

Roll No: 45/A

Exam Seat No:

**VIVEKANAND EDUCATION SOCIETY'S INSTITUTE OF
TECHNOLOGY**

Hashu Advani Memorial Complex, Collector's Colony, R. C.
Marg, Chembur, Mumbai – 400074. Contact No. 02261532532



Since 1962

CERTIFICATE

Certified that Mr. Pushkar Prasad Sane of FYMCA-A has satisfactorily completed a course of the necessary experiments in AI Development Tool under my supervision in the Institute of Technology in the academic year 2023 -2024.

Principal

Head of Department
(Dr. Shivkumar Goel)

External Examiner

Subject Teacher
(Dr. Dhanamma Jagli)



INDEX

SR. NO	CONTENTS	DATE OF PREPARATI ON	DATE OF SUBMISSION	MARKS	FACULTY SIGN
1.	Assignment 1 1. What is AWS ? 2. What is Cloud Computing ? Who uses Cloud Computing ? 3. What are benefits of Cloud Computing ?	25-01-2024	31-01-2024		
2.	Practical 1.1 Aim- To set up various components of AI ML on AWS from IAM user.	25-01-2024	31-01-2024		
3.	Practical 1.2 Aim :- To create a free S3 and EC2 instance on AWS	09-02-2024	14-02-2024		
4.	Practical 2.1 Aim :- To study various components Amazon Sagemaker	15-02-2024	22-02-2024		
5.	Practical 2.2 Aim:- To study various components of AWS studio and implement model in Canvas	21-03-2024	21-03-2024		
6.	Practical 3 Aim:- To implement Data Labeling in AWS Sagemaker GroundTruth	21-03-2024	25-03-2024		

7.	<p>Practical 4</p> <p>Aim:- To implement EDA using Data Wrangler in AWS</p>	21-03-2024	28-03-2024		
8.	<p>Practical 5</p> <p>Aim:- To implement sagemaker XGBoost algorithm</p>	21-03-2024	30-03-2024		

Name of Student: Pushkar Sane		
Roll Number: 45		Lab Assignment Number: 0
Title of Lab Assignment: Study of Azure Cloud, AWS and Google Cloud.		
DOP: 31-01-2024		DOS: 01-02-2024
CO Mapped: CO1	PO Mapped: PO1, PO2, PO3, PSO1, PSO2	Signature:

Practical No. 0

Aim: Study of Azure Cloud, AWS and Google Cloud.

Description:

1. Azure Cloud

Azure Cloud, commonly referred to as Microsoft Azure, is a cloud computing platform and service provided by Microsoft. Launched in 2010, Azure offers a wide range of cloud services, including computing power, storage solutions, networking capabilities, databases, analytics, artificial intelligence, and more. It enables individuals and organizations to build, deploy, and manage applications and services through Microsoft's global network of data centers.

Key Features and Components of Azure:

- a. Compute Services:
 - i. Azure Virtual Machines (VMs): Allows users to deploy and run Windows or Linux virtual machines in the cloud.
 - ii. Azure App Service: A fully managed platform for building, deploying, and scaling web apps.
- b. Storage Services:
 - i. Azure Blob Storage: Object storage service for unstructured data.
 - ii. Azure Table Storage: NoSQL key-value store for semi-structured data.
 - iii. Azure File Storage: Fully managed file shares in the cloud.
- c. Database Services:
 - i. Azure SQL Database: A fully managed relational database service.
 - ii. Azure Cosmos DB: A globally distributed, multi-model database for various types of data.
- d. Networking Services:
 - i. Azure Virtual Network: Allows users to create private, isolated networks in the cloud.
 - ii. Azure Load Balancer: Distributes incoming network traffic across multiple servers.

- e. Identity and Access Management:
 - i. Azure Active Directory (AD): Provides identity and access management services for applications and services.
- f. AI and Machine Learning:
 - i. Azure Machine Learning: Enables building, training, and deploying machine learning models.
- g. Analytics and Big Data:
 - i. Azure Synapse Analytics (formerly SQL Data Warehouse): An analytics service for large volumes of data.
 - ii. Azure Data Lake Storage: Scalable and secure data lake for big data analytics.
- h. Internet of Things (IoT):
 - i. Azure IoT Hub: Connects, monitors, and manages IoT assets.
- i. Security and Compliance:
 - i. Azure Security Center: Provides advanced threat protection across hybrid cloud workloads.
 - ii. Azure Policy: Helps enforce organizational standards and compliance.
- j. DevOps and Developer Tools:
 - i. Azure DevOps Services: A set of development tools for planning, tracking, and delivering software.

2. Amazon Web Services:

Amazon Web Services (AWS) is a comprehensive and widely adopted cloud computing platform provided by Amazon. It offers a vast array of cloud services, allowing individuals, businesses, and organizations to access computing power, storage, databases, machine learning, analytics, and other resources over the internet. AWS is known for its flexibility, scalability, and reliability, making it a popular choice for various applications and workloads.

Here are some key aspects of AWS:

- a. Service Offering:
 - i. AWS provides a broad range of services, including computing power (Amazon EC2), storage (Amazon S3), databases (Amazon RDS), machine learning (Amazon SageMaker), serverless computing (AWS Lambda), and more.
- b. Global Infrastructure:
 - i. AWS operates a global network of data centers, referred to as Availability Zones, located in different regions around the world. This infrastructure allows users to deploy applications and services close to their end-users, reducing latency and improving performance.
- c. Pricing Model:
 - i. AWS follows a pay-as-you-go pricing model, where users are billed based on their actual usage of resources. This flexibility is beneficial for businesses as they only pay for the resources they consume.
- d. Security and Compliance:
 - i. AWS places a strong emphasis on security, offering various security services and features. It has a robust compliance program, meeting various industry-specific standards and certifications.
- e. Ecosystem and Marketplace:
 - i. AWS has a vast ecosystem of partners, third-party tools, and a marketplace where users can find and deploy additional services, solutions, and applications to enhance their cloud environment.

- f. Community and Support:
 - i. AWS has a large and active community of developers, architects, and businesses. It provides extensive documentation, training resources, and customer support to help users make the most of the platform.
- g. Innovation and Evolution:
 - i. AWS continually introduces new services and features, staying at the forefront of cloud innovation. This allows users to leverage cutting-edge technologies and stay competitive in their respective industries

3. Google Cloud:

Google Cloud Platform (GCP) is a suite of cloud computing services offered by Google. It provides a wide range of infrastructure and platform services for building, deploying, and scaling applications in the cloud. GCP is designed to offer powerful and flexible solutions for businesses, developers, and data scientists. Here are some key aspects of Google Cloud:

Key Services and Features:

- a. Compute Services:
 - i. Google Compute Engine (GCE): Infrastructure as a Service (IaaS) that allows users to run virtual machines in the Google Cloud.
- b. Storage Services:
 - i. Google Cloud Storage: Object storage service designed for secure and scalable storage of data.
- c. Databases:
 - i. Google Cloud SQL: Fully managed relational database service.
 - ii. Google Cloud Firestore: NoSQL document database for mobile and web application development.
- d. Big Data and Analytics:
 - i. BigQuery: Serverless, highly scalable, and cost-effective multi-cloud data warehouse for analytics.
 - ii. Dataproc: Fully managed Apache Spark and Hadoop service for big data processing.

- e. Machine Learning and AI:
 - i. AI Platform: End-to-end platform for building, testing, and deploying machine learning models.
 - ii. TensorFlow: An open-source machine learning framework developed by the Google Brain team.
- f. Networking:
 - i. Virtual Private Cloud (VPC): Provides networking functionality for GCP resources.
 - ii. Cloud Load Balancing: Distributes incoming network traffic across multiple instances.
- g. Security and Identity:
 - i. Identity and Access Management (IAM): Manages access control for GCP resources.
 - ii. Cloud Identity-Aware Proxy (IAP): Controls access to cloud applications.
- h. Serverless Computing:
 - i. Cloud Functions: Executes event-driven functions without the need to provision or manage servers.
- i. Containers and Kubernetes:
 - i. Google Kubernetes Engine (GKE): Managed Kubernetes service for deploying, managing, and scaling containerized applications.
- j. Internet of Things (IoT):
 - i. Cloud IoT Core: Fully managed service for connecting, managing, and ingesting data from IoT devices.

Unique Aspects:

- a. Data Centers and Global Network:
 - i. Google Cloud operates a global network of data centers to provide low-latency access to services.
- b. Emphasis on Data Analytics and Machine Learning:
 - i. GCP is known for its strengths in data analytics, machine learning, and artificial intelligence services.

- c. Containerization and Kubernetes:
 - i. Google has been a pioneer in containerization, and GCP offers strong support for Kubernetes.
- d. Integration with Google Workspace:
 - i. Integration with Google Workspace (formerly G Suite) for collaboration and productivity tools.
- e. Open Source Contributions:
 - i. Google has a history of contributing to open-source projects, and many of its cloud services are built on open-source technologies.

Azure, AWS, and Google Cloud Platform (GCP) are all cloud service providers, but they differ in various aspects, including their services, architecture, pricing models, and target audiences. Choosing between Azure, AWS, and GCP often depends on specific organizational requirements, existing technology stacks, and individual preferences. Many organizations adopt a multi-cloud strategy to leverage the strengths of each provider for different aspects of their infrastructure and applications.

Name of Student: Pushkar Sane		
Roll Number: 45	Lab Assignment Number: 1	
Title of Lab Assignment: Introduction to AWS signup and Service Legal Agreement, Introduction to AWS sagemaker (Signup, freetier, billing & Iam), AIML components in AWS, Simple Storage Service (S3) and Elastic Computing Cloud (EC2) Introduction.		
DOP: 25-01-2024	DOS: 06-02-2024	
CO Mapped: CO1	PO Mapped: PO1, PO2, PO3, PSO1, PSO2	Signature:

Practical No. 1

Aim: Introduction to AWS signup and Service Legal Agreement, Introduction to AWS sagemaker (Signup, freetier, billing & iam), AIML components in AWS, Simple Storage Service (S3) and Elastic Computing Cloud (EC2) Introduction.

