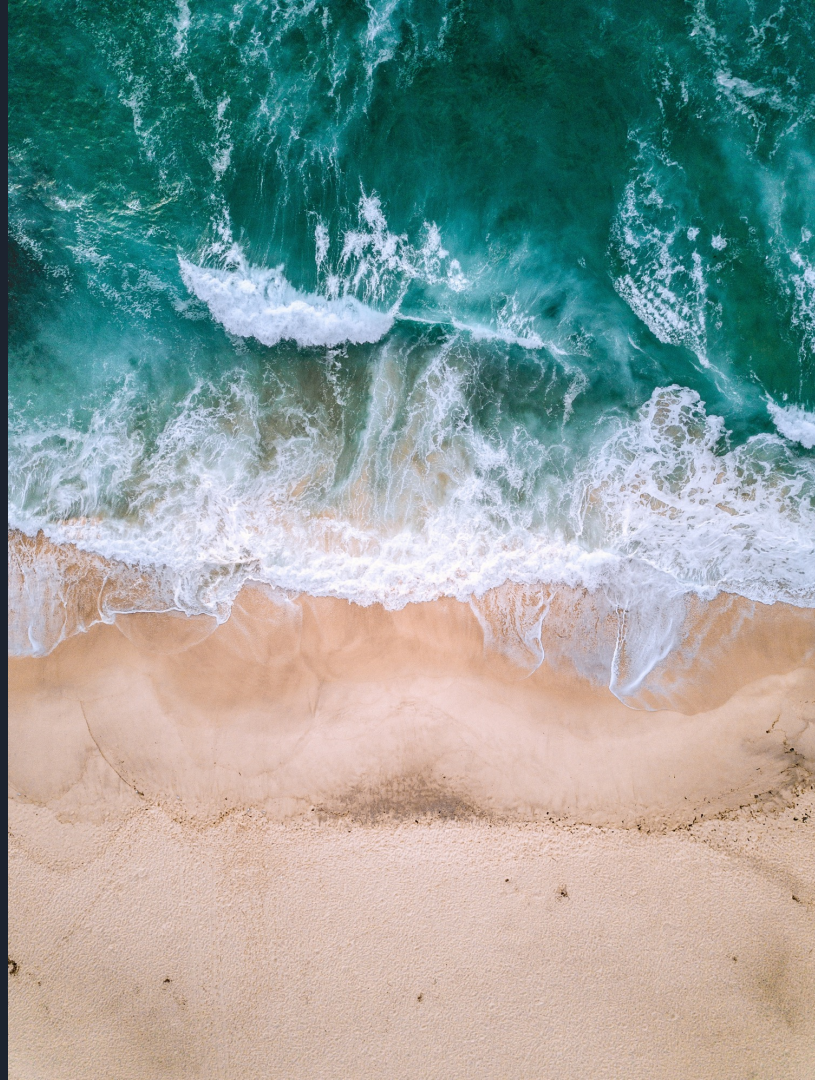


# Machine Learning on AWS with SageMaker

**Joe Wilson**  
Solutions Architecture, AWS  
Worldwide Public Sector

# Agenda

- Introduction + Recap
- AI Services
- SageMaker Deep Dive
  - Build / train / tune
- SageMaker Deep Dive
  - Deploy
  - Prepare
- Q/A



