

Open in app ↗

Medium

Search

Write

3



✦ Get unlimited access to the best of Medium for less than \$1/week. [Become a member](#)



How to Future-Proof Your Spark Jobs from Production Failures



Anurag Sharma · 4 min read · Jan 26, 2026



In this blog, we share how our Data Engineering team shifted from **reactive debugging** to a **predictive, future-proof approach** by embedding **Debug Mode** and **Analysis Mode** directly into our PySpark jobs. These modes continuously capture execution, resource, and data-growth signals from production workloads to help us **explain failures, predict risks, and prevent incidents before they happen.**

This blog walks through a **step-by-step implementation**, real-world failure stories, and practical PySpark code patterns that any team can adopt to reduce operational overhead and protect sprint velocity.

One or more Spark jobs would:

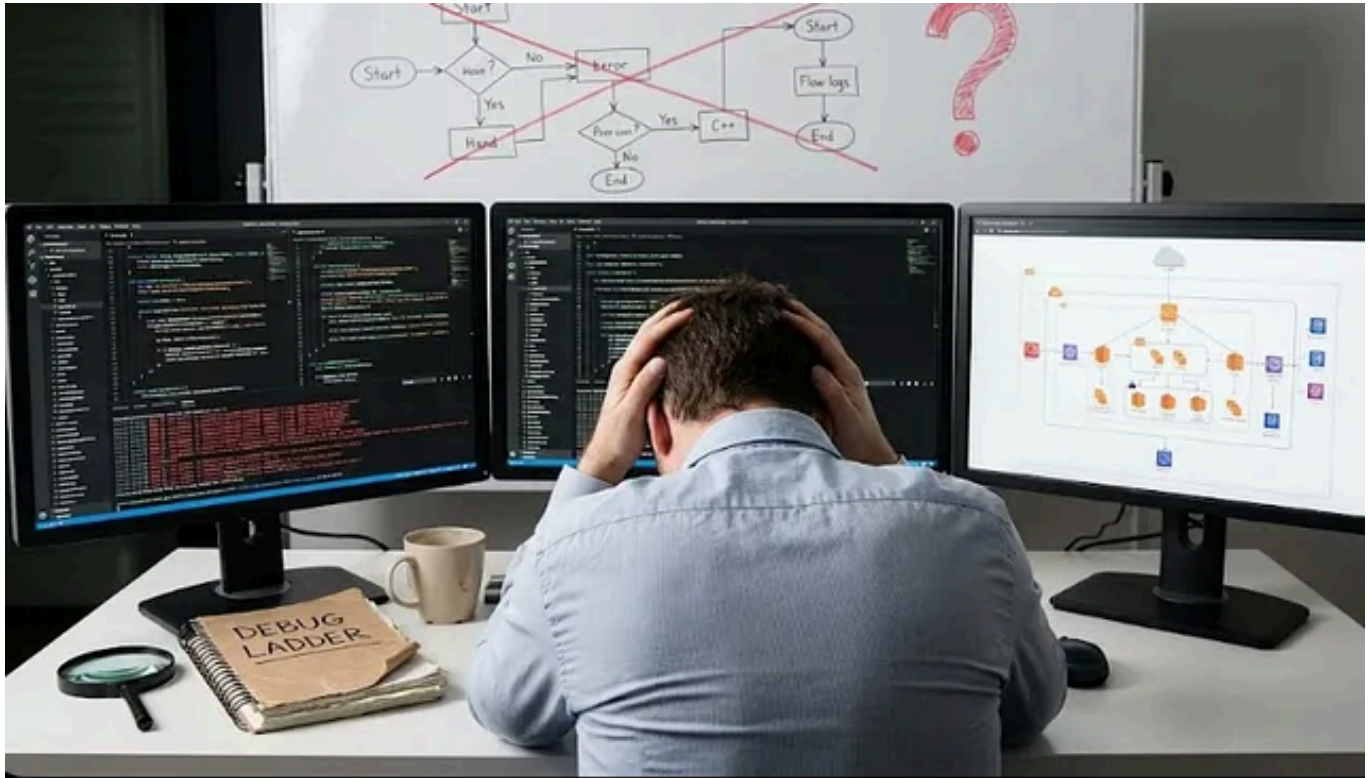
- Fail intermittently
- Fail once a week
- Fail randomly at peak load
- **Succeed immediately when re-run**

These failures were rarely code bugs. They were caused by **data growth, skew, memory pressure, or cluster saturation.**

Yet each incident pulled engineers away from roadmap work into:

- Log analysis
- Spark UI screenshots
- Emergency memory tuning
- Last-minute cluster resizing

Over time, **support and maintenance started consuming sprint capacity**, which is exactly what a healthy data platform should avoid.



Why Traditional Debugging Doesn't Work

Problem 1: You Can't Reproduce Prod Failures in UAT

- Smaller clusters
- Lower data volume
- Different data distribution

A job that fails in production due to memory pressure will often run flawlessly in UAT.

Problem 2: Production Data Is Always Different

- Sudden spikes in data
- Uneven partition growth
- Schema drift
- Unexpected nulls or outliers

By the time the job fails, you are already reacting.



