# Statistical ML for big data analysis with Spark

*Group Members:* Karthik Garimella, Sanjeev Raichur, Sandeep Alfred

**Introduction:**

Implementation of machine learning algorithms to predict the average temperature for each state in the United States over time. Considering this is a regression problem, four powerful regression methods were implemented namely, Random Forest Regression, Gradient Boosted Regression, Decision Tree Regression, and Elastic Net Linear Regression. The preprocessed data of 10,000 records was put through a pipeline of ML algorithms. Using metrics like Rsquared and RMSE, an appropriate model was determined. The entire task was performed in parallel using PySpark and different combinations of virtual machines. The computation speeds for each combination of VMs were recorded.

**Link to the code:** https://github.com/karthikgarimella37/SAT5165_PySpark_Project
**Data Downloaded from:** https://asmith.ucdavis.edu/

**ML Methods implemented by each group member:**
1. Karthik Garimella:  Gradient Boosted Regression and Decision Tree Regression.
2. Sanjeev Raichur: Elastic Net Linear Regression.
3. Sandeep Alfred: Random Forest Regression.

**Code Explanation and Results:**
- **Data Preprocessing:**
1. Dropped highly sparse columns like **stability,** separated the day, month and year columns from date.
2. Created dummy variables for all the categorical variables.
3. Centered and scaled the data using `StandardScaler()` from `pyspark.ml.feature.` The result is shown in figure 1.



Figure 1. Centered and Scaled data

- **Implementation of ML algorithms:**

1. Implemented **Pipeline** from **pyspark.ml** to run several ML models.
2. Used **VectorAssembler** from **pyspark.ml.feature** to merge all the columns of the data into a vector column.
3. Used **randomSplit** method with seed 123 to split the data into training and testing sets. 80% of the data was used to train the model and tested on the remaining 20% of the data.
4. From **pyspark.ml.regression** used **GBTRegressor** for Gradient Boosted Tree Regression with **maxIter** of 10.
5. Implemented Decision Tree Regression using **DecisionTreeRegressor** function.
6. Elastic Net Linear Regression was implemented with the **LinearRegression** function, using parameters of **maxIter** of 10.
7. **RandomForestRegressor** was used to implement Random Forest Regression. All the functions for the models were implemented from **pyspark.ml.regression**.
8. In the pipeline, initialized each model with **t_avg** as the response variable, which is the average temperature for each state.
9. The models were evaluated on the test set. Used **RegressionEvaluator** from **pyspark.ml.evaluation** to calculate RMSE and Rsquared for each model. The evaluation results are shown in figure 2.

```
RMSE for Gradient Boosted Tree Regression: 4.368317036438087
R2 for Gradient Boosted Tree Regression: 0.9641361790204939

RMSE for Decision Tree Regression: 3.2463124624985618
R2 for Decision Tree Regression: 0.99271165503378

RMSE for Elastic Net Linear Regression: 1.1654150696058916
R2 for Elastic Net Linear Regression: 0.967518117142999

RMSE for Random Forest Regression: 3.8560456735155126
R2 for Random Forest Regression: 0.8254559576937875
```

Figure 2. Model Evaluation Results on the Test Set.

From figure 2 we can see that Elastic Net Linear Regression has the lowest RMSE value of 1.16541, and a very good RMSE of 0.9675.

**Performance of PySpark:**

The performance of Spark to implement all the mentioned ML algorithms for each combination of virtual machines is mentioned below.

- The discussed code with preprocessing and ML algorithms using 10,000 rows of the data when run on 1 machine took 17 mins to run, as shown in figure 3.

**Spark** 3.5.0    **Spark Master at spark://hadoop1:7077**

**URL:** spark://hadoop1:7077
**Alive Workers:** 1
**Cores in use:** 8 Total, 0 Used
**Memory in use:** 1024.0 MiB Total, 0.0 B Used
**Resources in use:**
**Applications:** 0 Running, 9 Completed
**Drivers:** 0 Running, 0 Completed
**Status:** ALIVE

**▾ Workers (1)**

| Worker Id | Address | State | Cores | Memory | Resources |
|---|---|---|---|---|---|
| worker-20241105212827-192.168.13.153-35845 | 192.168.13.153:35845 | ALIVE | 8 (0 Used) | 1024.0 MiB (0.0 B Used) | |

**▾ Running Applications (0)**

| Application ID | Name | Cores | Memory per Executor | Resources Per Executor | Submitted Time | User | State | Duration |
|---|---|---|---|---|---|---|---|---|