# The AWS ML Stack

Broadest and most complete set of machine learning capabilities

## AI SERVICES



NEW

Amazon HealthLake

### HEALTH AI



Amazon Transcribe for Medical



Amazon Comprehend for Medical



NEW

AWS Panorama + Appliance



NEW

Amazon Monitron



NEW

Amazon Lookout for Equipment



NEW

Amazon Lookout for Vision

### ANOMALY DETECTION



NEW

Amazon Lookout for Metrics



NEW

Amazon DevOps Guru



Amazon CodeGuru

### CODE AND DEVOPS

### VISION



Amazon Rekognition

### SPEECH



Amazon Polly



Amazon Transcribe  
+Medical

### TEXT



Amazon Comprehend  
+Medical



Amazon Translate



Amazon Textract

### SEARCH



Amazon Kendra

### CHATBOTS



Amazon Lex

### PERSONALIZATION



Amazon Personalize

### FORECASTING



Amazon Forecast



Amazon Fraud Detector

### CONTACT CENTERS



Contact Lens

Voice ID

For Amazon Connect

## ML SERVICES



Amazon SageMaker

Label data

NEW

Aggregate & prepare data

NEW

Store & share features

Auto ML

Spark/R

NEW

Detect bias

Visualize in notebooks

Pick algorithm

Train models

Tune parameters

NEW

Debug & profile

Deploy in production

Manage & monitor

NEW

CI/CD

Human review

### SAGEMAKER STUDIO IDE

NEW: SageMaker JumpStart

NEW: Model management for edge devices

## FRAMEWORKS & INFRASTRUCTURE



DeepGraphLibrary

Deep Learning AMIs & Containers

GPUs & CPUs

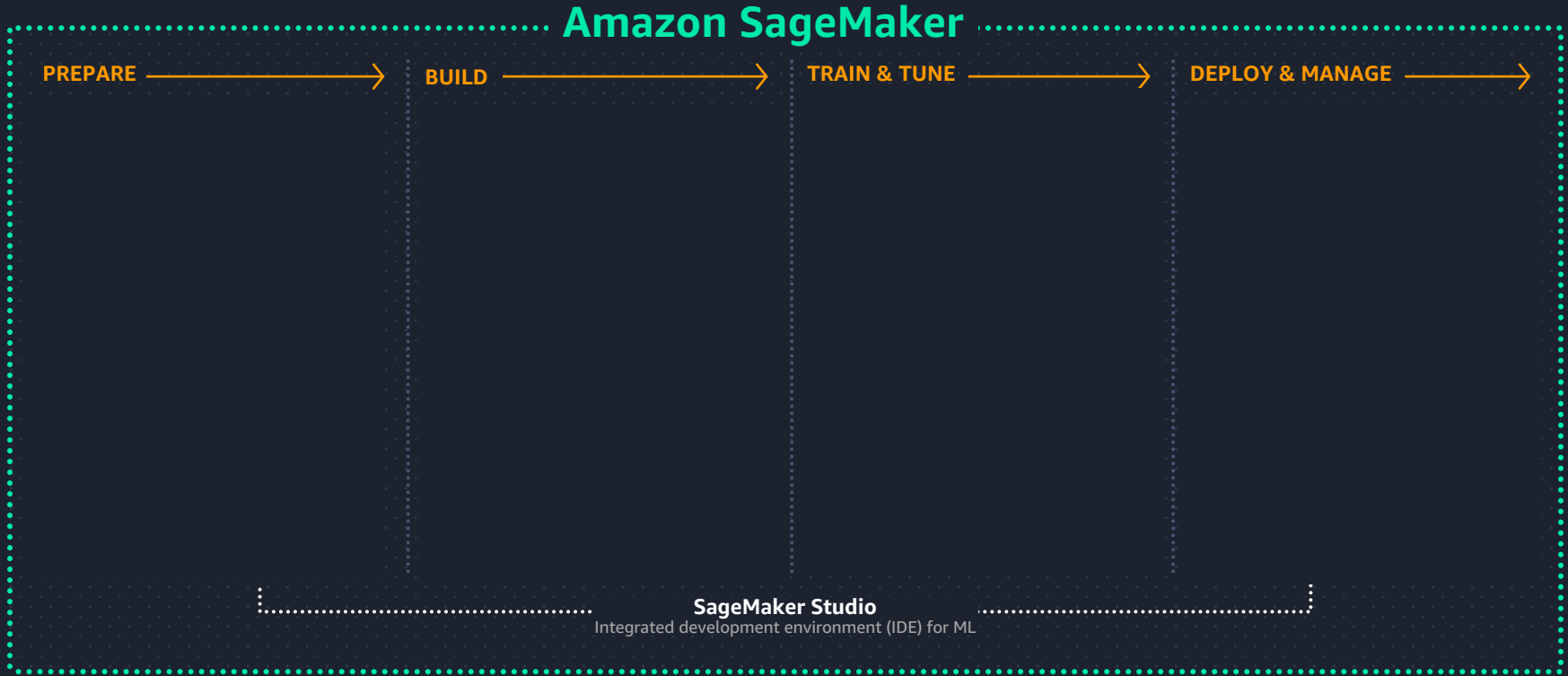
Elastic Inference

Habana Gaudi

Trainium Inferentia

FPGA

# Amazon SageMaker overview



# Developing on Amazon SageMaker

The “build” part of build, train, deploy

# Amazon SageMaker Studio

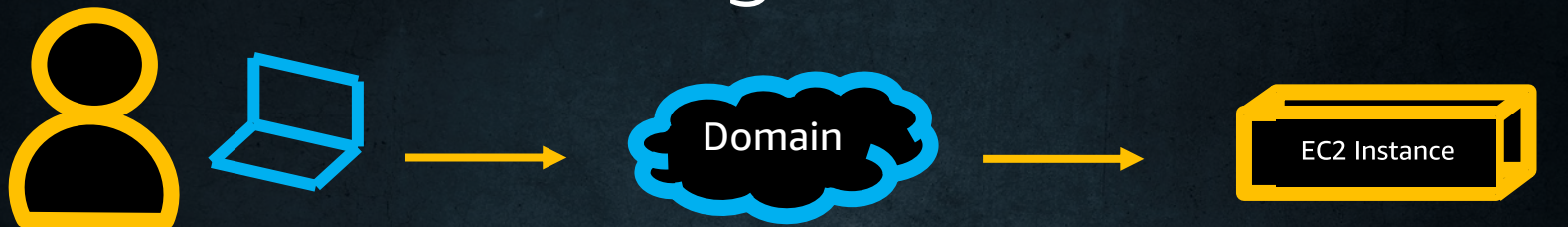
The screenshot displays the Amazon SageMaker Studio interface, which is a cloud-based environment for machine learning. The interface is divided into several panels:

