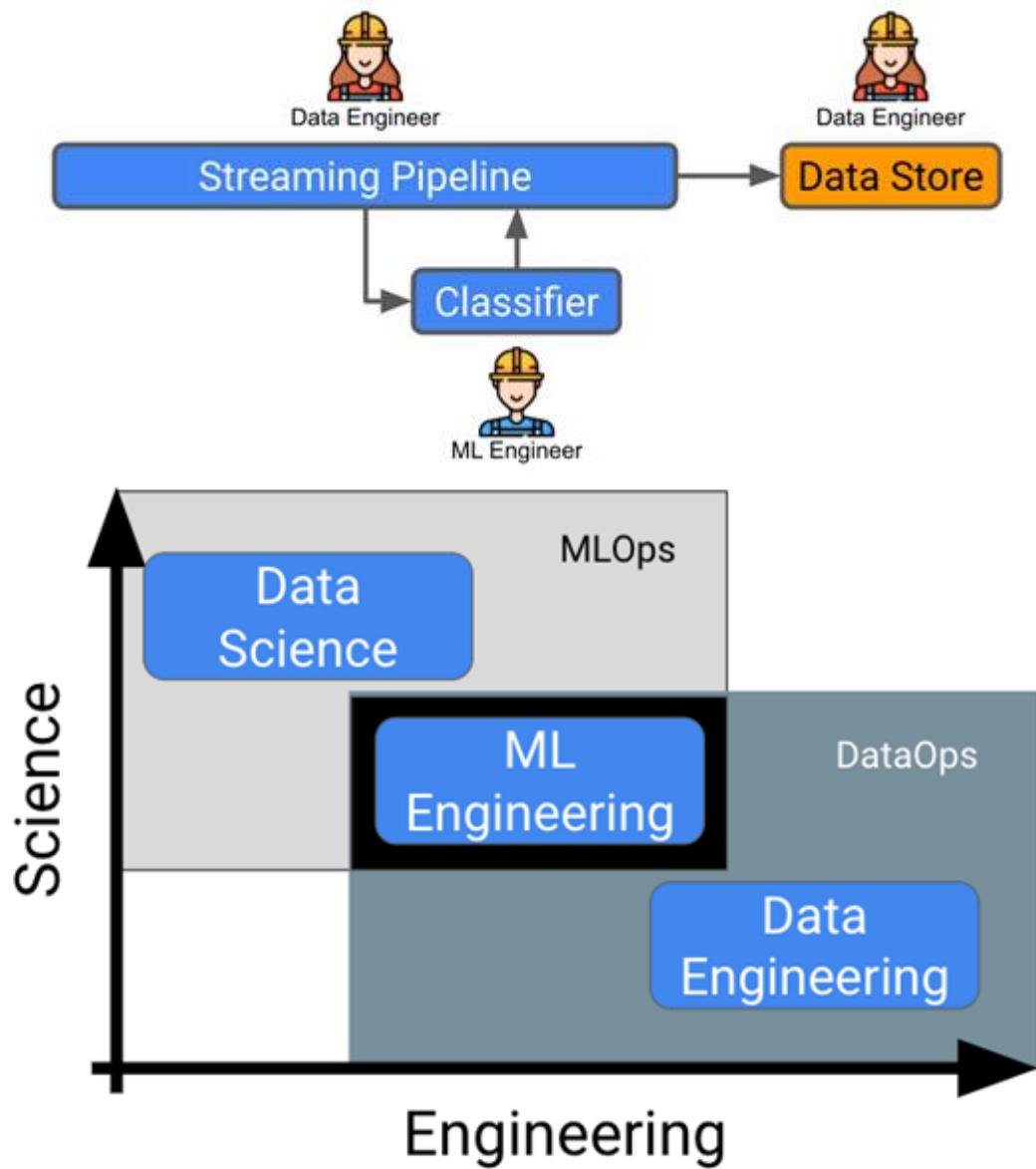


Chapter 1: Introduction to ML Engineering



Requirement	Is ML Appropriate?	Details
Anomaly detection of energy pricing signals	Yes	You will want to do this on large numbers of points on potentially varying time signals.
Improving data quality in an ERP system	No	This sounds more like a process problem. You can try and apply ML to this but often it is better to make the data entry process more automated or the process more robust.
Forecasting item consumption for a warehouse	Yes	ML will be able to do this more accurately than a human can, so this is a good area of application.
Summarizing data for business reviews	Maybe	This can be required at scale but is not an ML problem - simple queries against your data will do.

Surface

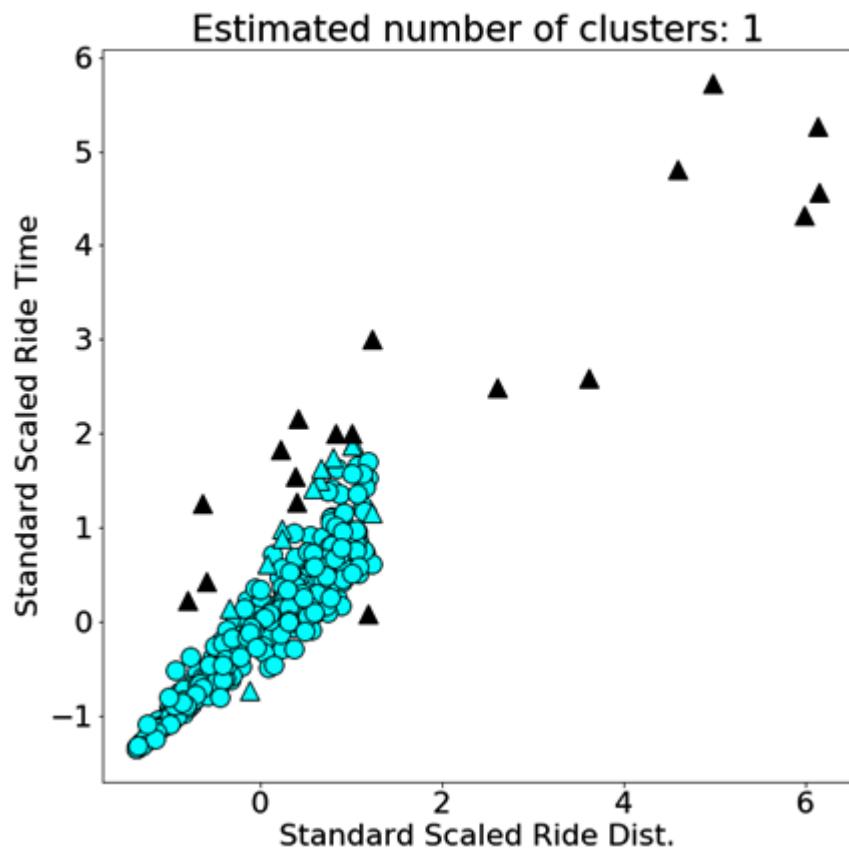
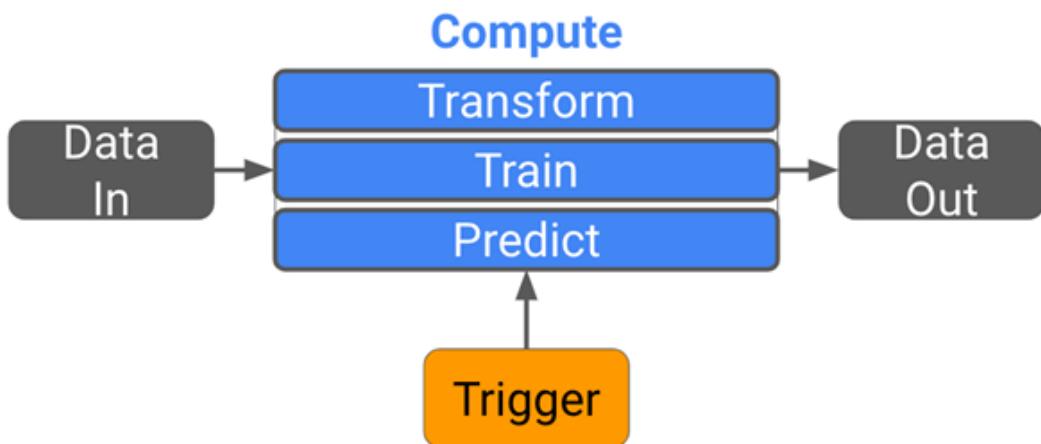
Visualisation, messaging ...

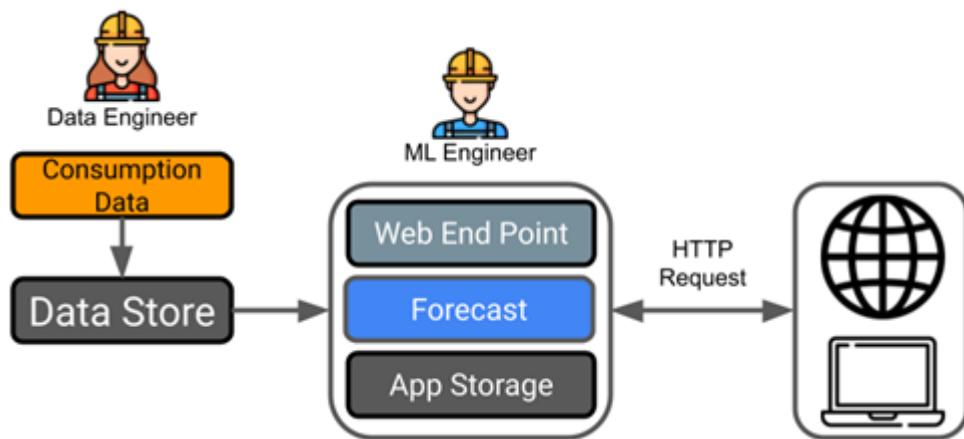
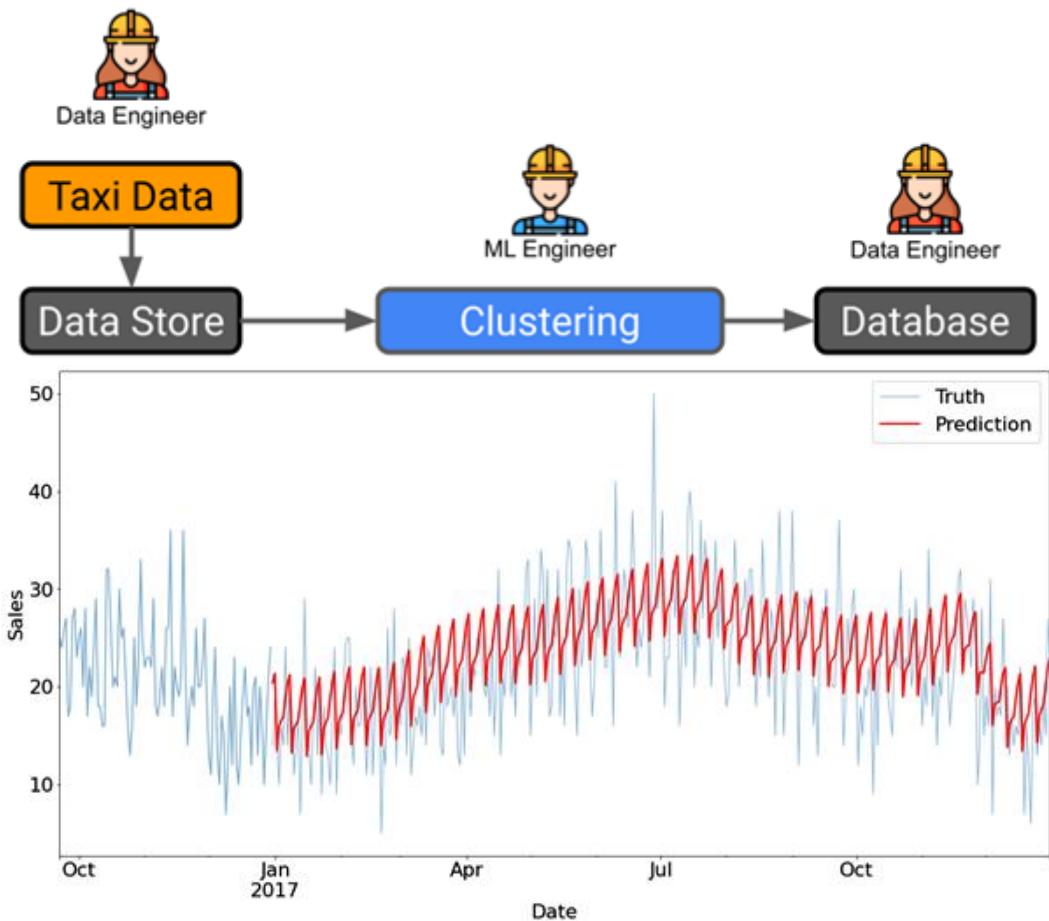
Compute

