



利用**NetApp**文件对象二元性和 **AWS SageMaker** 进行云数据管理

NetApp artificial intelligence solutions

NetApp
August 18, 2025

目录

利用NetApp文件对象二元性和 AWS SageMaker 进行云数据管理	1
TR-4967：使用NetApp文件对象二元性和 AWS SageMaker 进行云数据管理	1
解决方案技术	1
用例摘要	2
数据科学家和其他应用程序的数据二元性	2
技术要求	2
部署过程	2
通过 Jupyter Notebook 验证机器学习	16
结束语	28
在哪里可以找到更多信息	28

利用NetApp文件对象二元性和 AWS SageMaker 进行云数据管理

TR-4967: 使用NetApp文件对象二元性和 AWS SageMaker 进行云数据管理

Karthikeyan Nagalingam, NetApp

数据科学家和工程师经常需要访问以 NFS 格式存储的数据，但直接从 AWS SageMaker 中的 S3 协议访问这些数据可能具有挑战性，因为 AWS 仅支持 S3 存储桶访问。但是，NetApp ONTAP通过为 NFS 和 S3 启用双协议访问提供了解决方案。通过此解决方案，数据科学家和工程师可以通过NetApp Cloud Volumes ONTAP的 S3 存储桶访问来自 AWS SageMaker 笔记本的 NFS 数据。这种方法可以轻松访问和共享来自 NFS 和 S3 的相同数据，而无需额外的软件。

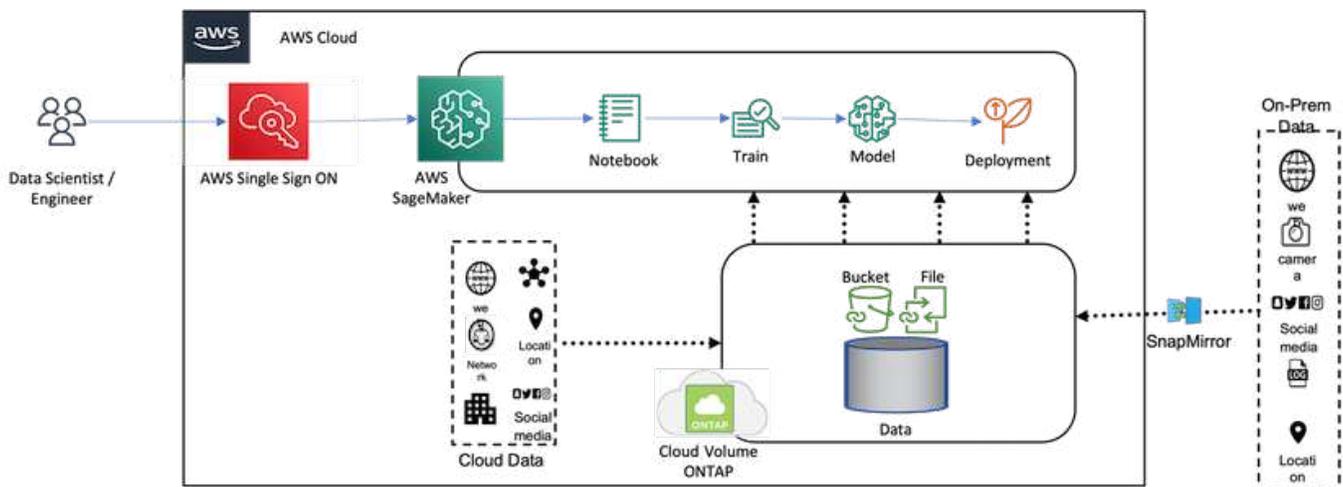
解决方案技术

该解决方案采用以下技术：

- AWS SageMaker 笔记本。为开发人员和数据科学家提供机器学习功能，以高效地创建、训练和部署高质量的 ML 模型。
- * NetApp BlueXP。*支持在本地以及 AWS、Azure 和 Google Cloud 上发现、部署和操作存储。它提供数据保护，防止数据丢失、网络威胁和意外中断，并优化数据存储和基础设施。
- * NetApp Cloud Volumes ONTAP。*在 AWS、Azure 和 Google Cloud 上提供具有 NFS、SMB/CIFS、iSCSI 和 S3 协议的企业级存储卷，让用户在访问和管理云中的数据时拥有更大的灵活性。

NetApp Cloud Volumes ONTAP由BlueXP创建，用于存储 ML 数据。

下图展示了该解决方案的技术组件。



用例摘要

NFS 和 S3 双协议访问的一个潜在用例是在机器学习和数据科学领域。例如，一个数据科学家团队可能正在使用 AWS SageMaker 开展机器学习项目，这需要访问以 NFS 格式存储的数据。但是，可能还需要通过 S3 存储桶访问和共享数据，以便与其他团队成员协作或与使用 S3 的其他应用程序集成。

通过利用 NetApp Cloud Volumes ONTAP，团队可以将其数据存储在单个位置，并可通过 NFS 和 S3 协议访问。数据科学家可以直接从 AWS SageMaker 访问 NFS 格式的数据，而其他团队成员或应用程序可以通过 S3 存储桶访问相同的数据。

这种方法可以轻松高效地访问和共享数据，而无需额外的软件或不同存储解决方案之间的数据迁移。它还允许团队成员之间更简化的工作流程和协作，从而更快、更有效地开发机器学习模型。

数据科学家和其他应用程序的数据二元性

数据可在 NFS 中使用，并可通过 AWS SageMaker 的 S3 访问。

技术要求

对于数据二元性用例，您需要 NetApp BlueXP、NetApp Cloud Volumes ONTAP 和 AWS SageMaker Notebooks。

软件要求

下表列出了实现用例所需的软件组件。

软件	数量
BlueXP	1
NetApp Cloud Volumes ONTAP	1
AWS SageMaker 笔记本	1

部署过程

部署数据二元性解决方案涉及以下任务：

- BlueXP 连接器
- NetApp Cloud Volumes ONTAP
- 机器学习数据
- AWS SageMaker
- 通过 Jupyter Notebook 验证机器学习

BlueXP 连接器

在本次验证中，我们使用了 AWS。它也适用于 Azure 和 Google Cloud。要在 AWS 中创建 BlueXP 连接器，请完成以下步骤：

1. 我们使用了基于 BlueXP 中的 mcarl-marketplace-subscription 的凭证。

2. 选择适合您环境的区域（例如，us-east-1 [N. Virginia]），并选择身份验证方法（例如，Assume Role 或 AWS keys）。在此验证中，我们使用 AWS 密钥。
3. 提供连接器的名称并创建角色。
4. 根据您是否需要公共 IP，提供网络详细信息，例如 VPC、子网或密钥对。
5. 提供安全组的详细信息，例如来自源类型的 HTTP、HTTPS 或 SSH 访问，例如任何地方和 IP 范围信息。
6. 审查并创建BlueXP连接器。
7. 验证BlueXP EC2 实例状态是否在 AWS 控制台中运行，并从 **Networking** 选项卡中检查 IP 地址。
8. 从BlueXP门户登录连接器用户界面，或者您可以使用 IP 地址从浏览器访问。

