

## ***Domain 4: Monitor Model***



## 4.1 Monitor Model Performance and Data Quality

### 4.1.1 Monitoring Machine Learning Solutions



#### Importance of Monitoring in ML

a) **Machine Learning Lens:** AWS Well-Architected Framework: Best practices and design principles

**Best practice:** When

|                                    |                                 |  |
|------------------------------------|---------------------------------|--|
|                                    |                                 |  |
| Optimize resources                 | <b>Resource pooling</b>         | Sharing compute, storage, and networking resources                       |
|                                    | <b>Caching</b>                  |  |
|                                    | <b>Data management</b>          | data compression, partitioning, and lifecycle management                 |
| Scale ML workloads based on demand | <b>AWS Auto Scaling</b>         | SageMaker built-in scaling. AWS Auto-Scaling                             |
|                                    | <b>Lambda</b>                   |  |
| Reduce Cost                        | <b>Monitor usage and costs</b>  | resource tagging   |
|                                    | <b>monitor ROI</b>              |  |
| Enable continuous improvement      | <b>Establish Feedback Loops</b> |  |
|                                    | <b>Monitor Performance</b>      | SageMaker <b>Model Monitor</b> (Drift)<br>CloudWatch alerts (deviations) |
|                                    | <b>Automate Retraining</b>      |  |

## Detecting Drift in Monitoring

### a) Drift Types

| Drift Type                       | Description   | Causes   | Implications  |
|----------------------------------|---|--|---|
| <b>Data Quality Drift</b>        | Production data distribution differs from training data distribution    | <ul style="list-style-type: none"><li>• Real-world data not as curated as training data</li><li>• Changes in data collection processes</li><li>• Shifts in real-world conditions</li></ul>   | <ul style="list-style-type: none"><li>• Model accuracy decreases</li><li>• Predictions become less reliable</li></ul>   |
| <b>Model Quality Drift</b>       | Model predictions differ from actual ground truth labels                | <ul style="list-style-type: none"><li>• Changes in the underlying relationship between features and target</li><li>• Model decay over time</li><li>• Concept drift</li><li>• Training data too small or not representative</li></ul> | <ul style="list-style-type: none"><li>• Decreased model performance</li><li>• Inaccurate predictions</li></ul>  |
| <b>Bias Drift</b>                | Increase in bias affecting model predictions over time                  | <ul style="list-style-type: none"><li>• Incorporation of societal assumptions in training data</li><li>• Exclusion of important data points</li><li>• Changes in real-world data distribution</li></ul>                              | <ul style="list-style-type: none"><li>• Model overgeneralization</li><li>• Unfair or discriminatory predictions</li><li>• Ethical concerns</li><li>• New groups in production</li></ul> |
| <b>Feature Attribution Drift</b> | Changes in the contribution of individual features to model predictions | <ul style="list-style-type: none"><li>• Shifts in feature importance over time</li><li>• Changes in the underlying problem domain</li><li>• Introduction of new, more predictive features</li></ul>                                  | <ul style="list-style-type: none"><li>• Model may rely on less relevant features</li><li>• Decreased interpretability</li><li>• Potentially reduced performance</li></ul>               |

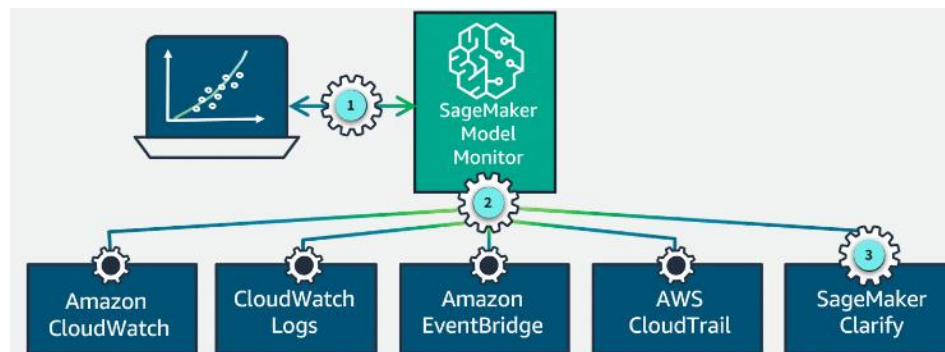
**Note:** Bias inverse of variance, which is the level of small fluctuations or noise common in complex data sets. Bias tends to cause model predictions to overgeneralize, and variance tends to cause models to undergeneralize. Increasing variance is one method for reducing the impact of bias.

b) *Monitoring Drift*

| Monitoring Type                             | What It Monitors  | How It Monitors   |
|---|---|---|
| <b>Data Quality Monitoring</b>              | <ul style="list-style-type: none"><li>• Missing values</li><li>• Outliers</li><li>• Data types</li><li>• Statistical metrics (mean, std dev, etc.)</li><li>• Data distribution</li></ul>                | <ul style="list-style-type: none"><li>• Implement data validation checks</li><li>• Calculate statistical metrics</li><li>• Compare metrics with baseline values</li><li>• Use data drift detection techniques (e.g., Kolmogorov-Smirnov tests, Maximum Mean Discrepancy)</li></ul>                                  |
| <b>Model Quality Monitoring</b>             | <ul style="list-style-type: none"><li>• Evaluation metrics (accuracy, precision, recall, F1, AUC, etc.)</li><li>• Prediction confidence</li><li>• Performance across different subpopulations</li></ul> | <ul style="list-style-type: none"><li>• Calculate evaluation metrics on held-out test set or production data sample</li><li>• Implement confidence thresholding or uncertainty estimation</li><li>• Flag low-confidence predictions</li><li>• Monitor performance on different data subsets</li></ul>               |
| <b>Model Bias Drift Monitoring</b>          | <ul style="list-style-type: none"><li>• Bias metrics (disparate impact, fairness, etc.)</li><li>• Performance across sensitive groups</li></ul>   | <ul style="list-style-type: none"><li>• Calculate bias metrics for different sensitive groups</li><li>• Compare bias metrics with baseline values or thresholds</li><li>• Implement bias mitigation techniques (e.g., adversarial debiasing, calibrated equalized odds)</li></ul>                                   |
| <b>Feature Attribution Drift Monitoring</b> | <ul style="list-style-type: none"><li>• Feature importance scores</li><li>• Statistical metrics of feature attributions</li></ul>   | <ul style="list-style-type: none"><li>• Use interpretability techniques (e.g., SHAP) to calculate feature attributions</li><li>• Calculate statistical metrics on feature attributions</li><li>• Compare metrics with baseline values</li><li>• Identify features with significantly changed attributions</li></ul> |