- Code Editor:** Shows a Python script in a file named `autopilot_customer_churn.ipynb`. The code defines a function to create an AutoML job. The output shows the job name `AutoMLJobName: automl-churn-16-18-04-54` and the ARN `AutoMLJobArn: 'arn:aws:sagemaker:us-east-1:865926247462:automl-16-18-04-54'`.
- AbalonePipeline:** A visualization of the pipeline execution. It shows a graph with a step named `AbaloneMSECond`. A `Start an execution` button is visible.
- Data flow:** A diagram showing the data flow process. It includes steps for `Source - sampled`, `Data types`, and `Custom Pandas`. The data is transformed from `S3: emissions_encoded.csv` and `S3: train_with_header.csv` into `Transform: emissions_encoded.csv` and `Transform: train_with_header.csv`.
- Experiment Summary:** A panel titled `EXPERIMENT: AUTOML-CHURN-16-18-04-54` showing a progress bar with stages: `Pre-processing`, `Candidate Definitions Generated`, `Feature Engineering`, and `Model Tuning`. The text indicates that Amazon SageMaker Autopilot is analyzing the input data and that the process can take hours to complete at the default setting of 250 trials.

The bottom status bar shows the current state: `0` files, `4` tabs, and `Git: idle`. The bottom right corner displays the job ID: `automl-churn-16-18-04-54`.



# Amazon SageMaker Studio



You!

Your IDE

Your Compute

- Can enable IAM or [single sign on](#)
- Fast notebook start-up
- Consolidated experience

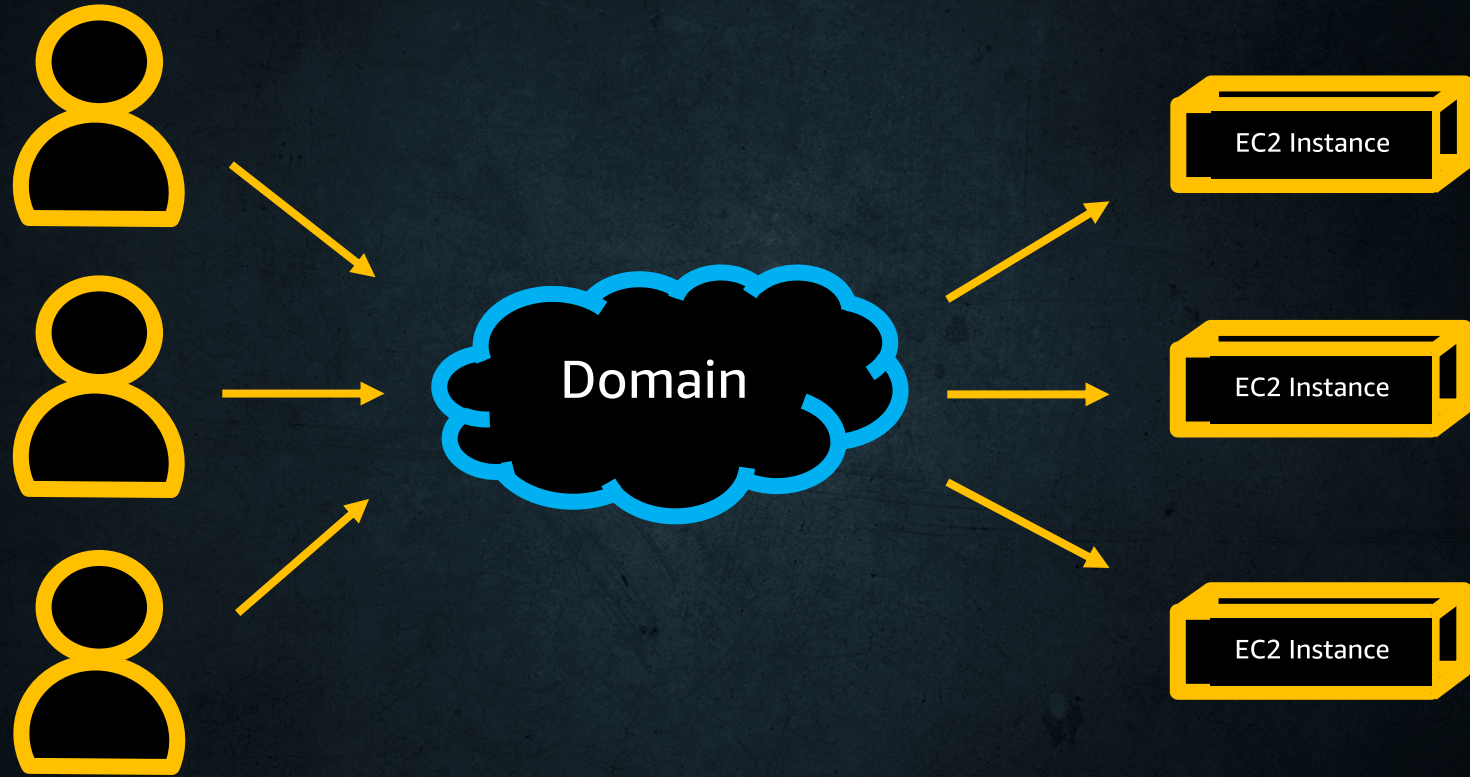
- Low-lift / automatic feature access
- Autopilot
  - Experiments
  - Debugger
  - Model Monitoring
  - Versioning
  - Endpoint management