NetApp Cloud Volumes ONTAP

要在BlueXP中创建Cloud Volumes ONTAP实例，请完成以下步骤：

1. 创建一个新的工作环境，选择云提供商，并选择Cloud Volumes ONTAP实例的类型（例如单 CVO、HA 或Amazon FSx ONTAP for ONTAP）。
2. 提供详细信息，例如Cloud Volumes ONTAP集群名称和凭据。在此验证中，我们创建了一个Cloud Volumes ONTAP svm_sagemaker_cvo_sn1。
3. 选择Cloud Volumes ONTAP所需的服务。在这次验证中，我们选择仅监控，因此我们禁用了*数据感知与合规性*和*备份到云服务*。
4. 在*位置和连接*部分中，选择 AWS 区域、VPC、子网、安全组、SSH 身份验证方法以及密码或密钥对。
5. 选择充电方式。我们使用*专业版*进行此验证。
6. 您可以选择预配置的包，例如*POC 和小型工作负载*、数据库和应用程序数据生产工作负载、经济高效的 **DR** 或 最高性能生产工作负载。在本次验证中，我们选择*Poc 和 Small Workloads*。
7. 创建具有特定大小、允许的协议和导出选项的卷。在此验证中，我们创建了一个名为 vol1。
8. 选择配置文件磁盘类型和分层策略。在本次验证中，我们禁用了*存储效率*和*通用 SSD - 动态性能*。
9. 最后，检查并创建Cloud Volumes ONTAP实例。然后等待 15-20 分钟让BlueXP创建Cloud Volumes ONTAP 工作环境。
10. 配置以下参数以启用 Duality 协议。从ONTAP 9 开始支持 Duality 协议 (NFS/S3)。 12.1 及更高版本。
 - a. 在此验证中，我们创建了一个名为 svm_sagemaker_cvo_sn1`和音量 `vol1。
 - b. 验证 SVM 是否支持 NFS 和 S3 协议。如果没有，请修改 SVM 以支持它们。

```

sagemaker_cvo_sn1::> vserver show -vserver svm_sagemaker_cvo_sn1
                                Vserver: svm_sagemaker_cvo_sn1
                                Vserver Type: data
                                Vserver Subtype: default
                                Vserver UUID: 911065dd-a8bc-11ed-bc24-
e1c0f00ad86b
                                Root Volume:
svm_sagemaker_cvo_sn1_root
                                Aggregate: aggr1
                                NIS Domain: -
                                Root Volume Security Style: unix
                                LDAP Client: -
                                Default Volume Language Code: C.UTF-8
                                Snapshot Policy: default
                                Data Services: data-cifs, data-
flexcache,
                                data-iscsi, data-nfs,
                                data-nvme-tcp
                                Comment:
                                Quota Policy: default
                                List of Aggregates Assigned: aggr1
                                Limit on Maximum Number of Volumes allowed: unlimited
                                Vserver Admin State: running
                                Vserver Operational State: running
                                Vserver Operational State Stopped Reason: -
                                Allowed Protocols: nfs, cifs, fcp, iscsi,
ndmp, s3
                                Disallowed Protocols: nvme
                                Is Vserver with Infinite Volume: false
                                QoS Policy Group: -
                                Caching Policy Name: -
                                Config Lock: false
                                IPspace Name: Default
                                Foreground Process: -
                                Logical Space Reporting: true
                                Logical Space Enforcement: false
                                Default Anti_ransomware State of the Vserver's Volumes: disabled
                                Enable Analytics on New Volumes: false
                                Enable Activity Tracking on New Volumes: false

sagemaker_cvo_sn1::>

```

11. 如果需要，创建并安装 CA 证书。

12. 创建服务数据策略。

```
sagemaker_cvo_sn1::*> network interface service-policy create -vserver
svm_sagemaker_cvo_sn1 -policy sagemaker_s3_nfs_policy -services data-
core,data-s3-server,data-nfs,data-flexcache
sagemaker_cvo_sn1::*> network interface create -vserver
svm_sagemaker_cvo_sn1 -lif svm_sagemaker_cvo_sn1_s3_lif -service-policy
sagemaker_s3_nfs_policy -home-node sagemaker_cvo_sn1-01 -address
172.30.10.41 -netmask 255.255.255.192
```

Warning: The configured failover-group has no valid failover targets for the LIF's failover-policy. To view the failover targets for a LIF, use the "network interface show -failover" command.

```
sagemaker_cvo_sn1::*>
```

```
sagemaker_cvo_sn1::*> network interface show
```

Logical Vserver Home	Status Interface	Network Admin/Oper	Current Address/Mask	Current Node	Is Port
sagemaker_cvo_sn1-01	cluster-mgmt	up/up	172.30.10.40/26	sagemaker_cvo_sn1-	e0a
true					
sagemaker_cvo_sn1-01	intercluster	up/up	172.30.10.48/26	sagemaker_cvo_sn1-	e0a
true					
sagemaker_cvo_sn1-01	sagemaker_cvo_sn1-01_mgmt1	up/up	172.30.10.58/26	sagemaker_cvo_sn1-	e0a
true					
svm_sagemaker_cvo_sn1-01	svm_sagemaker_cvo_sn1_data_lif	up/up	172.30.10.23/26	sagemaker_cvo_sn1-	e0a
true					
svm_sagemaker_cvo_sn1-01	svm_sagemaker_cvo_sn1_mgmt_lif	up/up	172.30.10.32/26	sagemaker_cvo_sn1-	e0a
true					
svm_sagemaker_cvo_sn1-01	svm_sagemaker_cvo_sn1_s3_lif	up/up	172.30.10.41/26	sagemaker_cvo_sn1-	

01

e0a

true

6 entries were displayed.

```
sagemaker_cvo_sn1::~*>
```

```
sagemaker_cvo_sn1::~*> vserver object-store-server create -vserver  
svm_sagemaker_cvo_sn1 -is-http-enabled true -object-store-server  
svm_sagemaker_cvo_s3_sn1 -is-https-enabled false  
sagemaker_cvo_sn1::~*> vserver object-store-server show
```

```
Vserver: svm_sagemaker_cvo_sn1
```

```
    Object Store Server Name: svm_sagemaker_cvo_s3_sn1
```

```
        Administrative State: up
```

```
            HTTP Enabled: true
```

```
    Listener Port For HTTP: 80
```

```
        HTTPS Enabled: false
```

```
    Secure Listener Port For HTTPS: 443
```

```
    Certificate for HTTPS Connections: -
```

```
        Default UNIX User: pcuser
```

```
    Default Windows User: -
```

```
        Comment:
```

```
sagemaker_cvo_sn1::~*>
```

13. 检查汇总详细信息。

```
sagemaker_cvo_sn1::*> aggr show
```

```
Aggregate      Size Available Used% State  #Vols  Nodes      RAID
Status
-----
-----
aggr0_sagemaker_cvo_sn1_01
      124.0GB   50.88GB   59% online    1 sagemaker_cvo_
raid0,
                                sn1-01
normal
aggr1      907.1GB   904.9GB   0% online    2 sagemaker_cvo_
raid0,
                                sn1-01
normal
2 entries were displayed.

sagemaker_cvo_sn1::*>
```

