# SageMaker MLOps

Paul Hargis
AWS Solutions Architect
AI/ML Specialist

Pranav Murthy
AWS Solutions Architect
AI/ML Specialist

# Agenda

Exploratory Data Analysis

Model Training and Model Tuning

Inference and Hosting

Operationalize Your Workflows

aws machine learning

# Flexible Exploratory Data Analysis

aws machine learning

# Exploratory Data Analysis

**Amazon SageMaker Data Wrangler**

A faster, visual way to aggregate and prepare data for machine learning

**SageMaker Studio Notebook EDA**

Locally Prepare and Analyze your data using SageMaker Studio Notebook

**SageMaker Processing Job**

To analyze data and evaluate machine learning models on Amazon SageMaker, use Amazon SageMaker Processing

aws machine learning

# Amazon SageMaker Studio to Quickly Analyze Data and Build Models

# Fast-start, shareable notebooks

- No wrangling of environments
- Start-up new notebooks in < 30 seconds
- One-click sharing of notebook and all its dependencies
- Easily switch instance types to match workload

- Exploratory data analysis

- Coding new algorithm

- Analyzing results

**M4** ⬅Switch➡ **P3**

Training and testing deep neural networks

Results➡ Enables notebook cost reduction of ~50%

aws machine learning

# Amazon SageMaker Processing

## Managed solution for data processing and model evaluation jobs

**Fully managed**

Fully managed clusters for distributed processing

**Bring your own custom processing script**

Bring your own script for feature engineering, data validation, model evaluation, and model interpretation

**Use SageMaker containers or BYO framework containers**

Use a SageMaker optimized framework container or bring your own

**Pay-as-you-go**

Infrastructure is created on-demand, configured, and terminated automatically

**Security and compliance**

All of the same security and compliance features

aws machine learning

# Flexible Model Training and Tuning At Scale

aws machine learning

# Amazon SageMaker – Built-in Algorithms

Rethinking algorithms design for large scale & streaming data

Data

Amazon SageMaker

Algorithms *as-a-service*

**K-Means** **XGBoost** **PCA**
**Linear Learner** **Seq2Seq**
**DeepAR**
**Factorization Machines**
**Object Detection**
**BlazingText**
**Image Classification**
**K-Nearest Neighbor**
**LDA** **Neural Topic Model**
**Random Cut Forest**
**Object2Vec** **IP Insight**
**Semantic Segmentation**