- Can create a shareable link with package dependences in single click
- Billing is per usage, not up time.
- Change environments per project

Remember! Features are backwards compatible, but differ in ease-of-use across instances vs Studio.



# 1-Click Notebook Sharing & Instance Upgrade



Studio defaults to 1:1 user to instance mapping, and you can share content on these via snapshots.

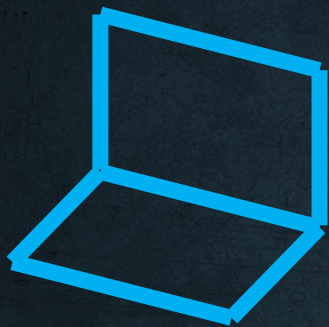




# Amazon SageMaker brings elastic, dedicated compute



Person



Project



Dataset



Step in your  
ML lifecycle

Allowing you to scale and manage each independently

# Training on Amazon SageMaker

The “train” part of build, train, deploy

# Many Ways to Train Models on Amazon SageMaker

1

Built-in  
Algorithms

2

Script Mode

3

Docker container

4

Locally\*

5

In a pipeline



# Train with your own XGBoost script

```
# Open Source distributed script mode
from sagemaker.session import Session
from sagemaker.inputs import TrainingInput
from sagemaker.xgboost.estimator import XGBoost

est = XGBoost(entry_point='abalone.py',
              framework_version='1.2-1', # Note: framework_version is mandatory
              hyperparameters=hyperparams,
              role=role,
              instance_count=2,
              instance_type=instance_type,
              output_path=output_path)

train_input = TrainingInput("s3://{}/{}/{}".format(bucket, prefix, "train"), content_type=content_type)
validation_input = TrainingInput("s3://{}/{}/{}".format(bucket, prefix, "validation"), content_type=content_type)

est.fit({'train': train_input, 'validation': validation_input})
```

2021-01-28 19:58:46 Starting - Starting the training job..

2021-01-28 19:59:10 Starting - Launching requested ML instancesProfilerReport-1611863926: InProgress

.....



# Train with your own deep learning model

```
from sagemaker.pytorch import PyTorch

estimator = PyTorch(entry_point='mnist.py',
                    role=role,
                    framework_version='1.8.0',
                    instance_count=2,
                    instance_type='ml.c4.xlarge',
                    hyperparameters={
                        'epochs': 6,
                        'backend': 'gloo'
                    })

inputs = sagemaker_session.upload_data(path='data', bucket=bucket, key_prefix=prefix)

estimator.fit({'training': inputs})
```

The logo for MXNet, featuring the letters 'mxnet' in a blue, lowercase, sans-serif font.The logo for PyTorch, consisting of a red circular icon with a white flame-like shape inside, followed by the text 'PyTorch' in a black, sans-serif font.The logo for scikit-learn, featuring a blue circle on the left and the text 'scikit learn' in a black, lowercase, sans-serif font on the right.The logo for Apache Spark ML, featuring the text 'APACHE Spark ML' in a black, sans-serif font, with a yellow star icon above the word 'Spark'.The logo for TensorFlow, featuring a yellow 3D 'T' icon above the text 'TensorFlow' in a black, sans-serif font.The logo for Chainer, featuring a red circular icon with a white neural network structure inside, followed by the text 'Chainer' in a red, sans-serif font.The logo for RAY, featuring a blue circular icon with a white neural network structure inside, followed by the text 'RAY' in a blue, sans-serif font.