## SageMaker Model Monitor

### Integration



### SageMaker - Monitoring for Data Quality Drift



#### STEPS

1. **Initiate data capture on the endpoint**
2. **Create a baseline**

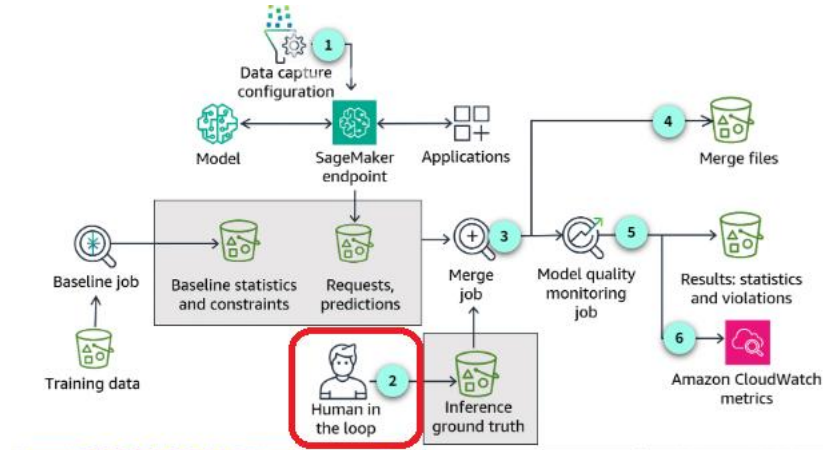
start a baseline processing job with the `suggest_baseline` method of the `ModelQualityMonitor` object using the SageMaker Python SDK.

3. **Schedule data quality monitoring jobs**
4. **Integrate data quality monitoring with Cloudwatch**
5. **Interpret results and analyze findings**

The report is generated as the `constraint_violations.json` file. The SageMaker Model Monitor prebuilt container provides the following violation checks.

- `data_type_check`
- `completeness_check`
- `baseline_drift_check`
- `missing_column_check`
- `extra_column_check`
- `categorical_values_check`

## SageMaker - Monitoring for Model Quality Drift using Model Monitor



To monitor model quality, SageMaker Model Monitor requires the following inputs:

1. Baseline data
2. Inference input and predictions made by the deployed model
3. Amazon SageMaker Ground Truth associated with the inputs to the model

## SageMaker - Monitoring for Bias using Clarify

Statistical bias drift occurs when the data used for training differs from the data encountered during prediction, leading to potentially biased outcomes. This is prominent when training data changes over time. In this lesson, you will learn about AWS services that help you monitor for statistical bias drift.

Post-training bias metrics in SageMaker Clarify help us answer two key questions:

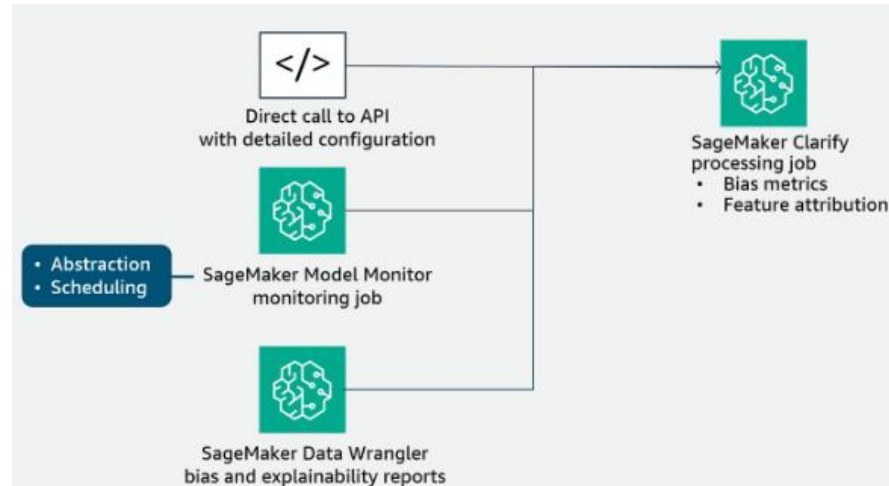
- Are all facet values represented at a similar rate in positive (favorable) model predictions?
- Does the model have similar predictive performance for all facet values?

**SageMaker Model Monitor automatically does the following:**

- **Merges** the prediction data with SageMaker Ground Truth labels
- **Computes** baseline statistics and constraints
- **Inspects** the merged data and generates bias metrics and violations
- **Emits** CloudWatch metrics to set up alerts and triggers
- **Reports and alerts** on bias drift detection
- **Provides** reports for visual analysis

**How it works:** It quantifies the contribution of each input feature (for example, audio characteristics) to the model's predictions, helping to explain how the model arrives at its decisions.