Description:

Introduction to AWS Sign Up and Service Legal Agreement: When signing up for Amazon Web Services (AWS), users embark on a streamlined process that involves creating an AWS account. This account provides access to a comprehensive suite of cloud services. During the signup, users must agree to AWS's Service Legal Agreement, which outlines terms and conditions governing the use of AWS services, ensuring a clear understanding of responsibilities and obligations.

Here's a step-by-step guide to help you with the AWS sign-up process:

1. Visit the AWS Website:

Open your web browser and go to the official AWS website at <https://aws.amazon.com/>

2. Choose "Sign in to the Console":

On the AWS homepage, click on the "Sign in to the Console" button located in the top right corner.

3. Create an AWS Account:

If you don't already have an AWS account, click on the "Create a new AWS account" button. If you have an existing AWS account, you can sign in with your credentials.

4. Provide Your Email Address:

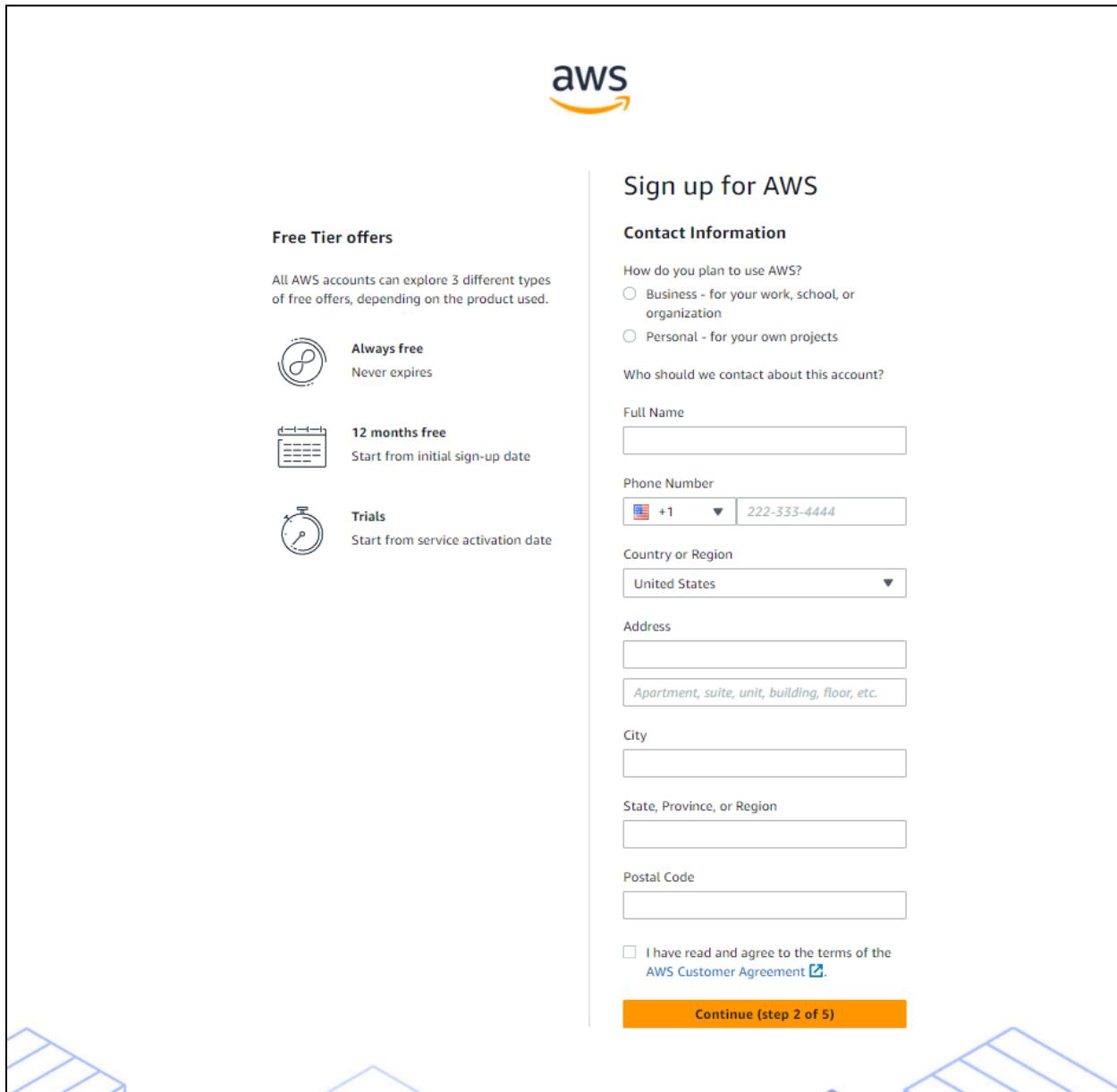
Enter the email address you want to associate with your AWS account. Make sure to use a valid email address as AWS will send important notifications and account-related information to this address.

5. Set up a Password:

Choose a secure password for your AWS account. AWS has specific requirements for passwords, so ensure that your chosen password meets the criteria.

6. Enter Your Account Information:

Fill in your account information, including your account name (which will be used as your AWS account alias), full name, and contact information.



The image shows the 'Sign up for AWS' page. At the top right, the AWS logo is displayed. Below it, the heading 'Sign up for AWS' is centered. To the left, there is a section titled 'Free Tier offers' with three options: 'Always free' (represented by a circular icon with a infinity symbol), '12 months free' (represented by a calendar icon), and 'Trials' (represented by a stopwatch icon). The 'Always free' option is selected. To the right, the 'Contact Information' section begins. It includes fields for 'Full Name' (a text input field), 'Phone Number' (a dropdown with country code +1 and number 222-333-4444), 'Country or Region' (a dropdown set to 'United States'), 'Address' (a text input field), 'City' (a text input field), 'State, Province, or Region' (a text input field), and 'Postal Code' (a text input field). At the bottom of this section is a checkbox labeled 'I have read and agree to the terms of the AWS Customer Agreement' with a link to the agreement. A large orange button at the bottom right says 'Continue (step 2 of 5)'.

7. Payment Information:

Provide your payment information. AWS requires valid payment details even if you plan to use free-tier services. This is to verify your identity and prevent misuse of AWS resources.

8. Identity Verification:

AWS may perform an identity verification process to ensure the security of your account. This may involve a phone call or text message verification.

The screenshot shows the AWS sign-up process at step 4 of 5. The left side features a large placeholder for a profile picture, with a blue circle containing a checkmark overlaid on it. The right side contains the form fields:

- Sign up for AWS**
- Confirm your identity** (Info)
- Name** (Info)
Choose the name that you want to use for identity verification.
 Pushkar Sane
- Primary purpose of account registration**
Choose one that best applies to you. If your account is tied to a business, select the one that applies to your business.
- Ownership type**
- India document type** (Info)
To verify your identity, the name on the document must match the name that you chose.
- Permanent Account Number (PAN)**

The PAN is 10 alphanumeric characters without spaces or tabs. Example: AAAAAA111B
- I consent to allowing AWS to use and send the information above to a third-party service for identity verification purposes.
-

9. Choose a Support Plan:

AWS offers different support plans, including a free basic plan and various premium plans with additional features. Choose the plan that best fits your needs.

The screenshot shows the AWS sign-up process. At the top is the AWS logo. Below it, the heading "Sign up for AWS" is followed by "Select a support plan". A note says "Choose a support plan for your business or personal account. Compare plans and pricing examples" with a link. It also says "You can change your plan anytime in the AWS Management Console." Three support plan options are listed in boxes:

- Basic support - Free** (radio button selected):
 - Recommended for new users just getting started with AWS
 - 24x7 self-service access to AWS resources
 - For account and billing issues only
 - Access to Personal Health Dashboard & Trusted Advisor
- Developer support - From \$29/month** (radio button):
 - Recommended for developers experimenting with AWS
 - Email access to AWS Support during business hours
 - 12 (business)-hour response times
- Business support - From \$100/month** (radio button):
 - Recommended for running production workloads on AWS
 - 24x7 tech support via email, phone, and chat
 - 1-hour response times
 - Full set of Trusted Advisor best-practice recommendations

