

# Project Report: Cloud Cost Anomaly Detection & Forecasting Using Data Analytics

## 1. Background & Objective

Cloud infrastructure costs are one of the largest expenses for IT organizations. Unexpected spikes, anomalies, or inefficient usage can lead to budget overruns, SLA penalties, and reduced profitability.

This project simulates a real-world IT analytics scenario where we use data science + business intelligence to:

- Detect unusual patterns in daily cloud usage costs.
- Forecast future spend to enable proactive budget planning.
- Provide actionable insights through an interactive Power BI dashboard.

### Business Objective:

- Detect anomalies in daily spend at both global and service levels.
- Forecast 30-day (operational) and 6-month (strategic) costs.
- Empower IT managers with insights to optimize cloud usage and avoid financial risk.

## 2. Data Description

- Synthetic dataset generated with ~1 million rows to simulate cloud usage and cost logs.
- Fields included:
  - timestamp → event time
  - service\_tag → cloud service (RDS, EC2, S3, etc.)
  - team → owning team (analytics, ML, payments, etc.)
  - region → region of deployment
  - usage & usage\_unit → amount of compute/storage/network used
  - cost\_usd → corresponding cost in USD (with injected noise, spikes, and missing values)
  - owner, source\_ip → metadata

### Dataset properties:

- 12 months of daily data.
- Imperfections included: missing values, duplicates, inconsistent service tags, currency conversions, outliers.
- Exported in .csv and .parquet formats.

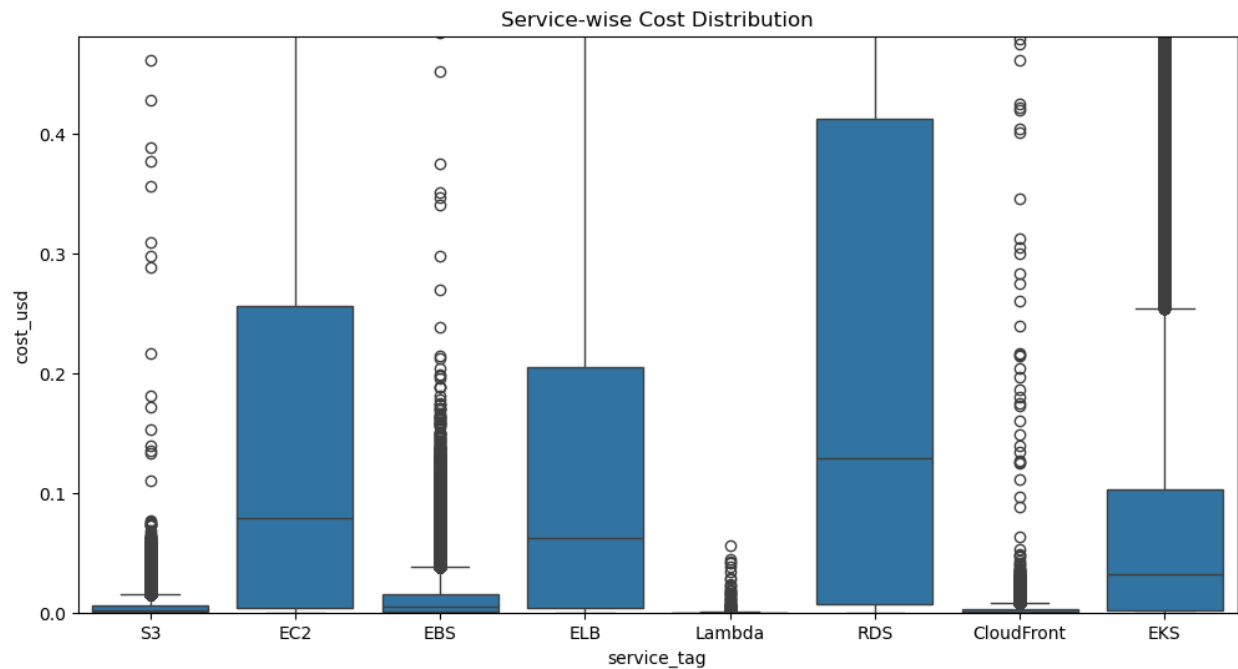
### 3. Methodology – PACE Framework

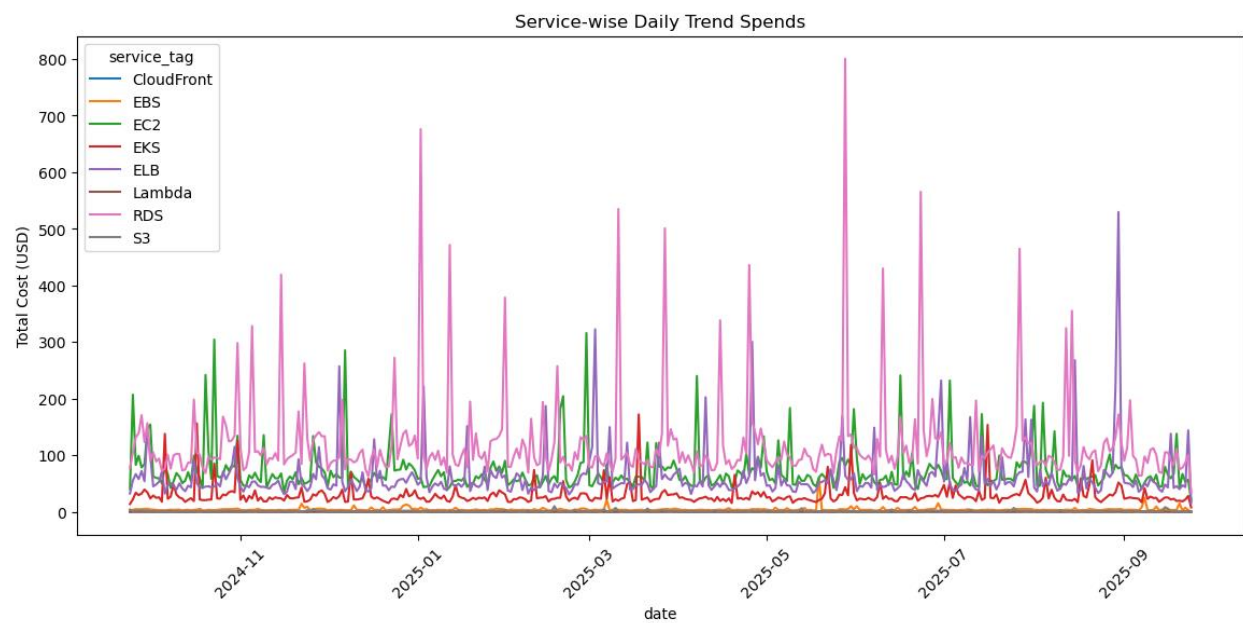
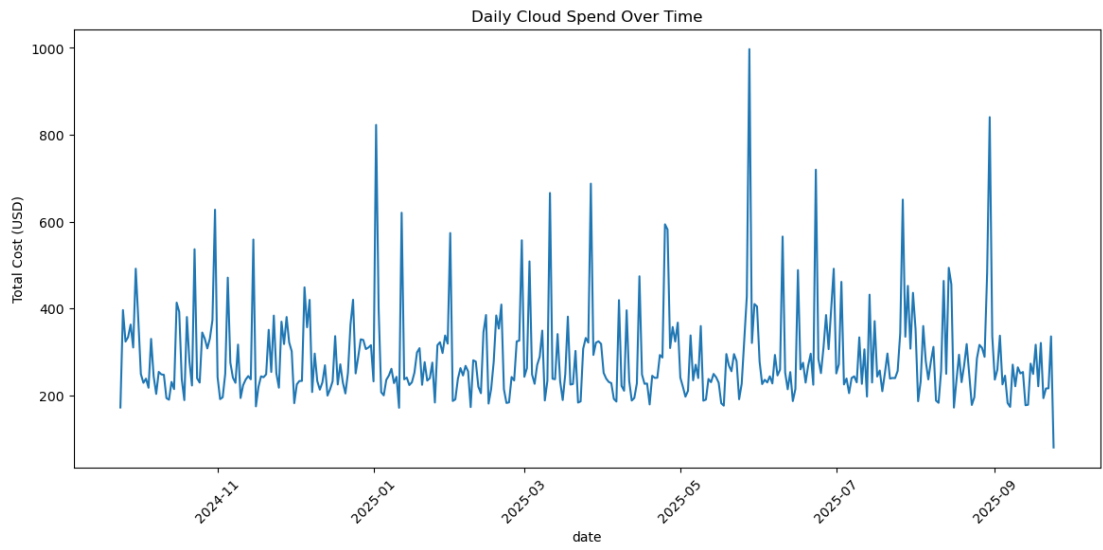
#### Phase 1: Plan