14. 创建用户和组。

```

sagemaker_cvo_sn1::*> vserver object-store-server user create -vserver
svm_sagemaker_cvo_sn1 -user s3user

sagemaker_cvo_sn1::*> vserver object-store-server user show
Vserver      User          ID          Access Key          Secret Key
-----
-----
svm_sagemaker_cvo_sn1
      root          0          -          -
      Comment: Root User
svm_sagemaker_cvo_sn1
      s3user        1          0ZNX21JW5Q8AP80CQ2E
PpLs4gA9K0_2gPhuykkp014gBjcC9Rbi3QDX_6rr
2 entries were displayed.

sagemaker_cvo_sn1::*>

sagemaker_cvo_sn1::*> vserver object-store-server group create -name
s3group -users s3user -comment ""

sagemaker_cvo_sn1::*>
sagemaker_cvo_sn1::*> vserver object-store-server group delete -gid 1
-vserver svm_sagemaker_cvo_sn1

sagemaker_cvo_sn1::*> vserver object-store-server group create -name
s3group -users s3user -comment "" -policies FullAccess

sagemaker_cvo_sn1::*>

```

15. 在 NFS 卷上创建一个存储桶。

```
sagemaker_cvo_sn1::~*> vservers object-store-server bucket create -bucket
ontapbucket1 -type nas -comment "" -vservers svm_sagemaker_cvo_sn1 -nas
-path /voll
sagemaker_cvo_sn1::~*> vservers object-store-server bucket show
Vserver      Bucket      Type      Volume      Size
Encryption  Role        NAS Path
-----
svm_sagemaker_cvo_sn1
                ontapbucket1    nas      voll        -        false
-            /voll
sagemaker_cvo_sn1::~*>
```

AWS SageMaker

要从 AWS SageMaker 创建 AWS Notebook，请完成以下步骤：

1. 确保创建 Notebook 实例的用户具有 AmazonSageMakerFullAccess IAM 策略或属于具有 AmazonSageMakerFullAccess 权限的现有组的一部分。在此验证中，用户是现有组的一部分。
2. 提供以下信息：
 - 笔记本实例名称。
 - 实例类型。
 - 平台标识符。
 - 选择具有 AmazonSageMakerFullAccess 权限的 IAM 角色。
 - 根访问 – 启用。
 - 加密密钥 - 选择无自定义加密。
 - 保留其余默认选项。
3. 本次验证中，SageMaker实例详情如下：

Amazon SageMaker > Notebook instances > nkarthiksagemaker

nkarthiksagemaker

Delete Stop Open Jupyter Open JupyterLab

Notebook instance settings Edit

Name	Status	Notebook instance type	Platform identifier
nkarthiksagemaker	✔ InService	ml.t2.medium	Amazon Linux 2, Jupyter Lab 3 (notebook-al2-v2)
ARN	Creation time	Elastic Inference	Minimum IMDS Version
arn:aws:sagemaker:us-east-1:210811600188:notebook-instance/nkarthiksagemaker	Feb 16, 2023 18:55 UTC	-	2
Lifecycle configuration	Last updated	Volume Size	
-	Mar 22, 2023 20:59 UTC	5GB EBS	

Permissions and encryption

IAM role ARN arn:aws:iam::210811600188:role/SageMakerFullRole	Root access Enabled	Encryption key
--	------------------------	----------------

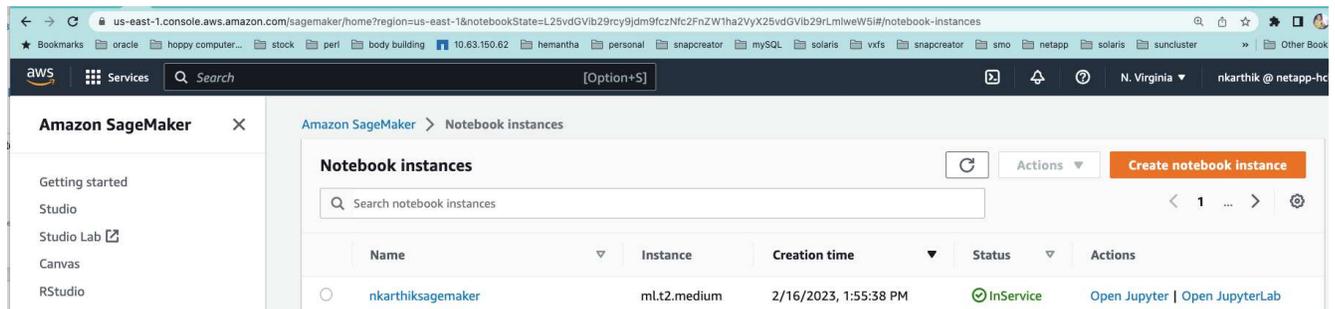
Network

Subnet(s)
[subnet-00f94558](#)

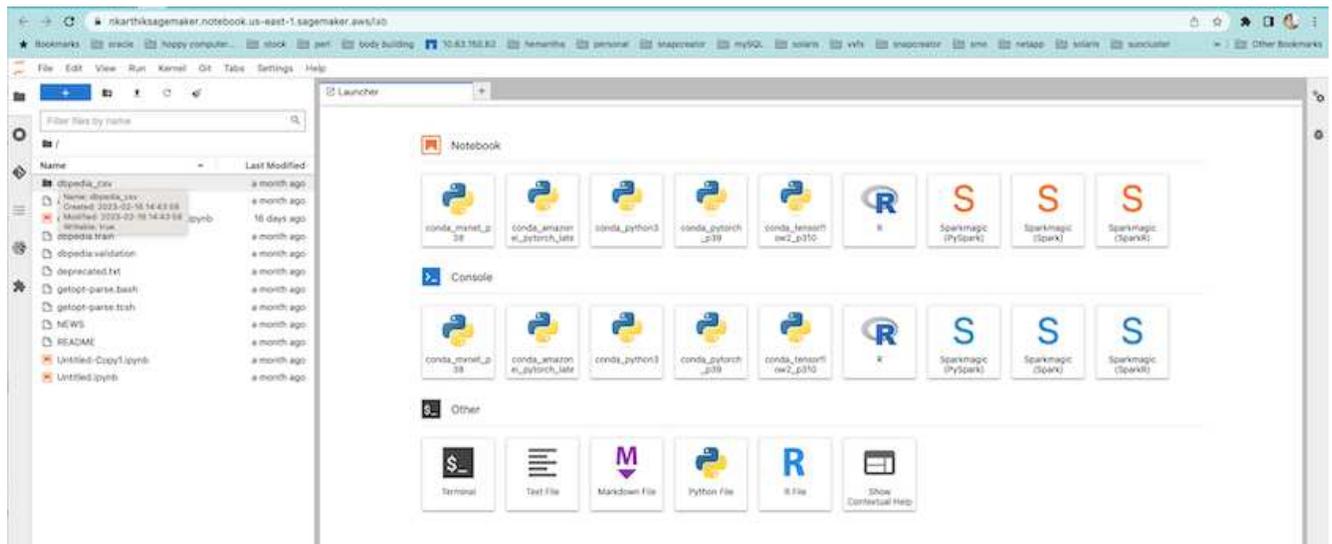
Security Group(s)
[sg-07111a8c16d67c81d](#)