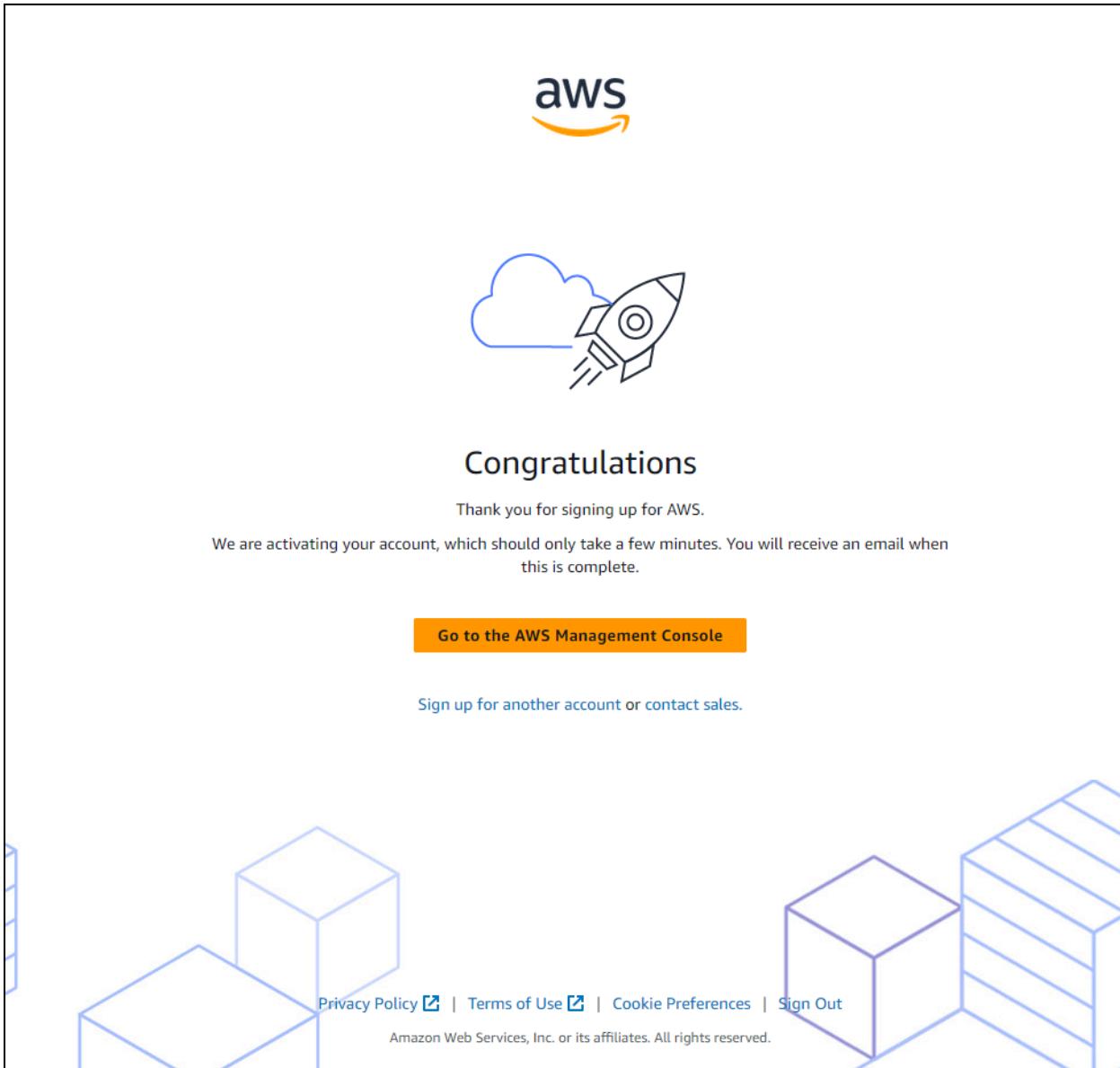
Below the plans, a section for "Need Enterprise level support?" shows an icon of buildings and text about receiving 15-minute response times and concierge-style experience with a Technical Account Manager, with a "Learn more" link. At the bottom is a large orange "Complete sign up" button.

10. Review and Confirm:

Review the information you provided, including your email address, account details, and payment information. Once you are satisfied, click on the "Create Account and Continue" button.

11. Confirmation Email:

AWS will send a confirmation email to the address you provided. Open the email and click on the confirmation link to verify your email address.

**12. Set up Multi-Factor Authentication (Optional):**

For added security, consider setting up Multi-Factor Authentication (MFA) for your AWS account. This can be done during the sign-up process or later within the AWS Management Console.

Once you've completed these steps, your AWS account will be created, and you can start using the AWS services. Keep your account credentials secure, and consider enabling additional security measures, such as IAM (Identity and Access Management) policies, for better control over your AWS resources.

Service Legal Agreement:



AWS Service Terms

Last Updated: January 25, 2024

1. Universal Service Terms (Applicable to All Services)

The Service Terms below govern your use of the Services. Capitalized terms used in these Service Terms but not defined below are defined in the [AWS Customer Agreement](#) or other agreement with us governing your use of the Services (the "Agreement"). For purposes of these Service Terms, "Your Content" includes any "Company Content" and any "Customer Content," and "AWS Content" includes "Amazon Properties."

1.1. You may not transfer outside the Services any software (including related documentation) you obtain from us or third party licensors in connection with the Services without specific authorization to do so.

1.2. You must comply with current technical documentation applicable to the Services (including applicable user, admin, and developer guides) posted on the AWS Site at <https://docs.aws.amazon.com/index.html> (and any successor or related locations designated by us).

1.3. You will provide information or other materials related to Your Content (including copies of any client-side applications) as reasonably requested by us to verify your compliance with the Agreement. You will reasonably cooperate with us to identify the source of any problem with the Services that we reasonably believe may be attributable to Your Content or any end user materials that you control.

1.4. In connection with your use of the Services, you are responsible for maintaining licenses and adhering to the license terms of any software you run. If we reasonably believe any of Your Content violates the law, infringes or misappropriates the rights of any third party, or otherwise violates a material term of the Agreement (including the Service Terms, or the Acceptable Use Policy) ("Prohibited Content"), we will notify you of the Prohibited Content and may request that such content be removed from the Services or access to it be disabled. If you do not remove or disable access to the Prohibited Content within 2 business days of our notice, we may remove or disable access to the Prohibited Content or suspend the Services to the extent we are not able to remove or disable access to the Prohibited Content. Notwithstanding the foregoing, we may remove or disable access to any Prohibited Content without prior notice in connection with illegal content, where the content may disrupt or threaten the Services or in accordance with applicable law or any judicial, regulatory



1.5. You will ensure that all information you provide to us via the AWS Site (e.g., information provided in connection with your registration for the Services, requests for increased usage limits) is accurate, complete, and not misleading.

1.6. From time to time, we may apply upgrades, patches, bug fixes, or other maintenance to the Services and AWS Content ("Maintenance"). We agree to use reasonable efforts to provide you with prior notice of any scheduled Maintenance (except for emergency Maintenance), and you agree to use reasonable efforts to comply with any Maintenance requirements that we notify you about.

1.7. If your Agreement does not include a provision on AWS Confidential Information, and you and AWS do not have an effective non-disclosure agreement in place, then you agree that you will not disclose AWS Confidential Information (as defined in the AWS Customer Agreement), except as required by law.

1.8. You may perform benchmarks or comparative tests or evaluations (each, a "Benchmark") of the Services. If you perform or disclose, or direct or permit any third party to perform or disclose, any Benchmark of any of the Services, you (i) will include in any disclosure, and will disclose to us, all information necessary to replicate such Benchmark, and (ii) agree that we may perform and disclose the results of Benchmarks of your products or services, irrespective of any restrictions on Benchmarks in the terms governing your products or services.

1.9. Only the applicable AWS Contracting Party (as defined in the AWS Customer Agreement) will have obligations with respect to each AWS account, and no other AWS Contracting Party has any obligation with respect to such account. The AWS Contracting Party for an account may change as described in the Agreement. invoices for each account will reflect the AWS Contracting Party that is responsible for that account during the applicable billing period.

If, as of the time of a change of the AWS Contracting Party responsible for your account, you have made an up-front payment for any Services under such account, then the AWS Contracting Party you paid such up-front payment to may remain the AWS Contracting Party for the applicable account only with respect to the Services related to such up-front payment.

1.10. When you use a Service, you may be able to use or be required to use one or more other Services (each, an "Associated Service"), and when you use an Associated Service, you are subject to the terms and fees that apply to that Associated Service.

1.11. If you process the personal data of End Users or other identifiable individuals in your use of a Service, you are responsible for providing legally adequate privacy notices and obtaining necessary consents for the processing of such data. You represent to us that you have provided all necessary privacy notices and obtained all necessary consents. You are responsible for processing such data in accordance with applicable law.



1.14.4 These Service Terms incorporate the [AWS UK GDPR Addendum](#) to the DPA, when the UK GDPR applies to your use of the AWS Services to process UK Customer Data (as defined in the AWS UK GDPR Addendum), and the [AWS Swiss Addendum](#) to the DPA, when the FOPA applies to your use of the AWS Services to process Swiss Customer Data (as defined in the AWS Swiss Addendum).

1.14.5 These Service Terms incorporate the [AWS CCPA Terms](#) ("CCPA Terms"), when the CCPA applies to your use of the AWS Services to process Personal Information (as defined in the CCPA Terms).

1.15. Following closure of your AWS account, we will delete Your Content in accordance with the technical documentation applicable to the Services.

1.16. Your receipt and use of any Promotional Credits is subject to the [AWS Promotional Credit Terms & Conditions](#).

1.17. Payment Currency

1.17.1 AWS provides a Service that enables payment in certain currencies ("Payment Currency") other than United States dollars when you purchase certain Services from AWS (the "Currency Service"). When you purchase Services in certain countries outside of the United States, we may require you, because of currency controls or other factors, to use the Currency Service. When using the Currency Service you are not tendering payment in one currency and receiving from us another currency.

1.17.2 When you use the Currency Service, Service fees and charges will automatically be invoiced in the Payment Currency. You must pay invoices in the currency specified on each invoice. but, for credit card or debit card purchases, you may only make payments in currencies supported by the issuer of your card. If the issuer of your credit card or debit card does not support the required Payment Currency, you must use a different payment method that does support paying in the Payment Currency.

1.17.3 Our fees and charges for your use of the Currency Service, if any, are included in the exchange rate applied to your invoice (the "Applicable Exchange Rate"). Third-parties, such as your bank, credit card issuer, debit card issuer, or card network, may charge you additional fees. The Applicable Exchange Rate is determined at the time your invoice is generated and, for invoices covering usage of Services over a period of time, will apply to all usage and Service charges listed on that invoice.

1.17.4 All refunds processed against an invoice will be provided in the currency in which the invoice was generated and reflected as a credit memo or a payment in your Payment Currency.

Introduction to AWS SageMaker:

Amazon SageMaker is a managed machine learning service by AWS. It streamlines the entire machine learning workflow, from data labeling and model training to deployment and monitoring. With built-in algorithms, hyperparameter optimization, and easy integration with AWS ecosystem, SageMaker makes machine learning accessible to developers and data scientists.