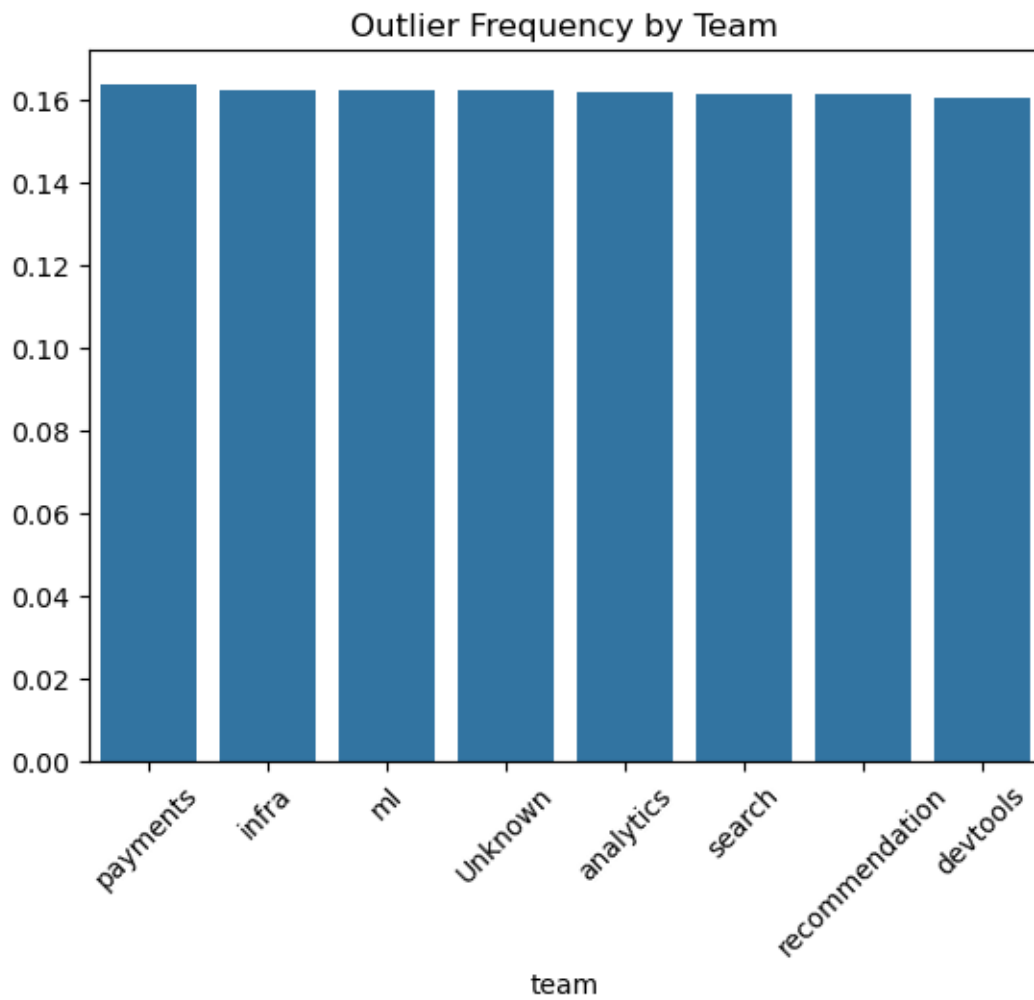
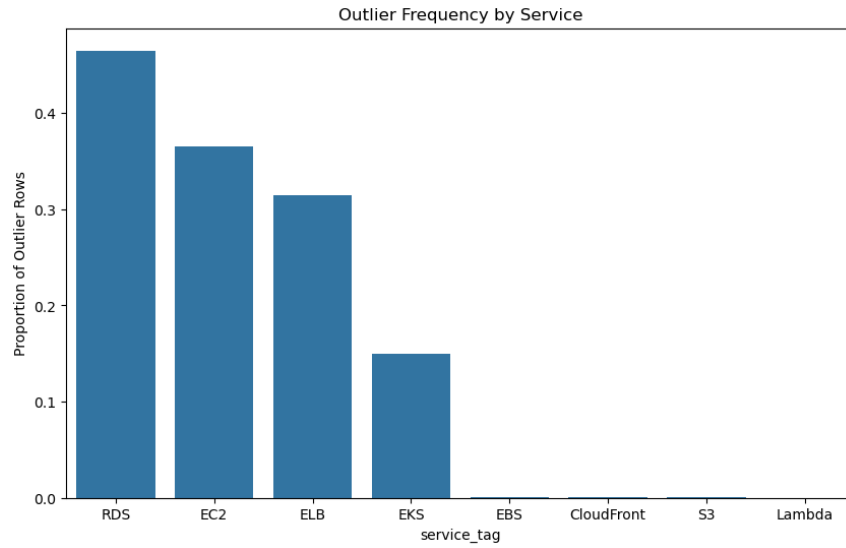
- Define the business problem: cost anomalies + forecasting.
- Identify methods: statistical + ML anomaly detection, Prophet forecasting.
- Design outputs: Jupyter notebooks + Power BI dashboard.