# Train using a Docker File

```
1 FROM ubuntu:20.04
2
3 ARG DEBIAN_FRONTEND=noninteractive
4
5 # Don't prompt for tzdata on new versions of Ubuntu:
6 ARG DEBIAN_FRONTEND=noninteractive
7
8 RUN apt-get -y update && apt-get install -y --no-install-recommends \
9     wget \
10    libcurl4-openssl-dev\
11    libsodium-dev \
12    r-base \
13    r-base-dev \
14    ca-certificates
15
16 RUN R -e "install.packages(c('mda', 'plumber'), repos='https://cloud.r-project.org')"
17
18 COPY mars.R /opt/ml/mars.R
19 COPY plumber.R /opt/ml/plumber.R
20
21 ENTRYPOINT ["/usr/bin/Rscript", "/opt/ml/mars.R", "--no-save"]
22
```

```
"AlgorithmSpecification": {
    "TrainingImage": '{}.dkr.ecr.{}.amazonaws.com/sagemaker-rmars:latest'.format(account,
region)
},
```

# Store, Search, Version, Track Model Artifacts by Default

Amazon SageMaker > Search

▼ Search

Property:  Operator:  Value:

Results: Training jobs

< 1 ... >

Name	Algorithm	Data sources
<a href="#">tuning-job-1-086396f1dda841258d-019-6c6367d0</a>	683313688378.dkr.ecr.us-east-1.amazonaws.com/sagemaker-xgboost:1.0-1-cpu-py3	train : <a href="#">s3://sagemaker-us-east-1-865926247462/automl-churn-sdk-18-19-41-09/transformed-data/dpp0/rpb/train</a> validation : <a href="#">s3://sagemaker-us-east-1-865926247462/automl-churn-sdk-18-19-41-09/transformed-data/dpp0/rpb/validation</a>
<a href="#">tuning-job-1-086396f1dda841258d-018-fbb3e0c3</a>	683313688378.dkr.ecr.us-east-1.amazonaws.com/sagemaker-xgboost:1.0-1-cpu-py3	train : <a href="#">s3://sagemaker-us-east-1-865926247462/automl-churn-sdk-18-19-41-09/transformed-data/dpp8/rpb/train</a> validation : <a href="#">s3://sagemaker-us-east-1-865926247462/automl-churn-sdk-18-19-41-09/transformed-data/dpp8/rpb/validation</a>
<a href="#">tuning-job-1-086396f1dda841258d-017-5c14c90a</a>	683313688378.dkr.ecr.us-east-1.amazonaws.com/sagemaker-xgboost:1.0-1-cpu-py3	train : <a href="#">s3://sagemaker-us-east-1-865926247462/automl-churn-sdk-18-19-41-09/transformed-data/dpp6/rpb/train</a> validation : <a href="#">s3://sagemaker-us-east-1-865926247462/automl-churn-sdk-18-19-41-09/transformed-data/dpp6/rpb/validation</a>

### Log events

You can use the filter bar below to search for and match terms, phrases, or values in your log events. [Learn more about filter patterns](#)

View as text  Actions

Clear 1m 30m 1h 12h Custom

Timestamp	Message
No older events at this moment. <a href="#">Retry</a>	
▶ 2021-02-18T14:49:15.656-06:00	Docker entrypoint called with argument(s): train
▶ 2021-02-18T14:49:15.656-06:00	Running default environment configuration script
▼ 2021-02-18T14:49:16.656-06:00	[02/18/2021 20:49:16 INFO 139632547493632] Input channel configuration: {'validation'...
	[02/18/2021 20:49:16 INFO 139632547493632] Input channel configuration: {'contentType': 'text/csv', 'TrainingInputMode': 'Pipe', 'S3DistributionType': 'FullyReplicated', 'RecordWrapperType': 'None'}, 'train': {'contentType': 'text/csv', 'TrainingInputMode': 'Pipe', 'S3DistributionType': 'FullyReplicated', 'RecordWrapperType': 'None'}}. <input type="button" value="Copy"/>
▶ 2021-02-18T14:49:18.657-06:00	[02/18/2021 20:49:18 INFO 139632547493632] nvidia-smi took: 0.025197267532348633 secs...
▶ 2021-02-18T14:49:18.657-06:00	[02/18/2021 20:49:18 INFO 139632547493632] Number of GPUs being used: 0
▶ 2021-02-18T14:49:18.657-06:00	[02/18/2021 20:49:18 INFO 139632547493632] Creating network in context cpu(0).
▶ 2021-02-18T14:49:18.657-06:00	[02/18/2021 20:49:18 INFO 139632547493632] Starting model trainer.
▼ 2021-02-18T14:49:20.658-06:00	[02/18/2021 20:49:19 INFO 139632547493632] #progress_notice: epoch=0, iterations=1, s...
	[02/18/2021 20:49:19 INFO 139632547493632] #progress_notice: epoch=0, iterations=1, speed=224.68171381600214 samples/sec <input type="button" value="Copy"/>
▼ 2021-02-18T14:49:20.658-06:00	[02/18/2021 20:49:19 INFO 139632547493632] #progress_notice: epoch=0, iterations=2, s...
	[02/18/2021 20:49:19 INFO 139632547493632] #progress_notice: epoch=0, iterations=2, speed=435.57591410342366 samples/sec <input type="button" value="Copy"/>
▼ 2021-02-18T14:49:20.658-06:00	[02/18/2021 20:49:19 INFO 139632547493632] #progress_notice: epoch=0, iterations=3, s...
	[02/18/2021 20:49:19 INFO 139632547493632] #progress_notice: epoch=0, iterations=3, speed=629.347907622412 samples/sec <input type="button" value="Copy"/>
▼ 2021-02-18T14:49:20.658-06:00	[02/18/2021 20:49:19 INFO 139632547493632] #progress_notice: epoch=0, iterations=4, s...
	[02/18/2021 20:49:19 INFO 139632547493632] #progress_notice: epoch=0, iterations=4, speed=816.3276545798652 samples/sec <input type="button" value="Copy"/>