## Options for using SageMaker Clarify

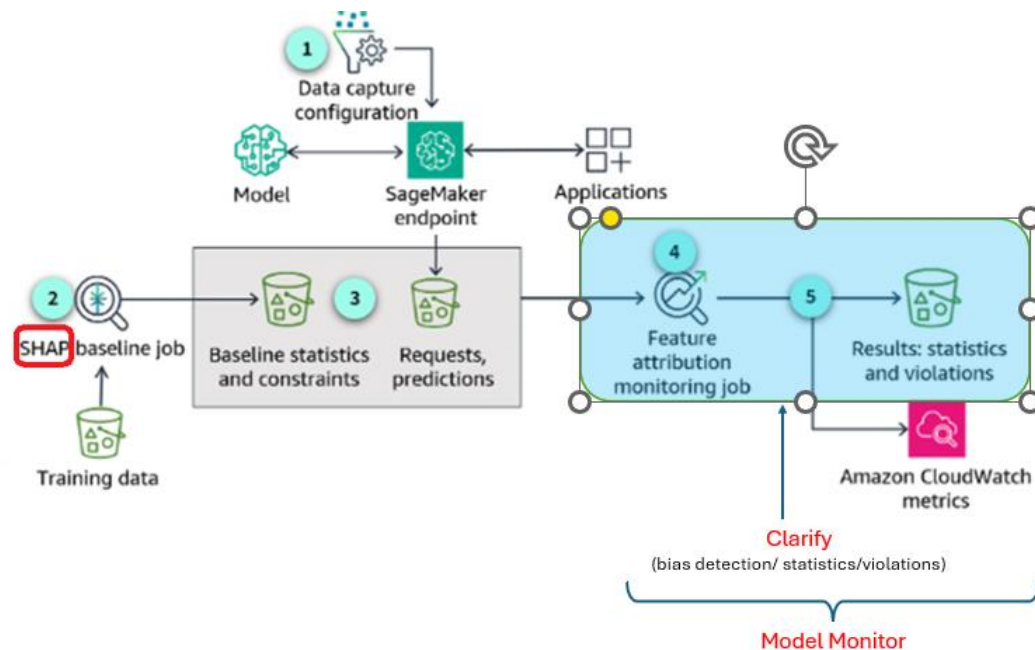


### When to use which

| Method                                   | Description   | When to Use  |
|--|---|--|
| <b>SageMaker Clarify Directly</b>        | Configure and run Clarify processing job using SageMaker Python SDK API               | <ul style="list-style-type: none"><li>• For <b>one-time or ad-hoc</b> bias analysis</li><li>• When you <b>need full control</b> over the analysis configuration</li><li>• For integrating bias analysis <b>into custom workflows</b></li></ul>       |
| <b>SageMaker Model Monitor + Clarify</b> | <b>Integrate Clarify with Model Monitor</b><br>Monitor for continuous bias monitoring | <ul style="list-style-type: none"><li>• When you want to <b>automate bias detection in production</b></li><li>• If you need to set up alerts for bias drift</li></ul>  |
| <b>SageMaker Data Wrangler</b>           | Utilize Clarify within Data Wrangler <b>during data preparation</b>                   | <ul style="list-style-type: none"><li>• <b>During the data preparation phase</b></li><li>• When <b>you want to identify potential bias early in the ML pipeline</b></li><li>• If you're already using Data Wrangler for data preprocessing</li></ul> |

## SageMaker - Monitoring for Feature Attribution Drift (**Model Monitor + Clarify**)

Feature attribution refers to understanding and quantifying the contribution or influence of each feature on the model's predictions or outputs. It helps to identify the most relevant features and their relative importance in the decision-making process of the model.



### Uses SHAP

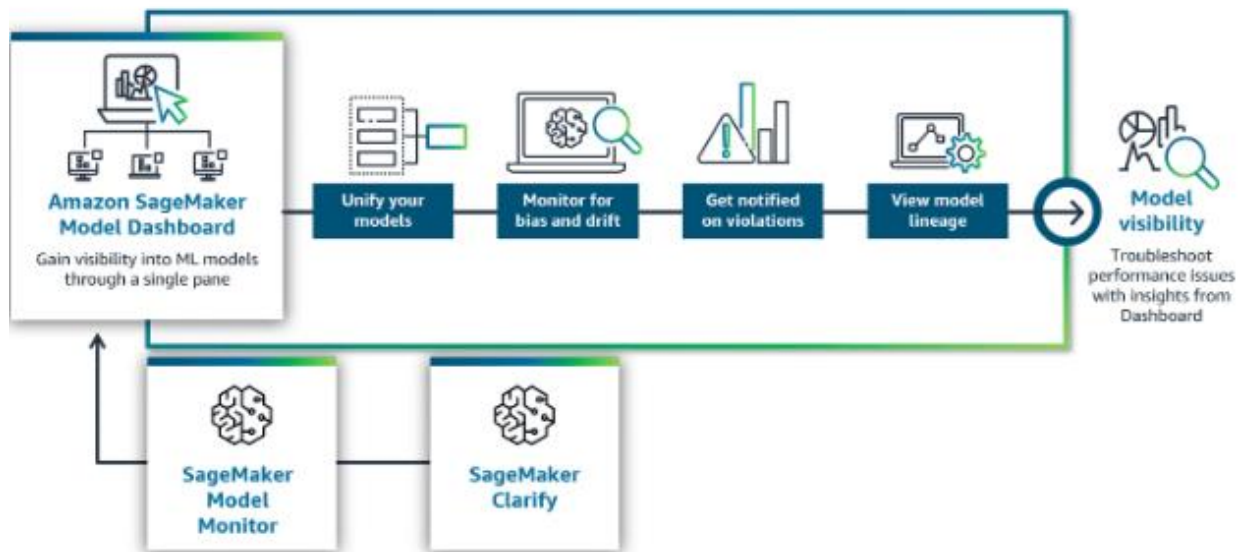
SageMaker Clarify provides feature attributions based on the concept of Shapley value. This is a game-theoretic approach that assigns an importance value (SHAP value) to each feature for a particular prediction.

### Here's how it works:

1. **SageMaker Clarify:** This is the core component that performs the actual bias detection and generate quality metrics and violations
2. **SageMaker Model Monitor:** This is the framework that can use Clarify's capabilities to perform continuous monitoring of deployed models.



## SageMaker Model Dashboard



### Features

#### 1. Alerts :

How it helps: The dashboard provides a record of all activated alerts, allowing the data scientist to review and analyze past issues.

Alert criteria depend upon two parameters:

- **Datapoints to alert:** Within the evaluation period, how many runtime failures raise an alert?
- **Evaluation period:** The # of most recent monitoring executions to consider when evaluating alert status.

#### 2. Risk rating

A user-specified parameter from the model card with a low, medium, or high value.

#### 3. Endpoint performance

You can select the endpoint column to view performance metrics, such as:

- **CpuUtilization:** The sum of each individual CPU core's utilization from 0%-100%.
- **MemoryUtilization:** The % of memory used by the containers on an instance, 0%-100%.
- **DiskUtilization:** The % of disk space used by the containers on an instance, 0%-100%.

#### 4. Most recent batch transform job

This information helps you determine if a model is actively used for batch inference.

#### 5. Model lineage graphs

When training a model, SageMaker creates a **model lineage graph**, a visualization of the entire ML workflow from data preparation to deployment.

#### 6. Links to model details

The dashboard links to a model details page where you can explore an individual model.