Billing Process:

The screenshot shows the AWS Billing Preferences page. At the top, there are three success notifications: "Introducing the new AWS Billing preferences page experience", "Your invoice delivery preferences were updated successfully.", and "Your alert preferences were updated successfully.". The main content area is titled "Billing preferences". It contains two sections: "Invoice delivery preferences" and "Alert preferences". Under "Invoice delivery preferences", it says "PDF invoices delivery by email" and "Activated". Under "Alert preferences", it lists "AWS Free Tier alerts" (Delivered to Root user email address) and "CloudWatch billing alerts" (Delivered). Below these sections is a box titled "Detailed billing reports (legacy)". It contains a note about the AWS Cost & Usage report and legacy report delivery to S3 (No S3 configured).

The screenshot shows the "Create alarm" wizard in the AWS CloudWatch Metrics console. The current step is "Step 1: Specify metric". A modal window titled "Select metric" is open, showing a graph of "EstimatedCharges" over time. The graph has a Y-axis from 0 to 1 and an X-axis from 10:00 to 16:00. Below the graph, the "Metrics (1)" section shows "Currency 1/1" and "USD" selected. The "EstimatedCharges" metric is listed under "Graphed metrics (1)". The "Source" dropdown shows "N. Virginia". The "Alarms" section indicates "No alarms". At the bottom right of the modal is a yellow "Select metric" button.

Conditions

Threshold type

Static
 Use a value as a threshold

Anomaly detection
 Use a band as a threshold

Whenever EstimatedCharges is...

Define the alarm condition.

Greater
 $>$ threshold

Greater/Equal
 \geq threshold

Lower/Equal
 \leq threshold

Lower
 $<$ threshold

than...

Define the threshold value.

USD

Must be a number

► Additional configuration

[Cancel](#)
Next

AWS Services Search [Alt+S]

N. Virginia Pushkar

CloudWatch > Alarms > Create alarm Step 1 Specify metric and conditions Step 2 Configure actions Step 3 Add name and description Step 4 Preview and create

Configure actions

Notification

Alarm state trigger Define the alarm state that will trigger this action.

In alarm
 The metric or expression is outside of the defined threshold.

OK
 The metric or expression is within the defined threshold.

Insufficient data
 The alarm has just started or not enough data is available.

Send a notification to the following SNS topic Define the SNS (Simple Notification Service) topic that will receive the notification.

Select an existing SNS topic
 Create new topic
 Use topic ARN to notify other accounts

Create a new topic... The topic name must be unique.

SNS topic names can contain only alphanumeric characters, hyphens (-) and underscores (_).

Email endpoints that will receive the notification... Add a comma-separated list of email addresses. Each address will be added as a subscription to the topic above.

user1@example.com, user2@example.com

[Create topic](#)
Add notification

Screenshot of the AWS CloudWatch Alarms 'Create alarm' wizard, Step 3: Add name and description.

The page title is "Add name and description". On the left sidebar, steps are listed: Step 1 (Specify metric and conditions), Step 2 (Configure actions), Step 3 (Add name and description), and Step 4 (Preview and create).

The main form fields include:

- Alarm name:** BillExceed
- Alarm description - optional:** [View formatting guidelines](#)
- Description content:**

```
# This is an H1
**double asterisks will produce strong character**
This is [an example](https://example.com/) inline link.
```
- Character limit:** Up to 1024 characters (0/1024)
- Note:** Markdown formatting is only applied when viewing your alarm in the console. The description will remain in plain text in the alarm notifications.

Buttons at the bottom right: Cancel, Previous, Next (highlighted in orange).

Screenshots of the AWS CloudWatch Alarms 'Create alarm' wizard, Steps 2 and 3.

Step 2: Configure actions

Actions

Notification: When In alarm, send a notification to "Default_CloudWatch_Alarms_Topic"

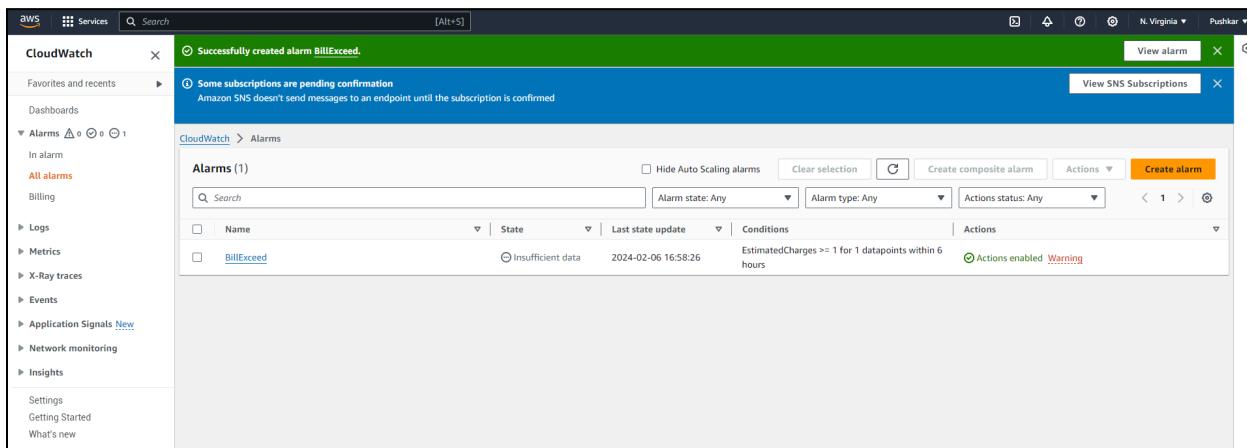
Step 3: Add name and description

Name and description

Name: BillExceed

Description: -

Buttons at the bottom right: Cancel, Previous, Create alarm (highlighted in orange).



IAMUser:

The screenshot shows the AWS IAM "Create user" wizard, Step 1: Specify user details. On the left, a sidebar lists steps: Step 1 (Specify user details), Step 2 (Set permissions), Step 3 (Review and create), and Step 4 (Retrieve password). The main form is titled "Specify user details" and contains a "User details" section. In the "User name" field, "Pushkar" is entered. Below it, a note says "The user name can have up to 64 characters. Valid characters: A-Z, a-z, 0-9, and + = . @ _ - (hyphen)". Under "Are you providing console access to a person?", the "I want to create an IAM user" option is selected. In the "Console password" section, "Autogenerated password" is chosen, and a note states "Users must create a new password at next sign-in - Recommended".

IAM > Users > Create user

Step 1
Specify user details

Step 2
Set permissions

Step 3
Review and create

Step 4
Retrieve password

Set permissions

Add user to an existing group or create a new one. Using groups is a best-practice way to manage user's permissions by job functions. [Learn more](#)

Permissions options

- Add user to group Add user to an existing group, or create a new group. We recommend using groups to manage user permissions by job function.
- Copy permissions Copy all group memberships, attached managed policies, and inline policies from an existing user.
- Attach policies directly Attach a managed policy directly to a user. As a best practice, we recommend attaching policies to a group instead. Then, add the user to the appropriate group.

Permissions policies (1/1175)

Choose one or more policies to attach to your new user.

Policy name	Type	Attached entities
AdministratorAccess	AWS managed - job function	0
AdministratorAccess-Amplify	AWS managed	0
AdministratorAccess-AWSElasticBeanstalk	AWS managed	0
AlexaForBusinessDeviceSetup	AWS managed	0
AlexaForBusinessFullAccess	AWS managed	0
AlexaForBusinessGatewayExecution	AWS managed	0
AlexaForBusinessLifesizeDelegatedAccessPolicy	AWS managed	0
AccessAnalyzerServiceRolePolicy	AWS managed	0

IAM > Users > Create user

Step 1
Specify user details

Step 2
Set permissions

Step 3
Review and create

Step 4
Retrieve password

Review and create

Review your choices. After you create the user, you can view and download the autogenerated password, if enabled.

User details

User name	Console password type	Require password reset
Pushkar	Autogenerated	Yes

Permissions summary

Name	Type	Used as
AdministratorAccess	AWS managed - job function	Permissions policy
IAMUserChangePassword	AWS managed	Permissions policy

Tags - optional

Tags are key-value pairs you can add to AWS resources to help identify, organize, or search for resources. Choose any tags you want to associate with this user.

No tags associated with the resource.

[Add new tag](#)

You can add up to 50 more tags.

[Cancel](#) [Previous](#) [Create user](#)

User created successfully

You can view and download the user's password and email instructions for signing in to the AWS Management Console.

[View user](#) [X](#)

IAM > Users > Create user

Step 1
Specify user details

Step 2
Set permissions

Step 3
Review and create

Step 4
Retrieve password

Retrieve password

You can view and download the user's password below or email users instructions for signing in to the AWS Management Console. This is the only time you can view and download this password.

Console sign-in details

Console sign-in URL	https://058264423218.signin.aws.amazon.com/console	Email sign-in instructions
---------------------	---	--

User name

Pushkar

Console password

***** [Show](#)

[Cancel](#) [Download .csv file](#) [Return to users list](#)

IAM > Users > Pushkar > Assign MFA device

Step 1
Select MFA device

Step 2
Set up device

Select MFA device Info

MFA device name

Device name
Enter a meaningful name to identify this device.

Maximum 128 characters. Use alphanumeric and '+ = , . @ - _' characters.

MFA device

Select an MFA device to use, in addition to your username and password, whenever you need to authenticate.

Authenticator app
Authenticate using a code generated by an app installed on your mobile device or computer.

Security Key
Authenticate using a code generated by touching a YubiKey or other supported FIDO security key.

Hardware TOTP token
Authenticate using a code displayed on a hardware Time-based one-time password (TOTP) token.

AWS Services Q mfa

Identity and Access Management (IAM)

MFA device assigned
You can register up to 8 MFA devices of any combination of the currently supported MFA types with your AWS account root and IAM user. With multiple MFA devices, you only need one MFA device to sign in to the AWS console or create a session through the AWS CLI with that user.