# Reduce training costs with fully-managed spot instances

```
from sagemaker.estimator import Estimator

est = Estimator(container,
                role,
                hyperparameters=hyperparameters,
                instance_count=1,
                instance_type=instance_type,
                volume_size=5, # 5 GB
                output_path=output_path,
                sagemaker_session=sagemaker.Session(),
                use_spot_instances=True,
                max_run=3600,
                max_wait=7200,
                checkpoint_s3_uri=checkpoint_s3_uri)

est.fit({'train': train_input}, job_name=job_name)
```

- Specify a max wait time
- SageMaker will default to giving you the lowest possible cost
- Store model checkpoints in S3 in case your job is interrupted for BYOM
- Many built-in algorithms automatically revert to a training job
- We have examples
- Save up to 90%!

**Amazon SageMaker price reductions: Up to 18% lower prices on ml.p3 and ml.p2 instances**

by Urvashi Chowdhary | on 07 OCT 2020 | in [Amazon SageMaker](#), [Artificial Intelligence](#) | [Permalink](#) | [Comments](#) | [Share](#)



# Deploying on Amazon SageMaker

The “deploy” part of build, train, deploy

# Deploy any open source model with Amazon SageMaker



# Enabling autoscaling on your endpoint is easy

1

```
client = boto3.client('sagemaker')

response = client.describe_endpoint(
    EndpointName=endpoint_name)

variant_name = response['ProductionVariants'][0]['VariantName']
resource_id = 'endpoint/{}/variant/{}'.format(endpoint_name, variant_name)
```

Use boto3's SageMaker client to get the *resource ID*

2

```
scaling_client = boto3.client('application-autoscaling')

response = scaling_client.register_scalable_target(
    ServiceNamespace='sagemaker',
    ResourceId=resource_id,
    ScalableDimension='sagemaker:variant:DesiredInstanceCount',
    MinCapacity=min_capacity,
    MaxCapacity=max_capacity,
    RoleARN=role)
```

Use boto3's application autoscaling client to register a scalable target

# Bring your own pretrained models to host on SageMaker using script mode

```
from sagemaker.tensorflow import TensorFlowModel

model = TensorFlowModel(model_data='s3://mybucket/model.tar.gz', role='MySageMakerRole')

predictor = model.deploy(initial_instance_count=1, instance_type='ml.c5.xlarge')
```