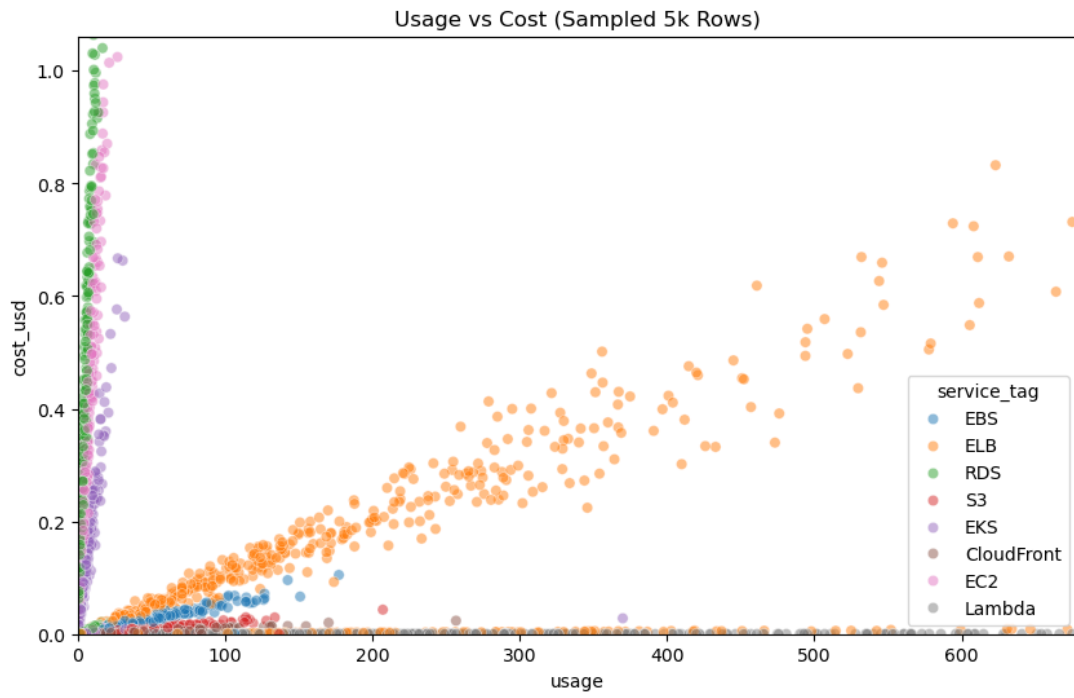
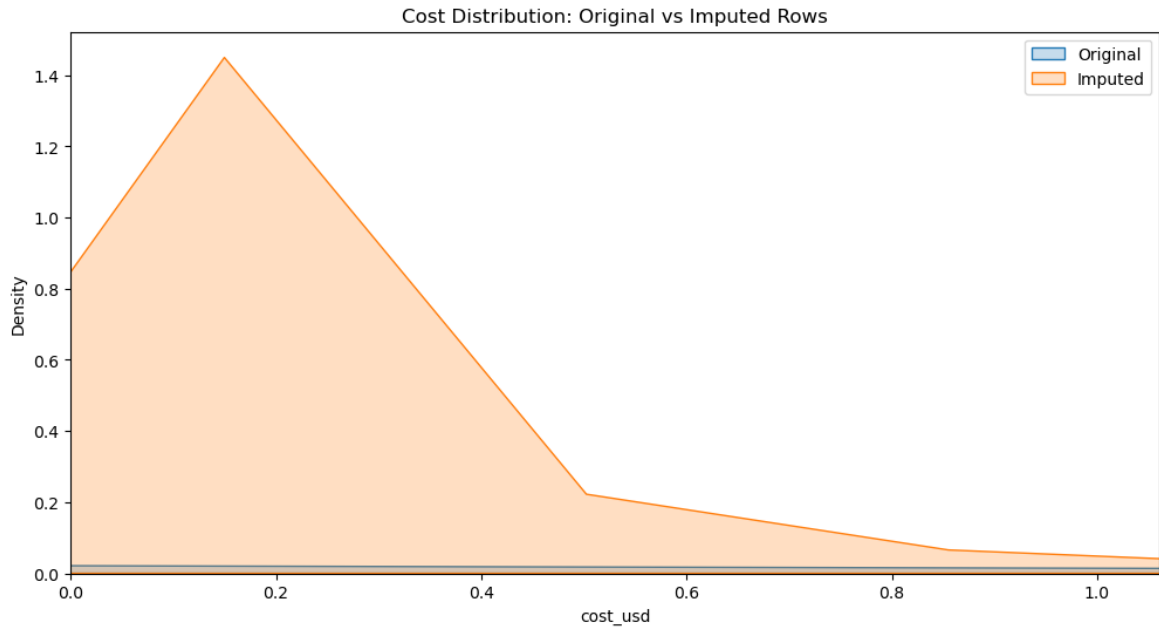
#### Phase 2: Analyze (EDA)