IAM > Users > Pushkar

Pushkar Info

Summary

ARN arn:aws:iam:058264423218:user/Pushkar	Console access Enabled with MFA	Access key 1 Create access key
Created February 06, 2024, 22:56 (UTC+05:30)	Last console sign-in Never	

Permissions Groups Tags Security credentials Access Advisor

Console sign-in

Console sign-in link
<https://058264423218.signin.aws.amazon.com/console>

Console password
Updated 2 minutes ago (2024-02-06 22:36 GMT+5:30)

Last console sign-in
Never

Multi-factor authentication (MFA) (1)

Use MFA to increase the security of your AWS environment. Signing in with MFA requires an authentication code from an MFA device. Each user can have a maximum of 8 MFA devices assigned. [Learn more](#)

Remove Resync Assign MFA device

Device type	Identifier	Certifications	Created on
-------------	------------	----------------	------------

AIML Components in AWS:

AWS provides a robust set of Artificial Intelligence and Machine Learning (AIML) services. These components include Amazon SageMaker for end-to-end machine learning workflows, Amazon Comprehend for natural language processing, Amazon Rekognition for image and video analysis, and more. These AIML services empower users to integrate powerful machine learning capabilities into their applications without the need for extensive expertise in the field.

Simple Storage Service (S3) Introduction:

Amazon S3 (Simple Storage Service) is a scalable object storage service in AWS, designed to store and retrieve any amount of data from anywhere on the web. S3 offers industry-leading durability, availability, and scalability, making it ideal for diverse storage needs. Users can organize data in buckets, control access through flexible permissions, and leverage features like versioning and lifecycle policies.

Elastic Computing Cloud (EC2) Introduction:

Amazon EC2 (Elastic Compute Cloud) forms the backbone of AWS's compute services, offering resizable virtual servers in the cloud. Users can choose from a variety of instance types based on their specific computational requirements. EC2 instances provide flexibility, scalability, and control over computing resources, enabling users to run applications efficiently and cost-effectively in the AWS cloud.

Conclusion: In conclusion, this practical provided an introduction to essential aspects of Amazon Web Services (AWS), beginning with the sign-up process and familiarization with service legal agreements. It then delved into Amazon SageMaker, emphasizing its ease of access, free tier availability, billing structure, and integration with AWS Identity and Access Management (IAM). Additionally, the practical covered key components of artificial intelligence and machine learning (AIML) within AWS, followed by an overview of fundamental services like Simple Storage Service (S3) and Elastic Compute Cloud (EC2). These foundational concepts equip users with the knowledge necessary to leverage AWS effectively for their computing and machine learning needs.

Name of Student: Pushkar Sane		
Roll Number: 45	Lab Assignment Number: 3, 4	
Title of Lab Assignment: Introduction to AWS SageMaker, Data Labeling in AWS Sage maker Ground Truth, Labeling Text, Bounding Boxes and Semantic Segmentation in Ground Truth.		
DOP: 14-02-2024	DOS: 22-02-2024	
CO Mapped:	PO Mapped:	Signature:

Practical No. 3 & 4

Aim: Introduction to AWS SageMaker, Data Labeling in AWS Sage maker Ground Truth, Labeling Text, Bounding Boxes and Semantic Segmentation in Ground Truth.

Theory:

Amazon SageMaker is a fully managed service by AWS that simplifies the machine learning workflow, from data preparation to model deployment. It provides a range of tools and services for building, training, and deploying machine learning models at scale. With SageMaker, developers and data scientists can efficiently experiment with different algorithms, automate model training with Autopilot, and deploy models with ease using managed endpoints. SageMaker also offers capabilities for data labeling, model monitoring, and governance, making it a comprehensive platform for AI development.

Components of AWS Sagemaker:

1. Jump Start:

- Amazon SageMaker JumpStart provides access to popular machine learning algorithms and model templates.
- Amazon SageMaker JumpStart provides access to pre-built machine learning models and solutions, accelerating the model development process.
- It offers a curated selection of algorithms and model templates, allowing users to quickly get started with common machine learning tasks without the need for extensive setup or configuration.

To use JumpStart:

- Log in to the AWS Management Console.
- Navigate to Amazon SageMaker.
- Choose "JumpStart" from the left-hand menu.
- Select the algorithm or model template that fits your needs.
- Follow the provided instructions to configure and deploy the chosen model.

2. Ground Truth:

- Amazon SageMaker Ground Truth helps label data for machine learning models efficiently.

- Amazon SageMaker Ground Truth is a data labeling service that makes it easy to annotate datasets for training machine learning models.
- Ground Truth offers both human annotation and automated data labeling capabilities, helping users generate high-quality labeled datasets efficiently.

Steps to use Ground Truth:

- Log in to the AWS Management Console.
- Navigate to Amazon SageMaker.
- Choose "Ground Truth" from the left-hand menu.
- Create a labeling job by specifying the input data, labeling categories, and workforce.
- Launch the labeling job and monitor its progress.
- Once labeling is complete, use the labeled data for training your model.

3. Governance:

- SageMaker provides governance features to manage access, security, and compliance of machine learning workflows.
- SageMaker Governance features enable users to manage access, security, and compliance for machine learning workflows.
- It provides tools for defining IAM policies, auditing model activities, implementing encryption, and enforcing data privacy and compliance requirements.

Steps for governance:

- Navigate to SageMaker in the AWS Management Console.
- Set up IAM roles and policies to control access to SageMaker resources.
- Use AWS CloudTrail to monitor API activity.
- Implement encryption for data at rest and in transit.
- Configure VPC settings for SageMaker resources to control network access.

4. Notebook:

- Amazon SageMaker Notebooks provide a fully managed Jupyter notebook instance for building and testing models.
- SageMaker Notebook Instances provide a fully managed Jupyter notebook environment for building, experimenting, and collaborating on machine learning projects.

- Users can leverage built-in libraries and frameworks, scale compute resources as needed, and integrate with other SageMaker components seamlessly.

To use SageMaker Notebooks:

- Navigate to SageMaker in the AWS Management Console.
- Choose "Notebook instances" from the left-hand menu.
- Create a new notebook instance, selecting the instance type and IAM role.
- Open the Jupyter notebook interface and start writing code.
- Save your work to Amazon S3 or GitHub for version control.

5. Processing:

- SageMaker Processing allows you to preprocess data and run custom data transformations at scale.
- SageMaker Processing allows users to preprocess and analyze data at scale before training machine learning models.
- It supports custom data transformations, parallel execution of processing tasks, and integration with other AWS services for data storage and retrieval.

Steps for processing:

- Navigate to SageMaker in the AWS Management Console.
- Choose "Processing jobs" from the left-hand menu.
- Create a processing job, specifying the input data, processing script, and output location.
- Monitor the status of the processing job and view logs for debugging if needed.

6. Training:

- SageMaker Training enables you to train machine learning models using built-in algorithms or custom scripts.
- SageMaker Training enables users to train machine learning models using built-in algorithms, custom code, or frameworks like TensorFlow and PyTorch.
- It offers distributed training capabilities, automatic model tuning, and support for training on large datasets stored in Amazon S3.

Steps for training:

- Navigate to SageMaker in the AWS Management Console.
- Choose "Training jobs" from the left-hand menu.

- Create a training job, specifying the algorithm or script, input data location, and output location.
- Configure hyperparameters and instance types for training.
- Monitor the training job's progress and view logs for insights.

7. Inference:

- Amazon SageMaker Inference allows you to deploy trained models for real-time or batch inference.
- SageMaker Inference enables users to deploy trained models for real-time or batch inference, serving predictions at scale with low latency.
- It provides managed endpoints, automatic scaling, and integration with AWS Lambda and API Gateway for building scalable inference pipelines.

To deploy a model for inference:

- Navigate to SageMaker in the AWS Management Console.
- Choose "Endpoints" from the left-hand menu.
- Create a new endpoint, specifying the trained model, instance type, and number of instances.
- Once the endpoint is created, you can use it to make predictions by sending HTTP requests.

8. Augmented AI:

- SageMaker Augmented AI (A2I) helps you build human review workflows for ML predictions to improve model accuracy.
- SageMaker Augmented AI (A2I) helps improve the accuracy of machine learning models by integrating human review into the prediction process.
- It allows users to create custom review workflows, route predictions to human reviewers, and incorporate feedback to continuously improve model performance

Steps for AI:

- Navigate to SageMaker in the AWS Management Console.
- Choose "Ground Truth" from the left-hand menu.
- Create a human review workflow by specifying the conditions for sending predictions to human reviewers.
- Configure the workforce and define the instructions for reviewers.

- Monitor the human review process and incorporate feedback to improve model performance.