**Model Monitor vs SageMaker Dashboard vs Clarify: When to use which one**

| Tool                       | Description                                      | Why to use   | When to Use  |
|----------------------------|--|--|--|
| <b>Model Monitor</b>       | Continuous monitoring of ML models in production | <ul style="list-style-type: none"><li>• data and model quality issues</li><li>• model drift</li></ul>                | <ul style="list-style-type: none"><li>• To set up automated alerts for performance degradation</li><li>• When you need to monitor resource utilization</li><li>• Monitor real-time endpoints, batch transform, On-demand monitoring job</li></ul>                                  |
| <b>SageMaker Dashboard</b> | Centralized view of SageMaker resources and jobs | <ul style="list-style-type: none"><li>•</li></ul>  | <ul style="list-style-type: none"><li>• For a high-level overview of all SageMaker activities</li><li>• To track training jobs, endpoints, and notebook instances</li></ul>  |
| <b>SageMaker Clarify</b>   | Bias detection and model explainability tool     | <ul style="list-style-type: none"><li>• Detecting Bias</li><li>• Triggers statistics and Violations report</li></ul> | <ul style="list-style-type: none"><li>• To detect bias in training data and model predictions</li><li>• When you need to explain model decisions</li><li>• For regulatory compliance requiring model transparency</li><li>• To improve model fairness and accountability</li></ul> |

## 4.1.2 Remediating Problems Identified by Monitoring

### Automated remediations and notifications

- **Stakeholder notifications:** When monitoring metrics indicate **changes that impact business KPIs** or the **underlying problem**
- **Data Scientist notification:** You can use automated notifications to **data scientists** when your monitoring **detects data drift** or when **expected data is missing**.
- **Model retraining:** Configure your model training pipeline to **automatically retrain models when monitoring detects drift, bias, or performance degradation**.
- **Autoscaling:** You use resource utilization metrics gathered by infrastructure monitoring to initiate autoscaling actions.

### Model retraining strategies

| Strategy     | When to Use  | Advantages   | Considerations  |
|--------------|--|--|---|
| Event-driven | <ul style="list-style-type: none"><li>• When drift is detected above a certain threshold</li><li>• In response to significant changes in data or performance</li></ul> | <ul style="list-style-type: none"><li>• Timely response to changes</li><li>• Efficient use of resources</li></ul>                              | <ul style="list-style-type: none"><li>• May be frequent if thresholds are too sensitive</li><li>• Retraining can be expensive and time-consuming</li></ul>  |
| On-demand    | <ul style="list-style-type: none"><li>• When market conditions change significantly</li><li>• In response to new competitors or strategies</li></ul>                   | <ul style="list-style-type: none"><li>• Allows for human judgment in decision-making</li><li>• Can incorporate business context</li></ul>      | <ul style="list-style-type: none"><li>• Requires constant monitoring by data scientists or stakeholders</li><li>• May lead to delayed responses</li></ul>   |
| Scheduled    | <ul style="list-style-type: none"><li>• When there are known seasonal patterns</li><li>• For maintaining model accuracy over time</li></ul>                            | <ul style="list-style-type: none"><li>• Predictable maintenance schedule</li><li>• Can anticipate and prepare for retraining periods</li></ul> | <ul style="list-style-type: none"><li>• May retrain unnecessarily if no significant changes occur</li><li>• Might miss sudden, unexpected changes</li></ul> |

## 4.2 Monitor and Optimize Infrastructure and Costs

### 4.2.1 Monitor Infrastructure

#### Monitor Performance Metrics - CloudWatch vs Model Monitor

| Feature                 | SageMaker Model Monitor   | CloudWatch Logs   |
|-------------------------|---|---|
| <b>Purpose</b>          | Continuous <i>monitoring of ML models</i> in production   | Monitoring, storing, and accessing log files  |
| <b>Key Capabilities</b> | (all 4 ML monitoring types) <ul style="list-style-type: none"><li>• Data quality monitoring</li><li>• Model quality monitoring</li><li>• Bias drift monitoring</li><li>• Feature attribution drift monitoring</li></ul> | <ul style="list-style-type: none"><li>• Log collection from various sources</li><li>• Log storage in S3</li><li>• Pattern recognition</li><li>• Log anomaly detection</li></ul> |
| <b>Monitoring Types</b> | <ul style="list-style-type: none"><li>• Real-time endpoint monitoring</li><li>• Batch transform job monitoring</li><li>• On-schedule monitoring for async batch jobs</li></ul>  | <ul style="list-style-type: none"><li>• EC2 instances</li><li>• CloudTrail</li><li>• Amazon Route 53</li><li>• Other sources</li></ul>  |
| <b>Alert System</b>     | Set alerts for deviations in model quality  | Notifications based on preset thresholds  |
| <b>Customization</b>    | <ul style="list-style-type: none"><li>• Pre-built monitoring capabilities (no coding)</li><li>• Custom analysis options</li></ul>   | Customizable log patterns and anomaly detection   |

#### Monitoring vs. Observability

|                       | Monitoring  | Observability   |
|-----------------------|---|---|
| <b>Definition</b>     | Continuous collection and analysis of metrics   | Deep insights into internal state and behavior of ML systems  |
| <b>Focus</b>          | Detecting anomalies and deviations  | Understanding complex interactions and dependencies   |
| <b>Key Activities</b> | <ul style="list-style-type: none"><li>• Collecting metrics</li><li>• Logging</li><li>• Alerting</li></ul>                           | <ul style="list-style-type: none"><li>• Analyzing system behavior</li><li>• Identifying root causes</li><li>• Reasoning about system health</li></ul> |
| <b>Techniques</b>     | <ul style="list-style-type: none"><li>• Metric collection</li><li>• Threshold-based alerting</li><li>• Basic log analysis</li></ul> | <ul style="list-style-type: none"><li>• Distributed tracing</li><li>• Structured logging</li><li>• Advanced data visualization</li></ul>              |
| <b>Outcome</b>        | Detect issues and invoke alerts or automated actions  | Provide deeper insights for troubleshooting and optimization  |
| <b>Scope</b>          | Primarily focused on predefined metrics and thresholds  | Enables asking and answering questions about system behavior  |