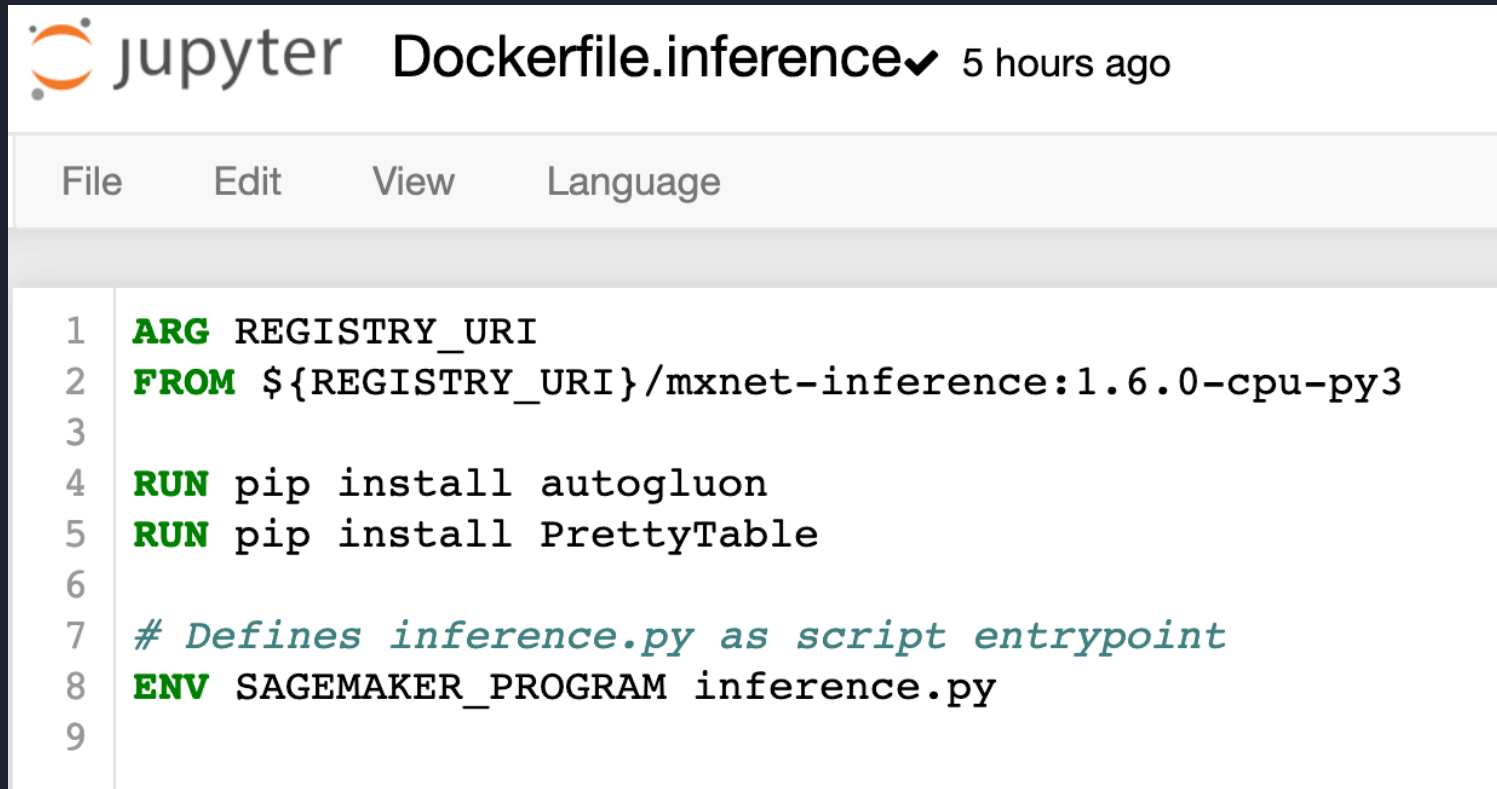
You can use the built-in inference script with the deep learning container, *or you can bring your own inference script.*

```
pytorch_model = PyTorchModel(model_data='s3://my-bucket/my-path/model.tar.gz', role=role,
                              entry_point='inference.py')

predictor = pytorch_model.deploy(instance_type='ml.c4.xlarge', initial_instance_count=1)
```



# Deploy *any open source model* on SageMaker with Docker



The screenshot shows a Jupyter interface for editing a Dockerfile. The title bar reads "jupyter Dockerfile.inference ✓ 5 hours ago". Below the title bar is a menu bar with "File", "Edit", "View", and "Language". The main area contains a Dockerfile with the following content:

```
1 ARG REGISTRY_URI
2 FROM ${REGISTRY_URI}/mxnet-inference:1.6.0-cpu-py3
3
4 RUN pip install autogluon
5 RUN pip install PrettyTable
6
7 # Defines inference.py as script entrypoint
8 ENV SAGEMAKER_PROGRAM inference.py
9
```

Remember, you need to build this image first, then push to ECR

# More features in SageMaker Studio

# Amazon SageMaker Autopilot

---

Automatic model creation with full visibility & control



Quick to start

Provide your data in a tabular form & specify target prediction



Automatic model creation

Get ML models with feature engineering & model tuning automatically done



Visibility & control

Get notebooks for your models with source code



Recommendations & Optimization

Get a leaderboard & continue to improve your model

# Use Amazon SageMaker Autopilot to create and review top performing regression and classification models

The screenshot displays the Amazon SageMaker Autopilot interface. The left sidebar shows a file explorer with a folder named 'automl-preview' containing subfolders 'bank-additional', 'model', and 'sagemaker\_auto...'. The main content area is titled 'EXPERIMENT: MY-SAGEMAKER-AUTOPILOT' and shows a table of trials. The table has columns for 'Trial name', 'Status', 'Start time', 'End time', and 'Objective'. All trials listed are in a 'Completed' status. There are buttons for 'Open candidate generation notebook' and 'Open data exploration notebook' at the top right, and a 'Deploy model' button at the bottom right of the table.