Modelling, training, transformation

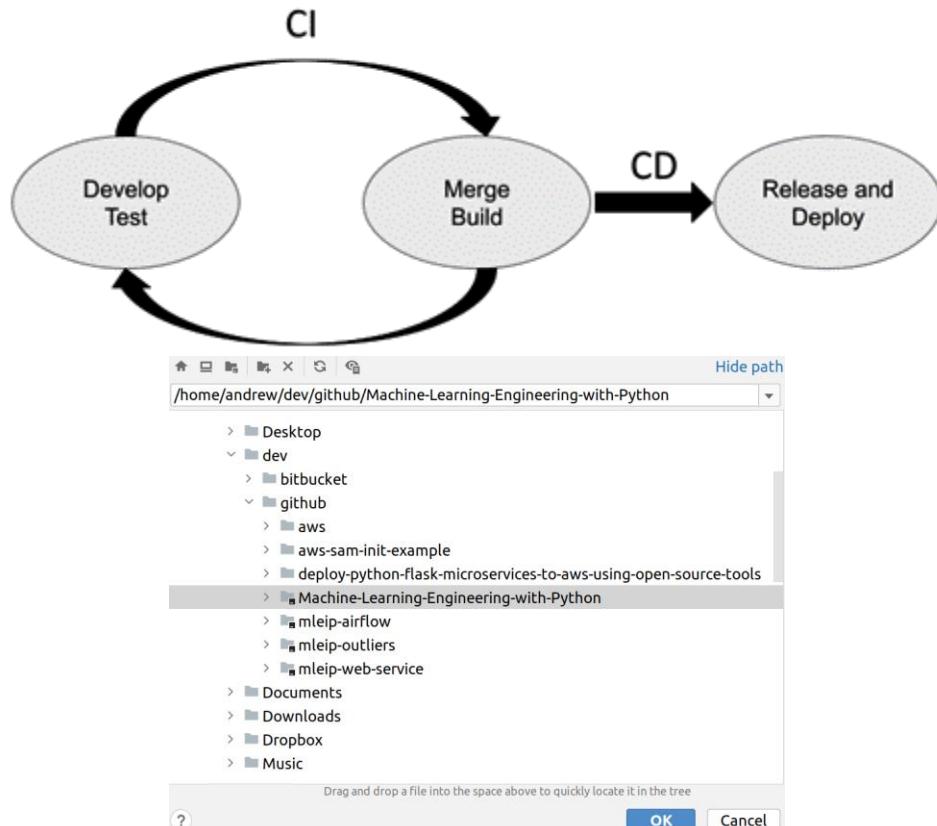
Storage

Data, metadata, artefacts ...





Chapter 2: The Machine Learning Development Process



Create a new repository

A repository contains all project files, including the revision history. Already have a project repository elsewhere? [Import a repository](#).

Owner * Repository name *

 AndyMc629 /

Great repository names are short and memorable. Need inspiration? How about [bookish-goggles](#)?

Description (optional)

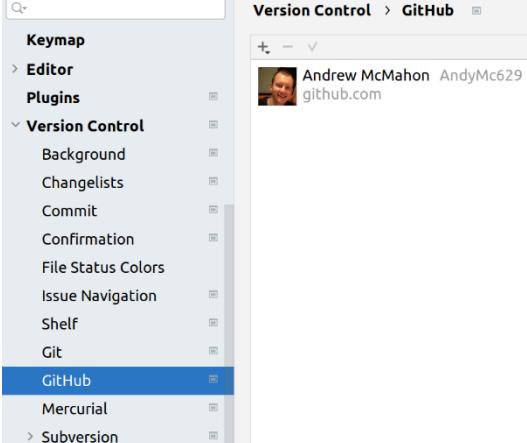
 Public
Anyone on the internet can see this repository. You choose who can commit.

 Private
You choose who can see and commit to this repository.

Initialize this repository with:
Skip this step if you're importing an existing repository.

Add a README file
This is where you can write a long description for your project. [Learn more](#).

Add .gitignore
Choose which files not to track from a list of templates. [Learn more](#).



The screenshot shows the IntelliJ IDEA settings interface with the 'Version Control' section selected. Under 'GitHub', it lists 'Background', 'Changelists', 'Commit', 'Confirmation', 'File Status Colors', 'Issue Navigation', 'Shelf', 'Git', and 'GitHub'. The 'GitHub' item is highlighted with a blue selection bar. On the right side of the interface, there is a preview window showing a GitHub profile for 'Andrew McMahon' (AndyMc629) with the URL 'github.com'.

Projects / ml-engineering-in-python

MEIP board

The board has the following issues:

- To Do (4 issues):**
 - As a store demand planner, I want to see forecasts for items split by region to anticipate extra orders that need to be made (MEIP-8)
 - Add DBSCAN functionality to DetectionModels in outliers package (MEIP-26)
- In Progress (5 issues):**
 - As a store demand planner, I want to be able to trigger retraining of models I think are out of date to improve forecast performance (MEIP-9)
 - Build forecasting algorithm: build basic prophet algorithm (use code already developed) (MEIP-6)
- Done (✓)**
 - See all Done issues

AWS Management Console

AWS services

- Recently visited services:
 - Billing
 - Managed Apache Airflow
 - EC2
 - Elastic Container Service
 - VPC
- All services

Stay connected to your AWS resources on-the-go

AWS Console Mobile App now supports four additional regions. Download the AWS Console Mobile App to your iOS or Android mobile device.

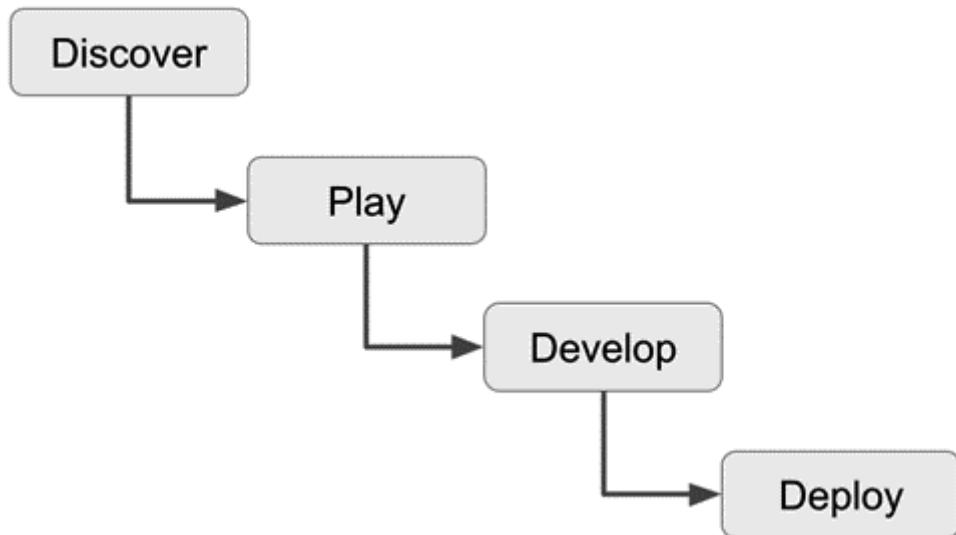
[Learn more](#)

Explore AWS

Introducing AWS Backup Audit Manager

Maintain and demonstrate your data backup and compliance posture at scale. [Learn more](#)

AWS Certification



Stage	Outputs
Discover	Clarity on the business question Clear arguments for ML over another approach Definition of the KPIs and metrics you want to optimize A sketch of the route to value
Play	Detailed understanding of the data Working proof of concept Agreement on the model/algorithm/logic that will solve the problem Evidence that a solution is doable within realistic resource scenarios Evidence that good ROI can be achieved
Develop	A working solution that can be hosted on appropriate and available infrastructure Thorough test results and performance metrics (for algorithms and software) An agreed retraining and model deployment strategy Unit tests, integration tests, and regression tests Solution packaging and pipelines
Deploy	A working and tested deployment process Provisioned infrastructure with appropriate security and performance characteristics Mode retraining and management processes An end-to-end working solution!