## Monitoring Tools (for Performance and Latency)

| Feature                    | AWS X-Ray  | CloudWatch Lambda Insights   | CloudWatch Logs Insights  | QuickSight  |
|----------------------------|--|--|---|---|
| <b>Purpose</b>             | Trace information about responses and calls in applications  | In-depth performance monitoring for <b>Lambda fns only</b>   | Interactive log analytics service   | BI and data visualization service   |
| <b>Key Features</b>        | <ul style="list-style-type: none"> <li>• Works across AWS and third-party services</li> <li>• Generates detailed service graphs</li> <li>• Identifies performance bottlenecks</li> </ul> | <ul style="list-style-type: none"> <li>• Monitors metrics (memory, duration, invocation count)</li> <li>• Provides detailed logs and traces</li> <li>• Helps identify bottlenecks in Lambda functions</li> </ul> | <ul style="list-style-type: none"> <li>• Interactive querying and analysis of log data</li> <li>• Correlates log data from different sources</li> <li>• Visualizes time series data</li> <li>• Supports aggregations, filters, and regex</li> </ul> | <ul style="list-style-type: none"> <li>• Interactive dashboards</li> <li>• ML-powered insights</li> <li>• Supports various data sources</li> </ul>          |
| <b>Compatible Services</b> | <b>EC2, ECS, Lambda, Elastic Beanstalk</b>   | <b>Lambda</b>  | Any service that generates logs in CloudWatch   | Various AWS services and external data sources  |
| <b>ML Use Cases</b>        | <ul style="list-style-type: none"> <li>• Analyze bottlenecks in ML systems</li> <li>• Trace requests in ML applications (e.g., chatbot inference)</li> </ul>                             | <ul style="list-style-type: none"> <li>• Monitor and optimize ML models deployed as Lambda functions</li> <li>• Identify root causes of Lambda function issues</li> </ul>  | <ul style="list-style-type: none"> <li>• Analyze logs from ML workloads</li> <li>• Identify patterns and anomalies in ML system behavior</li> </ul>   | <ul style="list-style-type: none"> <li>• Create dashboards for ML experiment results</li> <li>• Analyze and present insights from ML predictions</li> </ul> |
| <b>Visualization</b>       | Service maps, Trace views  | Performance dashboards, Trace details  | Time series graphs, Log event views   | Interactive dashboards, Charts, Graphs  |
| <b>Primary Benefit</b>     | <b>End-to-end</b> request tracing and bottleneck identification  | Detailed Lambda function performance insights  | Flexible, interactive log analysis and visualization  | Comprehensive data visualization and business intelligence  |

## SageMaker w/ EventBridge

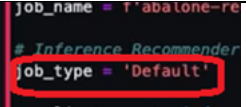
### Actions that can be automatically invoked using EventBridge:

- Invoking an AWS Lambda function
- Invoking Amazon **EC2 run** command (**not create or deploy**)
- Relaying event to **Kinesis Data Streams**
- Activating an **AWS Step Functions** state machine.
- Notifying SNS topic** or an **AWS Server Migration Service (AWS SMS) queue**.

## 4.2.2 Optimize Infrastructure

### Inference Recommender types

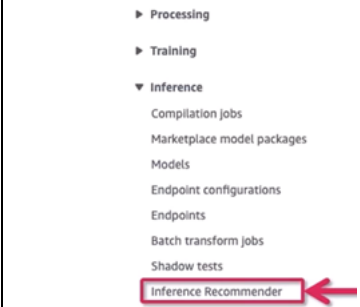
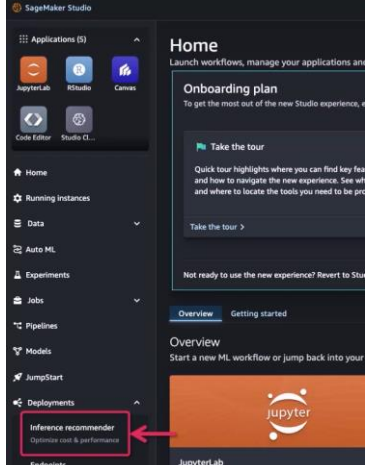
#### a) Inference Recommendation Types

| Default   | Advanced                                     |
|---|--|
| Endpoint Recommender  | Endpoint Recommender + Inference Recommender |
| 45 mins   | 2 hours                                      |
|  |  |
|   |  |

#### b) Endpoint Recommender vs Inference Recommender

|                     | Endpoint Recommender                              | Inference Recommender   |
|---------------------|---|---|
| Output              | list (or ranking) of <b>prospective</b> instances | <b>Same</b>   |
|                     | run a set of load tests.                          | based on a <b>custom</b> load test.   |
| What you need to do | <b>-N/A -</b>                                     | your desired ML instances or a serverless endpoint, provide a custom traffic pattern, and provide requirements for latency and throughput |

#### c) How to start

| SageMaker   | Inference Recommender  |
|---|--|
|  |  |

#### d) Sample Recommender output

| Inference recommendations  |           |                |               |               |                      |  |
|--|-----------|----------------|---------------|---------------|----------------------|--|
| Inference recommendations help you select the best instance type and configuration (such as instance count, container parameters, and model optimizations) for your ML models and workloads. |           |                |               |               |                      |  |
| Instance type  | Status    | Instance count | Model Latency | Cost per hour | Cost per millisecond |  |
| ml.c5.xlarge   | Completed | 1              | 135ms         | \$0.20        | \$0.85               |  |
| ml.g4dn.xlarge   | Completed | 1              | 69ms          | \$0.94        | \$6.86               |  |
| ml.c5.large  | Completed | 1              | 102ms         | \$0.10        | \$0.66               |  |
|  | Failed    | -              | -             | -             | -                    |  |
|  | Failed    | -              | -             | -             | -                    |  |
| ml.c5.2xlarge  | Completed | 1              | 198ms         | \$0.41        | \$0.64               |  |
|  | Failed    | -              | -             | -             | -                    |  |
| ml.c5.2xlarge (c)  | Completed | 1              | 56ms          | \$0.41        | \$0.48               |  |
| ml.g4dn.xlarge   | Completed | 1              | 86ms          | \$0.74        | \$5.78               |  |

### 4.2.3 Optimize Costs

#### Inference Recommender types

| Option                             | Description   | Best For                                    | Cost Savings                  | Example Use Case                                |
|------------------------------------|---|---|-------------------------------|---|
| <b>Spot Instances</b>              | Spare EC2 capacity at lower prices; <b>can be interrupted</b> | Interruptible workloads                     | <b>Up to 90% vs On-Demand</b> | Data preprocessing or batch processing          |
| <b>On-Demand Instances</b>         | Pay-per-use with no long-term commitment                      | Short-term, unpredictable workloads         | <b>None</b> (baseline)        | Real-time inference services                    |
| <b>Reserved Instances</b>          | Discounted rates for <b>1 or 3-year commitments</b>           | Steady-state, predictable workloads         | <b>Up to 72% vs On-Demand</b> | Long-running ML training jobs                   |
| <b>Capacity Blocks</b>             | Reserved capacity for AWS Outposts or Wavelength Zones        | <b>Ensuring capacity during peak demand</b> | Varies                        | ML workloads in <b>on-premises environments</b> |
| <b>Savings Plans for SageMaker</b> | Commit to a specific <b>compute</b> usage for 1 or 3 years    | Flexible, recurring SageMaker usage         | Up to <b>64% vs On-Demand</b> | <b>Regular model training and deployment</b>    |