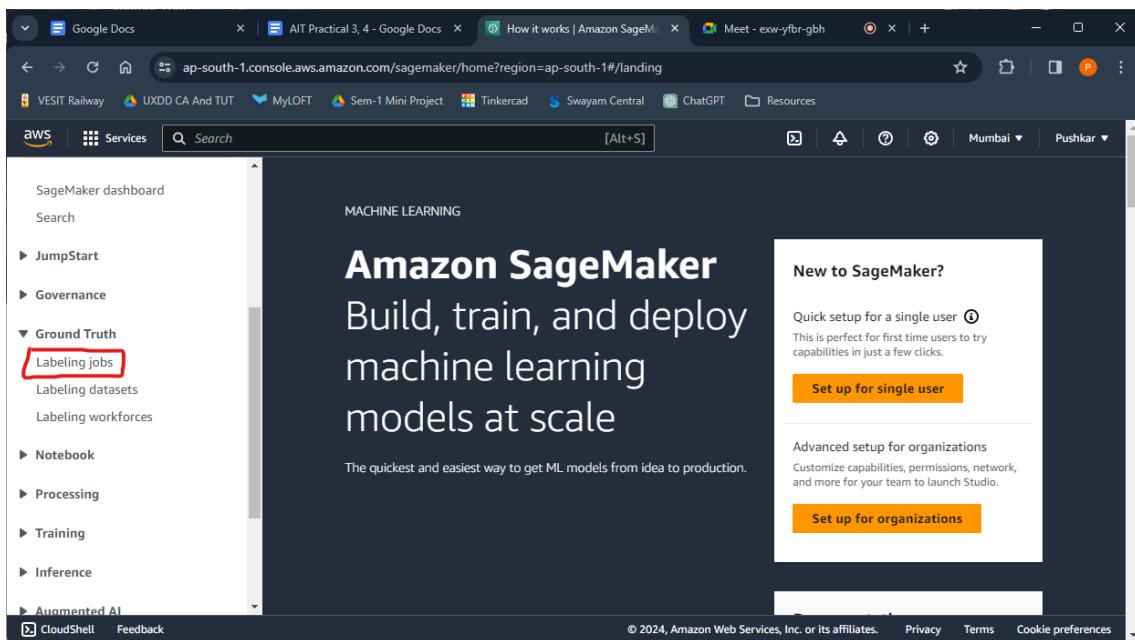
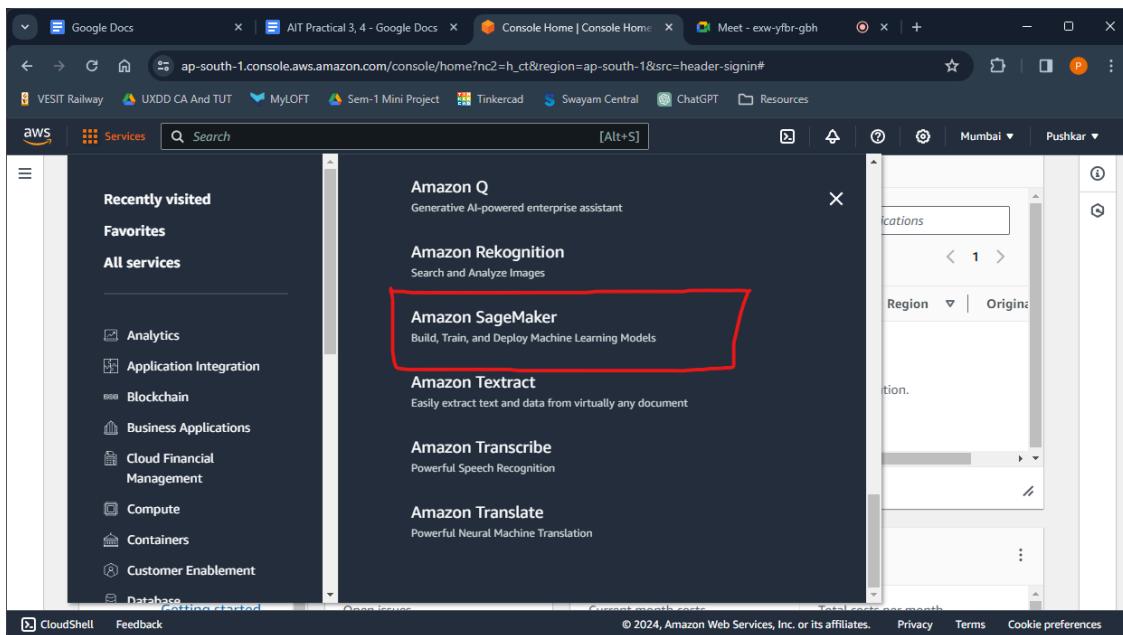
9. AWS Marketplace (GluonCV YOLO v3 object detector):

- The AWS Marketplace offers pre-built machine learning models and algorithms for easy integration into SageMaker.
- The AWS Marketplace offers a wide range of pre-built machine learning models and algorithms, including the GluonCV YOLO v3 object detector.
- Users can easily deploy these models in SageMaker for tasks like object detection, image classification, and natural language processing, accelerating the development of AI applications.

To use GluonCV YOLO v3 object detector:

- Navigate to the AWS Marketplace in the AWS Management Console.
- Search for "GluonCV YOLO v3 object detector" and subscribe to the listing.
- Follow the provided instructions to deploy the model in SageMaker.
- Once deployed, you can use the model for object detection tasks by sending input images to the endpoint.

Remember to manage costs effectively by monitoring resource usage and shutting down instances when not in use. Additionally, ensure compliance with data privacy regulations when handling sensitive data.

Steps to explore the AWS sagemaker:

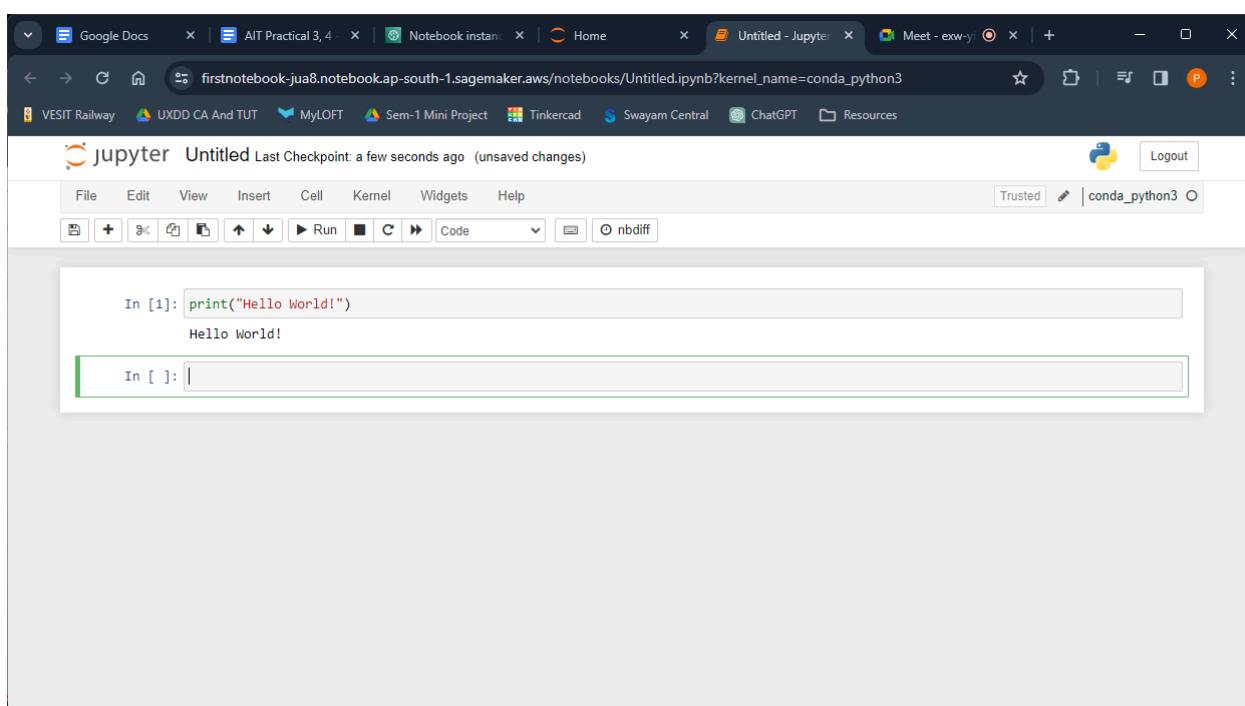
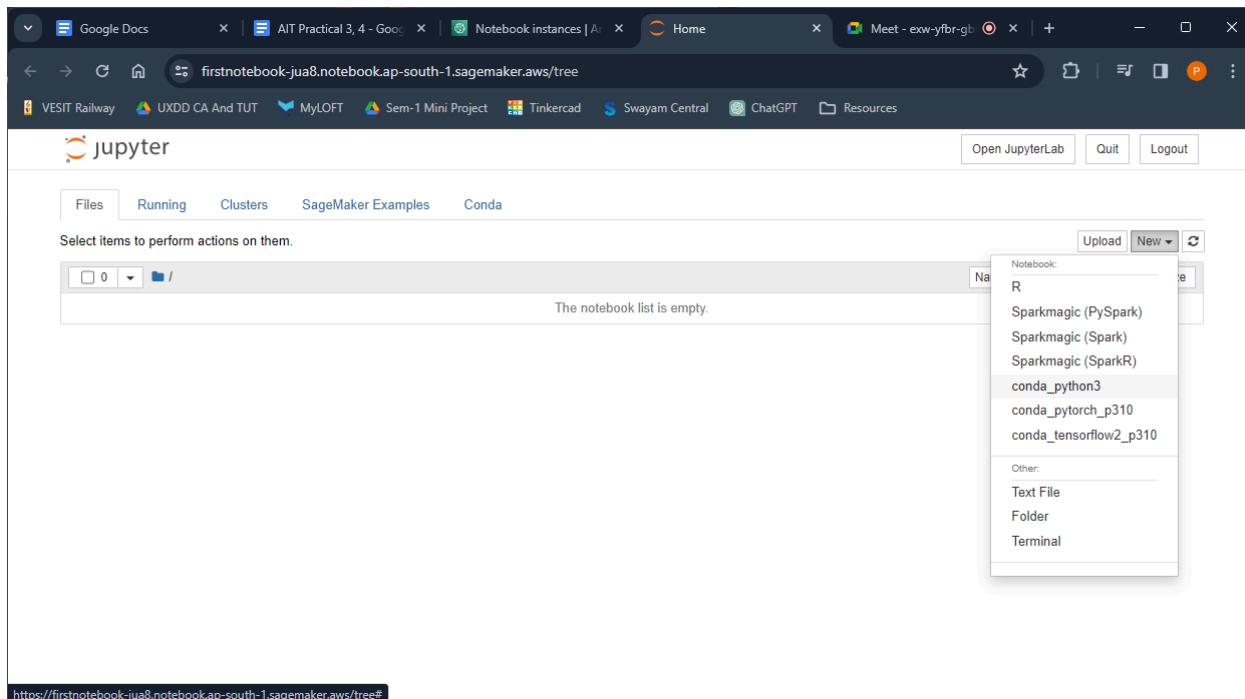
The screenshot shows the Amazon SageMaker service in the AWS Management Console. The left sidebar navigation includes 'Role manager', 'Images', 'Lifecycle configurations', 'SageMaker dashboard', 'Search', 'JumpStart', 'Governance', 'Ground Truth', 'Notebook' (selected), 'Notebook instances' (selected), 'Git repositories', 'Processing', and 'Training'. The main content area is titled 'Notebook instances' and shows a table with columns: Name, Instance, Creation time, Status, and Actions. A message at the bottom states 'There are currently no resources.' The top navigation bar shows tabs for 'Google Docs', 'AIT Practical 3, 4 - Google Docs', 'Notebook instances | Amazon S', and 'Meet - exw-yfbr-gbh'. The top right corner shows account information for 'Pushkar'.