Projects / 🎯 ml-engineering-in-py... / ⚒ Add epic / MEIP-14

As a taxi ride analyst, I want to see anomalous journeys so that I can tailor our offers to customers



Description

Acceptance criteria (scenario):

Given I have access to data and/or visualizations of ML results,
where there are anomalous taxi rides,
then the system accurately identifies and labels these rides and I can see them

GitHub integration for Jira

Branches Pull Requests Commits Tags

Prep for Prophet

```
df.rename(columns= {'Datetime': 'ds', 'AEP_MW': 'y'}, inplace=True)
df['ds']=df['ds'].astype('datetime64[ns]')
df.dtypes
#Initialize Split Class, we'll split our data 5 times for cv
ts_splits = TimeSeriesSplit(n_splits=5)
```

Train and Forecast

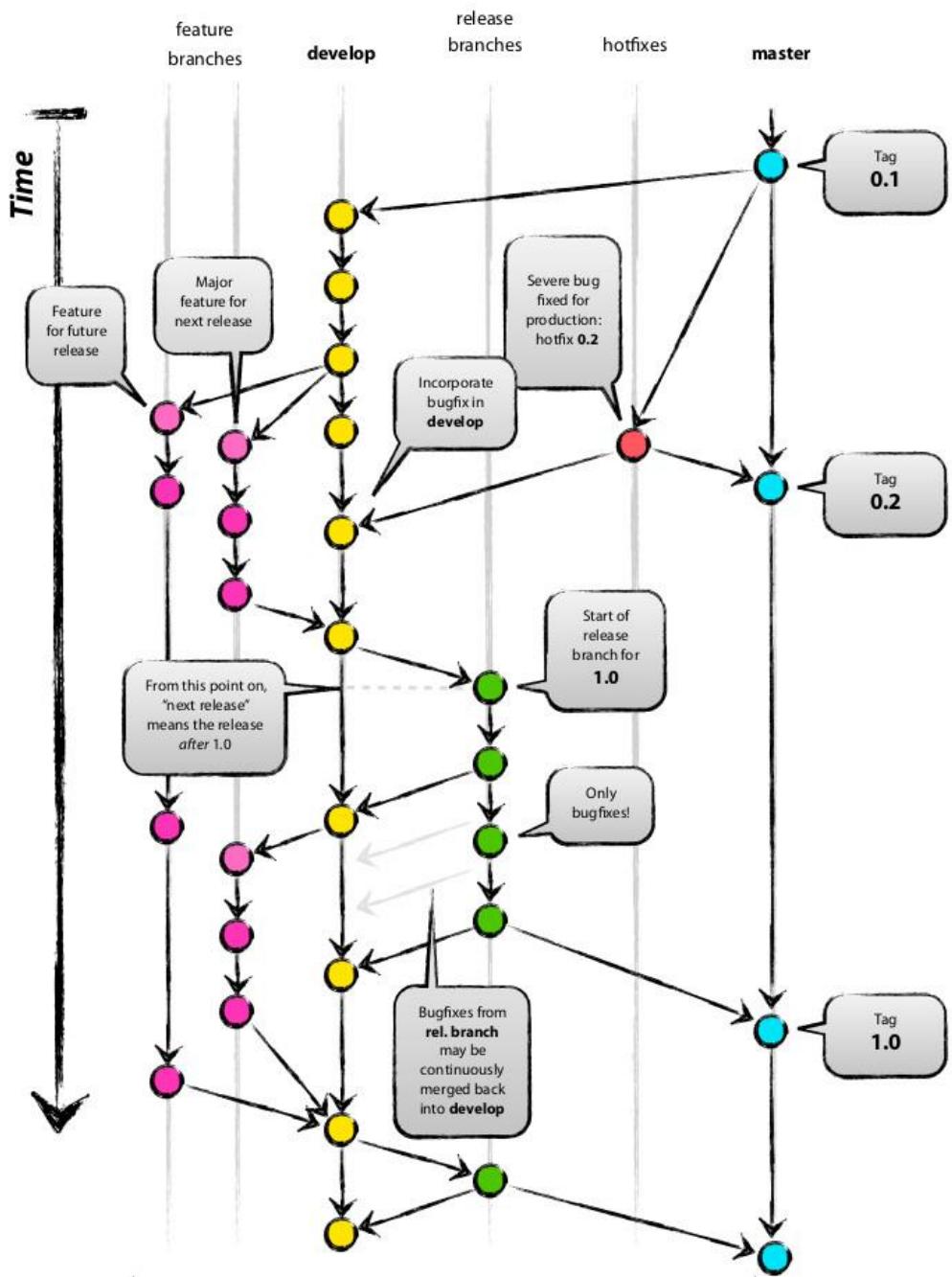
```
tmp = time_split_train_test(df.sort_values('ds', ascending=True).iloc[-1000:], ts_splits)
tmp.head()
```

Plot

```
nrow = 5; ncol = 1;
fig, axs = plt.subplots(nrows=nrow, ncols=ncol, figsize=(20,30))
fig.subplots_adjust(hspace=0.4, wspace=0.4)
for i, ax in enumerate(fig.axes):
    split_rmse = tmp[(tmp['split']==i) & (tmp['train']==False)]['rmse'].iloc[0]
    ax.set_title('Split '+str(i+1) + ' - RMSE: '+'{:.2f}'.format(split_rmse))

    tmp[(tmp['split']==i) & (tmp['train']==True)].plot(x='ds', y='y', ax=ax, color='blue', marker='o')
    tmp[(tmp['split']==i) & (tmp['train']==False)].plot(x='ds', y='y', ax=ax, color='red', marker='o')
    tmp[(tmp['split']==i) & (tmp['train']==False)].plot(x='ds', y='yhat', ax=ax, color='orange', marker='^')
```

Methodology	Pros	Cons
Agile	Flexibility is expected. Faster dev to deploy cycles.	If not well managed, can easily have scope drift. Kanban or Sprints may not work well for some projects.
Waterfall	Clearer path to deployment. Clear staging and ownership of tasks.	Lack of flexibility. Higher admin overheads.



```

v ~ 9 chapter1/bad-git/pipeline.py ⌂
...
... @@ -1,9 +1,14 @@
1 - # EXAMPLE BELOW TAKEN FROM THE SPARK API DOCS, BEFORE BEING UPDATED FOR THE BOOK
2 + # EXAMPLE BELOW ADAPTED FROM THE SPARK DOCS
3   # https://spark.apache.org/docs/latest/ml-pipeline.html#pipeline
4   from pyspark.ml import Pipeline
5   from pyspark.ml.classification import LogisticRegression
6   from pyspark.ml.feature import HashingTF, Tokenizer
7
8 +
9 + with open("model_config.json") as f:
10 +     model_config = json.load(f)
11 +
12   # Prepare training documents from a list of (id, text, label) tuples.
13   training = spark.createDataFrame([
14       (0, "a b c d e spark", 1.0),
15   @@ -15,7 +20,7 @@
16   # Configure an ML pipeline, which consists of three stages: tokenizer, hashingTF, and lr.
17   tokenizer = Tokenizer(inputCol="text", outputCol="words")
18   hashingTF = HashingTF(inputCol=tokenizer.getOutputCol(), outputCol="features")
19 - lr = LogisticRegression(maxIter=10, regParam=0.001)
20 + lr = LogisticRegression(maxIter=model_config['maxIter'], regParam=model_config['regParam'])
21
22 pipeline = Pipeline(stages=[tokenizer, hashingTF, lr])

```

mlflow Experiments Models GitHub Docs

Experiments

Default

Search Experiments

Experiment ID: 0

Artifact Location: file:///home/andrew/dev/github/Machine-Learning-Engineering-with-Python/chapter1/mlflow/mlruns/0

Notes

None

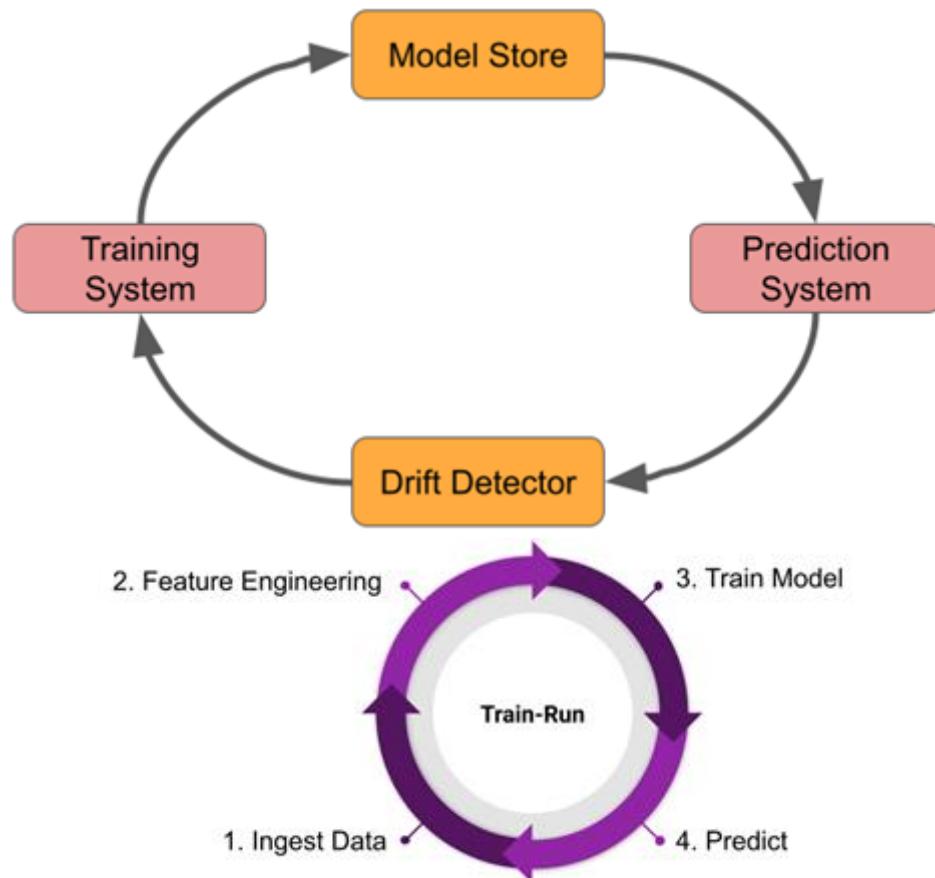
Search Runs:

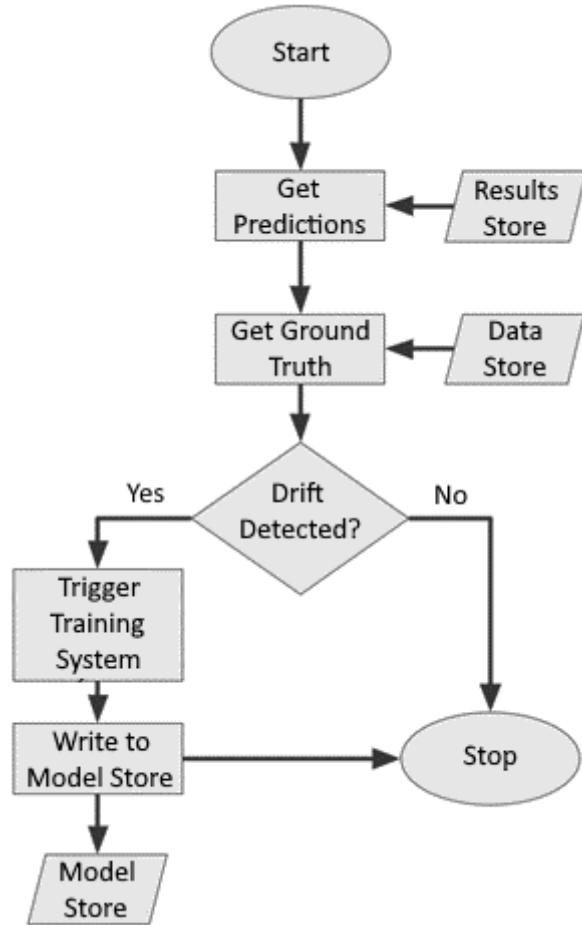
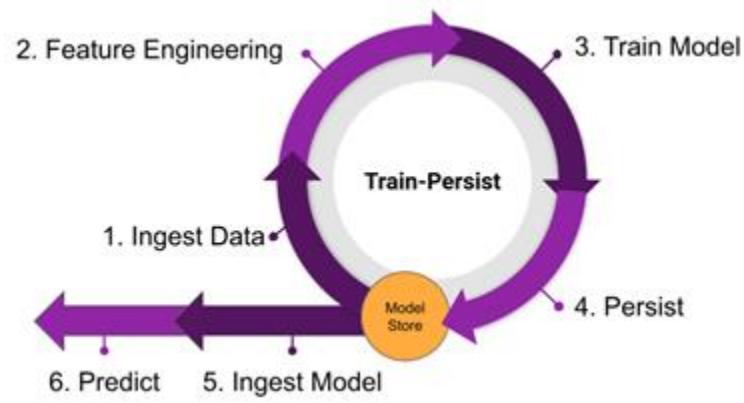
Showing 3 matching runs