Direct internet access
Enabled: [Learn more](#)

4. 启动 AWS Notebook。



5. 打开 Jupyter 实验室。



6. 登录终端并挂载Cloud Volumes ONTAP卷。

```
sh-4.2$ sudo mkdir /vol1; sudo mount -t nfs 172.30.10.41:/vol1 /vol1
sh-4.2$ df -h
```

Filesystem	Size	Used	Avail	Use%	Mounted on
devtmpfs	2.0G	0	2.0G	0%	/dev
tmpfs	2.0G	0	2.0G	0%	/dev/shm
tmpfs	2.0G	624K	2.0G	1%	/run
tmpfs	2.0G	0	2.0G	0%	/sys/fs/cgroup
/dev/xvda1	140G	114G	27G	82%	/
/dev/xvdf	4.8G	72K	4.6G	1%	/home/ec2-user/SageMaker
tmpfs	393M	0	393M	0%	/run/user/1001
tmpfs	393M	0	393M	0%	/run/user/1002
tmpfs	393M	0	393M	0%	/run/user/1000
172.30.10.41:/vol1	973M	189M	785M	20%	/vol1

```
sh-4.2$
```

7. 使用 AWS CLI 命令检查在Cloud Volumes ONTAP卷上创建的存储桶。

```
sh-4.2$ aws configure --profile netapp
AWS Access Key ID [None]: 0ZNAX21JW5Q8AP80CQ2E
AWS Secret Access Key [None]: PpLs4gA9K0_2gPhuykkp014gBjcC9Rbi3QDX_6rr
Default region name [None]: us-east-1
Default output format [None]:
sh-4.2$

sh-4.2$ aws s3 ls --profile netapp --endpoint-url
2023-02-10 17:59:48 ontapbucket1

sh-4.2$ aws s3 ls --profile netapp --endpoint-url s3://ontapbucket1/

2023-02-10 18:46:44          4747 1
2023-02-10 18:48:32          96 setup.cfg

sh-4.2$
```

机器学习数据

在这次验证中，我们使用了来自众包社区努力的 DBpedia 的数据集，从各种维基媒体项目创建的信息中提取结构化内容。

1. 从 DBpedia GitHub 位置下载数据并提取。使用与上一节相同的终端。

```
sh-4.2$ wget
--2023-02-14 23:12:11--
Resolving github.com (github.com)... 140.82.113.3
Connecting to github.com (github.com)|140.82.113.3|:443... connected.
HTTP request sent, awaiting response... 302 Found
Location: [following]
--2023-02-14 23:12:11--
Resolving raw.githubusercontent.com (raw.githubusercontent.com)...
185.199.109.133, 185.199.110.133, 185.199.111.133, ...
Connecting to raw.githubusercontent.com
(raw.githubusercontent.com)|185.199.109.133|:443... connected.
HTTP request sent, awaiting response... 200 OK
Length: 68431223 (65M) [application/octet-stream]
Saving to: 'dbpedia_csv.tar.gz'

100%[=====
=====
=====>] 68,431,223  56.2MB/s   in 1.2s

2023-02-14 23:12:13 (56.2 MB/s) - 'dbpedia_csv.tar.gz' saved
[68431223/68431223]

sh-4.2$ tar -zxvf dbpedia_csv.tar.gz
dbpedia_csv/
dbpedia_csv/test.csv
dbpedia_csv/classes.txt
dbpedia_csv/train.csv
dbpedia_csv/readme.txt
sh-4.2$
```

2. 将数据复制到Cloud Volumes ONTAP位置并使用 AWS CLI 从 S3 存储桶中进行检查。

```

sh-4.2$ df -h
Filesystem                Size      Used Avail Use% Mounted on
devtmpfs                  2.0G         0   2.0G   0% /dev
tmpfs                     2.0G         0   2.0G   0% /dev/shm
tmpfs                     2.0G   628K   2.0G   1% /run
tmpfs                     2.0G         0   2.0G   0% /sys/fs/cgroup
/dev/xvda1                140G   114G    27G  82% /
/dev/xvdf                  4.8G     52K   4.6G   1% /home/ec2-user/SageMaker
tmpfs                    393M         0   393M   0% /run/user/1002
tmpfs                    393M         0   393M   0% /run/user/1001
tmpfs                    393M         0   393M   0% /run/user/1000
172.30.10.41:/vol1       973M   384K   973M   1% /vol1
sh-4.2$ pwd
/home/ec2-user
sh-4.2$ cp -ra dbpedia_csv /vol1
sh-4.2$ aws s3 ls --profile netapp --endpoint-url s3://ontapbucket1/
                PRE dbpedia_csv/
2023-02-10 18:46:44          4747 1
2023-02-10 18:48:32           96 setup.cfg
sh-4.2$

```

3. 执行基本验证以确保读/写功能在 S3 存储桶上正常运行。

```

sh-4.2$ aws s3 cp --profile netapp --endpoint-url /usr/share/doc/util-
linux-2.30.2 s3://ontapbucket1/ --recursive
upload: ../../usr/share/doc/util-linux-2.30.2/deprecated.txt to
s3://ontapbucket1/deprecated.txt
upload: ../../usr/share/doc/util-linux-2.30.2/getopt-parse.bash to
s3://ontapbucket1/getopt-parse.bash
upload: ../../usr/share/doc/util-linux-2.30.2/README to
s3://ontapbucket1/README
upload: ../../usr/share/doc/util-linux-2.30.2/getopt-parse.tcsh to
s3://ontapbucket1/getopt-parse.tcsh
upload: ../../usr/share/doc/util-linux-2.30.2/AUTHORS to
s3://ontapbucket1/AUTHORS
upload: ../../usr/share/doc/util-linux-2.30.2/NEWS to
s3://ontapbucket1/NEWS
sh-4.2$ aws s3 ls --profile netapp --endpoint-url
s3://ontapbucket1/s3://ontapbucket1/

An error occurred (InternalError) when calling the ListObjectsV2
operation: We encountered an internal error. Please try again.
sh-4.2$ aws s3 ls --profile netapp --endpoint-url s3://ontapbucket1/
                PRE dbpedia_csv/

