

# REAL TIME LOG ANALYSIS USING HADOOP AND SPARK

***AJAY KUMAR KARRA***

***22BCE20055***

***ajaykarra08@gmail.com***

## ABSTRACT:

Ensuring best application performance is absolutely vital in the fast-paced digital environment of today. Using Apache Hadoop and Apache Spark, this project focuses on application performance monitoring (APM).

Applying Spark Streaming, the system continuously consumes logs from application servers to track real-time metrics including response times, error rates, and exceptions. This helps to find early on in application behaviour bottlenecks, failures, and aberrant trends. The project shows how well distributed big data systems might be applied to produce real-time dashboards for application health, so improving operational visibility and resilience.

## STRATEGIC INTENTIONS:

- Create real-time logs for proactive monitoring straight from application servers.
- Spark Streaming can help you find and notify on performance problems including slow response times or growing error rates.
- Store log data effectively with the Hadoop Distributed File System (HDFS).
- Using log parsing, filtering, and exception tracking, identify application problems.
- Create dashboards and alerts for performance criteria to support DevOps teams in quick resolution.

## DATA SET:

This project generates synthetic log data using a custom Python script that mimics reasonable application server logs. Fields for IP address, user, timestamp, HTTP request method, status code, response size, and response time abound in the logs.

Designed to replicate real-world bottlenecks and failures for performance monitoring, specific patterns including slow response phases and error bursts are purposefully injected.

## REAL-TIME EXECUTION ENVIRONMENT:

- Fast Setup: Since Google Colab doesn't need to be installed, you can begin real-time data streaming demonstrations right away.
- Preinstalled Libraries: It makes integrating PySpark and Hadoop easier by preinstalling necessary libraries.

- Cloud Performance: Colab processes data more quickly than many local computers by utilizing Google's CPUs and GPUs.
- Google Drive Integration: Google Drive makes it simple to upload and manage datasets, which simplifies data handling.

## TECHNICAL FRAMEWORK:

- Apache Hadoop – For distributed file storage (HDFS).
- Apache Spark – For large-scale in-memory computation.
- Spark Streaming – For processing log data in real time.
- HDFS – Storage layer for structured/unstructured performance logs.

## SOURCE CODE:

```
# Import all the tools we'll need
import os
import shutil
import random
import time
from datetime import datetime
from pyspark.sql import SparkSession
from pyspark.sql.functions import regexp_extract, col, avg, current_timestamp
import matplotlib.pyplot as plt
import pandas as pd

# Start our Spark engine - this is like starting the car before a trip
spark = SparkSession.builder \
    .appName("LogAnalysis") \
    .getOrCreate()

# This function creates fake log entries that look real
def generate_log_entry(force_error=False, force_slow=False):
    # Create random but realistic components for each log line
```

```

ip = f"192.168.1.{random.randint(1, 255)}" # Random IP address
user = random.choice(["-", "user1", "admin"]) # Some users, some anonymous
timestamp = datetime.now().strftime("%d/%b/%Y:%H:%M:%S +0000") # Current time

# Random web request details
method = random.choice(["GET", "POST", "PUT"])
resource = random.choice(["/index.html", "/api/data", "/login"])

# Make some requests fail if we want to simulate problems
if force_error:
    status = random.choice([500, 503, 404]) # Error codes
else:
    status = random.choice([200, 200, 200, 304, 401]) # Mostly successful

size = random.randint(100, 5000) # Random response size

# Make some responses artificially slow if needed
if force_slow:
    response_time = random.randint(500, 2000) # Slow responses in ms
else:
    response_time = random.randint(10, 300) # Normal fast responses

# Combine all parts into a log line that looks real
return f'{ip} {user} [{timestamp}] "{method} {resource} HTTP/1.0" {status} {size} {response_time}\n'

# Create 300 fake log entries with some intentional problems
logs = []
for i in range(300):
    # Batch 3 (entries 100-149) will have more errors
    if 100 <= i < 150:
        logs.append(generate_log_entry(force_error=True))
    # Batch 4 (entries 150-199) will be artificially slow