Trial name	Status	Start time	End time	Objective
my-sagemaker-tuning-job-...	Completed	9 hours ago		0.9206119775772095
my-sagemaker-tuning-job-...	Completed	9 hours ago		0.9202479720115662
my-sagemaker-tuning-job-...	Completed	7 hours ago		0.9200050234794617
my-sagemaker-tuning-job-...	Completed	7 hours ago		0.9195190072059631
my-sagemaker-tuning-job-...	Completed	9 hours ago		0.9191550016403198
my-sagemaker-tuning-job-...	Completed	7 hours ago		0.9190340042114258
my-sagemaker-tuning-job-...	Completed	8 hours ago		0.9189119935035706
my-sagemaker-tuning-job-...	Completed	8 hours ago		0.9186699986457825
my-sagemaker-tuning-job-...	Completed	8 hours ago		0.9186699986457825



# Amazon SageMaker Experiments

Organize, track, and compare training experiments



Tracking at scale

Track parameters & metrics across experiments & users



Custom organization

Organize experiments by teams, goals, & hypotheses



Visualization

Easily visualize experiments and compare



Metrics and logging

Log custom metrics using the Python SDK & APIs



Fast Iteration

Quickly go back & forth & maintain high-quality

# Use Amazon SageMaker Experiments to track and manage thousands of experiments

Amazon SageMaker Studio File Edit View Run Kernel Git Tabs Settings Help

Trial Component List Trial Component Chart

**TRIAL COMPONENTS** 5 rows selected. Select rows to toggle chart visibility. Add Chart

	Experiment	Trial	Trial Component	Type	Training time	Actions
<input type="checkbox"/>	customer-churn-predi...	Trial-9	Training-Run-9-aws-training-job	Training job	~6 minutes	Remove
<input type="checkbox"/>	customer-churn-predi...	Trial-8	Training-Run-8-aws-training-job	Training job	~5 minutes	Remove
<input type="checkbox"/>	customer-churn-predi...	Trial-7	Training-Run-7-aws-training-job	Training job	~6 minutes	Remove
<input type="checkbox"/>	customer-churn-predi...	Trial-6	Training-Run-6-aws-training-job	Training job	~4 minutes	Remove
<input type="checkbox"/>	customer-churn-predi...	Trial-5	Training-Run-5-aws-training-job	Training job	~4 minutes	Remove

**1 CHART**

train:loss\_last with 1-minute aggregation

train:loss\_last

period

**CHART PROPERTIES**

Data type

- Time series
- Summary statistics

Chart type

- Histogram
- Line

X-axis dimension

- Epoch
- Time
- Periods from start

X-axis aggregation

- 1-minute
- 5-minute
- 60-minute

Y-axis

train:loss\_last

Trial Component Chart

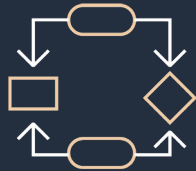
# Amazon SageMaker Debugger

Analysis and debugging, explainability, and alert generation



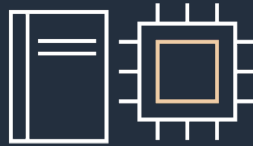
Relevant data capture

Data is automatically captured for analysis



Data analysis & debugging

Analyze & debug data with no code changes



Automatic error detection

Errors are automatically detected based on rules



Improved productivity with alerts

Take corrective action based on alerts



Visual analysis and debugging

Visually analyze & debug from SageMaker Studio

# Amazon SageMaker is fully managed

---

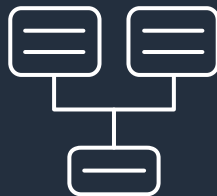
## One click model deployment



Auto-scaling



Low latency and  
high throughput



Bring your  
own model



Python SDK



Deploy multiple  
models on an  
endpoint



# Amazon SageMaker Model Monitor

Continuous monitoring of models in production



Automatic data collection

Data is automatically collected from your endpoints



Continuous Monitoring

Define a monitoring schedule and detect changes in quality against a pre-defined baseline



Flexibility with rules

Use built-in rules to detect data drift or write your own rules for custom analysis



Visual data analysis

See monitoring results, data statistics, and violation reports in SageMaker Studio



CloudWatch Integration

Automate corrective actions based on Amazon CloudWatch alerts

# Use Amazon SageMaker Model Monitor to identify model drift and take action

Amazon SageMaker Studio File Edit View Run Kernel Git Tabs Settings Help

Amazon S x UC-DEMO x UC-DEMO x model-mo x SageMaker x model-mo x DEMO-xgt x model-mo x

Monitoring results Monitoring Job history AWS settings

### AMAZON SAGEMAKER MODEL MONITOR

Amazon SageMaker Model Monitor detects data drift and other issues that can affect models in production and alerts you so you can take corrective action.

2 CHARTS [Add chart](#)

#### Int'l Plan\_yes: Sum

Date and time	Baseline	Current
01 PM	250	0
03 PM	250	550
05 PM	250	550
07 PM	250	550
09 PM	250	550
11 PM	250	550
01 AM	250	550
03 AM	250	550
05 AM	250	550
07 AM	250	550
09 AM	250	550
11 AM	250	550
01 PM	250	550

source: Baseline (blue line), Current (orange line)

### CHART PROPERTIES

Timeline

- 1 week
- 1 day
- 12 hours
- 6 hours
- 1 hour

Statistic

- Average
- SampleCount
- Sum
- Minimum
- Maximum

Feature

feature\_data\_Int'l Plan\_yes ▼

© 2021, Amazon Web Services, Inc. or its Affiliates.

# Q&A!