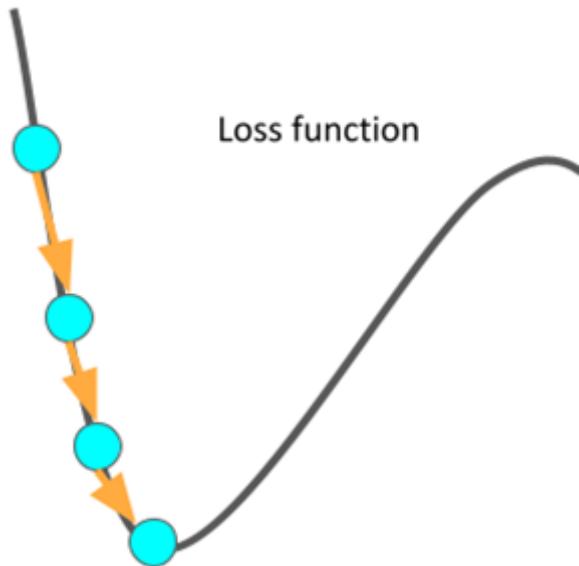
						Parameters	Metrics		
	Start Time	Run Name	User	Source	Version	Models	param1	foo	rmse
<input type="checkbox"/>	2021-02-22 13:12:57	-	andrew	mlflow_fore_c6790e	-	-	-	-	4.853
<input type="checkbox"/>	2021-02-22 10:23:15	-	andrew	mlflow_fore_c6790e	-	-	-	-	4.853
<input type="checkbox"/>	2021-02-17 19:23:53	-	andrew	mlflow_base_ed9a02	-	20	2.762	-	

Chapter 3: From Model to Model Factory

Algorithm	Hyperparameters	What This Controls
K-Nearest Neighbors	<ul style="list-style-type: none">KDistance metric	<ul style="list-style-type: none">The number of clustersHow to define the distance between points
DBSCAN	<ul style="list-style-type: none">EpsilonMinimum number of samplesDistance metric	<ul style="list-style-type: none">The max distance to be considered neighborsHow many neighbors are required to be considered coreHow to define the distance between points

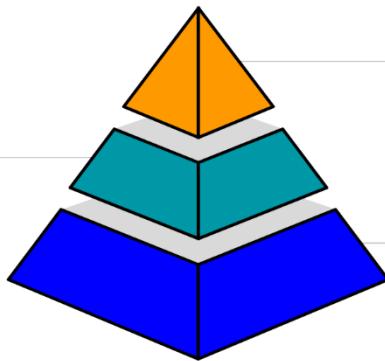






Hyperparameter Optimization

Tune the parameters that tell your ML algorithm how to learn. Use tools that allow for efficient search over hyperparameter spaces.



AutoML

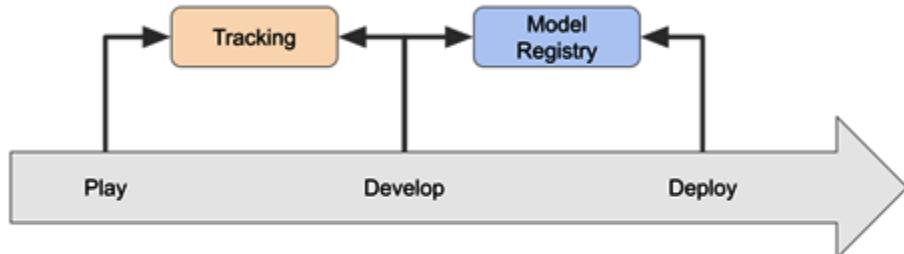
Perform automated selection of the best hyperparameters and the best ML algorithm. Provide only a high level steer of this process.

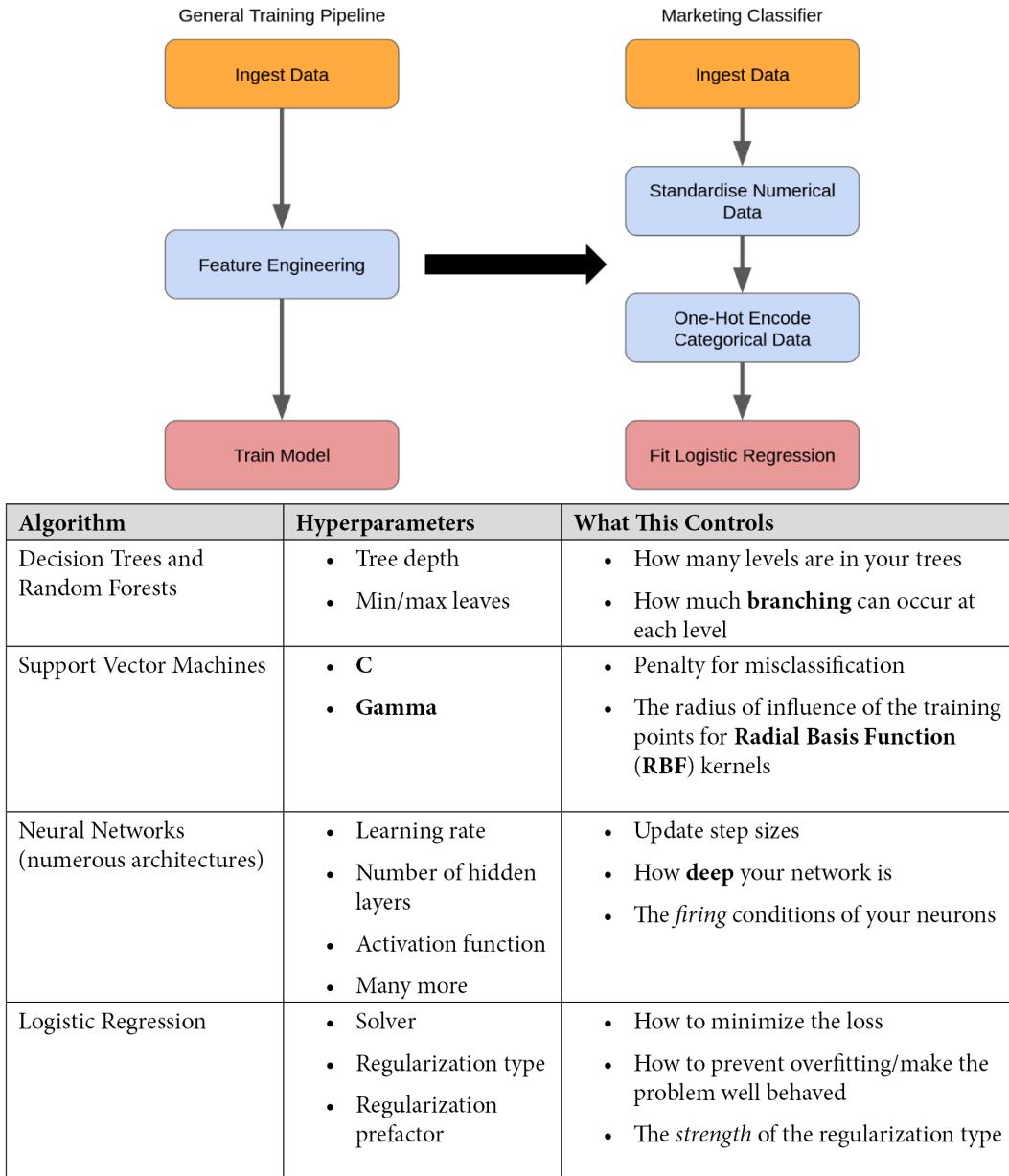
Hand Cranking

Define the hyperparameters for learning explicitly and select one ML algorithm to try at a time.

Algorithm	Hyperparameters	What This Controls
Decision Trees & Random Forests	<ul style="list-style-type: none"> • Tree depth • Min/max leaves 	<ul style="list-style-type: none"> • How many levels in your trees • How much ‘branching’ can occur at each level
Support Vector Machines	<ul style="list-style-type: none"> • ‘C’ • ‘Gamma’ 	<ul style="list-style-type: none"> • Penalty for misclassification • Radius of influence of training points for Radial Basis Function (RBF) kernels
Neural Networks (numerous architectures)	<ul style="list-style-type: none"> • Learning rate • Number of hidden layers • Activation function • Many more 	<ul style="list-style-type: none"> • Update step sizes • How ‘deep’ your net is • The ‘firing’ conditions of your neurons
Logistic Regression	<ul style="list-style-type: none"> • Solver • Regularization type • Regularization prefactor 	<ul style="list-style-type: none"> • How to minimise the loss • How to prevent overfitting/make problem well behaved • The amount to do this

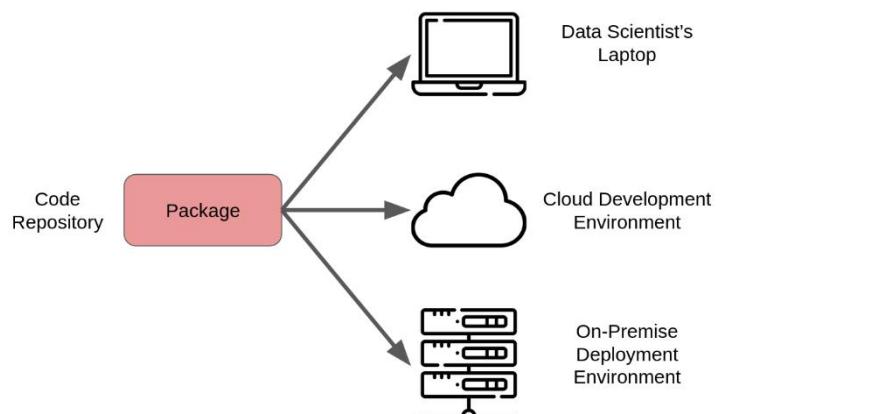
Algorithm	Hyperparameters	What This Controls
K-Nearest Neighbors	<ul style="list-style-type: none"> • “K” • Distance metric 	<ul style="list-style-type: none"> • The number of clusters • How to define distance between points
DBSCAN	<ul style="list-style-type: none"> • “Epsilon” • Minimum number of samples • Distance metric 	<ul style="list-style-type: none"> • The max distance to be considered ‘neighbors’ • How many neighbors required to be considered ‘core’ • How to define distance between points





Chapter 4: Packaging Up

Container	Description
deque	This is a double-ended queue and allows you to add and remove elements to either end of the object in a scalable way. It's useful if you want to add to the beginning or end of large data lists or if you want to search for the last occurrences of X in your data.
Counter	Counters take in iterables such as dicts or lists and return the count of each of the elements. They're really useful to get quick summaries of the content of these objects.
OrderedDict	The standard dict object does not maintain order, so OrderedDict introduces this functionality. This can be really useful if you need to loop back over a dictionary you have created in the same order as it was created for new processing.