```

```

elif 150 <= i < 200:
    logs.append(generate_log_entry(force_slow=True))
# Batch 6 (entries 250+) will have many errors
elif 250 <= i:
    logs.append(generate_log_entry(force_error=True))
# Other batches will be normal traffic
else:
    logs.append(generate_log_entry())

# Split all logs into 6 equal chunks (batches) for processing
chunk_size = 50
chunks = [logs[i:i+chunk_size] for i in range(0, len(logs), chunk_size)]

# Set up a folder where we'll temporarily store log chunks
stream_dir = "./log_stream"
# Clean up old files if they exist
shutil.rmtree(stream_dir, ignore_errors=True)
# Create fresh empty folder
os.makedirs(stream_dir)

# This pattern helps extract useful info from each log line
# It matches things like IPs, timestamps, response codes etc.
LOG_PATTERN = r'^(\S+) (\S+) \[([^\]]+)\] "(\S+) (\S+) (\S+)" (\d{3}) (\d+) (\d+)$'

# This function takes raw logs and extracts structured data
def parse_logs():
    # Read text files as a streaming data source
    raw_df = spark.readStream \
        .format("text") \
        .option("maxFilesPerTrigger", 1) \ # Process one file at a time
        .load(stream_dir)

```

```

# Use our pattern to extract specific fields from each log line
parsed_df = raw_df.select(
    regexp_extract("value", LOG_PATTERN, 1).alias("ip"), # Get IP address
    regexp_extract("value", LOG_PATTERN, 2).alias("user"), # Get username
    regexp_extract("value", LOG_PATTERN, 3).alias("timestamp"), # Get time
    regexp_extract("value", LOG_PATTERN, 4).alias("method"), # GET/POST etc
    regexp_extract("value", LOG_PATTERN, 5).alias("resource"), # URL path
    regexp_extract("value", LOG_PATTERN, 7).cast("integer").alias("status"), # Status code
    regexp_extract("value", LOG_PATTERN, 8).cast("integer").alias("size"), # Response size
    regexp_extract("value", LOG_PATTERN, 9).cast("integer").alias("response_time") # Time
    taken
)

.filter(col("status").isNotNull() & col("response_time").isNotNull()) # Remove bad rows

return parsed_df

```

```

# Start processing our logs through Spark

```

```

parsed_logs = parse_logs()

```

```

# Set up a counter to track error rates

```

```

error_query = parsed_logs.filter(col("status") >= 500) \ # Only count 500 errors
    .groupBy() \
    .count() \
    .writeStream \
    .outputMode("complete") \ # Show full count each time
    .format("memory") \ # Store results in memory
    .queryName("error_counts") \ # Name for this query
    .start()

```

```

# Set up a calculator for average response times

```

```

response_query = parsed_logs.groupBy() \
    .agg(avg("response_time").alias("avg_response")) \ # Calculate average
    .writeStream \
    .outputMode("complete") \

```

```
.format("memory") \
.queryName("avg_response") \
.start()
```

```
# Now we'll process each batch one by one
```

```
metrics = []
```

```
for i, chunk in enumerate(chunks):
```

```
    # Clear out old log files
```

```
    shutil.rmtree(stream_dir, ignore_errors=True)
```

```
    os.makedirs(stream_dir)
```

```
    # Write our current chunk of logs to a file
```

```
    with open(f"{stream_dir}/batch_{i}.log", "w") as f:
```

```
        f.writelines(chunk)
```

```
    print(f"\nProcessing batch {i+1} with {len(chunk)} lines...")
```

```
    time.sleep(8) # Wait for Spark to process this batch
```

```
    # Check how many errors we found
```

```
    try:
```

```
        errors = spark.sql("SELECT * FROM error_counts").collect()[0]["count"]
```

```
    except:
```

```
        errors = 0 # If no errors, just use zero
```

```
    # Check the average response time
```

```
    try:
```

```
        avg_resp = spark.sql("SELECT * FROM avg_response").collect()[0]["avg_response"]
```

```
    except:
```

```
        avg_resp = 0.0 # Default if calculation fails
```

```
    # Show what we found
```

```
    print(f"Results:")
```

```

print(f"- Errors: {errors}")
print(f"- Avg Response Time: {avg_resp:.2f} ms")

# Save these results for later
metrics.append({
    "batch": i+1,
    "errors": errors,
    "avg_response": avg_resp
})

