Activities Terminal Nov 9 21:53 hdoop@srn-pc ~/sparkpro

hdoop@srn-pc:~$ echo "export SPARK_HOME=/opt/spark" >> ~/.profile
hdoop@srn-pc:~$ echo "export PATH=$PATH:$SPARK_HOME/bin:$SPARK_HOME/sbin" >> ~/.profile
hdoop@srn-pc:~$ echo "export PYTHONPATH=$PYTHONPATH:/usr/bin/python3" >> ~/.profile
hdoop@srn-pc:~$ nano ~/.profile
hdoop@srn-pc:~$ source ~/.profile
hdoop@srn-pc:~$ start-master.sh
starting org.apache.spark.deploy.master.Master, logging to /opt/spark/logs/spark-hadoop-org.apache.spark.deploy.master.Master-1-srn-pc.out
hdoop@srn-pc:~$ start-worker.sh spark://master:8080
starting org.apache.spark.deploy.worker.Worker, logging to /opt/spark/logs/spark-hadoop-org.apache.spark.deploy.worker.Worker-1-srn-pc.out
hdoop@srn-pc:~$ start-worker.sh spark://srn-pc:7077
starting org.apache.spark.deploy.worker.Worker, logging to /opt/spark/logs/spark-hadoop-org.apache.spark.deploy.worker.Worker-1-srn-pc.out
hdoop@srn-pc:~$ spark-shell
Setting default log level to "WARN".
To adjust logging level use sc.setLogLevel(newLevel). For SparkR, use setLogLevel(newLevel).
22/11/05 13:45:55 WARN NativeCodeLoader: Unable to load native-hadoop library for your platform... using builtin-java classes where applicable
Spark context Web UI available at http://srn-pc:4040
Spark context available as 'sc' (master = local[*]), app id = local-1667636150755).
Spark session available as 'spark'.
Welcome to

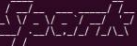
Spark version 3.3.1

Using Scala version 2.12.15 (OpenJDK 64-Bit Server VM, Java 1.8.0_342)
Type in expressions to have them evaluated.
Type :help for more information.

scala> :q
hdoop@srn-pc:~$ pyspark
Python 3.8.10 (default, Jun 22 2022, 20:18:18)
[GCC 9.4.0] on linux
Type "help", "copyright", "credits" or "license()" for more information.
Setting default log level to "WARN".
To adjust logging level use sc.setLogLevel(newLevel). For SparkR, use setLogLevel(newLevel).
22/11/05 13:46:14 WARN HiveConf: HiveConf of name hive.metastore.wm.default.pool.size does not exist
22/11/05 13:46:14 WARN HiveConf: HiveConf of name hive.llap.task.scheduler.preempt.independent does not exist
22/11/05 13:46:14 WARN HiveConf: HiveConf of name hive.llap.output.format.arrow does not exist
22/11/05 13:46:14 WARN HiveConf: HiveConf of name hive.tez.llap.mn.reducer.per.executor does not exist
22/11/05 13:46:14 WARN HiveConf: HiveConf of name hive.arrow.root.allocator.limit does not exist
22/11/05 13:46:14 WARN HiveConf: HiveConf of name hive.vectorized.use.checked.expressions does not exist
22/11/05 13:46:14 WARN HiveConf: HiveConf of name hive.tez.dynamic.senjoin.reduction.for.mapjoin does not exist
22/11/05 13:46:14 WARN HiveConf: HiveConf of name hive.vectorized.complex.types.enabled does not exist
22/11/05 13:46:14 WARN HiveConf: HiveConf of name hive.server2.um.worker.threads does not exist
22/11/05 13:46:14 WARN HiveConf: HiveConf of name hive.repl.partitions.dump.parallelism does not exist
22/11/05 13:46:14 WARN HiveConf: HiveConf of name hive.metastore.url.selection does not exist
22/11/05 13:46:14 WARN HiveConf: HiveConf of name hive.strict.checks.no.partition.filter does not exist
```

```
Activities Terminal Nov 9 21:55
hadoop@srn-pc: ~/sparkpro

hadoop@srn-pc:~$ echo "export SPARK_HOME=/opt/spark" >> ~/.profile
hadoop@srn-pc:~$ echo "export PATH=$PATH:$SPARK_HOME/bin:$SPARK_HOME/sbin" >> ~/.profile
hadoop@srn-pc:~$ echo "export PYSPARK_PYTHON=/usr/bin/python3" >> ~/.profile
hadoop@srn-pc:~$ nano .profile
hadoop@srn-pc:~$ source ~/.profile
hadoop@srn-pc:~$ start-master.sh
starting org.apache.spark.deploy.master.Master, logging to /opt/spark/logs/spark-hadoop.org.apache.spark.deploy.master.Master-1-srn-pc.out
hadoop@srn-pc:~$ start-worker.sh spark://master:port
starting org.apache.spark.deploy.worker.Worker, logging to /opt/spark/logs/spark-hadoop.org.apache.spark.deploy.worker.Worker-1-srn-pc.out
hadoop@srn-pc:~$ start-worker.sh spark://srn-pc:7077
starting org.apache.spark.deploy.worker.Worker, logging to /opt/spark/logs/spark-hadoop.org.apache.spark.deploy.worker.Worker-1-srn-pc.out
hadoop@srn-pc:~$ spark-shell
Setting default log level to "WARN".
To adjust logging level use sc.setLogLevel(newLevel). For SparkR, use setLogLevel(newLevel).
22/11/05 13:45:55 WARN NativeCodeLoader: Unable to load native-hadoop library for your platform... using builtin-java classes where applicable
Spark context Web UI available at http://srn-pc:4040
Spark context available as 'sc' (master = local[*]), app id = local-1667636156755).
Spark session available as 'spark'.
Welcome to


 version 3.3.1

Using Scala version 2.12.15 (OpenJDK 64-Bit Server VM, Java 1.8.0_342)
Type in expressions to have them evaluated.
Type :help for more information.