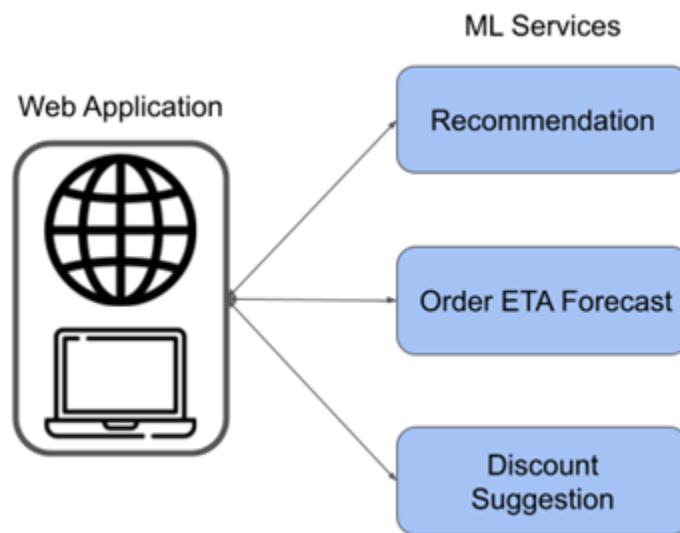
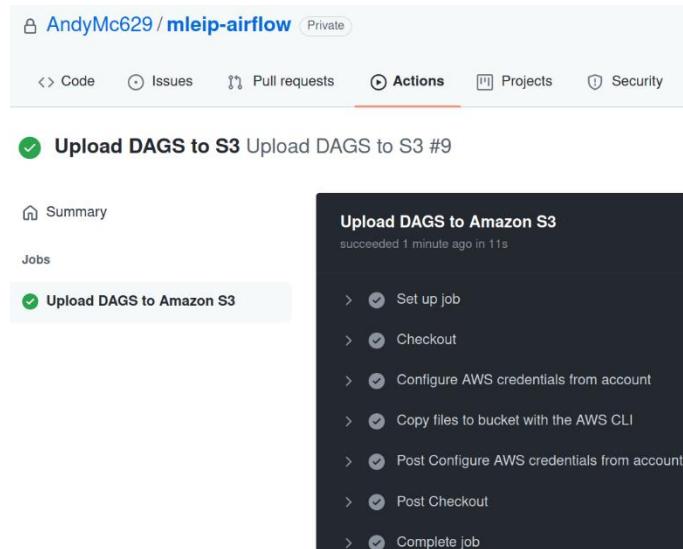
```
===== test session starts =====
platform linux -- Python 3.8.5, pytest-6.1.1, py-1.9.0, pluggy-0.13.1 -- /home/andrew/anaconda3/envs/mleng/bin/python
cachedir: .pytest_cache
rootdir: /home/andrew/dev/github/Machine-Learning-Engineering-with-Python/chapter4/outlier_package
collected 2 items

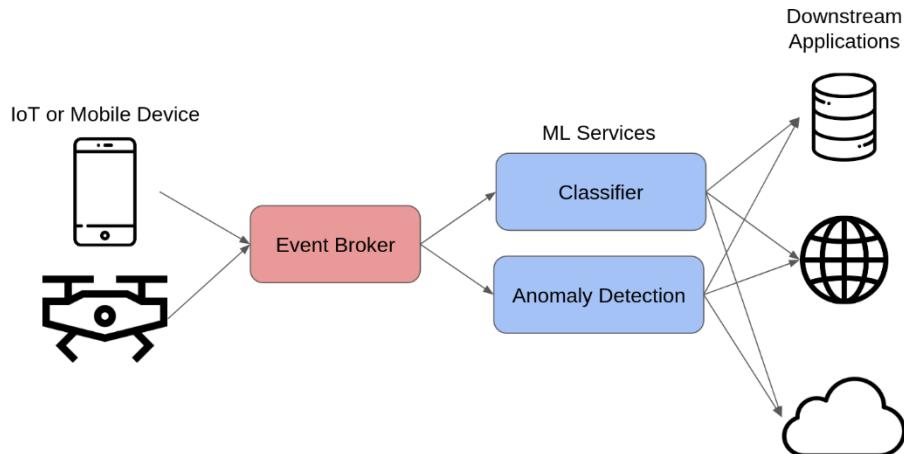
outliers/tests/test_create_data.py::test_data_is_numpy PASSED
outliers/tests/test_create_data.py::test_data_is_large PASSED

===== 2 passed in 0.45s =====
===== test session starts =====
platform linux -- Python 3.8.5, pytest-6.1.1, py-1.9.0, pluggy-0.13.1 -- /home/andrew/anaconda3/envs/mleng/bin/python
cachedir: .pytest_cache
rootdir: /home/andrew/dev/github/Machine-Learning-Engineering-with-Python/chapter4/outlier_package
collected 5 items

outliers/tests/test_create_data.py::test_data_is_numpy PASSED
outliers/tests/test_create_data.py::test_data_is_large PASSED
outliers/tests/test_detectors.py::test_model_creation PASSED
outliers/tests/test_detectors.py::test_model_get_models PASSED
outliers/tests/test_detectors.py::test_model_evaluation PASSED
```

Chapter 5: Deployment Patterns and Tools





MLEIP / MLEIP-web-service-basic

POST ▼

Params	Authorization	Headers (8)	Body	Pre-request Script	Tests	Settings
<input type="radio"/> none	<input type="radio"/> form-data	<input type="radio"/> x-www-form-urlencoded	<input checked="" type="radio"/> raw	<input type="radio"/> binary	<input type="radio"/> GraphQL	JSON ▼

```

1  [
2    "store_number": 55,
3    "forecast_start_date": "2021-10-10T00:00:00"
4  ]

```

Body Cookies Headers (4) Test Results ▼

Pretty	Raw	Preview	Visualize	JSON	▼	=
--------	-----	---------	-----------	------	---	---

```

1  [
2    "result": [
3      0.5119601362414734,
4      0.048732590944116194,
5      0.6843997801791688,
6      0.8694554790764257,
7      0.6549083843755388,
8      0.7591360074029415,
9      0.02741600455310511,
10     0.9855584281616281,
11     0.054214771538483975,
12     0.3851831190248214
13   ],
14   "store_number": 55
15 ]

```

IMAGE ID	CREATED AT	REPOSITORY
cfc60d3d8055	2021-08-22 21:28:21 +0100 BST	basic-ml-webservice
cfc60d3d8055	2021-08-22 21:28:21 +0100 BST	508972911348.dkr.ecr.eu-west-1.amazonaws.com/basic-ml-microservice
cfc60d3d8055	2021-08-22 21:28:21 +0100 BST	basic-ml-microservice
c369d19a3e8f	2021-08-22 21:28:03 +0100 BST	<none>
4e8f4d160116	2021-08-07 19:40:11 +0100 BST	ml-microservice
06f5a53dccfb	2021-08-07 19:39:52 +0100 BST	<none>
9d753150d71c	2021-08-07 19:32:29 +0100 BST	<none>
f68dc7fd2a66	2021-08-07 19:32:16 +0100 BST	<none>
dd8edfcc5a84	2021-08-07 19:09:58 +0100 BST	<none>
0589878a491c	2021-08-07 19:09:45 +0100 BST	<none>
a27d063bdac	2021-07-13 16:09:49 +0100 BST	inferencefunction
b5712d3d8d03	2021-07-13 16:05:55 +0100 BST	508972911348.dkr.ecr.eu-west-2.amazonaws.com/mleip-lambda-example-repo
b5712d3d8d03	2021-07-13 16:05:55 +0100 BST	inferencefunction
e6f9b123f9a0	2021-07-13 13:52:37 +0100 BST	public.ecr.aws/lambda/python

Linux x86 (64-bit) | Linux ARM

For the latest version of the AWS CLI, use the following command block:

```
$ curl "https://awscli.amazonaws.com/awscli-exe-linux-x86_64.zip" -o "awscliv2.zip"
unzip awscliv2.zip
sudo ./aws/install
```

For a specific version of the AWS CLI, append a hyphen and the version number to the filename. For this example the filename for version [2.0.30](#) would be `awscli-exe-linux-x86_64-2.0.30.zip` resulting in the following command:

```
$ curl "https://awscli.amazonaws.com/awscli-exe-linux-x86_64-2.0.30.zip" -o "awscliv2.zip"
unzip awscliv2.zip
sudo ./aws/install
```

For a list of versions, see the [AWS CLI version 2 changelog](#) on GitHub.

The screenshot shows the AWS Management Console interface for the Amazon ECS service. The top navigation bar includes the AWS logo, a search bar, and a 'Services' dropdown. Below the navigation, there's a promotional message for the New ECS Experience. The left sidebar has a tree view with 'Amazon ECS' expanded, and 'Clusters' is currently selected. The main content area is titled 'Clusters' and contains a brief description of what an ECS cluster is. It features two prominent buttons: 'Create Cluster' and 'Get Started'. At the bottom of the content area, there are 'View' options for 'list' and 'card'.

<p>Networking only</p> <p><u>Resources to be created:</u></p> <ul style="list-style-type: none"> Cluster VPC (optional) Subnets (optional) <p>For use with either AWS Fargate or External instance capacity.</p>	<p>EC2 Linux + Networking</p> <p><u>Resources to be created:</u></p> <ul style="list-style-type: none"> Cluster VPC Subnets <p>Auto Scaling group with Linux AMI</p>
---	--

<p>EC2 Windows + Networking</p> <p><u>Resources to be created:</u></p> <ul style="list-style-type: none"> Cluster VPC Subnets <p>Auto Scaling group with Windows AMI</p>
--

Instance configuration

Provisioning Model On-Demand Instance

With On-Demand Instances, you pay for compute capacity by the hour, with no long-term commitments or upfront payments.

Spot

Amazon EC2 Spot Instances let you take advantage of unused EC2 capacity in the AWS cloud. Spot Instances are available at up to a 90% discount compared to On-Demand prices.

[Learn more](#)

EC2 instance type*  
 Manually enter desired instance type

Number of instances* 

Launch status

Your container instances are launching, and it may take a few minutes until they are immediately available and continue to accrue until you stop or terminate them.

ECS status - 2 of 3 complete mleip-web-app-demos

- ECS cluster**
ECS Cluster mleip-web-app-demos successfully created
- ECS Instance IAM Policy**
IAM Policy for the role ecsInstanceRole successfully attached
- CloudFormation Stack**
Creating CloudFormation stack resources

New ECS Experience Tell us what you think

Clusters > mleip-web-app-demos

Cluster : mleip-web-app-demos Update Cluster Delete Cluster

Get a detailed view of the resources on your cluster.