aws machine learning

# Amazon SageMaker – Bring-Your-Own-Algorithm

Build your own algorithms, SageMaker handles the rest

Amazon SageMaker Training Service Supported Containers

GLUON

Keras

PYTORCH
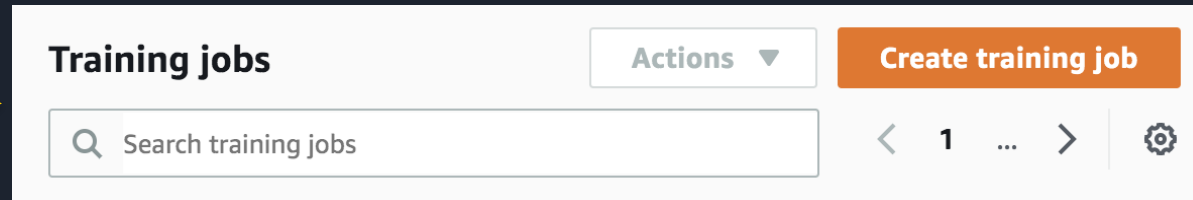
scikit learn

mxnet

TensorFlow

Chainer

Data

Custom algorithms

aws machine learning

# Amazon SageMaker – Bring-Your-Own-Container

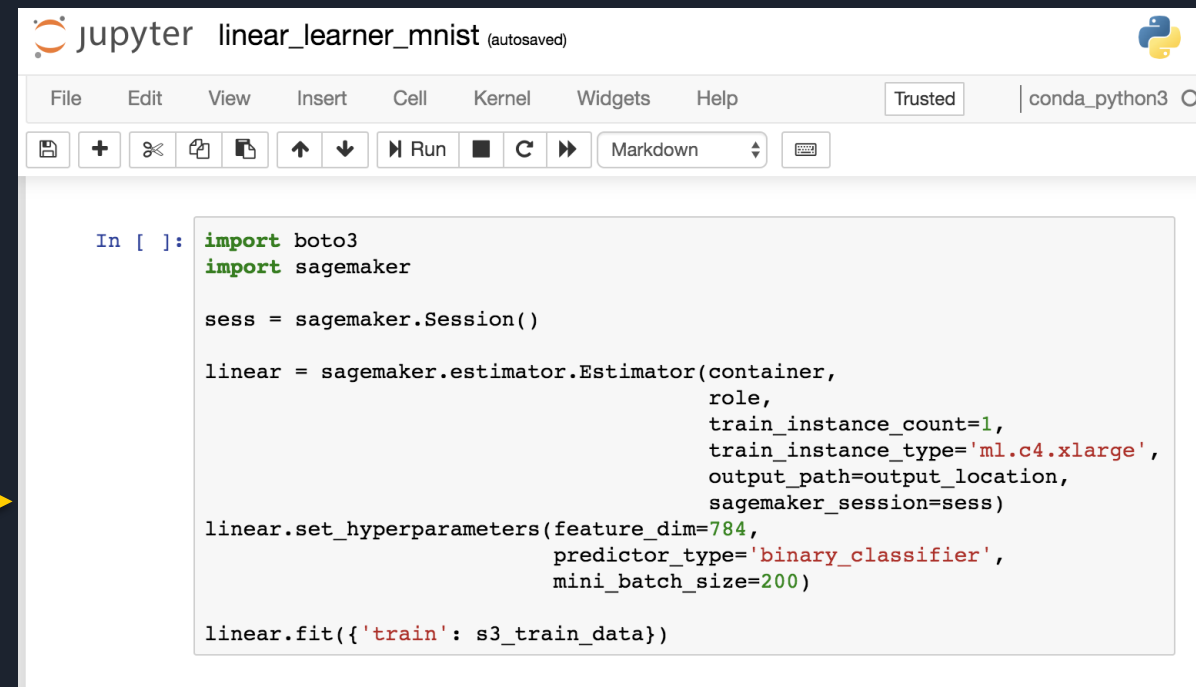Bring your custom code and container, train at scale in SageMaker

Data

Docker container with custom algorithms (e.g. R, Java)

Amazon SageMaker Training Service

Fully managed training service

aws machine learning

# Launching Training Jobs from Jupyter



Launch Training

# Model Registry

aws machine learning

# SageMaker Model Registry

# Speedy Deployment at Scale

aws machine learning
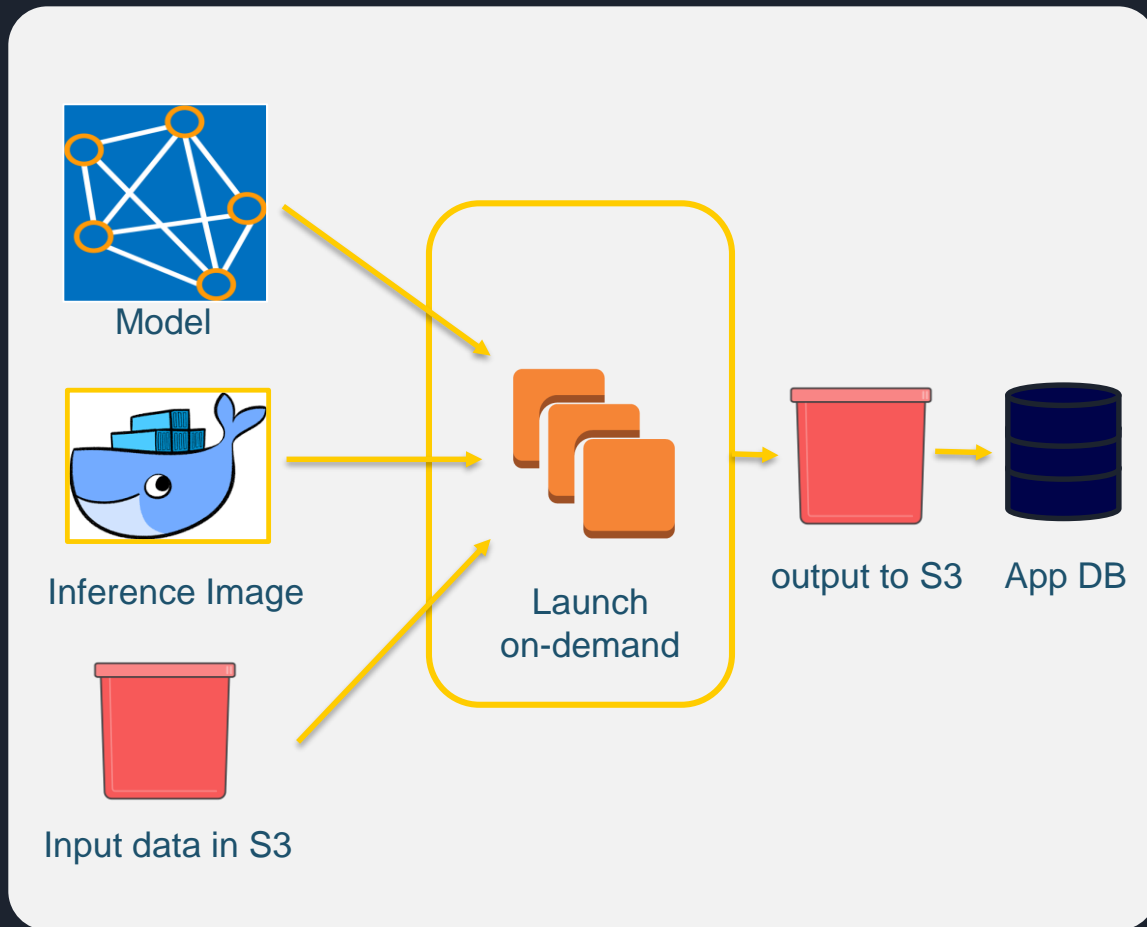
# Amazon SageMaker – Real time inference