scala> :q
hadoop@srn-pc:~$ pyspark
python 3.8.10 (default, Jun 22 2022, 20:18:18)
[GCC 9.4.0] on linux
Type "help", "copyright", "credits" or "license()" for more information.
To adjust logging level use sc.setLogLevel(newLevel). For SparkR, use setLogLevel(newLevel).
22/11/05 13:46:14 WARN HiveConf: Hiveconf of name hive.metastore.wm.default.pool.size does not exist
22/11/05 13:46:14 WARN HiveConf: Hiveconf of name hive.llap.task.scheduler.preempt.independent does not exist
22/11/05 13:46:14 WARN HiveConf: Hiveconf of name hive.llap.output.format.arrow does not exist
22/11/05 13:46:14 WARN HiveConf: Hiveconf of name hive.tez.llap.min.reducer.per.executor does not exist
22/11/05 13:46:14 WARN HiveConf: Hiveconf of name hive.arrow.root.allocator.limit does not exist
22/11/05 13:46:14 WARN HiveConf: Hiveconf of name hive.vectorized.use.checked.expressions does not exist
22/11/05 13:46:14 WARN HiveConf: Hiveconf of name hive.tez.dynamic.sen.join.reduction.for.mapjoin does not exist
22/11/05 13:46:14 WARN HiveConf: Hiveconf of name hive.vectorized.complex.types.enabled does not exist
22/11/05 13:46:14 WARN HiveConf: Hiveconf of name hive.server2.worker.threads does not exist
22/11/05 13:46:14 WARN HiveConf: Hiveconf of name hive.repl.partitions.dump.parallelism does not exist
22/11/05 13:46:14 WARN HiveConf: Hiveconf of name hive.metastore.url.selection does not exist
22/11/05 13:46:14 WARN HiveConf: Hiveconf of name hive.strict.checks.no.partition.filter does not exist
```

```
Activities Terminal Nov 9 21:55
hadoop@srn-pc: ~/sparkpro

Welcome to

 version 3.3.1

Using Python version 3.8.10 (default, Jun 22 2022 20:18:18)
Spark context Web UI available at http://srn-pc:4040
Spark context available as 'sc' (master = local[*]), app id = local-1667636175818).
Spark session available as 'spark'.
>>> quit()
hadoop@srn-pc:~$ mkdir sparkpro
hadoop@srn-pc:~$ cd sparkpro
hadoop@srn-pc:~/sparkpro$ jupyter notebook
[I 13:51:36.957 NotebookApp] Writing notebook server cookie secret to /home/hadoop/.local/share/jupyter/runtime/notebook_cookie_secret
[I 13:51:38.375 NotebookApp] Serving notebooks from local directory: /home/hadoop/sparkpro
[I 13:51:38.375 NotebookApp] 0.0.0.0:8888 -> http://localhost:8888/?token=87888ca2470f6ce3d66d4fc2406f327413bc3e567d63a298
[I 13:51:38.375 NotebookApp] or http://127.0.0.1:8888/?token=87888ca2470f6ce3d66d4fc2406f327413bc3e567d63a298
[I 13:51:38.375 NotebookApp] Use Control-C to stop this server and shut down all kernels (twice to skip confirmation).

[C 13:51:38.442 NotebookApp]

To access the notebook, open this file in a browser:
file:///home/hadoop/.local/share/jupyter/runtime/nbserver-61786-open.html
Or copy and paste one of these URLs:
http://localhost:8888/?token=87888ca2470f6ce3d66d4fc2406f327413bc3e567d63a298
or http://127.0.0.1:8888/?token=87888ca2470f6ce3d66d4fc2406f327413bc3e567d63a298
[I 13:51:39.012 NotebookApp] Creating new notebook in
[I 13:51:39.088 NotebookApp] Writing notebook-signing key to /home/hadoop/.local/share/jupyter/notebook_secret
[I 13:51:41.141 NotebookApp] Kernel started: e2f14c06-4e08-4d0b-bc08-746918bc3453, name: python3
[IPKernelApp] ERROR | No such comm target registered: jupyter.widget.control
[IPKernelApp] WARNING | No such comm: 5fb1bd3a-12bf-4a05-b18e-44e02ec073ea
[I 13:51:41.128 NotebookApp] Saving file at /Untitled.ipynb
/home/hadoop/.local/lib/python3.8/site-packages/nbformat/_util.py:120: MissingIDFieldWarning: Code cell is missing an id field, this will become a hard error in future nbformat versions. You may want to use 'normalize()' on your notebooks before validations (available since nbformat 5.1.4). Previous versions of nbformat are fixing this issue transparently, and will stop doing so in the future.
/home/hadoop/.local/lib/python3.8/site-packages/nbformat/_util.py:133: MissingIDFieldWarning: Code cell is missing an id field, this will become a hard error in future nbformat versions. You may want to use 'normalize()' on your notebooks before validations (available since nbformat 5.1.4). Previous versions of nbformat are fixing this issue transparently, and will stop doing so in the future.
validate(nb)
validate(nb.model['content'])
[I 13:57:41.093 NotebookApp] Saving file at /Untitled.ipynb
[I 14:19:41.119 NotebookApp] Saving file at /Untitled.ipynb
[I 14:23:41.115 NotebookApp] Saving file at /Untitled.ipynb
[I 14:23:41.131 NotebookApp] Saving file at /Untitled.ipynb
[I 14:47:41.886 NotebookApp] Saving file at /Untitled.ipynb
Setting default log level to "WARN".
To adjust logging level use sc.setLogLevel(newLevel). For SparkR, use setLogLevel(newLevel).
```

CODE:

## WATER QUALITY PREDICTION

## Importing libraries