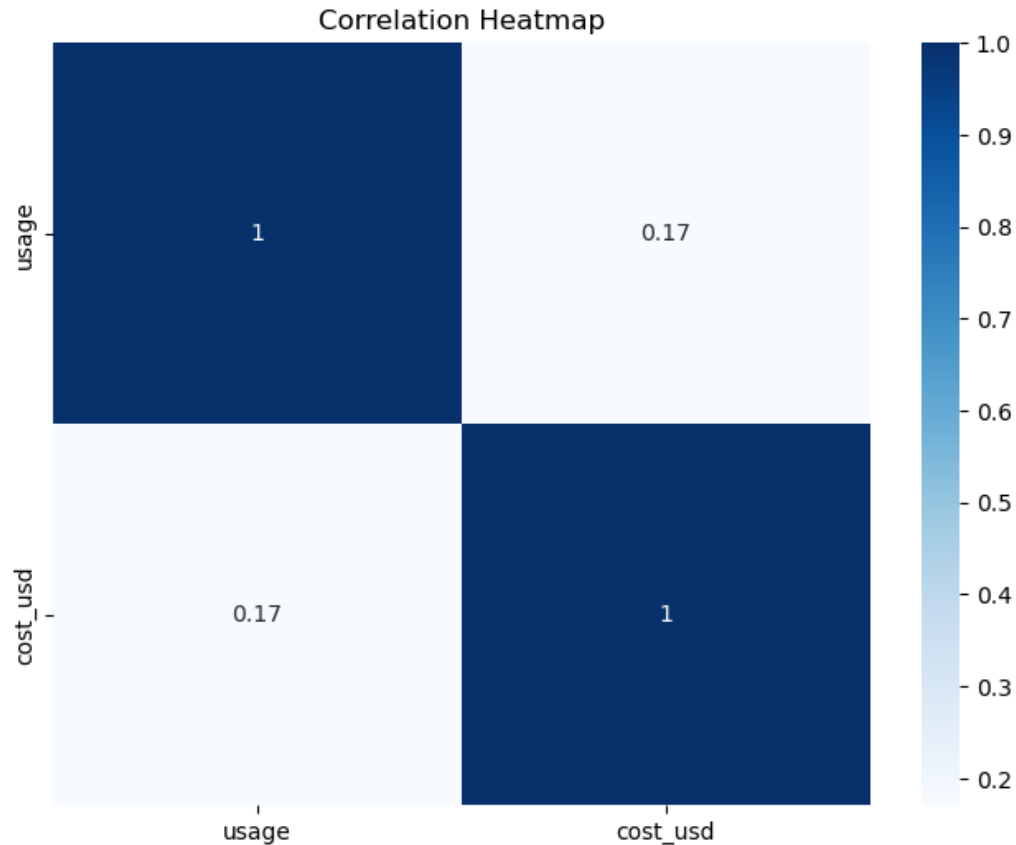
- Data cleaning: handled missing values, duplicates, invalid timestamps, and currency conversion.
- Outlier detection: identified ~16% anomalies in raw usage logs.
- Exploratory analysis:
  - RDS & ELB showed major spikes.
  - Teams contributed ~equally to anomalies, with Payments slightly higher.
  - Daily costs mostly between 200–400 USD, with extreme spikes up to 1000 USD.
- Visualized trends: daily costs, service-level patterns, anomalies by team/region.





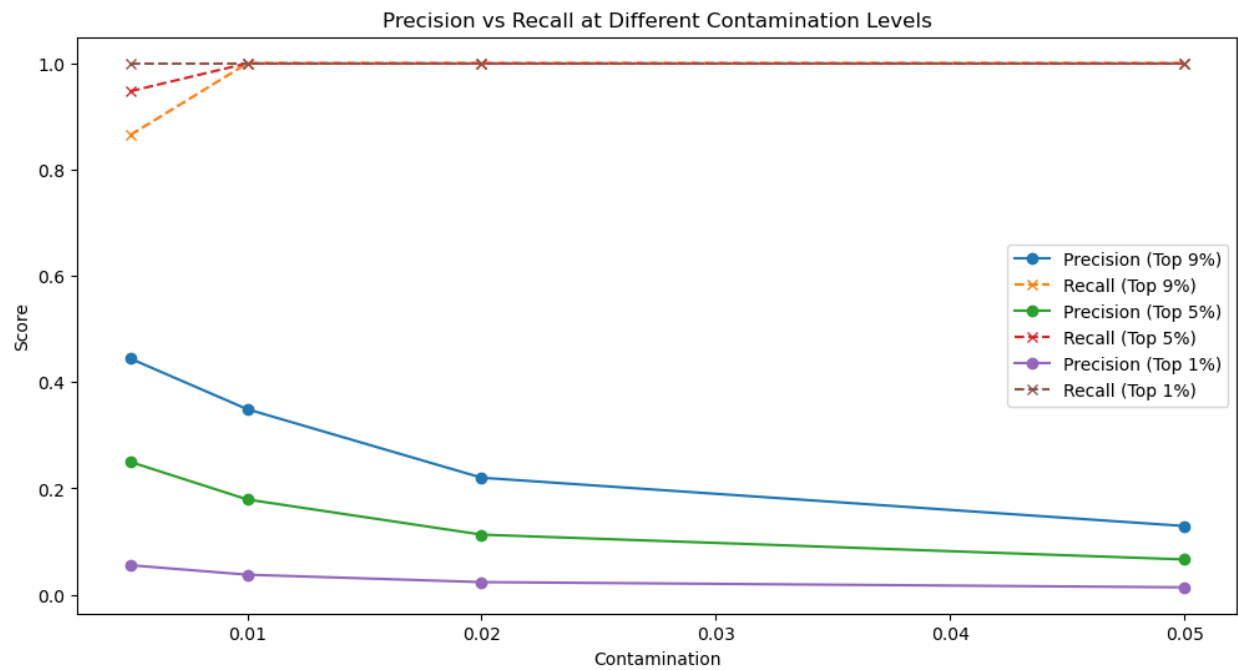
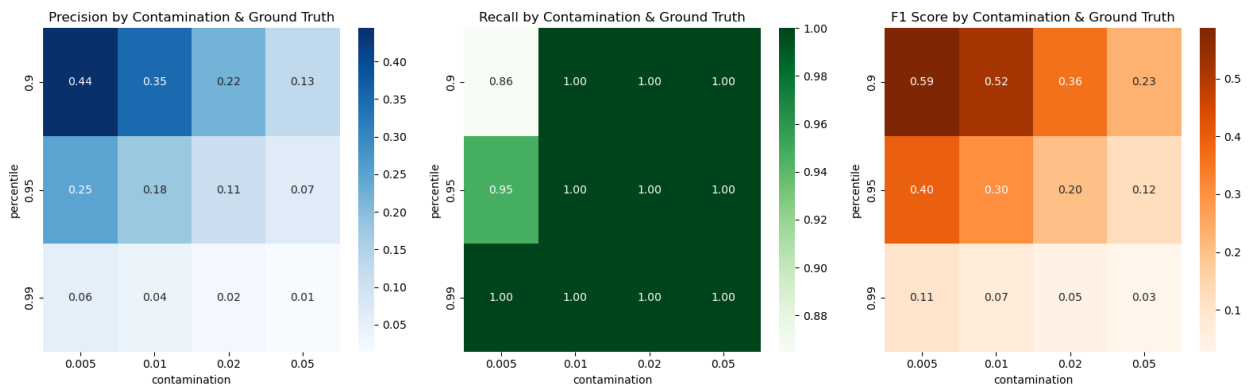
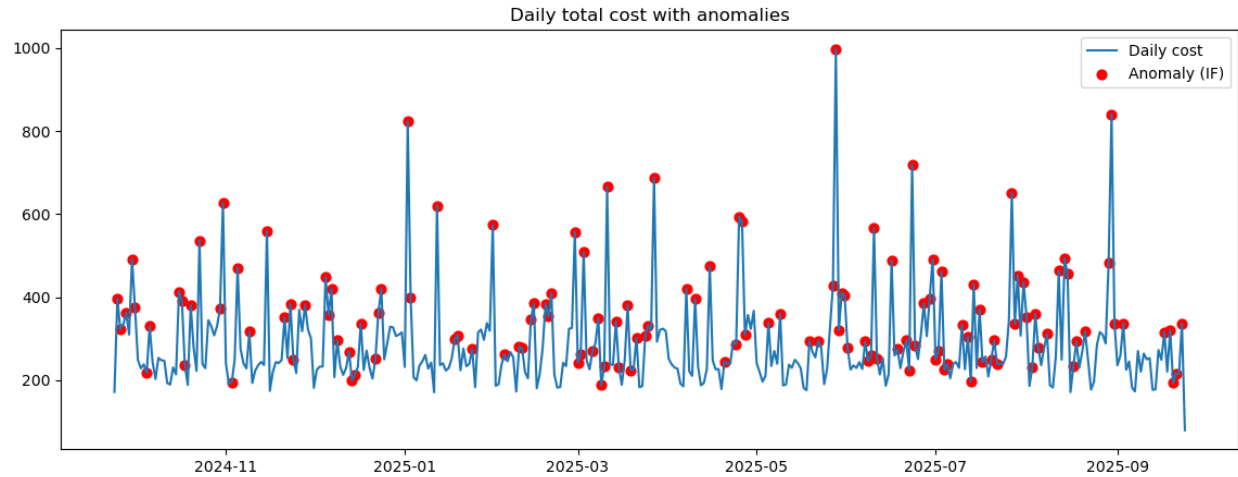