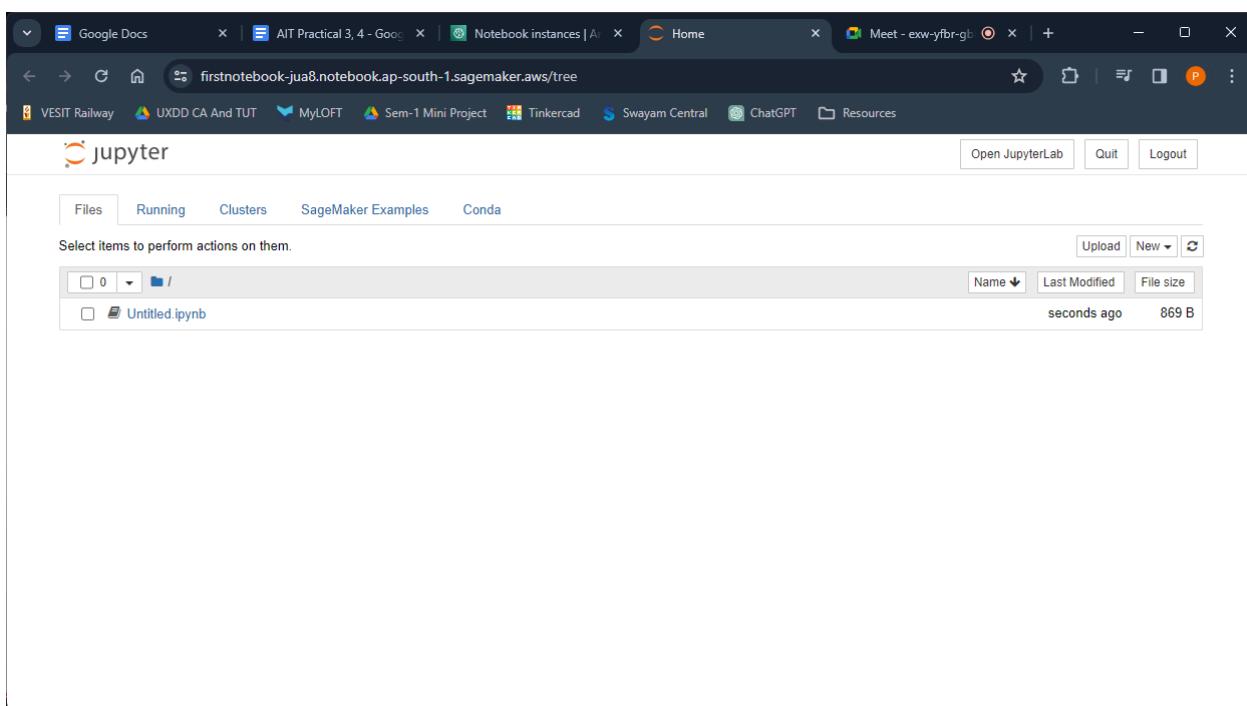
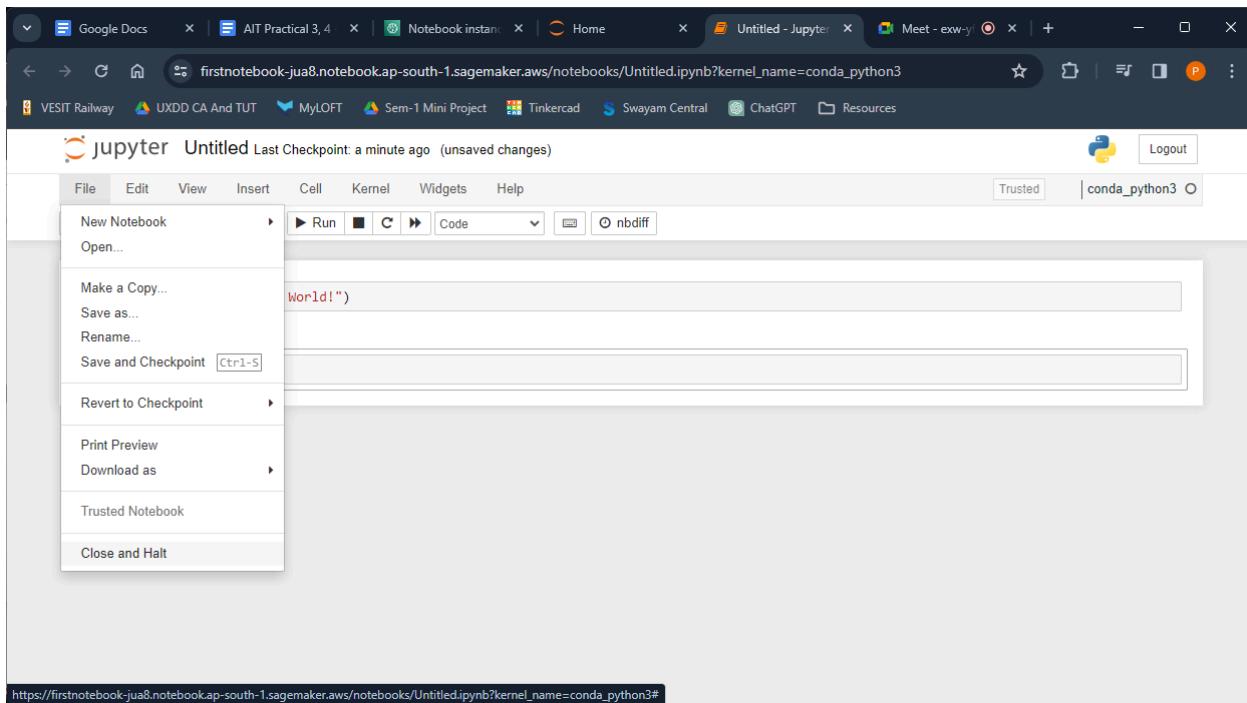
The screenshot shows the 'Create notebook instance' wizard. Step 1: Notebook instance settings. It asks for a 'Notebook instance name' (entered as 'firstnotebook'), 'Notebook instance type' (selected as 'ml.t3.medium'), and 'Platform identifier' (selected as 'Amazon Linux 2, Jupyter Lab 3'). Below these fields is a 'Additional configuration' section with a '▶' icon. At the bottom is a 'Permissions and encryption' section. The top navigation bar shows tabs for 'Google Docs', 'AIT Practical 3, 4 - Google Docs', 'Create notebook instance | Note...', and 'Meet - exw-yfbr-gbh'. The top right corner shows account information for 'Pushkar'.

Screenshot of the AWS SageMaker IAM role creation wizard. The page shows a success message: "Success! You created an IAM role." followed by the ARN "AmazonSageMaker-ExecutionRole-20240222T195292". Below this, there are sections for "Create role using the role creation wizard" and "Encryption key - optional". The "Root access - optional" section has a radio button for "Enable - Give users root access to the notebook" selected. The "Network - optional" section is collapsed.

Screenshot of the AWS SageMaker Notebook instances list. The page shows a table of existing notebook instances:

Name	Instance	Creation time	Status	Actions
firstnotebook	mlt2.medium	2/22/2024, 7:53:54 PM	InService	Open Jupyter Open JupyterLab





The screenshot shows the 'Notebook instance settings' page for a notebook named 'firstnotebook'. The left sidebar lists various options like 'Getting started', 'Studio', 'Canvas', 'RStudio', and 'Admin configurations' (Domains, Role manager, Images, Lifecycle configurations). The main panel displays the following details:

Name	Notebook instance type
firstnotebook	ml.t2.medium
ARN	Volume Size
arn:aws:sagemaker:ap-south-1:058264423218:notebook-instance/firstnotebook	5GB EBS
Lifecycle configuration	Platform identifier
-	Amazon Linux 2, Jupyter Lab 3 (notebook-al2-v2)
Status	Minimum IMDS Version
InService	2
Creation time	

At the bottom right of the main panel is an 'Edit' button. At the very bottom of the screen are standard browser navigation buttons and a footer bar.

The screenshot shows the 'Notebook instances' list page. The left sidebar is identical to the previous screenshot. The main panel displays a table of notebook instances:

Name	Instance	Creation time
firstnotebook	ml.t2.medium	2/22/2024, 7:53:54 PM

To the right of the table is a context menu for the 'firstnotebook' row, with 'Stop' highlighted. Other options in the menu include 'Start', 'Update settings', 'Add/Edit tags', and 'Delete'. The footer of the screen includes standard browser navigation and a footer bar.

The screenshot shows the AWS Marketplace search interface. On the left, there is a sidebar with various navigation options like Canvas, RStudio, Admin configurations, SageMaker dashboard, Search, JumpStart, Governance, Ground Truth, Notebook, Processing, and Training. The main area displays a search bar with the placeholder "Search" and a message: "Need help creating a custom machine learning solution? Connect with an AWS IQ expert to train a custom model in Amazon SageMaker." Below this, there is a section titled "Featured model packages" with three items: "Face and License Plate Anonymizer" by NavInfo Europe B.V., "GluonCV YOLOv3 Object Detector" by Amazon Web Services, and "Passport Data Page Detection" by Gtrip. The "GluonCV YOLOv3 Object Detector" item has a "Free Trial" button.