```
In [2]: 1 import os
2 import pandas as pd
3 import numpy as np
4 %matplotlib inline
5 import matplotlib.pyplot as plt
6 import seaborn as sns
7 import geopandas as gpd
8 import warnings
9 warnings.filterwarnings("ignore")
10
11 from pylab import *
12 from pyspark.sql.functions import udf, concat, col, lit
13 from pyspark import SparkConf, SparkContext
14 from pyspark.sql import SparkSession, SQLContext
15
16 from pyspark.sql.types import *
17 import pyspark.sql.functions as F
18 sc = SparkContext.getOrCreate(SparkConf().setMaster("local[*]"))
19 from pyspark.sql import SparkSession
20 spark = SparkSession \
21     .builder \
22     .getOrCreate()
23 sqlContext = SQLContext(sc)
```

## Reading the dataset and describing

```
In [4]: 1 df = spark.read.format("csv").option("header", "true").load('/home/hadoop/sparkpro/waterquality.csv')
```

```
In [5]: 1 df.show(5)
```

Water Quality Data - Godavari River											
STATION CODE	LOCATIONS	STATE	TEMP	DO	pH	CONDUCTIVITY	BOD	NITRATE	N_NITRITE	N_FECAL_COLIFORM	TOTAL
73	1312	GODAVARI AT JAYAK...	MAHARASHTRA	29.2	6.4	8.1	735	3.4	2	3	
182	2177	GODAVARI RIVER NE...	MAHARASHTRA	24.5	6	8	270	3.1	2	72	
133	2182	GODAVARI RIVER AT...	MAHARASHTRA	25.8	5.5	7.8	355	4.2	9	59	
283	2179	GODAVARI RIVER AT...	MAHARASHTRA	24.8	5.5	7.8	371	5.6	3.55	90	
132	2183	GODAVARI RIVER AT...	MAHARASHTRA	25.7	5.7	7.9	294	3.2	2.69	45	
only showing top 5 rows											

```

In [6]:      dl.dtypes

Out[6]: [('STATION CODE', 'string'),
          ('LOCATIONS', 'string'),
          ('STATE', 'string'),
          ('TEMP', 'string'),
          ('DO', 'string'),
          ('pH', 'string'),
          ('CONDUCTIVITY', 'string'),
          ('BOD', 'string'),
          ('NITRATE_N_NITRITE_N', 'string'),
          ('FECAL COLIFORM', 'string'),
          ('TOTAL COLIFORM', 'string')]

In [7]:      dt.head()

      ? + | " ROW( STATION CODE= '1312', LDCA+IDNS= 'GODAVARI AT 3AYAKNADI DAM, AURNAGABAD MAHARASHTRA', STATE= 'MAHARASHTRA', TEMP= '29.2', DO= '6.4', pH= '8.1', CONDUCTIVITY= '735', BOD= '3.4', NITRATE_N_NITRITE_N= '2', FECAL COLIFORM= '3', TOTAL COLIFORM= '731'

In [10]:      d.describe()

Elit "ld" Data range [ s unrrary: string, STATION CODE: string, LOCATIONS: string, STATE: string, TEMP: string, DD: string, pH: string, CONDUCTIVITY: string, BOD: string, NITRATE_N_NITRITE_N: string, FECAL COLIFORM: string, +OTAL_CDLIFORM: string]

```

## Data cleaning

```

In [11]:      from pyspark.sql.types import FloatType

```

As we can see above that all the columns have string data types, but for the calculation of water quality index we need to convert them in float data type. So we will convert the required columns to float data type

```

In [12]:      d = df.withColumn('TEMP', df['TEMP'].cast(FloatType()))
      d = df.withColumn('pH', df['pH'].cast(FloatType(111)))
      d = df.withColumn('DO', df['DO'].cast(FloatType(111)))
      d = df.withColumn('CONDUCTIVITY', df['CONDUCTIVITY'].cast(FloatType(1)))
      d = df.withColumn('BOD', df['BOD'].cast(FloatType(1)))
      d = df.withColumn('NITRATE_N_NITRITE_N', df['NITRATE_N_NITRITE_N'].cast(FloatType(1)))
      d = df.withColumn('FECAL_COLIFORM', df['FECAL_COLIFORM'].cast(FloatType(1)))

In [13]:      d.dtypes

```

```

Out[13]: [('STATION CODE', 'string'),
          ('LOCATIONS', 'string'),
          ('STATE', 'string'),
          ('TEMP', 'float'),
          ('DO', 'float'),
          ('CONDUCTIVITY', 'float'),
          ('BOD', 'float'),
          ('NITRATE_N_NITRITE_N', 'float'),
          ('FECAL_COLIFORM', 'float'),
          ('TOTAL_COLIFORM', 'float')]

```

Now the column TOTAL COLIFORM is not required so we will drop this column.