```

```

2023-02-16 19:19:27      26774 AUTHORS
2023-02-16 19:19:27      72727 NEWS
2023-02-16 19:19:27       4493 README
2023-02-16 19:19:27       2825 deprecated.txt
2023-02-16 19:19:27       1590 getopt-parse.bash
2023-02-16 19:19:27       2245 getopt-parse.tcsh
sh-4.2$ ls -ltr /voll
total 132
drwxrwxr-x 2 ec2-user ec2-user  4096 Mar 29  2015 dbpedia_csv
-rw-r--r-- 1 nobody  nobody   2245 Apr 10 17:37 getopt-parse.tcsh
-rw-r--r-- 1 nobody  nobody   2825 Apr 10 17:37 deprecated.txt
-rw-r--r-- 1 nobody  nobody   4493 Apr 10 17:37 README
-rw-r--r-- 1 nobody  nobody   1590 Apr 10 17:37 getopt-parse.bash
-rw-r--r-- 1 nobody  nobody  26774 Apr 10 17:37 AUTHORS
-rw-r--r-- 1 nobody  nobody  72727 Apr 10 17:37 NEWS
sh-4.2$ ls -ltr /voll/dbpedia_csv/
total 192104
-rw----- 1 ec2-user ec2-user 174148970 Mar 28  2015 train.csv
-rw----- 1 ec2-user ec2-user  21775285 Mar 28  2015 test.csv
-rw----- 1 ec2-user ec2-user    146 Mar 28  2015 classes.txt
-rw-rw-r-- 1 ec2-user ec2-user   1758 Mar 29  2015 readme.txt
sh-4.2$ chmod -R 777 /voll/dbpedia_csv
sh-4.2$ ls -ltr /voll/dbpedia_csv/
total 192104
-rwxrwxrwx 1 ec2-user ec2-user 174148970 Mar 28  2015 train.csv
-rwxrwxrwx 1 ec2-user ec2-user  21775285 Mar 28  2015 test.csv
-rwxrwxrwx 1 ec2-user ec2-user    146 Mar 28  2015 classes.txt
-rwxrwxrwx 1 ec2-user ec2-user   1758 Mar 29  2015 readme.txt
sh-4.2$ aws s3 cp --profile netapp --endpoint-url http://172.30.2.248/
s3://ontapbucket1/ /tmp --recursive
download: s3://ontapbucket1/AUTHORS to ../../tmp/AUTHORS
download: s3://ontapbucket1/README to ../../tmp/README
download: s3://ontapbucket1/NEWS to ../../tmp/NEWS
download: s3://ontapbucket1/dbpedia_csv/classes.txt to
../../tmp/dbpedia_csv/classes.txt
download: s3://ontapbucket1/dbpedia_csv/readme.txt to
../../tmp/dbpedia_csv/readme.txt
download: s3://ontapbucket1/deprecated.txt to ../../tmp/deprecated.txt
download: s3://ontapbucket1/getopt-parse.bash to ../../tmp/getopt-
parse.bash
download: s3://ontapbucket1/getopt-parse.tcsh to ../../tmp/getopt-
parse.tcsh
download: s3://ontapbucket1/dbpedia_csv/test.csv to
../../tmp/dbpedia_csv/test.csv
download: s3://ontapbucket1/dbpedia_csv/train.csv to
../../tmp/dbpedia_csv/train.csv

```

```
sh-4.2$  
sh-4.2$ aws s3 ls --profile netapp --endpoint-url s3://ontapbucket1/  
                PRE dbpedia_csv/  
2023-02-16 19:19:27      26774 AUTHORS  
2023-02-16 19:19:27      72727 NEWS  
2023-02-16 19:19:27      4493 README  
2023-02-16 19:19:27      2825 deprecated.txt  
2023-02-16 19:19:27      1590 getopt-parse.bash  
2023-02-16 19:19:27      2245 getopt-parse.tcsh  
sh-4.2$
```

通过 Jupyter Notebook 验证机器学习

以下验证通过使用以下 SageMaker BlazingText 示例通过文本分类提供机器学习构建、训练和部署模型：

1. 安装 boto3 和 SageMaker 包。

```
In [1]: pip install --upgrade boto3 sagemaker
```

输出：

```
Looking in indexes: https://pypi.org/simple,  
https://pip.repos.neuron.amazonaws.com  
Requirement already satisfied: boto3 in /home/ec2-  
user/anaconda3/envs/python3/lib/python3.10/site-packages (1.26.44)  
Collecting boto3  
  Downloading boto3-1.26.72-py3-none-any.whl (132 kB)  
-----  
132.7/132.7 kB 14.6 MB/s eta 0: 00:00  
Requirement already satisfied: sagemaker in /home/ec2-  
user/anaconda3/envs/python3/lib/python3.10/site-packages (2.127.0)  
Collecting sagemaker  
  Downloading sagemaker-2.132.0.tar.gz (668 kB)  
-----  
668.0/668.0 kB 12.3 MB/s eta 0:  
00:0000:01  
  Preparing metadata (setup.py) ... done  
Collecting botocore<1.30.0,>=1.29.72  
  Downloading botocore-1.29.72-py3-none-any.whl (10.4 MB)  
-----  
10.4/10.4 MB 44.3 MB/s eta 0: 00:0000:010:01  
Requirement already satisfied: s3transfer<0.7.0,>=0.6.0 in /home/ec2-  
user/anaconda3/envs/python3/lib/python3.10/site-packages (from boto3)  
(0.6.0)
```

Requirement already satisfied: jmespath<2.0.0,>=0.7.1 in /home/ec2-user/anaconda3/envs/python3/lib/python3.10/site-packages (from boto3) (0.10.0)

Requirement already satisfied: attrs<23,>=20.3.0 in /home/ec2-user/anaconda3/envs/python3/lib/python3.10/site-packages (from sagemaker) (22.1.0)

Requirement already satisfied: google-pasta in /home/ec2-user/anaconda3/envs/python3/lib/python3.10/site-packages (from sagemaker) (0.2.0)

Requirement already satisfied: numpy<2.0,>=1.9.0 in /home/ec2-user/anaconda3/envs/python3/lib/python3.10/site-packages (from sagemaker) (1.22.4)

Requirement already satisfied: protobuf<4.0,>=3.1 in /home/ec2-user/anaconda3/envs/python3/lib/python3.10/site-packages (from sagemaker) (3.20.3)

Requirement already satisfied: protobuf3-to-dict<1.0,>=0.1.5 in /home/ec2-user/anaconda3/envs/python3/lib/python3.10/site-packages (from sagemaker) (0.1.5)

Requirement already satisfied: smdebug_rulesconfig==1.0.1 in /home/ec2-user/anaconda3/envs/python3/lib/python3.10/site-packages (from sagemaker) (1.0.1)

Requirement already satisfied: importlib-metadata<5.0,>=1.4.0 in /home/ec2-user/anaconda3/envs/python3/lib/python3.10/site-packages (from sagemaker) (4.13.0)

Requirement already satisfied: packaging>=20.0 in /home/ec2-user/anaconda3/envs/python3/lib/python3.10/site-packages (from sagemaker) (21.3)

Requirement already satisfied: pandas in /home/ec2-user/anaconda3/envs/python3/lib/python3.10/site-packages (from sagemaker) (1.5.1)

Requirement already satisfied: pathos in /home/ec2-user/anaconda3/envs/python3/lib/python3.10/site-packages (from sagemaker) (0.3.0)