**▾ Completed Applications (9)**

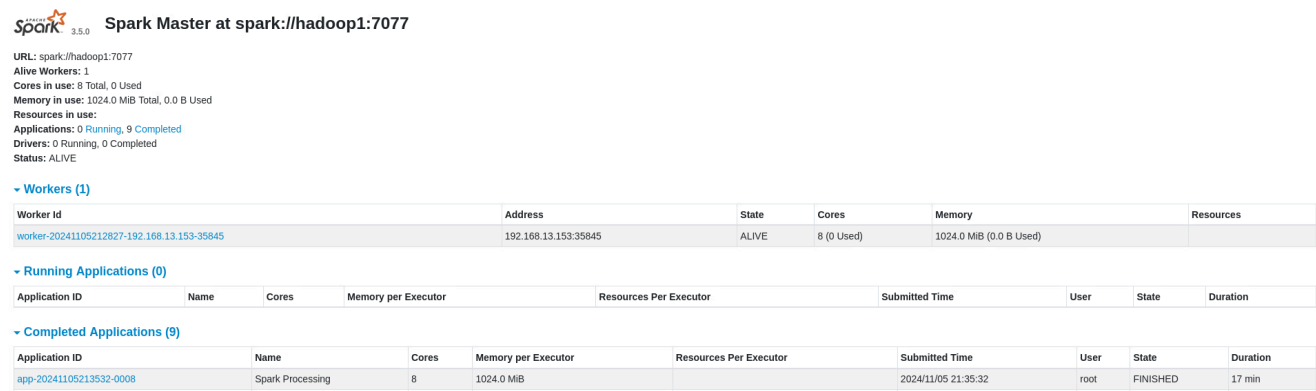| Application ID | Name | Cores | Memory per Executor | Resources Per Executor | Submitted Time | User | State | Duration |
|---|---|---|---|---|---|---|---|---|
| app-20241105213532-0008 | Spark Processing | 8 | 1024.0 MiB | | 2024/11/05 21:35:32 | root | FINISHED | 17 min |

Figure 3. Execution time of 1 VM, 8 cores.

- The discussed code with preprocessing and ML algorithms using 10,000 rows of the data when run on 2 machines took 27 mins to run, as shown in figure 4. This was not expected with two VMs.
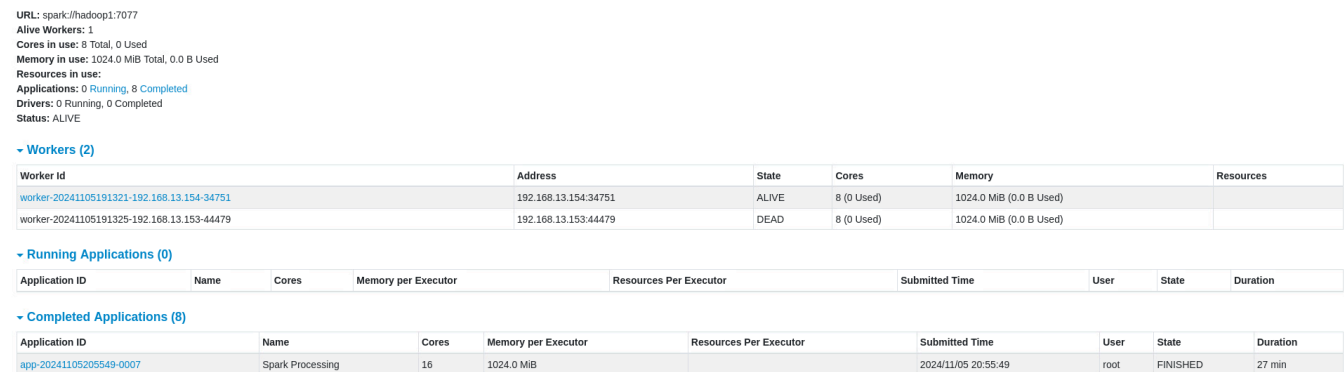
**URL:** spark://hadoop1:7077
**Alive Workers:** 1
**Cores in use:** 8 Total, 0 Used
**Memory in use:** 1024.0 MiB Total, 0.0 B Used
**Resources in use:**
**Applications:** 0 Running, 8 Completed
**Drivers:** 0 Running, 0 Completed
**Status:** ALIVE

**▾ Workers (2)**

| Worker Id | Address | State | Cores | Memory | Resources |
|---|---|---|---|---|---|
| worker-20241105191321-192.168.13.154-34751 | 192.168.13.154:34751 | ALIVE | 8 (0 Used) | 1024.0 MiB (0.0 B Used) | |
| worker-20241105191325-192.168.13.153-44479 | 192.168.13.153:44479 | DEAD | 8 (0 Used) | 1024.0 MiB (0.0 B Used) | |

**▾ Running Applications (0)**

| Application ID | Name | Cores | Memory per Executor | Resources Per Executor | Submitted Time | User | State | Duration |
|---|---|---|---|---|---|---|---|---|

**▾ Completed Applications (8)**

| Application ID | Name | Cores | Memory per Executor | Resources Per Executor | Submitted Time | User | State | Duration |
|---|---|---|---|---|---|---|---|---|
| app-20241105205549-0007 | Spark Processing | 16 | 1024.0 MiB | | 2024/11/05 20:55:49 | root | FINISHED | 27 min |

Figure 4. Execution time of 2 VMs, 16 cores.