```
[ 14 ]: df=dl.dropna(subset=['TOTAL COLIFORM'])
```

Now we have to check all the rDws for null values and remove all the rDws which contain any null value in it.

So for applying a SOL query we first have to register it as a virtual temporary table and then we will issue SOL query.

We are doing this because it is important to perform data cleaning as it will make our model work with better prediction and accuracy.

```
In [ 15 ]: df.createOrReplaceTempView('d')
```

```
In [ 16 ]: df_clean = spark.sql("""SELECT * FROM d WHERE DO IS NOT NULL AND pH IS NOT NULL AND BOD IS NOT NULL AND CONDUCTIVITY IS NOT NULL AND NITRATE_N IS NOT NULL AND NITRITE_N IS NOT NULL AND FECAL_COLIFORM IS NOT NULL""")
```

```
[ 17 ]: df_clean.show(5)
```

STATION CODE	LOCATIONS	STATE	TEMP	DO	pH	CONDUCTIVITY	BOD	NITRATE_N	NITRITE_N	FECAL_COLIFORM
1312	GODAVARI AT AYAK...	NAHARASHTRA	29.2	5.4	8.1	735.0	3.4	2.0	3.01	
2177	GODAVARI RIVER NE...	NAHARASHTRA	24.3	6.0	8.0	270.0	3.1	2.0	72.01	
2182	GODAVARI RIVER AT...	NAHARASHTRA	23.8	3.5	7.8	355.0	4.2	9.0	69.01	
2179	GODAVARI RIVER AT...	NAHARASHTRA	24.8	3.5	7.8	371.0	5.6	3.55	90.01	
2183	GODAVARI RIVER AT...	NAHARASHTRA	25.7	5.7	7.9	294.0	3.2	2.69	45.01	

only showing top 5 rows

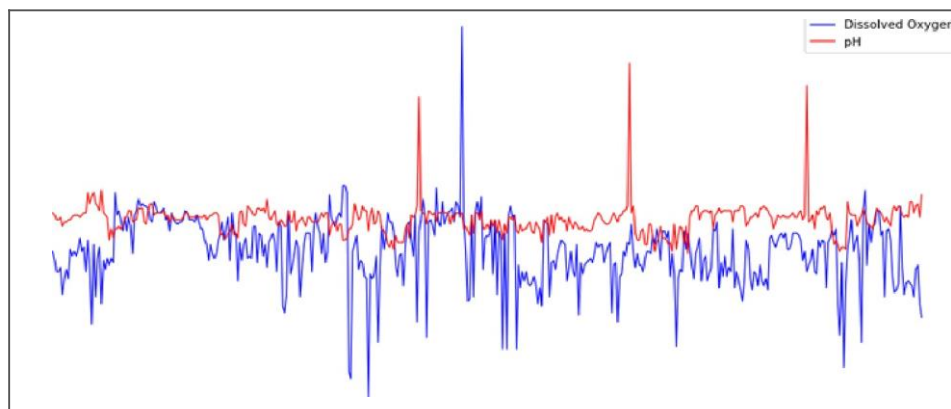
## EDA - Exploratory data analysis

```
In [ 18 ]: df_clean.createOrReplaceTempView('af')
```

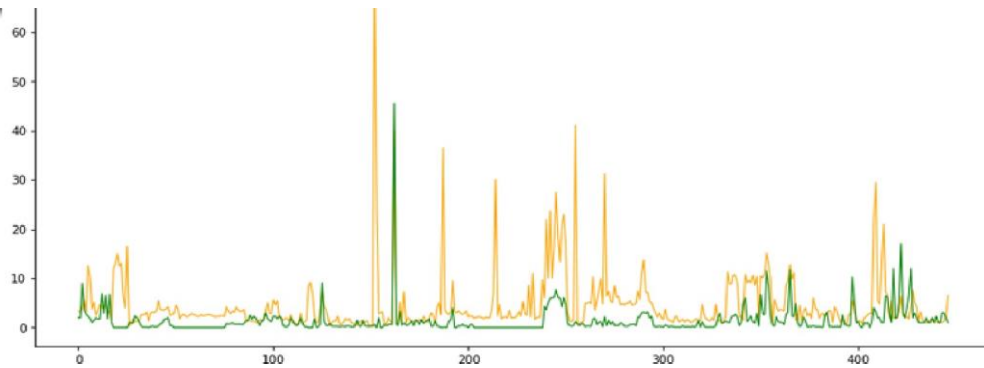
```
In [ 20 ]: do = spark.sql("""SELECT * FROM af WHERE DO < 5""")
do = do.rdd.map(lambda row: row.DO).collect()
ph = spark.sql("""SELECT * FROM af WHERE pH < 7""")
ph = ph.rdd.map(lambda row: row.pH).collect()
bod = spark.sql("""SELECT * FROM af WHERE BOD > 5""")
bod = bod.rdd.map(lambda row: row.BOD).collect()
nn = spark.sql("""SELECT * FROM af WHERE NITRATE_N > 5 OR NITRITE_N > 5""")
nn = nn.rdd.map(lambda row: row.NITRATE_N + row.NITRITE_N).collect()
```

Data visualization

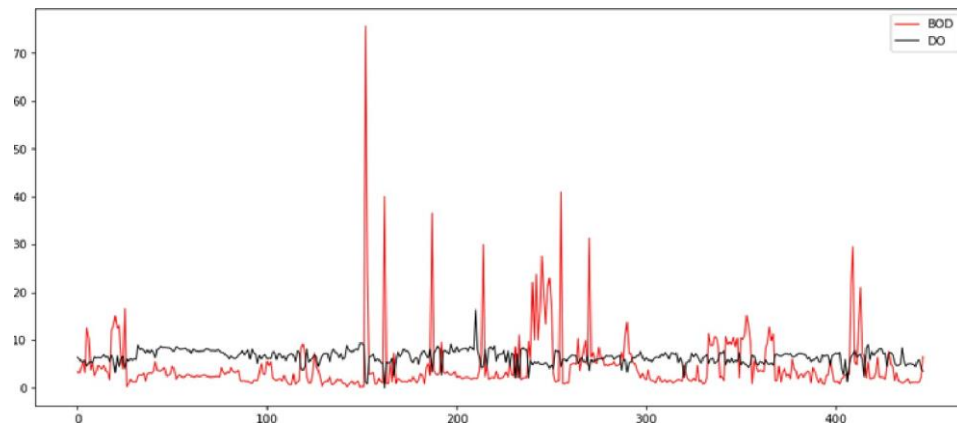
```
In [ 21 ]: fig, ax = plt.subplots(figsize=(14, 10), dpi=80, facecolor='w', edgecolor='k')
size = len(df)
ax.plot(range(0, size), do, color='b', label='Dissolved Oxygen')
ax.plot(range(0, size), ph, color='r', label='pH')
legend = ax.legend()
```



