



```

73 // Let's fix this by converting the timestamp column to the format spark knows
74 val month_map = Map("Jan" -> 1, "Feb" -> 2, "Mar" -> 3, "Apr" -> 4, "May" -> 5, "Jun" -> 6, "Jul" -> 7, "Aug" -> 8
75 | , "Sep" -> 9, "Oct" -> 10, "Nov" -> 11, "Dec" -> 12)
76 def parse_clf_time(s: String) = {
77   "%3$s-%2$s-%1$s %4$s:%5$s:%6$s".format(s.substring(0,2),month_map(s.substring(3,6)),s.substring(7,11)
78   | ,s.substring(12,14),s.substring(15,17),s.substring(18))
79 }
80 val totimestamp = udf[String, String](parse_clf_time(_))
81 val logs_df = cleaned_df.select($"*", to_timestamp(to_timestamp($"timestamp")).alias("time")).drop("timestamp")
82 logs_df.printSchema()
83 logs_df.show(2)
84 // We cache the dataset so the next action would be faster
85 logs_df.cache()
86
87 // =====< Analysis walk-through >=====
88
89 /*
90 | | status column statistics
91 */
92 logs_df.describe("status").show()
93
94 /*
95 | | HTTP status analysis
96 */
97 logs_df.groupBy("status").count().sort("status").show()
98
99 /*
100 | | Frequent Hosts
101 */
102 logs_df.groupBy("host").count().filter($"count" > 10).show()
103
104 /*
105 | | Visualizing Paths
106 */
107 logs_df.groupBy("path").count().sort(desc("count")).show()
108
109 /*
110 | | Top Paths
111 */

```

```

106 /*
107 | | Top Paths
108 */
109 logs_df.groupBy("path").count().sort(desc("count")).show(10)
110
111 // =====< Analyzing Web Server Log File >=====
112
113 /*
114 | | Top Ten Error Paths
115 */
116 logs_df.filter($"status" != 200).groupBy("path").count().sort(desc("count")).show(10)
117
118 /*
119 | | Number of unique Hosts
120 */
121 val unique_host_count = logs_df.select("host").distinct().count()
122 println("Unique hosts : %d".format(unique_host_count))
123
124 /*
125 | | Number of Unique Daily Hosts :
126 */
127 val daily_hosts_df = logs_df.withColumn("day", dayofyear($"time")).withColumn("year", year($"time")).select("host", "day", "year").distinct().groupBy("day", "year").count().show(5)
128
129 /*
130 | | Average Number of Daily Requests per Host
131 */
132 val total_req_per_day_df = logs_df.withColumn("day", dayofyear($"time")).withColumn("year", year($"time")).groupBy("day", "year").count()
133 val avg_daily_request_per_host_df = total_req_per_day_df.join(daily_hosts_df, total_req_per_day_df("day") === daily_hosts_df("day") && total_req_per_day_df("year") === daily_hosts_df("year")).show(5)
134
135 // =====< Exploring 404 status codes >=====
136
137 /*
138 | | Let's drill down and explore the error 404 status records, We've all seen those "404 Not Found" web pages.
139 | | 404 errors are returned when the server cannot find the resource (page or object) the browser client requested.
140 */
141
142 // Counting 404 Response Codes
143 val not_found_df = logs_df.where($"status" === 404).cache()

```

```

140 // Counting 404 Response Codes
141 val not_found_df = logs_df.where($"status" === 404).cache()
142 println("found %d 404 Urls".format(not_found_df.count()))
143
144 // Listing 404 Status Code Records
145 not_found_df.select("path").distinct().show(40,false)
146
147 // Listing The Top Twenty 404 Response Code Paths :
148 not_found_df.groupBy("path").count().sort("count").show(20,false)
149 not_found_df.groupBy("path").agg("host" -> "collect_list", "status" -> "count").sort("count(status)").show(20)
150 not_found_df.groupBy("path").agg("host" -> "collect_set", "status" -> "count").sort("count(status)").show(20)
151
152 // Listing the Top Twenty-five 404 Response Code Hosts
153 not_found_df.groupBy("host").count().sort(desc("count")).show(truncate = false)
154
155 // Listing 404 Errors per Day
156 val errors_by_date_pair = not_found_df.withColumn("day", dayofyear($"time")).withColumn("year", year($"time")).groupBy("day", "year").count()
157 not_found_df.withColumn("day", dayofyear($"time")).withColumn("year", year($"time")).groupBy("day", "year").count().sort($"year", $"day").show(10)
158
159
160
161 /* To run the program
162 scala> :load WebLog_Processing.scala
163 */
164
165

```