## 4.3 Secure AWS ML Resources

### 4.3.1 Securing ML Resources

#### Access Control using IAM

##### a) Roles vs Policies

| Category      | Type                        | Description  | Key Responsibilities/Features   |
|---------------|-----------------------------|--|---|
| User Roles    | Data Scientist/ ML Engineer | Provides access for experimentation                  | Access to S3, Athena , SageMaker Studio   |
|               | Data Engineer               | Provides access for data management                  | Access to S3, Athena, AWS Glue, EMR   |
|               | MLOps Engineer              | Provides access for ML operations                    | Access to SageMaker, CodePipeline, CodeBuild, CloudFormation, ECR, Lambda, Step Functions |
| Service Roles | SageMaker Execution         | Allows SageMaker to perform tasks on behalf of users | General SageMaker operations  |
|               | Processing Job              | Specific to SageMaker processing jobs                | Data processing tasks   |
|               | Training Job                | Specific to SageMaker training jobs                  | Model training tasks  |
|               | Model                       | Specific to SageMaker model deployment               | Model deployment and hosting  |
| IAM Policies  | Identity-based              | Attached to IAM users, groups, or roles              | Define actions allowed on specific resources  |
|               | Resource-based              | Attached to resources (e.g., S3 buckets)             | Control who can access specific resources   |



## IAM Policy – Examples for ML workflows

| ID | Purpose                                    | Key Permissions  | Resource Scope  | Notes   |
|----|--|--|---|---|
| 1  | Least privilege access for ML workflow     | <ul style="list-style-type: none"> <li>• SageMaker: <b>CreateTrainingJob</b>, <b>CreateModel</b></li> <li>• S3: GetObject, PutObject</li> <li>• ECR: BatchGetImage</li> <li>• CloudWatch: PutMetricData</li> </ul> <pre> {   "Effect": "Allow",   "Action": [     "sagemaker:CreateTrainingJob",     "sagemaker:DescribeTrainingJob",     "sagemaker:StopTrainingJob"   ],   "Resource": "*" } </pre>  | Specific ARNs for each service                        | Adheres to principle of least privilege             |
| 2  | Read metadata of ML resources              | <ul style="list-style-type: none"> <li>• <b>machinelearning:Get*</b></li> <li>• <b>machinelearning:Describe*</b></li> </ul> <pre> {   "Version": "2012-10-1/",   "Statement": [     {       "Effect": "Allow",       "Action": [         "machinelearning:Get*"       ],       "Resource": [         "arn:aws:machinelearning:us-east-1:555555555555:mlmodel/-ML-5555",         "arn:aws:machinelearning:us-east-1:666666666666:mlmodel/-ML-6666",         "arn:aws:machinelearning:us-east-1:777777777777:mlmodel/-ML-7777",         "arn:aws:machinelearning:us-east-1:888888888888:mlmodel/-ML-8888",         "arn:aws:machinelearning:us-east-1:555555555555:mlmodel/-ML-5555"       ]     },     {       "Effect": "Allow",       "Action": [         "machinelearning:Describe*"       ],       "Resource": [         "*"       ]     }   ] } </pre> | Specific MLModel ARNs for Get*<br>(all) for Describe* | Allows reading metadata but not modifying resources |
| 3  | Create ML resources                        | <ul style="list-style-type: none"> <li>• machinelearning:<b>CreateDataSourceFrom*</b></li> <li>• machinelearning:<b>CreateMLModel</b></li> <li>• machinelearning:<b>CreateBatchPrediction</b></li> <li>• machinelearning:<b>CreateEvaluation</b></li> </ul>  | * (all)   | Cannot be restricted to specific resources          |
| 4  | Manage real-time endpoints and predictions | <ul style="list-style-type: none"> <li>• machinelearning:<b>CreateRealtimeEndpoint</b></li> <li>• machinelearning:<b>DeleteRealtimeEndpoint</b></li> <li>• machinelearning:<b>Predict</b></li> </ul>   | Specific MLModel ARN                                  | Allows management of endpoints for a specific model |

## Detailed examples

### 1. identity-based policy used in a machine learning use case

|  |   |  |
|--|---|--|
| <pre>{   "Version": "2012-10-17",   "Statement": [     {       "Effect": "Allow",       "Action": [         "s3:GetObject",         "s3:PutObject"       ],       "Resource": [         "arn:aws:s3:::ml-data-bucket/*"       ]     },     {       "Effect": "Allow",       "Action": [         "sagemaker:CreateTrainingJob",         "sagemaker:DescribeTrainingJob",         "sagemaker:StopTrainingJob"       ],       "Resource": [         "*"       ]     }   ] }</pre> | <pre>{   "Effect": "Allow",   "Action": [     "cloudwatch:PutMetricData"   ],   "Resource": "*",   "Condition": {     "StringEquals": {       "cloudwatch:namespace": "SageMaker"     }   } }, {   "Effect": "Allow",   "Action": [     "logs:CreateLogGroup",     "logs:CreateLogStream",     "logs:PutLogEvents"   ],   "Resource": [     "arn:aws:logs:*:*:log-group:/aws/sagemaker/*"   ] }</pre> |  |
|--|---|--|

### 2. Allow users to read machine learning resources metadata

```
{
  "Version": "2012-10-17",
  "Statement": [
    {
      "Effect": "Allow",
      "Action": [
        "machinelearning:Get*"
      ],
      "Resource": [
        "arn:aws:machinelearning:us-east-1:555555555555:mlmodel/-ML-5555",
        "arn:aws:machinelearning:us-east-1:666666666666:mlmodel/-ML-6666",
        "arn:aws:machinelearning:us-east-1:777777777777:mlmodel/-ML-7777",
        "arn:aws:machinelearning:us-east-1:888888888888:mlmodel/-ML-8888",
        "arn:aws:machinelearning:us-east-1:555555555555:mlmodel/-ML-5555"
      ]
    },
    {
      "Effect": "Allow",
      "Action": [
        "machinelearning:Describe*"
      ],
      "Resource": [
        "*"
      ]
    }
  ]
}
```