```
In [22]: fig, ax2 = plt.subplots(num=None, figsize=(L4, G), dpi=80, facecolor='w', edgecolor='k')
ax2.plot(iangle\8, size, bad, color='aiznge', animated=True, linewidth=l, label='BOX')
ax2.plot(iangle\8, size, nn, color='gre:n' animated=True, linewidth=l, label='Jf.')
legend=ax2.legend()
```



```
In [23]: fig, ax3 = plt.subplots(num=None, figsize=(L4, G), dpi=80, facecolor='w', edgecolor='k')
ax3.plot(iangle, size).bad, color='reo', animated=True, linewidth=l, label='BOD')
ax3.plot(iangle, size, do, color='lack'a animated=True, linewidth=l, label='DC')
legend=ax3.legend()
```



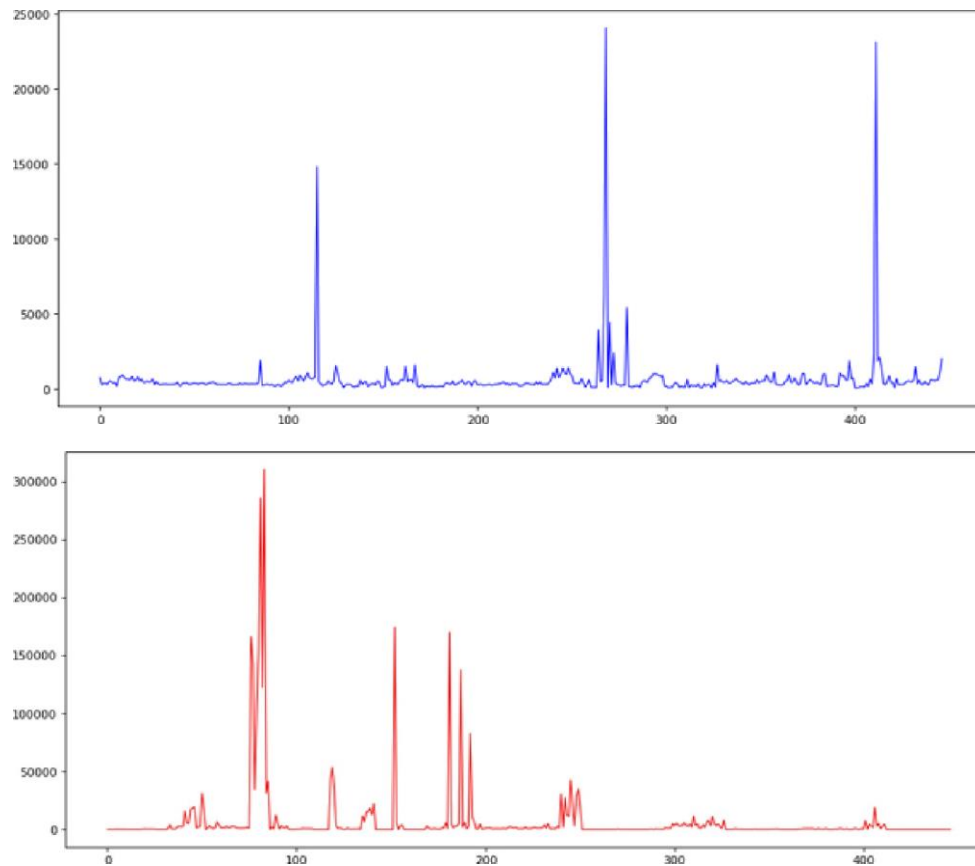
```

In [2d1]: 1 com = spar1c.sql-I 5e\ec t CO/@UCTIVITY f rd d f sql ")
          2 com = con.rdd.eap (£ aaAda rou: re.CB@U£ TIYITY ).collect ()
          fee = spar-k.sql-f Setec I FECAL COLIF0£¥I Iron d f sql ")
          4 fee = fee.rdd.eap (£ aaAda rou: re. FECAL_C6LI F0£¥¥ ).collect ()