Requirement already satisfied: schema in /home/ec2-user/anaconda3/envs/python3/lib/python3.10/site-packages (from sagemaker) (0.7.5)

Requirement already satisfied: python-dateutil<3.0.0,>=2.1 in /home/ec2-user/anaconda3/envs/python3/lib/python3.10/site-packages (from botocore<1.30.0,>=1.29.72->boto3) (2.8.2)

Requirement already satisfied: urllib3<1.27,>=1.25.4 in /home/ec2-user/anaconda3/envs/python3/lib/python3.10/site-packages (from botocore<1.30.0,>=1.29.72->boto3) (1.26.8)

Requirement already satisfied: zipp>=0.5 in

```

/home/ec2-user/anaconda3/envs/python3/lib/python3.10/site-packages
(from importlib-metadata<5.0,>=1.4.0->sagemaker) (3.10.0)
Requirement already satisfied: pyparsing!=3.0.5,>=2.0.2 in /home/ec2-
user/anaconda3/envs/python3/lib/python3.10/site-packages (from
packaging>=20.0->sagemaker) (3.0.9)
Requirement already satisfied: six in /home/ec2-
user/anaconda3/envs/python
3/lib/python3.10/site-packages (from protobuf3-to-dict<1.0,>=0.1.5-
>sagemaker) (1.16.0)
Requirement already satisfied: pytz>=2020.1 in /home/ec2-
user/anaconda3/envs/python3/lib/python3.10/site-packages (from pandas-
>sagemaker) (2022.5)
Requirement already satisfied: ppft>=1.7.6.6 in /home/ec2-
user/anaconda3/envs/python3/lib/python3.10/site-packages (from pathos-
>sagemaker) (1.7.6.6) Requirement already satisfied:
multiprocess>=0.70.14 in /home/ec2-user/anac
onda3/envs/python3/lib/python3.10/site-packages (from pathos->sagemaker)
(0.70.14)
Requirement already satisfied: dill>=0.3.6 in /home/ec2-
user/anaconda3/envs/python3/lib/python3.10/site-packages (from pathos-
>sagemaker) (0.3.6)
Requirement already satisfied: pox>=0.3.2 in /home/ec2-
user/anaconda3/envs/python3/lib/python3.10/site-packages (from pathos-
>sagemaker) (0.3.2) Requirement already satisfied: contextlib2>=0.5.5 in
/home/ec2-user/anaconda3/envs/python3/lib/python3.10/site-packages
(from schema->sagemaker) (21.
6.0) Building wheels for collected packages: sagemaker
  Building wheel for sagemaker (setup.py) ... done
  Created wheel for sagemaker: filename=sagemaker-2.132.0-py2.py3-none-
any.whl size=905449
sha256=f6100a5dc95627f2e2a49824e38f0481459a27805ee19b5a06ec
83db0252fd41
  Stored in directory: /home/ec2-
user/.cache/pip/wheels/60/41/b6/482e7ab096
520df034fbf2d44a1d7ba0681b27ef45aa61
Successfully built sagemaker
Installing collected packages: botocore, boto3, sagemaker
  Attempting uninstall: botocore      Found existing installation:
botocore 1.24.19
    Uninstalling botocore-1.24.19:      Successfully uninstalled
botocore-1.24.19
  Attempting uninstall: boto3      Found existing installation: boto3
1.26.44
    Uninstalling boto3-1.26.44:
      Successfully uninstalled boto3-1.26.44
  Attempting uninstall: sagemaker      Found existing installation:

```

```
sagemaker 2.127.0
```

```
Uninstalling sagemaker-2.127.0:
```

```
Successfully uninstalled sagemaker-2.127.0
```

```
ERROR: pip's dependency resolver does not currently take into account  
all the packages that are installed. This behaviour is the source of  
the following dependency conflicts.
```

```
awscli 1.27.44 requires botocore==1.29.44, but you have botocore 1.29.72  
which is incompatible.
```

```
aiobotocore 2.0.1 requires botocore<1.22.9,>=1.22.8, but you have  
botocore 1.29.72 which is incompatible. Successfully installed boto3-
```

```
1.26.72 botocore-1.29.72 sagemaker-2.132.0 Note: you may need to restart  
the kernel to use updated packages.
```

2. 在下一步中，数据(dbpedia_csv)从 s3 bucket 下载 `ontapbucket1` 到机器学习中使用的 Jupyter Notebook 实例。

```

In [2]: import sagemaker
In [3]: from sagemaker import get_execution_role
In [4]:
import json
import boto3
sess = sagemaker.Session()
role = get_execution_role()
print(role)
bucket = "ontapbucket1"
print(bucket)
sess.s3_client = boto3.client('s3',region_name='',aws_access_key_id =
'0ZNAX21JW5Q8AP80CQ2E', aws_secret_access_key =
'PpLs4gA9K0_2gPhuykkp014gBjcC9Rbi3QDX_6rr',
                                use_ssl = False, endpoint_url =
'http://172.30.10.41',

config=boto3.session.Config(signature_version='s3v4',
s3={'addressing_style':'path'}) )
sess.s3_resource = boto3.resource('s3',region_name='',aws_access_key_id
= '0ZNAX21JW5Q8AP80CQ2E', aws_secret_access_key =
'PpLs4gA9K0_2gPhuykkp014gBjcC9Rbi3QDX_6rr',
                                use_ssl = False, endpoint_url =
'http://172.30.10.41',

config=boto3.session.Config(signature_version='s3v4',
s3={'addressing_style':'path'}) )
prefix = "blazingtext/supervised"
import os
my_bucket = sess.s3_resource.Bucket(bucket)
my_bucket = sess.s3_resource.Bucket(bucket)
#os.mkdir('dbpedia_csv')
for s3_object in my_bucket.objects.all():
    filename = s3_object.key
    # print(filename)
    # print(s3_object.key)
    my_bucket.download_file(s3_object.key, filename)

```

3. 以下代码创建从整数索引到类标签的映射，用于在推理期间检索实际的类名。

```

index_to_label = {}
with open("dbpedia_csv/classes.txt") as f:
    for i,label in enumerate(f.readlines()):
        index_to_label[str(i + 1)] = label.strip()

```

输出列出了 `ontapbucket1` 存储桶用作 AWS SageMaker 机器学习验证的数据。

```
arn:aws:iam::210811600188:role/SageMakerFullRole ontapbucket1
AUTHORS
AUTHORS
NEWS
NEWS
README README
dbpedia_csv/classes.txt dbpedia_csv/classes.txt dbpedia_csv/readme.txt
dbpedia_csv/readme.txt dbpedia_csv/test.csv dbpedia_csv/test.csv
dbpedia_csv/train.csv dbpedia_csv/train.csv deprecated.txt
deprecated.txt getopt-parse.bash getopt-parse.bash getopt-parse.tcsh
getopt-parse.tcsh
In [5]: ls
AUTHORS          deprecated.txt    getopt-parse.tcsh NEWS
Untitled.ipynb dbpedia_csv/    getopt-parse.bash lost+found/
README
In [6]: ls -l dbpedia_csv
total 191344
-rw-rw-r-- 1 ec2-user ec2-user      146 Feb 16 19:43 classes.txt
-rw-rw-r-- 1 ec2-user ec2-user     1758 Feb 16 19:43 readme.txt
-rw-rw-r-- 1 ec2-user ec2-user  21775285 Feb 16 19:43 test.csv
-rw-rw-r-- 1 ec2-user ec2-user 174148970 Feb 16 19:43 train.csv
```