- The discussed code with preprocessing and ML algorithms using 10,000 rows of the data when run on 4 machines took 15 mins to run, as shown in figure 5.



**Spark Master at spark://hadoop1:7077**

**URL:** spark://hadoop1:7077
**Alive Workers:** 4
**Cores in use:** 32 Total, 0 Used
**Memory in use:** 4.0 GiB Total, 0.0 B Used
**Resources in use:**
**Applications:** 0 Running, 7 Completed
**Drivers:** 0 Running, 0 Completed
**Status:** ALIVE

**Workers (4)**

| Worker Id | Address | State | Cores | Memory | Resources |
|---|---|---|---|---|---|
| worker-20241105191321-192.168.13.154-34751 | 192.168.13.154:34751 | ALIVE | 8 (0 Used) | 1024.0 MiB (0.0 B Used) | |
| worker-20241105191325-192.168.13.153-44479 | 192.168.13.153:44479 | ALIVE | 8 (0 Used) | 1024.0 MiB (0.0 B Used) | |
| worker-20241105191326-192.168.13.121-37139 | 192.168.13.121:37139 | ALIVE | 8 (0 Used) | 1024.0 MiB (0.0 B Used) | |
| worker-20241105203343-192.168.13.119-35355 | 192.168.13.119:35355 | ALIVE | 8 (0 Used) | 1024.0 MiB (0.0 B Used) | |

**Running Applications (0)**

| Application ID | Name | Cores | Memory per Executor | Resources Per Executor | Submitted Time | User | State | Duration |
|---|---|---|---|---|---|---|---|---|

**Completed Applications (7)**

| Application ID | Name | Cores | Memory per Executor | Resources Per Executor | Submitted Time | User | State | Duration |
|---|---|---|---|---|---|---|---|---|
| app-20241105203620-0006 | Spark Processing | 32 | 1024.0 MiB | | 2024/11/05 20:36:20 | root | FINISHED | 15 min |

Figure 5. Execution time of 4 VMs, 32 cores

- The discussed code with preprocessing and ML algorithms using 10,000 rows of the data when run on all 6 machines took 12 mins to run, as shown in figure 6.



**URL:** spark://hadoop1:7077
**Alive Workers:** 6
**Cores in use:** 48 Total, 0 Used
**Memory in use:** 6.0 GiB Total, 0.0 B Used
**Resources in use:**
**Applications:** 0 Running, 3 Completed
**Drivers:** 0 Running, 0 Completed
**Status:** ALIVE

**Workers (6)**

| Worker Id | Address | State | Cores | Memory | Resources |
|---|---|---|---|---|---|
| worker-20241105191321-192.168.13.154-34751 | 192.168.13.154:34751 | ALIVE | 8 (0 Used) | 1024.0 MiB (0.0 B Used) | |
| worker-20241105191321-192.168.13.174-33187 | 192.168.13.174:33187 | ALIVE | 8 (0 Used) | 1024.0 MiB (0.0 B Used) | |
| worker-20241105191322-192.168.13.173-46149 | 192.168.13.173:46149 | ALIVE | 8 (0 Used) | 1024.0 MiB (0.0 B Used) | |
| worker-20241105191323-192.168.13.119-38971 | 192.168.13.119:38971 | ALIVE | 8 (0 Used) | 1024.0 MiB (0.0 B Used) | |
| worker-20241105191325-192.168.13.153-44479 | 192.168.13.153:44479 | ALIVE | 8 (0 Used) | 1024.0 MiB (0.0 B Used) | |
| worker-20241105191326-192.168.13.121-37139 | 192.168.13.121:37139 | ALIVE | 8 (0 Used) | 1024.0 MiB (0.0 B Used) | |

**Running Applications (0)**

**Completed Applications (3)**

| Application ID | Name | Cores | Memory per Executor | Resources Per Executor | Submitted Time | User | State | Duration |
|---|---|---|---|---|---|---|---|---|
| app-20241105191404-0002 | Spark Processing | 48 | 1024.0 MiB | | 2024/11/05 19:14:04 | root | FINISHED | 12 min |

Figure 6. Execution time of 6 VMs, 48 cores.

Performance summary to run the ML methods on 10,000 records of data:

| Number of VMs Used | Time Taken | Cores Used |
|---|---|---|
| One | 17 minutes | 8 |
| Two | 27 minutes | 16 |
| Four | 15 minutes | 32 |
| Six | 12 minutes | 48 |