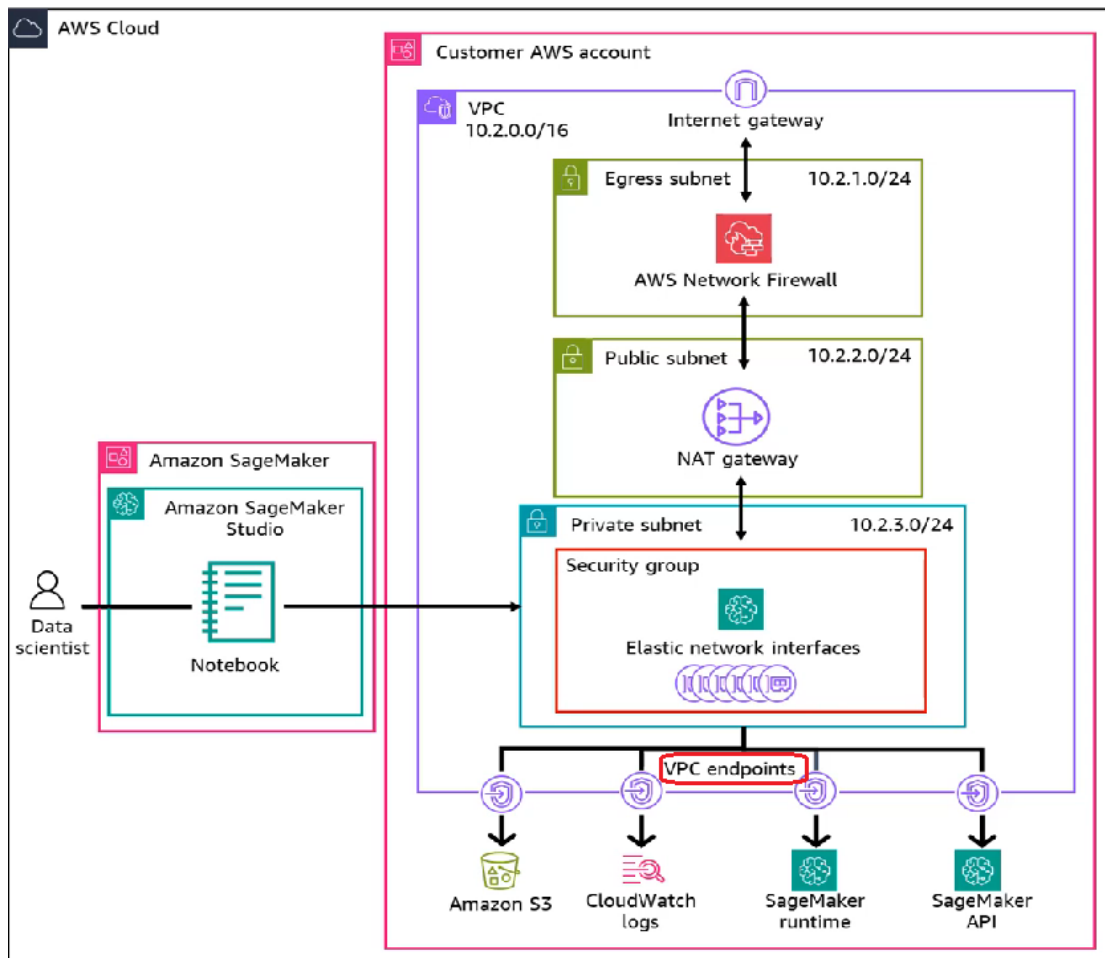
### 3. Allow users to create machine learning resources

```
{
  "Version": "2012-10-17",
  "Statement": [
    {
      "Effect": "Allow",
      "Action": [
        "machinelearning:CreateDataSourceFrom*",
        "machinelearning:CreateMLModel",
        "machinelearning:CreateBatchPrediction",
        "machinelearning:CreateEvaluation"
      ],
      "Resource": [
        "*"
      ]
    }
  ]
}
```

### 4. Allow users to create /delete real-time endpoints and perform real-time predictions on an ML model

```
{
  "Version": "2012-10-17",
  "Statement": [
    {
      "Effect": "Allow",
      "Action": [
        "machinelearning:CreateRealtimeEndpoint",
        "machinelearning>DeleteRealtimeEndpoint",
        "machinelearning:Predict"
      ],
      "Resource": [
        "arn:aws:machinelearning:us-east-1:555555555555:mlmodel/-ML-5555"
      ]
    }
  ]
}
```

## Detailed examples



To ensure access only from VPC, use **VPC Endpoints** for:

- S3
- CloudWatch Logs
- SageMaker runtime
- SageMaker API

### 4.3.3 SageMaker Compliance & Governance

#### AWS Services for Compliance and Governance

| Service                    | Purpose   | Key Features   | ML-Related Use Case   |
|----------------------------|---|--|---|
| <b>AWS Artifact</b>        | Provide on-demand access to AWS compliance reports and agreements | <ul style="list-style-type: none"><li>Self-service portal</li><li>Access to compliance documentation</li></ul>   | Access HIPAA compliance reports for healthcare ML projects                      |
| <b>AWS Config</b>          | Monitor & Evaluate AWS resource configurations                    | <ul style="list-style-type: none"><li>Continuous monitoring</li><li>Automated configuration evaluation</li></ul> | Monitor SageMaker resource configurations for compliance with security policies |
| <b>Audit Manager</b>       | Continuously audit AWS usage for risk and compliance assessment   | Streamlined auditing process, against regulations and standards  | Assess compliance of ML workflows with industry standards                       |
| <b>Security Hub</b>        | View of security alerts and posture                               | Centralized security alerts  | Monitor security posture across ML workflows and resources                      |
| <b>Amazon Inspector</b>    | Automated vulnerability management                                | Continuous scanning for vulnerabilities  | Scan container images in ECR for ML model deployments                           |
| <b>AWS Service Catalog</b> | Create and manage catalogs of pre-approved resources              | Governance-compliant resource catalogs   | Create catalogs of compliant SageMaker resources and ML models                  |