### Phase 3: Construct (Modeling) Anomaly Detection

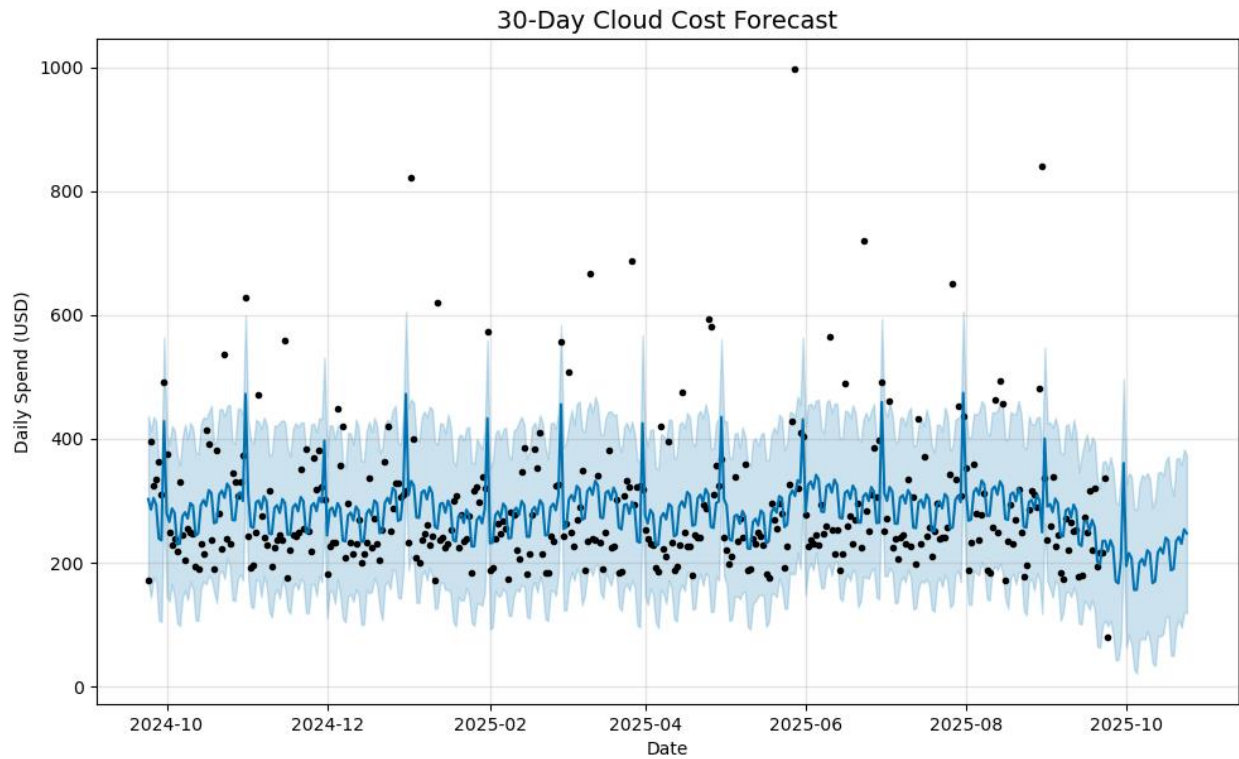
- Statistical Baseline: Z-score using rolling mean & std.
- Machine Learning: Isolation Forest with features:
  - lag costs, rolling mean/std, delta cost, pct change, cost-per-usage, weekend flag.
- Model Tuning: Tested contamination levels (0.5–5%) vs proxy ground truths (top 1%, 5%, 10% spend days).
  - Best trade-off: 0.5% contamination, top 10% ground truth → Precision 44%, Recall 87%, F1 ≈ 0.59.
- Deliverable: service\_level\_anomalies.csv (service-level flagged anomalies).