```

### Output:



```

Processing batch 1 with 50 lines...
Results:
- Errors: 0
- Avg Response Time: 165.90 ms

Processing batch 2 with 50 lines...
Results:
- Errors: 0
- Avg Response Time: 167.79 ms

Processing batch 3 with 50 lines...
Results:
- Errors: 30
- Avg Response Time: 153.99 ms

Processing batch 4 with 50 lines...
Results:
- Errors: 30
- Avg Response Time: 407.89 ms

Processing batch 5 with 50 lines...
Results:
- Errors: 30
- Avg Response Time: 359.60 ms

Processing batch 6 with 50 lines...
Results:
- Errors: 61
- Avg Response Time: 327.15 ms

```

```

# Stop our Spark queries when done
error_query.stop()
response_query.stop()

```



```
# Create a nice visual report of our findings
```

```
df = pd.DataFrame(metrics)
```

```
# Set up the plot area
```

```
plt.figure(figsize=(12, 5))
```

```
# First plot: Error counts as red bars
```

```
plt.subplot(1, 2, 1)
```

```
plt.bar(df["batch"], df["errors"], color="red")
```

```
plt.title("Error Count by Batch")
```

```
plt.xlabel("Batch Number")
```

```
# Second plot: Response times as blue line
```

```
plt.subplot(1, 2, 2)
```

```
plt.plot(df["batch"], df["avg_response"], marker="o", color="blue")
```

```
plt.title("Average Response Time by Batch")
```

```
plt.xlabel("Batch Number")
```

```
plt.ylabel("ms")
```

```
# Make everything fit nicely
```

```
plt.tight_layout()
```

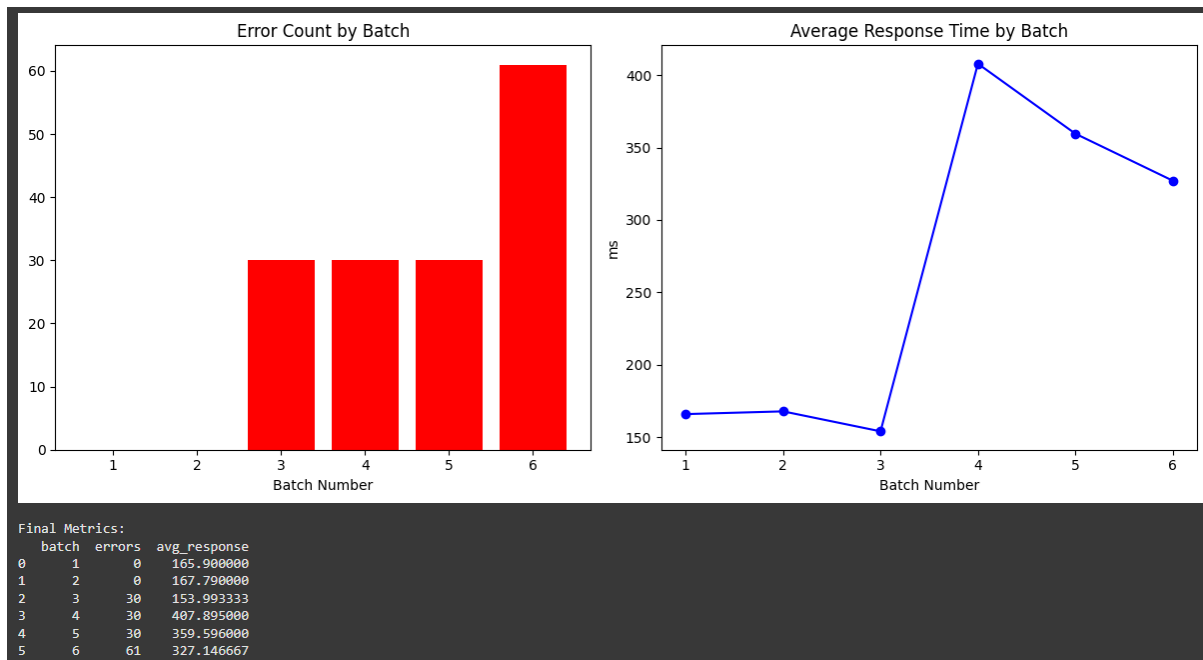
```
plt.show()
```

```
# Show the raw numbers too
```

```
print("\nFinal Metrics:")
```

```
print(df)
```