4. 开始数据预处理阶段，将训练数据预处理为空格分隔的标记化文本格式，BlazingText 算法和 nltk 库可以使用该格式对来自 DBPedia 数据集的输入句子进行标记化。下载 nltk 标记器和其他库。这 `transform_instance` 并行应用于每个数据实例使用 Python 多处理模块。

```
In [7]: from random import shuffle
import multiprocessing
from multiprocessing import Pool
import csv
import nltk
nltk.download("punkt")
def transform_instance(row):
    cur_row = []
    label = "__label__" + index_to_label [row[0]] # Prefix the index-ed
label with __label__
    cur_row.append (label)
    cur_row.extend(nltk.word_tokenize(row[1].lower ()))
    cur_row.extend(nltk.word_tokenize(row[2].lower ()))
    return cur_row
def preprocess(input_file, output_file, keep=1):
    all_rows = []
    with open(input_file,"r") as csvinfile:
```

```

        csv_reader = csv.reader(csvinfile, delimiter=",")
        for row in csv_reader:
            all_rows.append(row)
    shuffle(all_rows)
    all_rows = all_rows[: int(keep * len(all_rows))]
    pool = Pool(processes=multiprocessing.cpu_count())
    transformed_rows = pool.map(transform_instance, all_rows)
    pool.close()
    pool.join()
    with open(output_file, "w") as csvoutfile:
        csv_writer = csv.writer (csvoutfile, delimiter=" ",
lineterminator="\n")
        csv_writer.writerows (transformed_rows)

# Preparing the training dataset
# since preprocessing the whole dataset might take a couple of minutes,
# we keep 20% of the training dataset for this demo.
# Set keep to 1 if you want to use the complete dataset
preprocess("dbpedia_csv/train.csv","dbpedia.train", keep=0.2)
# Preparing the validation dataset
preprocess("dbpedia_csv/test.csv","dbpedia.validation")
sess = sagemaker.Session()
role = get_execution_role()
print (role) # This is the role that sageMaker would use to leverage Aws
resources (S3, Cloudwatch) on your behalf
bucket = sess.default_bucket() # Replace with your own bucket name if
needed
print("default Bucket::: ")
print(bucket)

```

输出:

```

[nltk_data] Downloading package punkt to /home/ec2-user/nltk_data...
[nltk_data]   Package punkt is already up-to-date!
arn:aws:iam::210811600188:role/SageMakerFullRole default Bucket:::
sagemaker-us-east-1-210811600188

```

5. 将格式化和训练数据集上传到 S3，以便 SageMaker 可以使用它来执行训练作业。然后使用 Python SDK 将两个文件上传到存储桶和前缀位置。

```
In [8]: %%time
train_channel = prefix + "/train"
validation_channel = prefix + "/validation"
sess.upload_data(path="dbpedia.train", bucket=bucket,
key_prefix=train_channel)
sess.upload_data(path="dbpedia.validation", bucket=bucket,
key_prefix=validation_channel)
s3_train_data = "s3://{}/{}".format(bucket, train_channel)
s3_validation_data = "s3://{}/{}".format(bucket, validation_channel)
```

输出:

```
CPU times: user 546 ms, sys: 163 ms, total: 709 ms
Wall time: 1.32 s
```

6. 在加载模型工件的 S3 处设置输出位置，以便工件可以作为算法训练作业的输出。创建一个 `sageMaker.estimator.Estimator` 对象来启动训练工作。

```
In [9]: s3_output_location = "s3://{}/{}/output".format(bucket, prefix)
In [10]: region_name = boto3.Session().region_name
In [11]: container =
sagemaker.amazon.amazon_estimator.get_image_uri(region_name,
"blazingtext", "latest")
print("Using SageMaker BlazingText container: {} ({})" .format(container,
region_name))
```

输出:

```
The method get_image_uri has been renamed in sagemaker>=2.
See: https://sagemaker.readthedocs.io/en/stable/v2.html for details.
Defaulting to the only supported framework/algorithm version: 1.
Ignoring f ramework/algorithm version: latest.
Using SageMaker BlazingText container: 811284229777.dkr.ecr.us-east-
1.amazo naws.com/blazingtext:1 (us-east-1)
```

7. 定义 SageMaker `Estimator` 使用资源配置和超参数在 c4.4xlarge 实例上使用监督模式在 DBPedia 数据集上训练文本分类。

```

In [12]: bt_model = sagemaker.estimator.Estimator(
    container,
    role,
    instance_count=1,
    instance_type="ml.c4.4xlarge",
    volume_size=30,
    max_run=360000,
    input_mode="File",
    output_path=s3_output_location,
    hyperparameters={
        "mode": "supervised",
        "epochs": 1,
        "min_count": 2,
        "learning_rate": 0.05,
        "vector_dim": 10,
        "early_stopping": True,
        "patience": 4,
        "min_epochs": 5,
        "word_ngrams": 2,
    },
)

```

8. 准备数据通道和算法之间的握手。为此，请创建 `sagemaker.session.s3_input` 来自数据通道的对象，并将它们保存在字典中以供算法使用。

```

In [13]: train_data = sagemaker.inputs.TrainingInput(
    s3_train_data,
    distribution="FullyReplicated",
    content_type="text/plain",
    s3_data_type="S3Prefix",
)
validation_data = sagemaker.inputs.TrainingInput(
    s3_validation_data,
    distribution="FullyReplicated",
    content_type="text/plain",
    s3_data_type="S3Prefix",
)
data_channels = {"train": train_data, "validation": validation_data}

```

9. 作业完成后，将出现“作业完成”消息。训练好的模型可以在设置为 `output_path` 在估算器中。

```

In [14]: bt_model.fit(inputs=data_channels, logs=True)