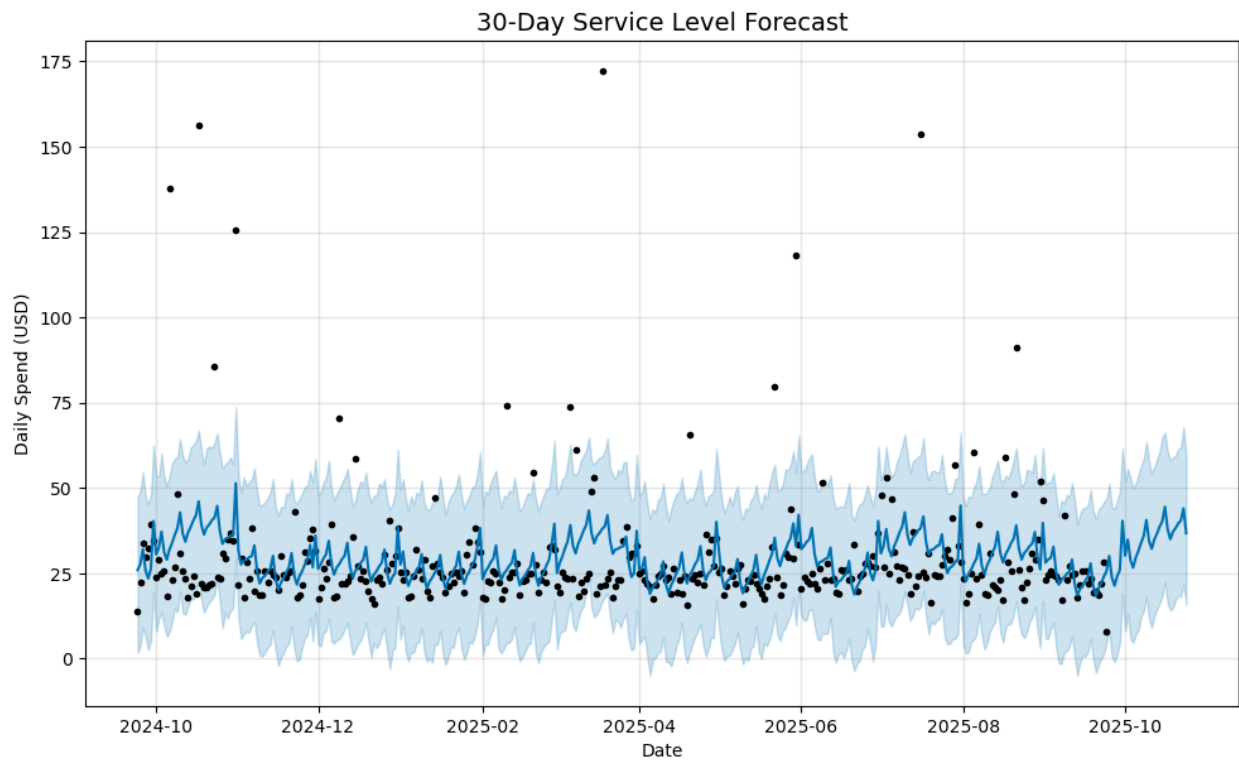
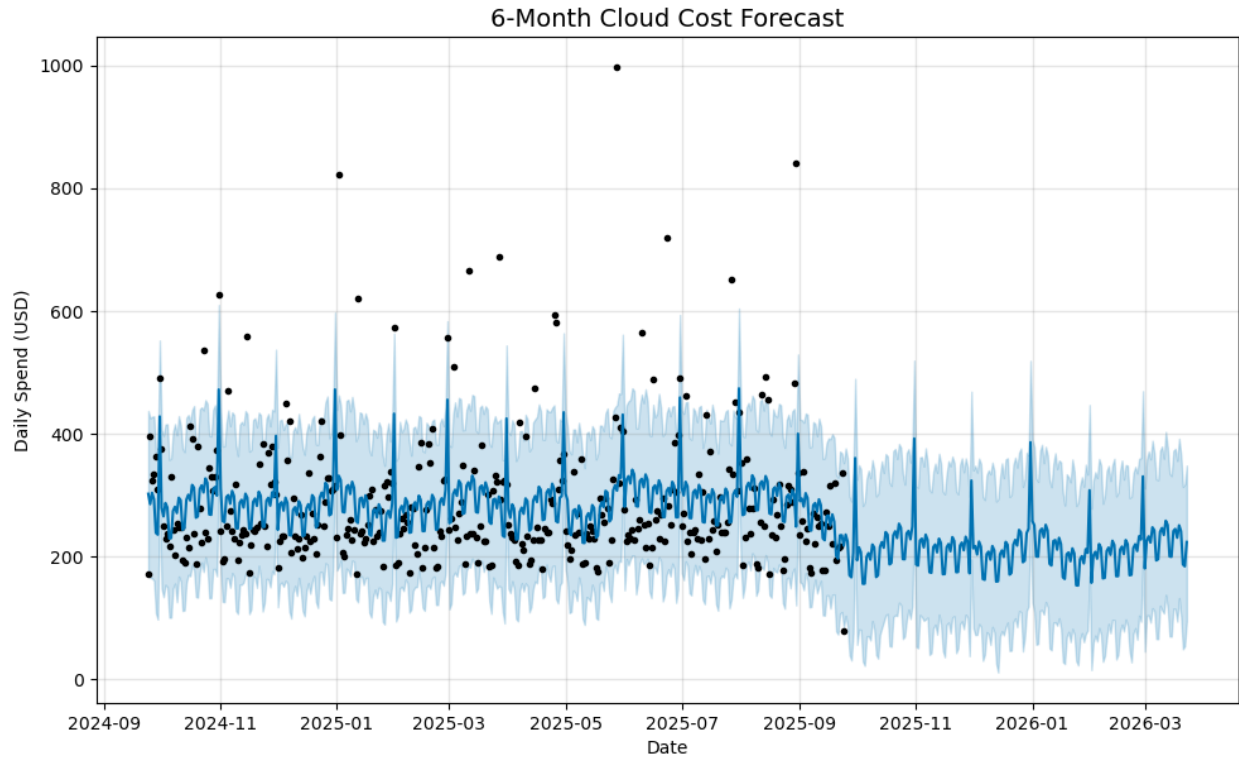
## Forecasting