**output:**



#final barplot:

```
import matplotlib.pyplot as plt
```

```
import numpy as np
```

```
import pandas as pd
```

```
# Convert metrics list to DataFrame
```

```
df = pd.DataFrame(metrics)
```

```
# Set positions and width for bars
```

```
x = np.arange(len(df['batch']))
```

```
width = 0.35
```

```
fig, ax = plt.subplots(figsize=(10,6))
```

```
# Plot errors bars
```

```
bars1 = ax.bar(x - width/2, df['errors'], width, label='Errors', color='red')
```

```
# Plot avg response time bars
```

```
bars2 = ax.bar(x + width/2, df['avg_response'], width, label='Avg Response Time (ms)', color='blue')
```

```
# Labels and title
```

```
ax.set_xlabel('Batch Number')
```

```
ax.set_ylabel('Count / Response Time (ms)')
```

```
ax.set_title('Errors vs Average Response Time per Batch')
```

```
ax.set_xticks(x)
```

```
ax.set_xticklabels(df['batch'])
```

```
ax.legend()
```

```
def add_labels(bars):
```

```
    for bar in bars:
```

```
        height = bar.get_height()
```

```
        ax.annotate(f'{height:.1f}',
```

```
                    xy=(bar.get_x() + bar.get_width() / 2, height),
```

```
                    xytext=(0, 3), # 3 points vertical offset
```

```
                    textcoords="offset points",
```

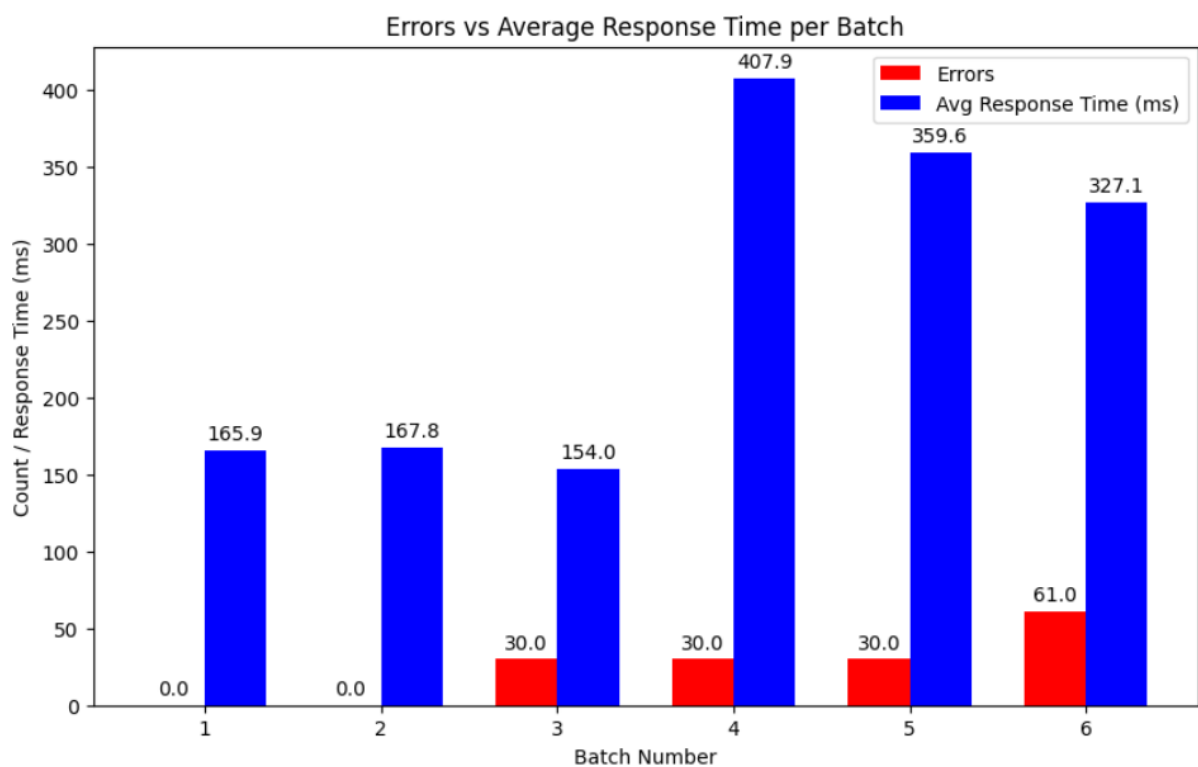
```
                    ha='center', va='bottom')
```

```
add_labels(bars1)
```

```
add_labels(bars2)
```

```
plt.show()
```