Cluster ARN	arn:aws:ecs:eu-west-1:508972911348:cluster/mleip-web-app-demos
Status	ACTIVE
Registered container instances	0
Pending tasks count	0 Fargate, 0 EC2, 0 External
Running tasks count	2 Fargate, 0 EC2, 0 External
Active service count	1 Fargate, 0 EC2, 0 External
Draining service count	0 Fargate, 0 EC2, 0 External

Task Definitions

Task definitions specify the container information for your application, such as the resources they will use, how they are linked together, and which host ports they

Create new Task Definition Create new revision Actions ▾

Last updated on September 20, 2018

Status: **ACTIVE** INACTIVE

Filter in this page

New ECS Experience Tell us what you think

Clusters

Task Definitions

Account Settings

Amazon EKS

Clusters

Amazon ECR

Repositories

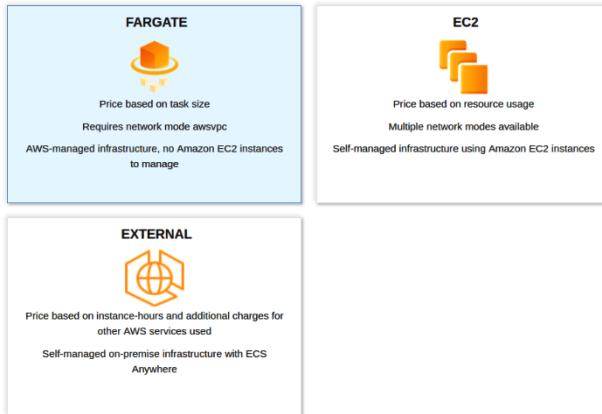
Create new Task Definition

Step 1: Select launch type compatibility

Step 2: Configure task and container definitions

Select launch type compatibility

Select which launch type you want your task definition to be compatible with based on where you want to launch your task.



Configure task and container definitions

A task definition specifies which containers are included in your task and how they interact with each other. You can also specify data volumes for your containers to use. [Learn more](#)

Task Definition Name* ⓘ

Requires Compatibilities* FARGATE

Task Role ⓘ

Optional IAM role that tasks can use to make API requests to authorized AWS services. Create an Amazon Elastic Container Service Task Role in the [IAM Console](#).

Network Mode ⓘ

If you choose <default>, ECS will start your container using Docker's default networking mode, which is Bridge on Linux and NAT on Windows. <default> is the only supported mode on Windows.

Task size

The task size allows you to specify a fixed size for your task. Task size is required for tasks using the Fargate launch type and is optional for the EC2 or External launch type. Container level memory settings are optional when task size is set. Task size is not supported for Windows containers.

Task memory (GB) ⓘ

The valid memory range for 0.25 vCPU is: 0.5GB - 2GB.

Task CPU (vCPU) ⓘ

The valid CPU for 0.5 GB memory is: 0.25 vCPU

Add container

Standard

Container name* basic-ml-microservice [i](#)

Image* 508972911348.dkr.ecr.eu-west-1.amazonaws.com/basic-ml-microservice:latest [i](#)

Private repository authentication* [i](#)

Memory Limits (MiB) 500 [i](#)

[+ Add Hard limit](#)

Define hard and/or soft memory limits in MiB for your container. Hard and soft limits correspond to the 'memory' and 'memoryReservation' parameters, respectively, in task definitions.
ECS recommends 300-500 MiB as a starting point for web applications.

Port mappings [Container port](#) [Protocol](#) [i](#)

5000	tcp	x
------	-----	---

[+ Add port mapping](#)

Task Definitions

Task definitions specify the container information for your application, such as how many containers are part of your task, what resources they more

Create new Task Definition		Create new revision	Actions ▾
Status: ACTIVE INACTIVE			i
Filter in this page			
<input type="checkbox"/> Task Definition			Latest revision status
<input type="checkbox"/> microservice-forecast-task			ACTIVE

Configure service

A service lets you specify how many copies of your task definition to run and maintain in a cluster. You can optionally use an Elastic Load Balancing load balancer to distribute incoming traffic to containers in your service. Amazon ECS maintains that number of tasks and coordinates task scheduling with the load balancer. You can also optionally use Service Auto Scaling to adjust the number of tasks in your service.

Launch type FARGATE i

EC2
 EXTERNAL

[Switch to capacity provider strategy](#) i

Task Definition Family microservice-forecast-task ▼

Revision 1 (latest) ▼

Platform version LATEST ▼ i

Cluster mleip-web-app-demos ▼ i

Service name micro-forecast-service i

Service type* REPLICA i

Number of tasks 2 i

Minimum healthy percent 50 i

Maximum percent 200 i

Deployment circuit breaker Disabled ▼ i

Deployments

Choose a deployment option for the service.

Deployment type* Rolling update i

Blue/green deployment (powered by AWS CodeDeploy) i

This sets AWS CodeDeploy as the deployment controller for the service. A CodeDeploy application and deployment group are created automatically with **default settings** for the service. To change to the rolling update deployment type after the service has been created, you must re-create the service and select the "rolling update" deployment type.

Service IAM role Task definitions that use the awsvpc network mode use the AWSServiceRoleForECS service-linked role, which is created for you automatically. [Learn more.](#)

Load balancer name 

Container to load balance

basic-ml-microservice : 5000

[Remove !\[\]\(f15da8627380db409bac161a6cb03047_img.jpg\)](#)

Production listener port* 

Production listener protocol* HTTP

Target group name  

Target group protocol 

Target type ip 

Path pattern Evaluation order

Path pattern: The first path pattern for a listener is the default path (/), which accepts all traffic that does not match another rule. You can later add additional patterns and priority values to this listener for other services.

Evaluation order: Rules are evaluated in priority order, from the lowest value to the highest value. Once a path pattern rule is matched, all other rules are ignored. You can route traffic from this listener to multiple services by creating a path for each service.

MLEIP / MLEIP-ch5-web-service1-ecs

POST

Params Authorization Headers (8) **Body** Pre-request Script Tests Settings

none form-data x-www-form-urlencoded raw binary GraphQL JSON

```
1 [ ]  
2 "store_number": 55,  
3 "forecast_start_date": "2021-10-10T00:00:00"  
4 [ ]
```

Body Cookies Headers (5) Test Results

Pretty Raw Preview Visualize JSON

```
1 [ ]  
2 "result": [  
3     0.5616842156165673,  
4     0.4780119433699437,  
5     0.2622048564650665,  
6     0.5268415039970527,  
7     0.45073700984020737,  
8     0.6778393849351809,  
9     0.7259141065542977,  
10    0.9281946888027568,  
11    0.46492673618700153,  
12    0.16346221334429245  
13 ],  
14 "store_number": 55  
15 [ ]
```

Select load balancer type

Elastic Load Balancing supports four types of load balancers: Application Load Balancer, Network Load Balancer, Regional Load Balancer, and Classic Load Balancer.

[Learn more about which load balancer is right for you](#)

Application Load Balancer



Create

Choose an Application Load Balancer when you need a flexible feature set for your web applications with HTTP and HTTPS traffic. Operating at the request level, Application Load Balancers provide advanced routing and visibility features targeted at application architectures, including microservices and containers.

[Learn more >](#)

Step 1: Configure Load Balancer

Basic Configuration

To configure your load balancer, provide a name, select a scheme, specify one or more listeners, and select a network. The default configuration is an Internet-facing load balancer.

Name	<input type="text" value="mleip-app-lb"/>
Scheme	<input checked="" type="radio"/> Internet-facing <input type="radio"/> Internal
IP address type	<input type="text" value="ipv4"/>

Listeners

A listener is a process that checks for connection requests, using the protocol and port that you configured.

Load Balancer Protocol	Load Balancer Port
<input type="text" value="HTTP"/>	<input type="text" value="80"/>
Add listener	

Availability Zones

Specify the Availability Zones to enable for your load balancer. The load balancer routes traffic to the targets in these Availability Zones only. You can specify only one VPC and up to five Availability Zones. This increases the availability of your load balancer.

VPC	<input type="text" value="vpc-422d8f3b (172.31.0.0/16) (default)"/>
Availability Zones	<input checked="" type="checkbox"/> eu-west-1a <input type="text" value="subnet-80223de6"/>
IPv4 address	<input type="text" value="Assigned by AWS"/>

Step 3: Configure Security Groups

A security group is a set of firewall rules that control the traffic to your load balancer. On this page, you can add rules to allow specific traffic to reach your load balancer. First, decide whether to create a new security group or select an existing one.

Assign a security group	<input checked="" type="radio"/> Create a new security group <input type="radio"/> Select an existing security group		
Security group name	<input type="text" value="mleip-app-lb-sg-wizard"/>		
Description	<input type="text" value="mleip-app-lb-sg-wizard-1 created on 2021-08-22T20:03:31.275+01:00"/>		
Type	Protocol	Port Range	Source
<input type="text" value="Custom TCP"/>	<input type="text" value="TCP"/>	<input type="text" value="80"/>	<input type="text" value="Custom"/> <input type="text" value="0.0.0.0/:/0"/> <input type="button" value="X"/>
Add Rule			

Step 4: Configure Routing

Your load balancer routes requests to the targets in this target group using the protocol and port that you specify here. It also performs health checks on this load balancer. You can edit or add listeners after the load balancer is created.

Target group

Target group (i)	<input type="text" value="New target group"/> ↳
Name (i)	<input type="text" value="mleip-app-lb-tg"/>
Target type	<input checked="" type="radio"/> Instance <input type="radio"/> IP <input type="radio"/> Lambda function
Protocol (i)	<input type="text" value="HTTP"/> ↳
Port (i)	<input type="text" value="5000"/>
Protocol version (i)	<input checked="" type="radio"/> HTTP1 Send requests to targets using HTTP/1.1. Supported when the request protocol is HTTP/1.1 or HTTP/2. <input type="radio"/> HTTP2 Send requests to targets using HTTP/2. Supported when the request protocol is HTTP/2 or gRPC, but gRPC-specific features are not available. <input type="radio"/> gRPC Send requests to targets using gRPC. Supported when the request protocol is gRPC.

Health checks

Protocol (i)	<input type="text" value="HTTP"/> ↳
Path (i)	<input type="text" value="/"/>

Step 6: Review

Please review the load balancer details before continuing

▼ Load balancer

Name mleip-app-lb
Scheme internet-facing
Listeners Port:80 - Protocol:HTTP
IP address type ipv4
VPC vpc-422d8f3b
Subnets subnet-80223de6, subnet-23806e68, subnet-f70159ad
Tags

▼ Security groups

Security groups mleip-app-lb-sg-wizard

▼ Routing

Target group New target group
Target group name mleip-app-lb-tg
Port 5000
Target type instance
Protocol HTTP
Protocol version HTTP1
Health check protocol HTTP
Path /
Health check port traffic port
Healthy threshold 5
Unhealthy threshold 2
Timeout 5
Interval 30
Success codes 200

Load Balancer Creation Status

- ✓ Successfully created load balancer

Load balancer mleip-app-lb was successfully created.