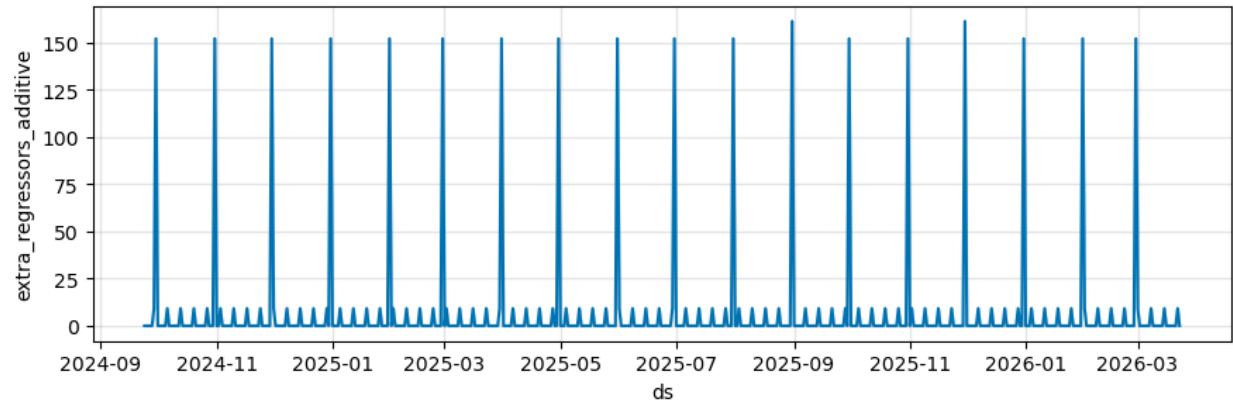
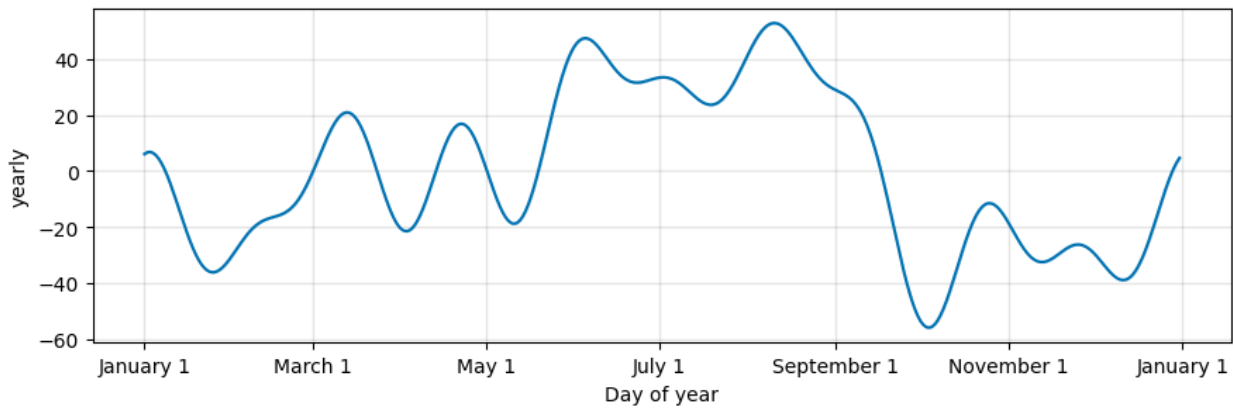
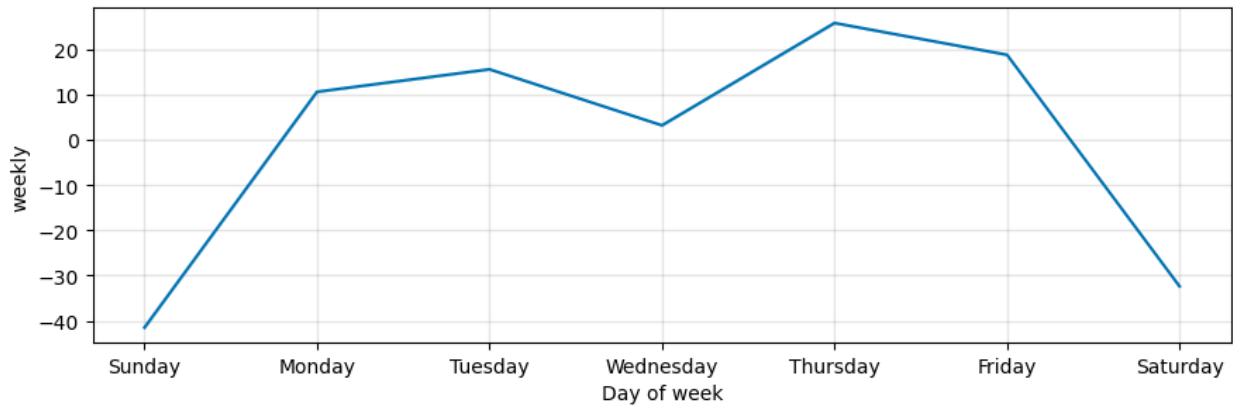
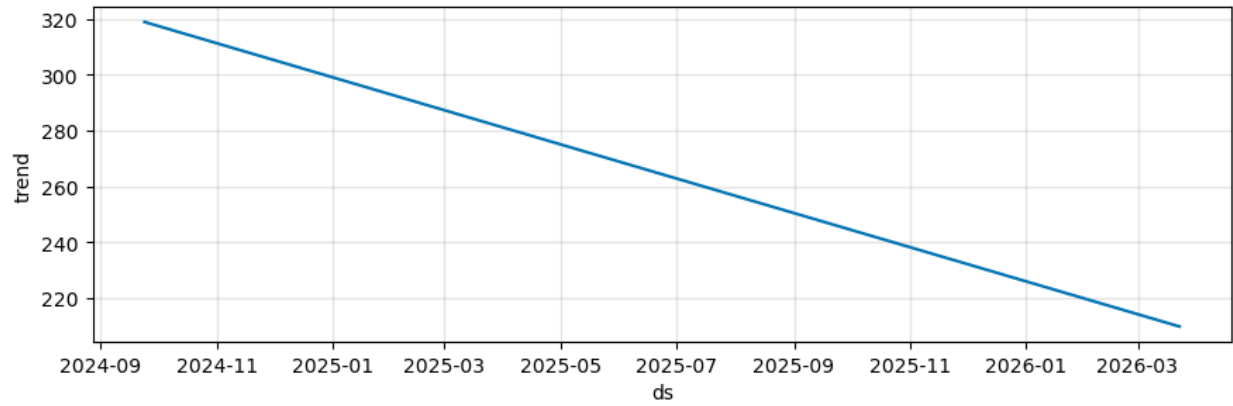
- Tool: Prophet (handles trend + seasonality + regressors).
- Global forecast:

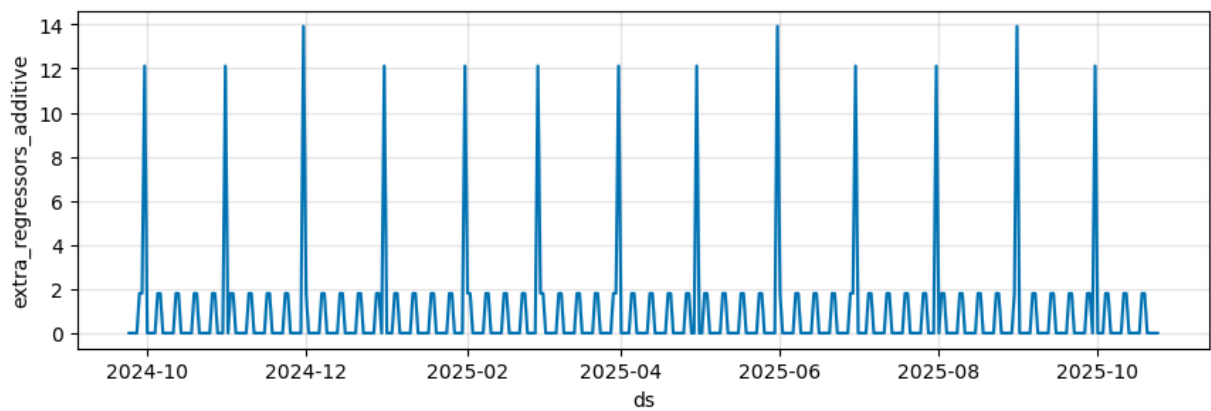
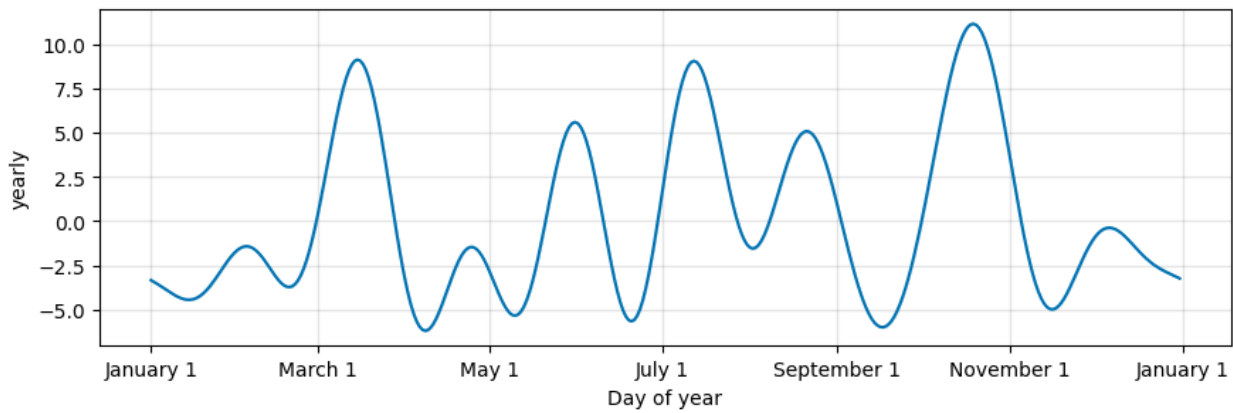
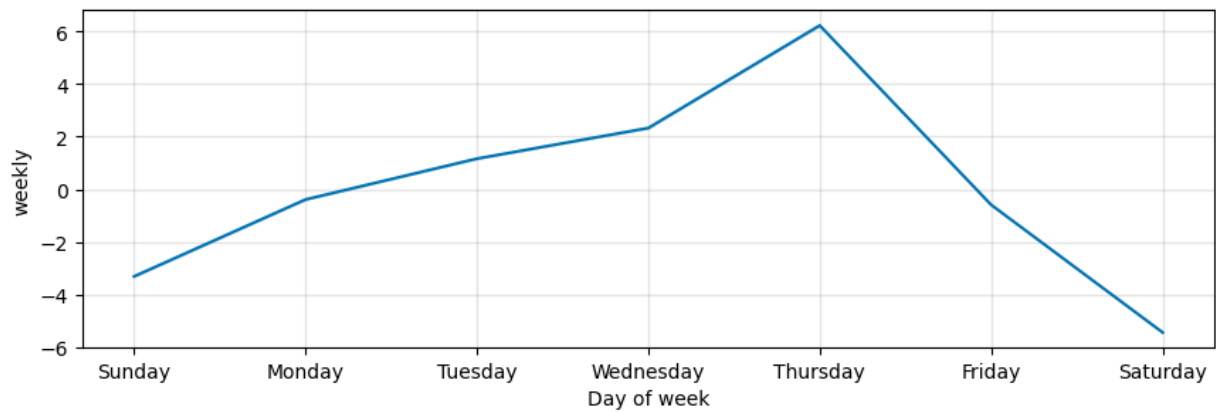
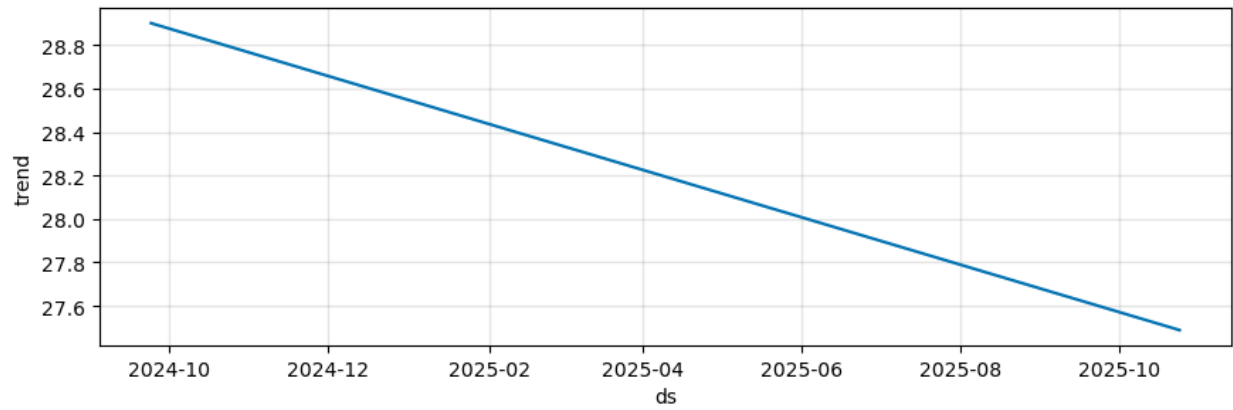
- 30 days: stable trend, costs ~200–400 USD/day.
- 6 months: seasonal peaks in May–July, dips in Sept–Nov.
- Regressors: weekends (lower cost), month-end (batch job spikes).
- Service-level forecasts:
  - RDS & ELB → high volatility.
  - S3 & Lambda → stable & predictable.
- Deliverables:
  - cloud\_cost\_forecast.csv (global).
  - service\_level\_forecasts.csv (per service).











## Phase 4: Execute (Power BI Dashboard)

- Designed interactive dashboard with 4 key pages:
  - Executive Overview: KPIs + global trends.
  - Anomaly Insights: anomalies by service, team, timeline.
  - Forecast Insights: global + per-service forecasts with seasonality patterns.
  - Team & Region Insights: spend by team, region maps, drill-through.
- Filters: service, team, region, date.
- Outputs: .pbix dashboard file + exports.

## 5. Key Insights

- **Anomalies:** RDS & ELB caused the largest spikes, often unrelated to usage (inefficient scaling).
- **Forecasting:** Costs are stable overall, but predictable cycles (weekends, month-end, seasonal peaks) matter for planning.
- **Business Value:**
  - Early detection of cost spikes.
  - Short-term planning for next 30 days.
  - Long-term visibility (6 months) for budgeting.
  - Accountability across services & teams.

## 6. Recommendations

- Set automated alerts for anomalies in RDS & ELB, with thresholds tied to historical rolling std.
- Implement cost governance policies for teams with repeated anomalies.
- Negotiate reserved capacity for predictable services (S3, Lambda) to reduce costs.
- Review weekend/month-end jobs to identify optimization opportunities.
- Use this dashboard as a cloud FinOps tool for continuous monitoring.

## 7. Deliverables

- Jupyter Notebooks (Data Cleaning, EDA, Anomaly Detection, Forecasting).
- Saved ML model (isolation\_forest\_final.joblib + scaler\_final.joblib).
- Forecast & anomaly CSVs for dashboarding.
- Power BI Dashboard (Cloud\_Cost\_Analytics.pbix).
- GitHub repo with code, data samples, and README.