

Predictive Auto Scaling using AWS CloudWatch and Flask

Author: Anurag Thakur

Abstract

This project demonstrates a predictive auto-scaling system built using AWS EC2, CloudWatch, and Flask. The system monitors CPU utilization metrics, predicts future usage trends, and provides scaling recommendations based on machine learning models. The approach helps optimize resource allocation while maintaining performance and cost efficiency.

Introduction

Auto-scaling in cloud computing allows infrastructure to dynamically adjust capacity based on workload demands. Traditional scaling methods react to real-time metrics, while predictive scaling uses historical data to anticipate future needs. In this project, AWS CloudWatch metrics are analyzed using Python and Flask to build a lightweight predictive auto-scaler.

Objective

The objective of this project is to design and implement a predictive auto-scaling system that analyzes AWS EC2 instance performance metrics and provides scaling recommendations based on forecasted resource utilization.

Methodology

1. Launch an AWS EC2 instance with the required permissions.
2. Configure AWS CLI and CloudWatch access.
3. Collect real-time CPU utilization metrics using CloudWatch APIs.
4. Use Python libraries such as NumPy and Flask to process and serve data.
5. Develop a Flask application to display current and predicted metrics.
6. Implement a simple prediction algorithm to recommend scaling actions.

Results

Below are the main outputs and screenshots from the project:

1. AWS EC2 Instance Screenshot:

```
[ec2-user@ip-172-31-22-184 predictive-autoscale]$ ls
README.md app.py metrics.json metrics.jsonsdcasvasvsgvasv requirements.txt
[ec2-user@ip-172-31-22-184 predictive-autoscale]$
```

```
[ec2-user@ip-172-31-22-184 predictive-autoscale]$ sudo python3 app.py
* Serving Flask app 'app'
* Debug mode: off
WARNING: This is a development server. Do not use it in a production deployment.
Use a production WSGI server instead.
* Running on all addresses (0.0.0.0)
* Running on http://127.0.0.1:5000
* Running on http://172.31.22.184:5000
Press CTRL+C to quit
106.51.86.205 - - [04/Nov/2025 12:40:11] "GET /predict HTTP/1.1" 200 -
```

i-07fc3035214d2e886 (Anurag1serverr)

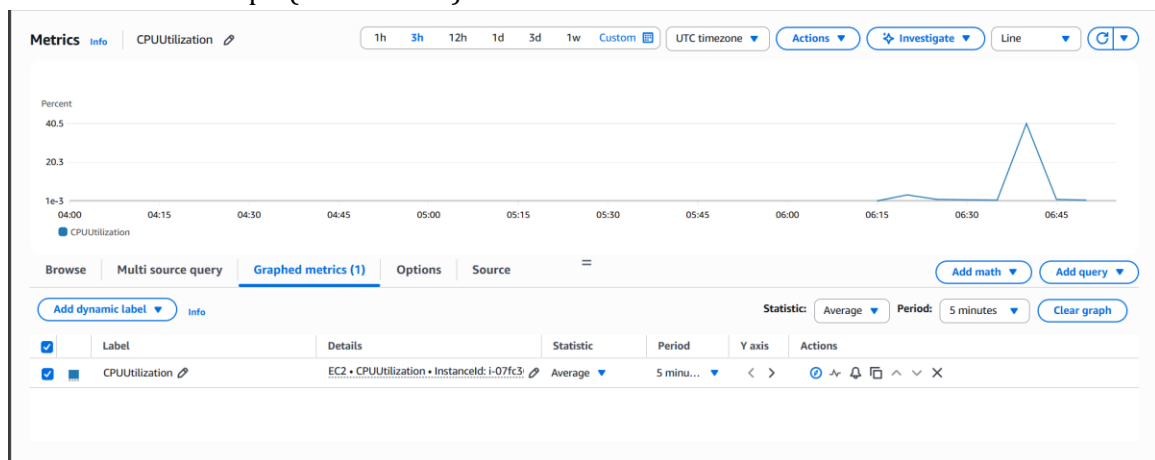
PublicIPs: 3.80.55.233 PrivateIPs: 172.31.22.184

2. Flask Application Prediction Output:

```
← ↻ ⚠ Not secure 3.80.55.233:5000/predict ☆
Pretty-print ✓

{
  "current_avg": 5.19085450100086,
  "predicted_next": 5.96948267615098,
  "recommendation": "Scale Down"
}
```

3. CPU Utilization Graph (CloudWatch):



Conclusion

The predictive auto-scaling system successfully demonstrates how machine learning and cloud monitoring tools can be integrated to improve cloud infrastructure efficiency. By

forecasting CPU utilization trends, it enables smarter scaling decisions and reduces unnecessary costs.

Future Work

Future improvements can include:

- Implementing advanced forecasting algorithms like ARIMA or LSTM.
- Automating scale-up/scale-down actions directly using AWS APIs.
- Integrating a user dashboard for real-time visualization and control.

References

- AWS CloudWatch Documentation
- Flask Web Framework
- NumPy and Pandas Python Libraries
- AWS EC2 User Guide