Note: It might take a few minutes for your load balancer to be fully set up and ready to route traffic, and for the targets to complete the registration process and pass the initial health checks.

Suggested next steps

- Discover other services that you can integrate with your load balancer. Visit the [Integrated services](#) tab within [mleip-app-lb](#)
- Consider using AWS Global Accelerator to further improve the availability and performance of your applications. [AWS Global Accelerator console](#)

Application Integration

Amazon Managed Workflows for Apache Airflow (MWAA)

Run Apache Airflow without provisioning or managing servers

Create an Airflow environment

Launch a complete, auto-scaling Airflow environment in minutes.

Create environment

▼ How Amazon MWAA works



Create an environment

An environment contains your Airflow cluster, including your scheduler, workers, and web server.

Upload your DAGs to Amazon S3

Package and upload your DAG (Directed Acyclic Graph) code to Amazon S3. Amazon MWAA loads the code into Airflow.

Run your DAGs in Airflow

Run your DAGs from the Airflow UI or CLI. Monitor your environment with Amazon CloudWatch.

Environment details [Info](#)

Name

mleip-airflow-dev-env

Use only letters, numbers, dashes, or underscores. Max 80 characters.

Airflow version

2.0.2 (Latest)

Buckets (6) [Info](#)

Buckets are containers for data stored in S3. [Learn more](#)



[Copy ARN](#)

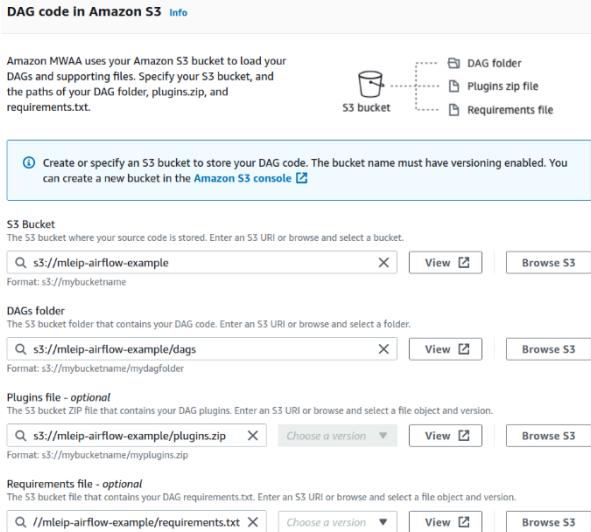
[Empty](#)

[Delete](#)

[Create bucket](#)

[Find buckets by name](#)

Name	AWS Region	Access	Creation date
mleip-airflow-example	EU (Ireland) eu-west-1	Bucket and objects not public	May 18, 2021, 19:26:16 (UTC+01:00)



Quick create stack

Template

Template URL

<https://mwaa-downloads.s3-us-west-2.amazonaws.com/mwaa-vpc-cfn-template.yaml>

Stack description

This template deploys a VPC, with a pair of public and private subnets spread across two Availability Zones. It deploys an internet gateway, with a default route on the public subnets. It deploys a pair of NAT gateways (one in each AZ), and default routes for them in the private subnets.

Stack name

Stack name

MWAA-VPC

Web server access

Private network (Recommended)

Additional setup required. Your Airflow UI can only be accessed by secure login behind your VPC. Choose this option if your Airflow UI is only accessed within a corporate network. IAM must be used to handle user authentication.

Public network (No additional setup)

Your Airflow UI can be accessed by secure login over the Internet. Choose this option if your Airflow UI is accessed outside of a corporate network. IAM must be used to handle user authentication.

 For private network access, the Airflow web server is reached via a VPC endpoint inside your VPC. Connecting to the endpoint requires additional setup. [Learn more about VPC endpoints](#)

Security group(s)

A VPC security group is required to allow traffic between your environment and your web server.

Create new security group

Allow MWAA to create a VPC security group with inbound and outbound rules based on your selection for web server access.

Existing security group(s)

You can choose 1 or more existing security groups to configure the inbound and outbound rules for your environment.

Max 5 security groups

DAG capacity*	Scheduler CPU	Worker CPU	Web server CPU
---------------	---------------	------------	----------------

<input checked="" type="radio"/> mw1.small	Up to 50	1 vCPU	1 vCPU	0.5 vCPU
<input type="radio"/> mw1.medium	Up to 250	2 vCPU	2 vCPU	1 vCPU
<input type="radio"/> mw1.large	Up to 1000	4 vCPU	4 vCPU	2 vCPU

*under typical usage

Maximum worker count

The maximum number of workers your environment is permitted to scale up to.

Must be between 1 and 25

Minimum worker count

The minimum number of workers always present in your environment.

Must be less than or equal to maximum workers. Minimum 1 worker

Permissions [Info](#)

Execution role
The IAM role used by your environment to access your DAG code, write logs, and perform other actions.

[Create a new role](#) [Edit](#) [Delete](#)

Role name

Use alphanumeric and '+,-,@,_-' characters. Maximum 64 characters.

Info Amazon MWAA will create and assume the execution role in IAM named **AmazonMWAA-mleip-airflow-dev-env-lyAb8s** on your behalf. This role is configured with permission to retrieve code from your Amazon S3 bucket, use your KMS key, and send data to Amazon CloudWatch. You must add permissions to your execution role if your Airflow DAGs require access to any other AWS services. [Info](#)

Airflow environments

Environments (1)					Edit	Delete	Actions	Create environment
<input type="text"/> Find environments					<	1	>	⚙️
Name	Status	Created date	Airflow version	Airflow UI				
mleip-airflow-dev-env	Creating	Sep 20, 2021 20:32:11 ...	2.0.2	Open Airflow ...				

Chapter 6: Scaling Up



pdays	previous	poutcome	y	month_as_int
-1	0	unknown	no	10
339	4	failure	no	5
330	1	failure	no	4
-1	0	unknown	no	6
-1	0	unknown	no	5
176	3	failure	no	2
330	2	other	no	5
-1	0	unknown	no	5
-1	0	unknown	no	5
147	2	failure	no	4
-1	0	unknown	no	5
-1	0	unknown	no	4
-1	0	unknown	no	8
-1	0	unknown	yes	4
241	1	failure	no	1
-1	0	unknown	no	8
-1	0	unknown	no	8
152	2	failure	no	4
-1	0	unknown	no	5
152	1	other	no	7

_1	_2	_3	_4	_5	_6	_7	_8	_9	_10	_11	_12	_13	class
14.23	1.71	2.43	15.6	127.0	2.8	3.06	0.28	2.29	5.64	1.04	3.92	1065.0	0
13.2	1.78	2.14	11.2	100.0	0.2	2.65	2.76	0.26	1.28	4.38	1.05	3.41	1050.0
13.16	2.36	2.67	18.6	101.0	2.8	3.24	0.31	2.81	5.68	1.03	3.17	1185.0	0
14.37	1.95	2.5	16.8	113.0	3.85	3.49	6.24	2.18	7.81	0.86	3.45	1488.0	0
13.24	2.59	2.87	21.0	118.0	2.8	2.69	0.39	1.02	4.32	1.04	2.93	735.0	2

Amazon EMR

EMR Studio

EMR on EC2

Clusters

Notebooks

Git repositories

Welcome to Amazon Elastic MapReduce

Amazon Elastic MapReduce (Amazon EMR) is a web service that enables businesses, researchers, data analysts, and developers to easily and cost-effectively process vast amounts of data.

You do not appear to have any clusters. Create one now:

[Create cluster](#)

General Configuration

Cluster name

Logging [i](#)

S3 folder [i](#)

Launch mode Cluster [i](#) Step execution [i](#)

Software configuration

Release [i](#)

Applications Core Hadoop: Hadoop 2.10.1, Hive 2.3.7, Hue 4.9.0, Mahout 0.13.0, Pig 0.17.0, and Tez 0.9.2
 HBase: HBase 1.4.13, Hadoop 2.10.1, Hive 2.3.7, Hue 4.9.0, Phoenix 4.14.3, and ZooKeeper 3.4.14
 Presto: Presto 0.245.1 with Hadoop 2.10.1 HDFS and Hive 2.3.7 Metastore
 Spark: Spark 2.4.7 on Hadoop 2.10.1 YARN and Zeppelin 0.9.0

Use AWS Glue Data Catalog for table metadata [i](#)

Hardware configuration

Instance type [i](#)

Number of instances (1 master and 1 core nodes)

Cluster scaling scale cluster nodes based on workload

Security and access

EC2 key pair [i](#) [Learn how to create an EC2 key pair.](#)

Permissions Default Custom
Use default IAM roles. If roles are not present, they will be automatically created for you with managed policies for automatic policy updates.

EMR role [EMR_DefaultRole](#) Use EMR_DefaultRole_V2 [i](#)

EC2 instance profile [EMR_EC2_DefaultRole](#) [i](#)

Cluster: My cluster Starting

Summary

ID: j-17UN3RZ4S7K7T
 Creation date: 2021-09-22 19:44 (UTC+1)
 Elapsed time: 1 minute
 After last step completes: Cluster waits
 Termination protection: Off [Change](#)
 Tags: -- [View All / Edit](#)

Master public DNS:
 ec2-34-243-3-14.eu-west-1.compute.amazonaws.com [Edit](#)
[Connect to the Master Node Using SSH](#)

Configuration details

Release label: emr-5.33.0
 Hadoop distribution: Amazon
 Applications: Spark 2.4.7, Zeppelin 0.9.0
 Log URI: s3://aws-logs-508972911348-eu-west-1
 /elasticmapreduce/ [Edit](#)
 EMRFS consistent view: Disabled

Which template source would you like to use? What package type would you like to use?

1 - AWS Quick Start Templates 2 - Custom Template Location Choice: <input type="text" value="1"/>	1 - Zip (artifact is a zip uploaded to S3) 2 - Image (artifact is an image uploaded to an ECR image repository) Package type: <input type="text" value="1"/>
---	--

Which base image would you like to use?

- 1 - amazon/nodejs14.x-base
- 2 - amazon/nodejs12.x-base
- 3 - amazon/nodejs10.x-base
- 4 - amazon/python3.8-base
- 5 - amazon/python3.7-base
- 6 - amazon/python3.6-base
- 7 - amazon/python2.7-base
- 8 - amazon/ruby2.7-base
- 9 - amazon/ruby2.5-base
- 10 - amazon/go1.x-base
- 11 - amazon/java11-base
- 12 - amazon/java8.al2-base
- 13 - amazon/java8-base
- 14 - amazon/dotnet5.0-base
- 15 - amazon/dotnetcore3.1-base
- 16 - amazon/dotnetcore2.1-base

Base image: Project name [sam-app]: mleip-lambda-example

AWS quick start application templates:

- 1 - Hello World Lambda Image Example
- 2 - PyTorch Machine Learning Inference API
- 3 - Scikit-learn Machine Learning Inference API
- 4 - Tensorflow Machine Learning Inference API
- 5 - XGBoost Machine Learning Inference API

Template selection:

 Generating application:

 Name: mleip-lambda-example
 Base Image: amazon/python3.8-base
 Dependency Manager: pip
 Output Directory: .