Our Shift in Thinking: Predict Failures Before They Happen

Instead of asking:

“Why did this Spark job fail?”

We started asking:

“What signals tell us this Spark job will fail in the future?”

That led us to embed **two execution modes** into every Spark job:

1. **Debug Mode** — explains *why* a job fails
2. **Analysis Mode** — predicts *when* a job will fail

Both are:

- Config-driven
- Optional
- Production-safe
- Zero overhead when disabled

Step-by-Step Implementation Guide (PySpark)

Step 1: Capture Query Execution Plans

```
if debug_enabled:
    print("==== LOGICAL PLAN ====")
    df.explain(mode="simple")

    print("==== PHYSICAL PLAN ====")
    df.explain(mode="extended")
```

Catches early:

- Accidental Cartesian joins

- Missing filters
- Unnecessary shuffles

Step 2: Detect Partition Skew

```
if debug_enabled:
    partition_sizes = (
        df.rdd
        .mapPartitions(lambda it: [sum(1 for _ in it)])
        .collect()

    print(f"Partition sizes: {partition_sizes}")
    print(f"Max: {max(partition_sizes)}")
    print(f"Min: {min(partition_sizes)}")
```

One partition 10x larger than others → future executor OOM.



Step 3: Capture Task and Executor Failures


```
from pyspark import SparkContext
from pyspark.sql import SparkSession
from pyspark.sql.utils import StreamingQueryException

if debug_enabled:
    sc = spark.sparkContext
```

```
def log_executor_info():  
    for executor_id, info in sc._jsc.sc().getExecutorMemoryStatus().items():  
        print(f"Executor: {executor_id}, Memory Status: {info}")  
  
log_executor_info()
```

This gives **live executor insights** during execution.

Step 4: Analyze Cluster Configuration



```
if analysis_enabled:  
    print("Executor Memory:", spark.conf.get("spark.executor.memory"))  
    print("Executor Cores:", spark.conf.get("spark.executor.cores"))  
    print("Shuffle Partitions:",  
          spark.conf.get("spark.sql.shuffle.partitions"))
```

Is the job still compatible with this cluster size?

Step 5: Memory Spill Analysis

```
if analysis_enabled:  
    status_store = spark.sparkContext.statusStore()  
    stages = status_store.stageList(None)  
  
    for stage in stages:  
        print(  
            f"Stage {stage.stageId()} | "  
            f"Memory Spill: {stage.memoryBytesSpilled()} | "  
            f"Disk Spill: {stage.diskBytesSpilled()}"  
        )
```

Increasing spill over weeks = silent failure loading

Step 6: Predict Data Growth: This was the biggest game-changer.

```
if analysis_enabled:
    row_count = input_df.count()
    avg_row_size_bytes = 500 # estimated
    current_size_gb = (row_count * avg_row_size_bytes) / (1024**3)

    previous_size_gb = load_metric("input_size_gb")

    if previous_size_gb:
        growth_pct = (
            (current_size_gb - previous_size_gb)
            / previous_size_gb
        ) * 100

        print(f>Data Growth: {growth_pct:.2f}%")

        if growth_pct > 20:
            print("High data growth detected. Failure risk ahead")

    save_metric("input_size_gb", current_size_gb)
```

This allowed us to **predict failures weeks in advance**

Step 7: Keep Normal Runs Clean

```
if not debug_enabled and not analysis_enabled:
    print("Running in normal execution mode")
```

No performance penalty & No extra logs.

Implementation Results

Our Story 1: The “Random” Weekly Failure

A Spark job failed every Sunday night.

Root cause (found via Analysis Mode):

- Weekly data load increased by ~25%
- One join key became highly skewed
- Executor memory unchanged for months

Fix:

- Salting + executor memory increase
- Failure eliminated permanently

Our Story 2: The Re-Run Mirage

A job failed at 2 AM but succeeded on a re-run.

Debug Mode showed:

- Heavy disk spill during shuffle
- The first run coincided with other peak jobs
- The second run had less cluster contention

Fix:

- Job rescheduled
- Shuffle partitions tuned

Our Story 3: The Silent Time Bomb

No failures. No alerts.

Analysis Mode detected:

- 18% week-over-week data growth
- Spill metrics are rising steadily

Action:

- Proactive cluster resize
- Zero production incidents

Final Takeaway

Spark jobs do not usually fail suddenly.

They age:

- Data grows
- Skew increases
- Memory assumptions break

By embedding **Debug Mode** and **Analysis Mode** directly into Spark jobs, we:

- Reduced firefighting
- Protected sprint velocity
- Shifted from reactive to predictive engineering

Production should never be your first alert system.

Spark

Data Engineering

Data Engineer

Optimization

Big Data



Written by Anurag Sharma

Edit profile

83 followers · 3 following

Data Engineering Specialist with 10+ exp. Passionate about optimizing pipelines, data lineage, and Spark performance and sharing insights to empower data pros!

No responses yet



Anurag Sharma him/he

What are your thoughts?