- Quick deployment of Restful API endpoint

- Autoscaling across multiple AZs for variable loads

- A/B testing support

- VPC Privatelink support

- Endpoint invocation and resource usage metrics

- No model re-engineering

aws machine learning

# Amazon SageMaker – Batch inference



Model

Inference Image

Input data in S3

Launch
on-demand

output to S3   App DB

- Use trained model to get inference on data in S3

- Launch fully managed transient inference resources

- Seamless integration with S3 for both input and output

aws machine learning

# Operationalize Your Workflows with SageMaker Pipelines and Projects

aws machine learning

# Challenges with creating a complete workflow for the ML lifecycle

1. Building, training, and deploying models is an iterative process

2. Getting models from concept to production involves multiple steps
   - Create standard code packages for each step of the ML lifecycle
   - Stitch them together into a structure called a workflow
   - Manage dependencies between steps
   - Execute the workflow as an orchestrated sequence

3. Track artifacts for each step of the workflow

4. Deploy and manage the right version of the model, across thousands of models

5. Automate and scale the complete workflow as part of ML Ops

aws machine learning

# Amazon SageMaker Pipelines

Build fully automated machine learning workflows at scale

## Accelerate machine learning development

With just a few clicks, you can create a fully automated ML workflow, reducing months of coding to a few hours

## Automatically track hundreds of model artifacts

Eliminate manual processes by automatically tracking model artifacts to keep a structured audit trail

## Scale to thousands of ML models in production

Choose from built-in templates to setup CI/CD pipelines to automate workflows and deploy ML models at scale

|    20

aws machine learning

## Key Features

**Amazon SageMaker Pipelines**

**Compose and manage ML workflows**
Create detailed workflows with an easy-to-use Python SDK and manage them visually

**Track model lineage for governance and audits**
Track code, datasets, and versions for each step of the ML lifecycle

**Replay and re-run workflows**
Keep models up to date by re-running every step based on custom schedules

**Visually compare, select, and deploy models**
Deploy and manage modules through the visual interface of SageMaker Studio

**Access a central registry of trained models**
Use the model registry to choose your best model for deployment to production

**Fully managed ML Ops with built-in CI/CD support**
Use CI/CD practices to build a fully automated machine learning workflow

|

aws machine learning

# Amazon SageMaker Pipelines - How it works



Automate ML Ops with CI/CD and end-to-end lineage tracking

**Model drift**

**Amazon SageMaker Pipelines**
Build fully automated machine learning workflows at scale

**Input data**

- Prepare or transform
- Explain
- Train
- Validate

**Model registry**
Catalog model versions, metrics, and approvals for model deployment

**Real-time inference**

**Batch scoring**

aws machine learning

# Compose & manage workflows

Create your ML workflows using the Python SDK

Visualize the workflows with SageMaker Studio

Manage every step including data transformations, training, debugging, and optimizing models

```python
class MyPipeline(dsl.Pipeline):
    # Pipeline parameter definitions
    training_instance = PipelineParam(string)
    num_trees = PipelineParam(int, default=100)
    tree_depth = PipelineParam(int, default=3)

# Example data query step
get_input_data_step = AthenaStep(query="SELECT a,b,c FROM table_xyz")

# Example training step
training_step = XGBoostStep(
    entry_point='scripts/xgb/train.py',
    source_dir='scripts/xgb',
    instance_type=training_instance,
    env_vars=[num_trees, tree_depth],
    inputs={'train': processing_step.outputs('clean_data')})
```

# Visualize status of model pipelines

Follow completed steps and monitor steps in progress

Understand the output from each step with the output logs

Monitor, change, and manage the parameters for each step

# Create workflows with built-in Steps

- Processing
- Training
- Condition
- BatchTransform
- RegisterModel
- CreateModel

# Choose models from the model registry

Model registry is a central repository of trained models

Register modes intended for production in the model registry

Access the model registry through SageMaker Studio or the Python SDK

# Execute CI/CD pipelines to update models

Use CI/CD practices to automate deployment

Operationalize machine learning at scale

Keep your models up to date and accurate

aws machine learning

# Execute CI/CD pipelines to update models
## Projects

**Project Templates**



**Built-In Templates**

**Custom Templates**

# Execute CI/CD pipelines to update models
## Projects

# Integrate with source control

# Thank you

aws machine learning