#### Amazon SageMaker Governance Tools Summary

| Tool                                       | Purpose                              | Key Features   |
|--|--------------------------------------|--|
| <b>SageMaker Role Manager</b>              | Simplify access control              | <ul style="list-style-type: none"><li>Define minimum permissions for ML activities</li><li>Quick setup &amp; Streamlined access management</li></ul>       |
| <b>SageMaker Model Cards</b>               | Document and share model information | <ul style="list-style-type: none"><li>Record intended uses</li><li>Document risk ratings</li></ul>   |
| <b>SageMaker Model Dashboard</b>           | Provide overview of models           | <ul style="list-style-type: none"><li>Unified view of all models in account</li><li>Monitor model behavior in production</li></ul>                         |
| <b>SageMaker Assets</b>                    | Streamline ML asset management       | <ul style="list-style-type: none"><li>Publish ML and data assets</li><li>Share assets across teams</li></ul>   |
| <b>Model Governance and Explainability</b> | Ensure compliance and transparency   | <ul style="list-style-type: none"><li>Protect data and workloads</li><li>Ensure compliance with standards</li><li>Enhance model interpretability</li></ul> |

## Compliance certifications and regulatory Frameworks

| Governance /Framework | Description  | AWS Services to Use   |
|-----------------------|--|---|
| <b>ISO 27001</b>      | Information Security Management System standard        | <ul style="list-style-type: none"><li>• AWS <a href="#">Config</a></li><li>• AWS <a href="#">Security Hub</a></li></ul>   |
| <b>SOC 2</b>          | Service Organization Control for service organizations | <ul style="list-style-type: none"><li>• AWS <a href="#">Artifact</a></li><li>• AWS <a href="#">Config</a></li><li>• SageMaker <a href="#">Model Cards</a></li></ul> |
| <b>PCI-DSS</b>        | Payment Card Industry Data Security Standard           | <ul style="list-style-type: none"><li>• AWS <a href="#">Config</a></li><li>• AWS <a href="#">WAF</a></li><li>• <a href="#">Amazon Inspector</a></li></ul>           |
| <b>HIPAA</b>          | Health Insurance Portability and Accountability Act    | <ul style="list-style-type: none"><li>• AWS <a href="#">Artifact</a></li><li>• AWS <a href="#">Security Hub</a></li><li>• AWS <a href="#">Config</a></li></ul>      |
| <b>FedRAMP</b>        | Federal Risk and Authorization Management Program      | <ul style="list-style-type: none"><li>• AWS CloudTrail</li><li>• AWS <a href="#">Config</a></li></ul>   |

**Note:** [AWS Config](#) common to all

### 4.3.3 Security Best Practices for CI/CD Pipelines

#### CI/CD pipeline stages

**Best practice:** When

,

| CI/CD Stage       | Security Tools/Practices   |
|-------------------|--|
| <b>Pre-Commit</b> | <ul style="list-style-type: none"><li>• pre-commit hooks (scripts)</li><li>• IDE plugins to<ul style="list-style-type: none"><li>○ analyze code, detect issues</li><li>○ provide recommendations for improvements.</li><li>○ handle linting, formatting, beautifying, and securing code.</li></ul></li></ul> |
| <b>Commit</b>     | <i>Static Application Security Testing (SAST),</i>   |
| <b>Build</b>      | <i>Software Composition Analysis (SCA)</i> <ul style="list-style-type: none"><li>○ identifies the open-source packages used in code</li><li>○ defining vulnerabilities and potential compliance-based issues</li><li>○ scan infrastructure as code (IaC) manifest files</li></ul>                            |
| <b>Test</b>       | <ul style="list-style-type: none"><li>• Dynamic Application Security Testing (DAST)</li><li>• Interactive Application Security Testing (IAST)<ul style="list-style-type: none"><li>○ Combine the advantages of SAST and DAST tools.</li></ul></li></ul>  |
| <b>Deploy</b>     | <ul style="list-style-type: none"><li>• Penetration testing</li></ul>  |
| <b>Monitor</b>    | <ul style="list-style-type: none"><li>• Red/Blue/Purple teaming</li></ul>  |

#### 4.3.4 Implement Security & Compliance w/ Monitoring, Logging and Auditing

##### CloudTrail for ML Resource Monitoring and Logging

| Use Case                     | Description  | Benefits  |
|------------------------------|--|---|
| <b>Compliance Auditing</b>   | Generate audit trails using CloudWatch Logs and CloudTrail | <ul style="list-style-type: none"><li>• Demonstrate compliance with regulations</li><li>• Meet internal policy requirements</li></ul> |
| <b>Resource Optimization</b> | Monitor resource utilization metrics                       | <ul style="list-style-type: none"><li>• Optimize ML workloads</li><li>• Prevent resource abuse and DoS attacks</li></ul>              |
| <b>Incident Response</b>     | Investigate and respond to security incidents              | <ul style="list-style-type: none"><li>• Identify unauthorized access attempts</li><li>• Detect and respond to data breaches</li></ul> |
| <b>Anomaly Detection</b>     | Implement ML models to detect unusual patterns             | <ul style="list-style-type: none"><li>• Identify potential security threats</li><li>• Detect deviations in monitoring data</li></ul>  |

##### SageMaker Security Troubleshooting and Debugging Summary

| Tool/Feature           | Purpose                       | Key Information Provided  | Use Case  |
|------------------------|-------------------------------|---|---|
| CloudTrail Logs        | Monitor API calls             | <ul style="list-style-type: none"><li>• Caller identity</li><li>• Timestamps</li><li>• API details</li></ul>                      | Identify unauthorized API calls to SageMaker resources            |
| <b>Data Event Logs</b> | Monitor data plane operations | Input/output data for training and inference  | <b>Verify if unauthorized entities accessed model data</b>        |
| IAM Policies           | Manage access control         | Permissions granted for SageMaker resources and operations  | Identify overly permissive policies, ensure least privilege       |
| <b>VPC Flow Logs</b>   | Monitor network traffic       | Network traffic to/from SageMaker resources   | Identify <b>suspicious IP addresses or communication patterns</b> |
| Encryption Settings    | Ensure data protection        | <ul style="list-style-type: none"><li>• Encryption status (at rest and in transit)</li><li>• AWS KMS key configurations</li></ul> | Verify proper data encryption and key management                  |
| AWS PrivateLink        | Enhance network security      | Private connections between VPC and SageMaker   | <b>Ensure traffic remains within AWS network</b>                  |