In [231]: 1 fig,ax = pit.subptots(nu==Me.i igsizet=14, 61,g i we. lacecofor='u' . ed9ecofor='k' l
          2 an.plot( r6ozj (B,size), con, color='blue', animated=7rue, lineuidh=i)
          2 fig,arZ = plt.subplots(num=6bne,figsize=(Lt,6), Opi=B0, facecoIor='w', edgecolor='k')
          J axe.plot(range18.size). 7ec. color='r*d', aninated=7me. Ii idtE=1)

```

Out [25 j]: [ zsat plott ib , lines. Li ne2D at 8x7fed75ae6ac• I



## Feature engineering

```
In [26]: df=df_clean.toPandas()
df.dtypes
```

```
Out[26]. STATION CODE      object
LOCATIONS                 object
STATE                     float32
TEMP                      float32
DO                        float32
pH                        float32
CONDUCTIVITY             float32
BOD                       float32
NITRATE N NITRITE N     float32
FETAL COLIFORM           float32
dtype: object
```

```
In [27]: sta rt=0
end=94B
station=df . i\oc [ s ta rt : end , 0]
\ocation=df . ioc [st a rt: end , I]
stat e=dT . il a c [ sta rt : end , 2]
do= df . iloc [start:end , 4].astype(np.float64)
value=0
ph -d I . hoc [ start : end , 3]
co -d I . i\oc [ start : end , 6] . as type ( np . I \o at 64
boo =dT . il a c [ sta rt : end , 7 ] . a stype ( np . I oat64)
na= df . itoc [start:end , 8] . astype ( np . I oat64)
I c=d f . i\oc [ 2 : end , 9 ] . a stype ( np . fi Dat64)

In [28]: df=pd.concat([station,location,state,do,ph,co,bod,na,fc],axis=1)
df. columns = ['station','location','state','do','ph','co','bod','na','fc']
```

the Water Quality Index is calculated by aggregating the quality with the weight linearly:

$$WQI = \sum (q_n \times W_n)$$

where  $q_n$  = Quality rating for the  $n$ th Water quality       $W_n$  = unit weight for the  $n$ th oarameters

```
In [29]: df[!EH] = d1.ph.apply({\anibda x: { 103 if (8.5>x>=7)
                                             el se (TO if (6.3>x>=8.5) a r (3.3—x>=6.8)
                                             else (63 if (B B>x>=8.6) a r (b.8>x>=6.T)
                                             et se (3 1+ (3>x>=0.8) a r (6.7>x>=6.3)
                                             e \ se C)))))
```

```
In [30]: df[inc a] = df.I.do . apply (\ lambda x: (f'00 If (x>=6)
                                         el se (a0 if (s>=x>=5. 1)
                                         else (6'3 if (5>=x>=4. 1)
                                         et se (=3 1+ (4>=x>=fi)
                                         else 0))))))
```

[illegible][illegible]

```
In [wee] = dI.co.apply(\ lambda x: (100 if *5>=x>=0)
                        else (0 if (15†>=x>= )
                                else (60 if (2<8>:x>=156)
                                        else( i* ‡t€>=x>=225)
                                              else 011))1
```

[illegible]

Now we apply the formula of wqi by first multiplying all the quality rating with its weight and then summed all the values

```
In [36]: df['n.th'] = df.npH * 0.165
df['v.da'] = df.ndo * 0.281
df['u.bd.a'] = df.nbdo * 0.234
df['sec'] = df.nec * 0.009

df['\-.c.D'] = df.nco * 0.281
df['\-.c.t'] = df.wph + df.wdo + df.wbdo + df.weo + df.wna + df.wco
```

station	location	stie	doph	oo	bod	ua	k	npH	.nboo	nec	uua	wph	wdo	wbdo	weo	wna	woo		
0	1342	GODAVARI AT JAYAKWADI AURMAGABAD MAHARASHTRA	MAHARASHTRA	6.8	1735.0	3.4	2.00	NaN	DO	80	0	10D	16.5	28.10	18.72	D.DO	2.8	0.00	
1		GODAVARI RIVER NEAR SOMESHWAR TEMPLE.	MAHARASHTRA	60	80	27D.0	31	200	HaN	80	40	10D	16.5	28.10	18.72	D.36	2.8	0.00	
2	B2	GODAVARI RIVER AT SAIKHEDA.	MAHARASHTRA	5.5	7.8	355.0	4.2	900	59.0	DO	80	0	10D	16.5	22.4B	18.72	DDO	2.8	16.86
3		GODAVARI RIVER AT HANUMAN GHAT, NASHIK CITY.	MAHARASHTRA	5.5	7.8	371.0	5.6	3.55	90.0	DO	80	0	10D	16.5	22.4B	18.72	D.DO	2.8	16.86
4	B3	GODAVARI RIVER AT NANDUR- MADAMESHWAR DAM.	MAHARASHTRA	5.7	7.4	294.0	3.2	2.69	45.0	DO	80	40	10D	16.5	22.4B	18.72	D.36	2.8	22.4B
442	2950	GAPE- SAGAR LAKE, DUNGARPUR, RAJASTHAN	RAJASTHAN	4.8	8.1	53B.0	1.2	1.00	5.0	DO	100	0	10D	16.5	16.86	23.40	D.DO	2.8	28.10
443	2941	LAKE JAISAMAND, SALUMBER, UDAIPUR, POINT	RAJASTHAN	5.6	8.4	591.0	1.1	3.00	4.0	DO	100	0	10D	16.5	22.4B	23.40	D.DO	2.8	28.10
444	2942	LAKE JAISAMAND, SALUMBER, UDAIPUR, POINT	RAJASTHAN	5.8	8.5	58B.0	1.2	3.00	4.0	DO	100	0	10D	16.5	22.4B	23.40	D.DO	2.8	28.10
445	2953	LODHA TALAB, BANSWARA- DUNGARPUR ROAD, BANSWARA...	RAJASTHAN	4.1	7.9	1133.0	2.3	2.00	4.0	DO	100	0	10D	16.5	0.00	23.40	D.DO	2.8	22.4B

Now we classify the water on the basis of their water quality index