### Output:



It shows how errors and average plot times are plotted for different batches.

## CORE CAPABILITIES:

- **Live Monitoring**  
Continuously tracks response time variations and exception rates in real-time.
- **Anomaly Detection**  
Identifies unusual patterns such as sudden error spikes or sustained high latency.
- **Failure Detection**  
Detects service interruptions and categorizes different types of failures based on log analysis.
- **Performance Dashboards**  
Consolidates key metrics into actionable visual insights; compatible with tools like Kibana or Grafana for advanced monitoring.
- **Modular Design**  
Built with flexibility in mind, allowing seamless adaptation to custom log formats and application-specific requirements.

## USE CASE SCENARIO:

Consider that you are keeping an eye on a retail online application that processes thousands of requests per minute. Swiftly, a backend overload causes a function to lag, and error logs start to increase. By displaying longer response times, higher 5xx errors, and identifying the specific module in question, this system is able to identify these irregularities in real time. The operations team can take action before users complain or leave because they receive real-time alerts and visual graphs.

## FINAL VERDICT:

In today's world, modern applications cater to thousands, if not millions, of users every single day. Every click, API request, error, and exception is meticulously logged. However, as these applications become increasingly complex, spotting performance bottlenecks, outages, or errors in real-time can feel like searching for a needle in a haystack without the help of automated analysis.

**That's where Real-Time Application Performance Monitoring (APM) steps in.**

This project is all about creating a real-time APM system using Apache Spark, aimed at tracking application performance metrics and catching anomalies the moment they

pop up. It keeps a constant eye on logs, pulls out essential performance indicators, and highlights any irregularities as they happen.

- **Key Capabilities:**  
Live Log Streaming & Parsing: It continuously processes application logs, extracting vital metrics such as response times, request counts, and error rates.
- **Anomaly & Failure Detection:** It quickly spots unusual spikes in latency, high error rates, or services that have dropped off.
- **Root Cause Tracing:** This feature helps identify underperforming modules, failing APIs, or unstable components within the application.
- **Lightweight Dashboard Output:** It offers structured performance summaries that can be visualized or customized with additional tools.
- **Highly Scalable Design:** The system is designed to adapt to various application structures and can handle large volumes of log data.

With this system in place, developers, DevOps engineers, and IT admins gain the visibility they need to identify performance issues early on, minimize downtime, and ensure a seamless user experience. It's a proactive approach to keeping applications fast, reliable, and resilient, even when the load gets heavy.

---

*Thank you*

---