```

输出:

```
INFO:sagemaker:Creating training-job with name: blazingtext-2023-02-16-20-3
7-30-748
2023-02-16 20:37:30 Starting - Starting the training job.....
2023-02-16 20:38:09 Starting - Preparing the instances for
training.....
2023-02-16 20:39:24 Downloading - Downloading input data
2023-02-16 20:39:24 Training - Training image download completed.
Training in progress... Arguments: train
[02/16/2023 20:39:41 WARNING 140279908747072] Loggers have already been
set up. [02/16/2023 20:39:41 WARNING 140279908747072] Loggers have
already been set up.
[02/16/2023 20:39:41 INFO 140279908747072] nvidia-smi took:
0.0251793861389
16016 secs to identify 0 gpus
[02/16/2023 20:39:41 INFO 140279908747072] Running single machine CPU
BlazingText training using supervised mode.
Number of CPU sockets found in instance is 1
[02/16/2023 20:39:41 INFO 140279908747072] Processing
/opt/ml/input/data/train/dbpedia.train . File size: 35.0693244934082 MB
[02/16/2023 20:39:41 INFO 140279908747072] Processing
/opt/ml/input/data/validation/dbpedia.validation . File size:
21.887572288513184 MB
Read 6M words
Number of words: 149301
Loading validation data from
/opt/ml/input/data/validation/dbpedia.validation
Loaded validation data.
----- End of epoch: 1 ##### Alpha: 0.0000 Progress: 100.00%
Million Words/sec: 10.39 ##### Training finished.
Average throughput in Million words/sec: 10.39
Total training time in seconds: 0.60
#train_accuracy: 0.7223
Number of train examples: 112000
#validation_accuracy: 0.7205
Number of validation examples: 70000
2023-02-16 20:39:55 Uploading - Uploading generated training model
2023-02-16 20:40:11 Completed - Training job completed
Training seconds: 68
Billable seconds: 68
```

10. 训练完成后，将训练好的模型部署为 Amazon SageMaker 实时托管终端节点以进行预测。

```
In [15]: from sagemaker.serializers import JSONSerializer
        text_classifier = bt_model.deploy(
            initial_instance_count=1, instance_type="ml.m4.xlarge",
            serializer=JSONS
        )
```

输出:

```
INFO:sagemaker:Creating model with name: blazingtext-2023-02-16-20-41-33-10
0
INFO:sagemaker:Creating endpoint-config with name blazingtext-2023-02-16-20-41-33-100
INFO:sagemaker:Creating endpoint with name blazingtext-2023-02-16-20-41-33-100
-----!
```

```
In [16]: sentences = [
        "Convair was an american aircraft manufacturing company which later expanded into rockets and spacecraft.",
        "Berwick secondary college is situated in the outer melbourne metropolitan suburb of berwick .",
    ]
# using the same nltk tokenizer that we used during data preparation for training
tokenized_sentences = [" ".join(nltk.word_tokenize(sent)) for sent in sentences]
payload = {"instances": tokenized_sentences} response = text_classifier.predict(payload)
predictions = json.loads(response)
print(json.dumps(predictions, indent=2))
```

```
[
  {
    "label": [
      "__label__Artist"
    ],
    "prob": [
      0.4090951681137085
    ]
  },
  {
    "label": [
      "__label__EducationalInstitution"
    ],
    "prob": [
      0.49466073513031006
    ]
  }
]
```

11. 默认情况下，模型返回一个概率最高的预测。检索顶部 `k` 预测，设置 `k` 在配置文件中。

```
In [17]: payload = {"instances": tokenized_sentences, "configuration":
{"k": 2}}
response = text_classifier.predict(payload)

predictions = json.loads(response)
print(json.dumps(predictions, indent=2))
```

```
[
  {
    "label": [
      "__label__Artist",
      "__label__MeanOfTransportation"
    ],
    "prob": [
      0.4090951681137085,
      0.26930734515190125
    ]
  },
  {
    "label": [
      "__label__EducationalInstitution",
      "__label__Building"
    ],
    "prob": [
      0.49466073513031006,
      0.15817692875862122
    ]
  }
]
```

12. 关闭笔记本之前删除端点。

```
In [18]: sess.delete_endpoint(text_classifier.endpoint)
WARNING:sagemaker.deprecations:The endpoint attribute has been renamed
in sagemaker>=2.
See: https://sagemaker.readthedocs.io/en/stable/v2.html for details.
INFO:sagemaker:Deleting endpoint with name: blazingtext-2023-02-16-20-
41-33
-100
```

结束语

基于此验证，数据科学家和工程师可以通过NetApp Cloud Volumes ONTAP的 S3 存储桶访问来自 AWS SageMaker Jupyter Notebooks 的 NFS 数据。这种方法可以轻松访问和共享来自 NFS 和 S3 的相同数据，而无需额外的软件。

在哪里可以找到更多信息

要了解有关本文档中描述的信息的更多信息，请查看以下文档和/或网站：

- 使用 SageMaker BlazingText 进行文本分类
- ONTAP版本对 S3 对象存储的支持

["https://docs.netapp.com/us-en/ontap/s3-config/ontap-version-support-s3-concept.html"](https://docs.netapp.com/us-en/ontap/s3-config/ontap-version-support-s3-concept.html)

版权信息

版权所有 © 2025 NetApp, Inc.。保留所有权利。中国印刷。未经版权所有者事先书面许可，本档中受版权保护的任何部分不得以任何形式或通过任何手段（图片、电子或机械方式，包括影印、录音、录像或存储在电子检索系统中）进行复制。

从受版权保护的 NetApp 资料派生的软件受以下许可和免责声明的约束：

本软件由 NetApp 按“原样”提供，不含任何明示或暗示担保，包括但不限于适销性以及针对特定用途的适用性的隐含担保，特此声明不承担任何责任。在任何情况下，对于因使用本软件而以任何方式造成的任何直接性、间接性、偶然性、特殊性、惩罚性或后果性损失（包括但不限于购买替代商品或服务；使用、数据或利润方面的损失；或者业务中断），无论原因如何以及基于何种责任理论，无论出于合同、严格责任或侵权行为（包括疏忽或其他行为），NetApp 均不承担责任，即使已被告知存在上述损失的可能性。

NetApp 保留在不另行通知的情况下随时对本文档所述的任何产品进行更改的权利。除非 NetApp 以书面形式明确同意，否则 NetApp 不承担因使用本文档所述产品而产生的任何责任或义务。使用或购买本产品不表示获得 NetApp 的任何专利权、商标权或任何其他知识产权许可。

本手册中描述的产品可能受一项或多项美国专利、外国专利或正在申请的专利的保护。

有限权利说明：政府使用、复制或公开本文档受 DFARS 252.227-7013（2014 年 2 月）和 FAR 52.227-19（2007 年 12 月）中“技术数据权利 — 非商用”条款第 (b)(3) 条规定的限制条件的约束。

本文档中所含数据与商业产品和/或商业服务（定义见 FAR 2.101）相关，属于 NetApp, Inc. 的专有信息。根据本协议提供的所有 NetApp 技术数据和计算机软件具有商业性质，并完全由私人出资开发。美国政府对这些数据的使用权具有非排他性、全球性、受限且不可撤销的许可，该许可既不可转让，也不可再许可，但仅限在与交付数据所依据的美国政府合同有关且受合同支持的情况下使用。除本文档规定的情形外，未经 NetApp, Inc. 事先书面批准，不得使用、披露、复制、修改、操作或显示这些数据。美国政府对国防部的授权仅限于 DFARS 的第 252.227-7015(b)（2014 年 2 月）条款中明确的权利。

商标信息

NetApp、NetApp 标识和 <http://www.netapp.com/TM> 上所列的商标是 NetApp, Inc. 的商标。其他公司和产品名称可能是其各自所有者的商标。