```
In [37]: df['at Qty'] = df.wqi.apply(lambda x: 'Excellent' if 125 >= x >= 0)
else ('Good' if 150 >= x >= 26)
else ('Poor' if 75 >= x >= 31)
else ('Very Poor' if 10 >= x >= 76)
else ('Unsuitable'))))
```

```
In [38]: spark_df = sqldf.createDataFrame(df)
```

```
In [39]: spark_df.show()
```

station	location	state	do	ph	o	ood	web
1312	GODAVARI AT AYAK...	MAHARASH RA	6.400090093367432	8.1000000381	69727	7J3.01s.4c0000	s5JO74316
2.0	NaN LOO 100 0 B01 0 LOO		16.3		28.1	18.72000000000002	0.0 2.800
00000000000003			66.12	Poor			
2177	GODAVARI RIDER NE...	MAHARASH+RA		6.0	8.0	1270.0	13.09999990 6323684
2.0	NaN UO JON 0 B01 40 U0		16.5		28.1	18.7200000000000002	0.36 2.800
00000000000003			66.4B	Poor			
2182	GODAVARI RIDER AT...	HAHARASH+RA		5.5 7.8000	0190734863	3^..^14	99999809265137
0.0	9.1 0 LOO B01 60 B01 0 LOO		T6.5 22.	480000000000004	18.72000000000002		0.0 2.800
00000000000003116.8600000000000003			77.36 Very	Poor			
2179	GODAVARI RIDER AT...	MAHARA9H*RA		5.5 7.800000190734863	371.1*	9999990463236B	3.24999999
32316284	90.01100	80	601	80 01T00		16.5122.48000000000000p	18.7200000000000002
2.8000000000000003			16.8600000000000003	77.36 Very	Poor		0.0
2183	GODAVARI RIDER AT...	MAHARASH*RA	5.699999809265137	7.900000095367432	294 "1	*000000^76B3716	2.69000000
37220459	45.01100	80	801	80 401T00		1c.5122.48000000000000p	18.7200000000000002
2.8000000000000003			22.480000000000004	1c.5122.48000000000000p	18.7200000000000002		0.36
0.0	2181 GODAVARI RIVER AT...	MAHARASHTRA		4.5	7.5 513.0	12.600000381469727	2.29999999
52316284	131.0 100	60	60 60 0 100		16.5 16.86000000000003	14.04900000000001	0.0
2.8000000000000003			16.86000000000003	67.96	Poor		

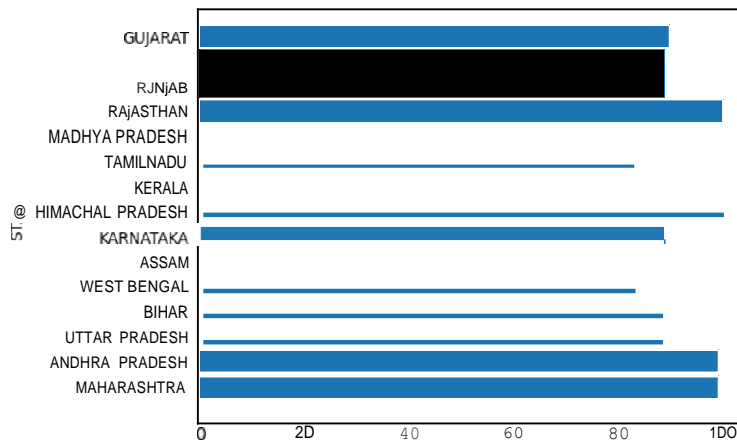
```
In [40]: spark_df.createOrReplaceTempView("dsql")

In [41]: StatB = spark.sql("Select state from dsql")
State = State.rdd.map(lambda row: row.state).collect()

In [42]: wqi = spark.sql("Select wqi from dsql")
Wqi = Wqi.rdd.map(lambda row: row.wqi).collect()

In [43]: pit.barh(State, wqi)
pit.xticks(["WOI"])
pit.yticks(["STATES"])
```

```
7 plt.show()
```



## Modeling

Now we apply machine learning algorithms to predict the data.

### Linear Regression Model

First we convey the required data to predict WQI into vector form by using VectorAssembler.

Then we normalize our data by using Normalizer.