Next steps can be found in the README file at ./mleip-lambda-example/README.md

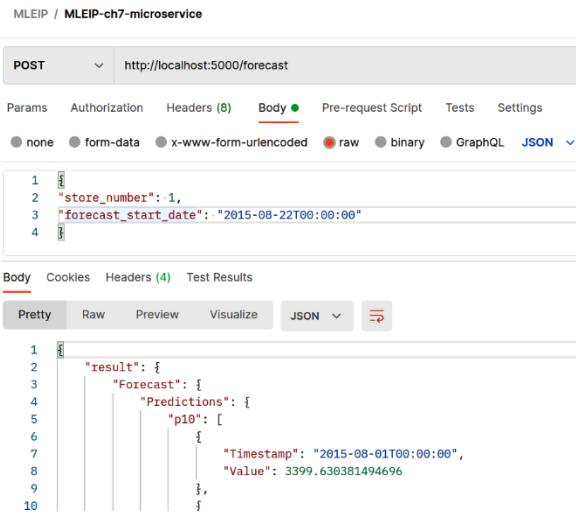
```
CloudFormation outputs from deployed stack
-----
Outputs
-----
Key           InferenceApi
Description   API Gateway endpoint URL for Prod stage for Inference function
Value         https://biwapzz3y5.execute-api.eu-west-2.amazonaws.com/Prod/classify_digit/
Key           InferenceFunctionIamRole
Description   Implicit IAM Role created for Inference function
Value         arn:aws:lambda:eu-west-2:508972911348:function:sam-app-InferenceFunction-WKDKOefcLIWn
Key           InferenceFunction
Description   Inference Lambda Function ARN
Value         arn:aws:lambda:eu-west-2:508972911348:function:sam-app-InferenceFunction-WKDKOefcLIWn
```

Successfully created/updated stack - sam-app in eu-west-2

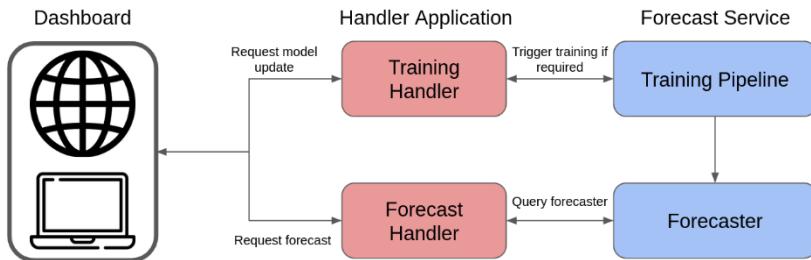
The screenshot shows the Postman interface with the following details:

- URL:** https://biwapzz3y5.execute-api.eu-west-2.amazonaws.com/Prod/classify_digit/
- Method:** POST
- Body:** form-data
- File:** ch6-flg15.png

Chapter 7: Building an Example ML Microservice



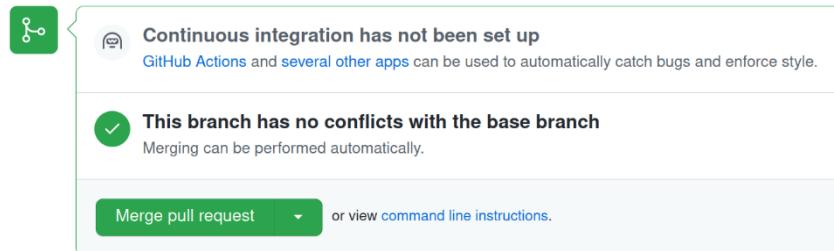
User story	Details	Technical requirements
1	As a local logistics planner, I want to log in to a dashboard in the morning and be able to see forecasts of item demand at store level for the next few days.	<ul style="list-style-type: none"> Target variable = item demand Forecast horizon = 1–7 days Interface with dashboard required
2	As a local logistics planner, I want to be able to request an update of my forecast if I see it is out of date. I want the new forecast to be retrieved in a reasonable time.	<ul style="list-style-type: none"> Lightweight retraining Model per store?
3	As a local logistics planner, I want to be able to filter for forecasts based on specific stores.	<ul style="list-style-type: none"> Model per store



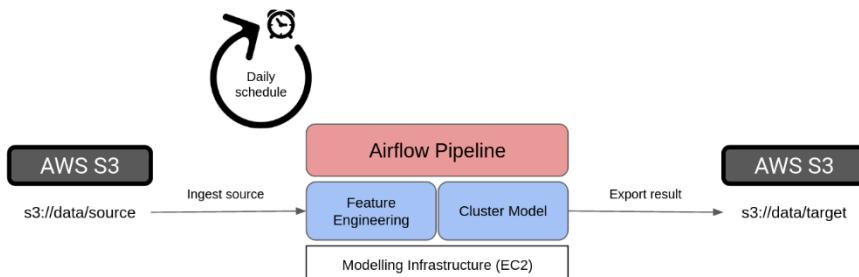
Tool/ framework	Pros	Cons
Sklearn	<ul style="list-style-type: none"> Already understood by almost all data scientists Very easy-to-use syntax Lots of great community support Good feature engineering and pipelining support 	<ul style="list-style-type: none"> No native time series modeling capabilities Will require some more feature engineering to apply to models to time series data More work and understanding required by engineer/scientist
Prophet	<ul style="list-style-type: none"> Purely focused on forecasting Has inbuilt hyperparameter optimization capabilities Provides a lot of easy-to-use functionality out of the box Often gives accurate results on a wide variety of problems Provides confidence intervals out of the box 	<ul style="list-style-type: none"> Not as commonly used as sklearn (but still relatively popular) Underlying methods are quite sophisticated – may lead to <i>black box</i> usage Not inherently scalable
Spark MLlib	<ul style="list-style-type: none"> Natively scalable to large volumes Good feature engineering and pipelining support 	<ul style="list-style-type: none"> No native time series modeling capabilities Algorithm options are relatively limited

Chapter 8: Building an Extract Transform Machine Learning Use Case

Add more commits by pushing to the [feature/MEIP-26-add-dbscan-functionality-to-detection-models-in-outliers-package](#) branch on [AndyMc629/mleip-outliers](#).



User Story	Details	Technical Requirements
1	As an operations analyst, I want to be given clear labels of rides that have anomalously long ride times or distances.	<ul style="list-style-type: none"> Algorithm type = anomaly detection/clustering/outlier detection Features = ride time and distance
2	As an internal application developer, I want to have a clear access point for data with anomalous labels. This data should be stored in the cloud.	<ul style="list-style-type: none"> System output destination = S3 on AWS
3	As an operations analyst, I would like to see labels for the previous day's rides every morning.	<ul style="list-style-type: none"> Batch frequency = daily



Solution Aspect	Potential Tools	Pros	Cons
Interfaces	AWS Command-Line Interface (CLI) and boto3	<ul style="list-style-type: none"> Simple to use Connects to a wide variety of other AWS tools and services 	<ul style="list-style-type: none"> Not cloud-agnostic Not applicable outside of AWS (on-premises systems, for example)

Potential Tools	Pros	Cons
Spark MLlib	Can scale to very large datasets	<ul style="list-style-type: none"> Requires cluster management Overkill for smaller datasets/processing requirements Limited algorithm set
Spark with pandas User-Defined Function (UDF)	<ul style="list-style-type: none"> Can scale to very large datasets Can use any Python-based algorithm 	Might not make sense for some problems where parallelization is not easily applicable
Scikit-learn	<ul style="list-style-type: none"> Well known by many data scientists Can run on many different types of infrastructure 	Not very scalable

Potential Tools	Pros	Cons
Apache Airflow	<ul style="list-style-type: none"> Good scheduling management Ability to build relatively complex pipelines Good documentation Cloud-hosted services available, such as AWS Managed Workflows for Apache Airflow (MWAA) 	<ul style="list-style-type: none"> Learning curve for usage Takes time to test pipelines and scheduling Cloud services (MWAA) can be expensive

Projects Branches Commits

Visibility	Project	Key	Repositories
<input checked="" type="checkbox"/>	 ml-engineering-in-python	MEIP	<input type="text" value="mleip-outlier"/> <div style="border: 1px solid #ccc; padding: 2px; width: 150px;"> AndyMc629/mleip-outliers </div>

[Add epic](#) / [MEIP-26](#)

[🔒](#) [1](#) [↑](#) [🔗](#) [...](#)

Add DBSCAN functionality to DetectionModels in outliers package



Description

As a machine learning engineer I want to be able to run and test DBSCAN functionality in the outlier package. I want to see successfully clustering using the same syntax as is currently employed but with this new algorithm.



Add a comment...

Pro tip: press [M](#) to comment

[Create](#) [Branches](#) [Commits](#) [PRs](#) [Tags](#)

Repository *

AndyMc629/mleip-outliers

Base branch *

↳ main

Branch name *

feature/MEIP-26-add-dbscan-functionality-to-det

[Branch Name Template](#)

[Create branch](#)

Overview

Yours

Active

Stale

All branches

Default branch

main

Updated last month by AndyMc629

[Default](#)



Your branches

feature/MEIP-26-add-dbscan-functionality-to-detection-models-in-outliers-package

Updated last month by AndyMc629

#1

[Open](#)



Active branches

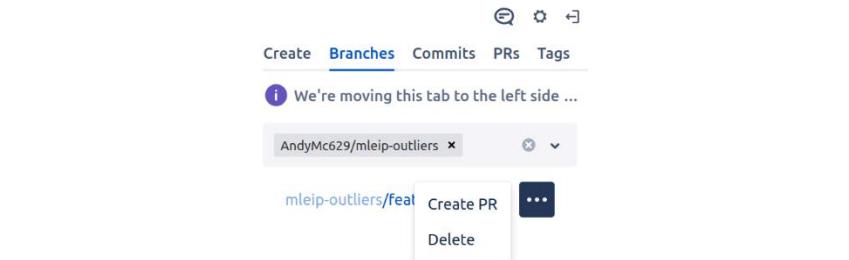
feature/MEIP-26-add-dbscan-functionality-to-detection-models-in-outliers-package

Updated last month by AndyMc629

#1

[Open](#)



A screenshot of a GitHub repository interface. At the top, there are navigation links: Create, Branches (which is underlined), Commits, PRs, and Tags. Below this, a message says "We're moving this tab to the left side ...". A dropdown menu is open for the branch "mleip-outliers/feat", showing options to "Create PR" and "Delete".

Feature/meip 26 add dbscan functionality to detection models in outliers package #1

Open

AndyMc629 wants to merge 3 commits into [main](#) from [feature/MEIP-26-add-dbscan-functionality-to-detection-mod](#)

Conversation 0 · Commits 3 · Checks 0 · Files changed 5

 AndyMc629 commented on 12 Aug

Main changes:

- Logging
- Exception handling
- Inclusion of DBSCAN functionality
- Taxi data simulation