```
In [44]: from pyspark.ml.feature import VectorAssembler
from pyspark.ml.feature import Normalizer

vectorAssembler = VectorAssembler(inputCols=["pH", "ndo", "bdo", "nec", "nna", "nco"], outputCol="features")
normalizer = Normalizer(inputCol="features", outputCol="features_norm")
```

Applying Linear Regression model

```
In [45]: from pyspark.ml.regression import LinearRegression
lr = LinearRegression(featuresCol="features_norm", labelCol="wqi", maxIter=10, regParam=0.3, elasticNetParam=0.2)
```

Now we streamline processes into

```
In [46]: from pyspark.ml import Pipeline

In [47]: pipeline = Pipeline(stages=[vectorAssembler, normalizer, lr])
```

Now we split the data into train and test data

```
In [48]: train_data, test_data = spark_df.randomSplit([0.8, 0.2])
```

Making predictions using the model

```
In [51]: predictions = model.transform(test_data)
```

```
In [52]: predictions.select('wqi', 'prediction').show(1)
```

wqi	prediction
82.03999999999999	B2.1238308668°088
82.4	B1.9274096394131B
66.12	67.34565604943359
66.12	67.3456b604943359
82.4	B1.92740963941318
77.72	77.78171638043395
77.36000000000001	77.89313076930208
66.12	67.3456b604943359
82.03999999999999	B2.1238308668p08 B
66.12	67.34565604943359
66.12	67.3456b604943359
82.03999999999999	B2.1238308668°088
82.03999999999999	B2.12383086684088
82.03999999999999	B2.1238308668q088
66.12	67.3456b604943359
77.36	78.5041637745495^
82.98	B3.11194223770477
82.4	B1.92740963941318
77.9	77.9108613083°056
93.28	91.66703037750938

Only showing top 20 rows

Now we check the performance of our model

```
In [53]: a = model.stages[2].summary.r2
```

```
Out[53]: 0.9753564789532392
```

## Logistic Regression Model

Here we are creating a logistic regression model because we don't have to predict a continuous value

```
In [54]: from pyspark.ml.feature import StringIndexer
```

As our quality column contains values in string format so first we indexed them using StringIndexer

Then data is converted which are required to predict water quality into vector form by using VectorAssembler

Then we normalize our data by using Normalizer.

```
In [55]: stringIndexer = StringIndexer(inputCol="quality", outputCol="label")
vectorAssembler = VectorAssembler(inputCols=["pH", "nd", "ndd", "red", "luna", "nd", "lci"], outputCol="features")
normalizer = Normalizer(inputCols=["features"], outputCol="features_norm")
```

Applying Logistic Regression Model

```
In [56]: from pyspark.ml.classification import LogisticRegression
logisticRegression = LogisticRegression(featuresCol="features_norm", labelCol="label", maxIter=10)
```

```
In [57]: pipeline = Pipeline(stages=[stringIndexer, vectorAssembler, normalizer, logisticRegression])
```

Now we split the data into train and test data

```
In [58]: (train_data, test_data) = spark_df.randomSplit([0.8, 0.2])
```

Fitting the data into the model

```
In [60]: model = pipeline.fit(train_data)
```

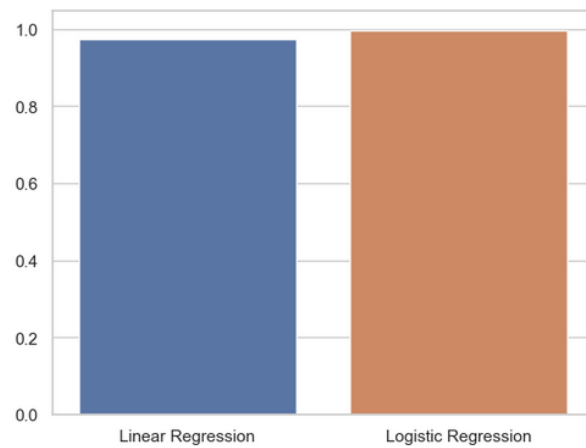
```
In [611: Dredictions2 = model2. I ransfo rm (t rain data )

In [621: Dredictions2. select ( 'label ", " p redict ion "). s how (I
```

PredIc ted : Ve ry Poo r Act uat : Very Poo r  
 PredIc ted : Ve ry Poo r Act uat : Very Pa a r  
 Pred be ted : PDo r Actual : PDo r  
 PredIc ted : Poo r Actual : Poo r  
 Pred be ted : Very Poo r Act uat : very Poo r  
 PredIc ted : Ve ry Poo r Act uat : Very Pa a r  
 Pred be ted : ve ry Poo r Act uat : very Poo r  
 PredIc ted : Poo r Actual : Poo r  
 Pred be ted : ve ry Poo r Act uat : very Poo r  
 PredIc ted : Ve ry Poo r Act uat : Very Pa a r  
 Pred 7c ted : Poo r Actual : Poo r  
 PredIc ted : Ve ry Poo r Act uat : Very Pa a r  
 Pred 7c ted : Ve ry Poo r Act uat : Very Pa a r  
 PredIc ted : Ve ry Poo r Act uat : Very Pa a r  
 PredIc ted : Ve ry Poo r Act uat : Very Pa a r  
 PredIc ted : Ve ry Poo r Act uat : Very Poo r  
 PredIc ted : Ve ry Poo r Act uat : Very Poo r  
 PredIc ted : Ve ry Poo r Act uat : Very Pa a r  
 Pred be ted : ve ry Poo r Act uat : very Poo r  
 PredIc ted : Ve ry Poo r Act uat : Very Pa a r  
 Pred be ted : ve ry Poo r Act uat : very Poo r  
 PredIc ted : Ve ry Poo r Act uat : Very Pa a r



```
In [70]: 1 sns.set_theme(style='whitegrid')
2 scores = pd.DataFrame(data=[[a,b]],columns=['Linear Regression', 'Logistic Regression'])
3 sns.barplot(data=scores);
```



**RESULT:**

The program is executed successfully.