The screenshot shows the product page for the "GluonCV YOLOv3 Object Detector" on the AWS Marketplace. At the top, there is a navigation bar with links for About, Categories, Delivery Methods, Solutions, AWS IQ, Resources, and Your Saved List. Below this, there is a "Continue to Subscribe" button, a "Try Product Demo" button, and a "Save To List" button. The main content area features the "mxnet" logo, the product name "GluonCV YOLOv3 Object Detector", and a brief description: "YOLOv3 is a powerful network for fast and accurate object detection, powered by GluonCV." There are tabs for Overview, Pricing, Usage, Support, and Reviews. The "Overview" tab is selected. In the "Key Data" section, there is a link: <https://aws.amazon.com/marketplace/ai/model-evaluation?productid=d9...>. On the right, there is a "Demo available" section with a "Try Product Demo" button and a small icon of a laptop.

Screenshot of the AWS Marketplace product demo page for GluonCV YOLOv3 Object Detector.

The page shows a living room scene with objects labeled: tv, couch, chair, vase, chair, chair, tv, dining table.

Object detection

Use your own image
Image must be in jpeg, png or bmp format and be no larger than 5 MB.
Do not upload any confidential or sensitive information. Use of this feature is for demonstration purposes only and is in a public environment.

Object details
Click on a row below or an object in the image to view details about the object.

Screenshot of the AWS Marketplace product demo page for GluonCV YOLOv3 Object Detector, showing the results of the object detection on a different image.

The image shows a living room scene with objects labeled: couch, chair, chair, vase, chair, chair, tv, dining table.

Object details
Label: couch
Confidence: 97.18%

Objects identified in image		
Tag	Confidence	Actions
couch	97.18%	Highlight in image
chair	96.99%	Highlight in image
vase	74.05%	Highlight in image
chair	73.53%	Highlight in image
chair	66.66%	Highlight in image
tv	58.90%	Highlight in image
dining table	31.67%	Highlight in image

This demo may not accurately represent the actual response times of the product.

Conclusion: Successfully explored all the components of the AWS Sagemaker.

Name of Student: Pushkar Sane

Roll Number: 45

Lab Assignment Number: 5 & 6

Title of Lab Assignment: Feature Engineering, one hot Encoding, Normalization, Standardization, EDA using SageMaker DataWrangler Linear and Multiple Linear Regression using SageMaker.

DOP: 16-03-2024

DOS: 16-03-2024

CO Mapped:

CO3

PO Mapped:

**PO2, PO3, PO4, PO5, PO6,
PO7, PSO1, PSO2**

Signature:

Practical 5 & 6

Aim: Feature Engineering, one hot Encoding, Normalization, Standardization, EDA using SageMaker DataWrangler Linear and Multiple Linear Regression using SageMaker.

Description:

Prerequisite

Create an S3 bucket and keep aside (it will be used later)

Click on “create” to create the bucket

Change the permissions

The screenshot shows the 'Edit bucket policy' page in the AWS S3 console. The breadcrumb path 'Amazon S3 > Buckets > alt-prac3 > Edit bucket policy' is highlighted with a red box. Below it, the 'Bucket policy' section is shown, with the word 'Policy' underlined by a red bracket. The JSON code for the policy is displayed:

```

1 {  
2   "Version": "2012-10-17",  
3   "Statement": [  
4     {  
5       "Sid": "Statement1",  
6       "Principal": "*",  
7       "Effect": "Allow",  
8       "Action": [  
9         "s3:GetObject"  
10      ],  
11      "Resource": "arn:aws:s3:::alt-prac3/*"  
12    }  
13  ]  
14 }

```

SageMaker Canvas

The screenshot shows the SageMaker Canvas landing page. The left sidebar menu has 'Canvas' underlined with a red box. The main content area features the heading 'SageMaker Canvas' and the subtext 'Generate accurate machine learning predictions — no code required'. To the right, there is a 'Get Started' box with the text 'Create an Amazon SageMaker domain to use Studio and Studio Notebooks.' and a prominent orange 'Create a SageMaker domain' button, which is also highlighted with a red arrow.

The screenshot shows the 'Set up SageMaker Domain' page. On the left, there's a sidebar with links like Getting started, Studio, Studio Lab, Canvas, RStudio, TensorBoard, Profiler, Admin configurations (Domains, Role manager, Images, Lifecycle configurations), SageMaker dashboard, and Search. The main content area has two sections: 'Set up for single user (Quick setup)' and 'Set up for organizations'. The 'Set up for single user' section is highlighted with a red arrow pointing to the 'Set up' button. It lists several configuration options with green checkmarks. Below this, a note says 'Perfect for single user domains and first time users looking to get started with SageMaker.' The 'Set up for organizations' section contains a note about controlling account configuration for large user groups. At the bottom right, there's a note about updating account settings later if desired.

The screenshot shows the 'Domain details' page for 'QuickSetupDomain-20240317T074006'. The sidebar is identical to the previous screenshot. The main content area shows a 'User profiles' tab selected. It displays a search bar and a table with columns for Name, Modified on, and Created on. A note at the bottom says 'No users' and 'To add a user, choose Add user and enter a user name.' The status bar at the bottom indicates the page was last updated on March 17, 2024, at 07:40:06.

The screenshot shows the 'Domain details' section of the Amazon SageMaker console. Under 'General settings', the domain name is 'QuickSetupDomain-20240317T074006', the status is 'Ready' (highlighted with a red arrow), and the VPC is 'vpc-02af8c3291e6078b9'. Under 'Storage configurations', the maximum space size is set to 5 (circled with a red oval). The sidebar on the left includes options like Getting started, Studio, Studio Lab, Canvas, RStudio, TensorBoard, Profiler, Admin configurations, JumpStart, and Governance.

Data Preparation using SageMaker DataWrangler

The screenshot shows the 'Amazon SageMaker' landing page with the 'Canvas' option selected in the sidebar. The main area features the heading 'SageMaker Canvas' and the subtext 'Generate accurate machine learning predictions — no code required'. A 'Get Started' button is visible, with a red arrow pointing to it. Below this, there is a section titled 'How it works' with a small preview image. The sidebar also lists 'Getting started', 'Studio', 'Studio Lab', 'Canvas' (selected), 'RStudio', 'TensorBoard', 'Profiler', 'Admin configurations', 'JumpStart', and 'Governance'.

Import Dataset

The screenshot shows the AWS SageMaker Data Wrangler interface. On the left, there is a sidebar with icons for Home, Data Wrangler (which is circled in red), My Models, ML Ops, Ready-to-use, Shared Models, Gen AI, and Help. The main area is titled "Data Wrangler" and has tabs for "Datasets" and "Data flows". Below the tabs, there is a flow diagram with four steps: "Import data" (database icon), "Prepare data" (lightbulb icon), "Scale data operations" (cogwheel and wrench icon), and "Build models" (handshake icon). A message says "You have not created any data flows yet." Below the message is a blue button with a plus sign and the text "Create a data flow". Red arrows point from the circled "Data Wrangler" icon in the sidebar to the "Create a data flow" button.

The screenshot shows the AWS SageMaker Data Wrangler interface for a specific data flow named "prac3-16Mar24.flow". The sidebar and main layout are identical to the first screenshot, but the main area now displays the details of the selected data flow. It shows the same four-step process: Import data, Prepare data, Scale data operations, and Build models. A message "To get started..." is displayed above two buttons: "Import data" and "Select existing dataset". A red arrow points from the "Select existing dataset" button towards the "Import data" button.

The screenshot shows the AWS SageMaker Canvas Data Wrangler interface. On the left is a sidebar with icons for Home, Data Wrangler (selected), My Models, ML Ops, Ready-to-use, Shared Models, Gen AI, and Help. The main area displays a workflow: Import data → Prepare data → Scale data operations → Build models. A tooltip for 'Import data' is open, showing 'Dataset type' with 'Tabular' (CSV or Parquet files) selected. Other options include 'Image' (PNG and JPG files).

Reference for sample datasets explanation:

<https://docs.aws.amazon.com/sagemaker/latest/dg/canvas-sample-datasets.html>

canvas-sample-housing.csv: This dataset contains data on the characteristics tied to a given housing price. You can use this dataset to predict housing prices. Use the **median_house_value** column as the target column, and use the numeric prediction model type with this dataset. To learn more about building a model with this dataset, see the [SageMaker Canvas workshop page](#). This is the California housing dataset obtained from the [StatLib repository](#).

The screenshot shows the 'Select existing dataset' step in the AWS SageMaker Canvas Data Wrangler interface. The sidebar is the same as the previous screenshot. The main area shows a list of datasets: canvas-sample-retail-electronics-forecasting.csv, canvas-sample-shipping-logs.csv, canvas-sample-product-descriptions.csv, canvas-sample-databricks-dolly-15k.csv, canvas-sample-housing.csv (selected and highlighted with a red circle), and canvas-sample-diabetic-readmission.csv. A red arrow points to the 'Add dataset' button at the bottom right.

Delete irrelevant columns

The screenshot shows the AWS Data Wrangler interface. On the left, there's a sidebar with icons for Home, Data Wrangler (selected), My Models, ML Ops, Ready-to-use, Shared Models, Gen AI, Help, and a feedback icon. The main area has tabs for Data and Analyses. Under Data, it shows 'Step 2. Data types' with three columns: longitude (float), latitude (float), and housing_median_age. Each column has a histogram and some numerical values below it. To the right, there's a 'Steps' sidebar with a '+ Add step' button highlighted by a red arrow. Below it, there are two steps listed: '1. 53 Source' and '2. Data types'. At the bottom, there's a checkbox for visualizing the first 2,000 rows.

This screenshot is similar to the one above, but the sidebar options are expanded. The 'Manage columns' option is highlighted with a red arrow. Other visible options include Handle missing, Handle outliers, Handle structured column, Manage rows, Manage vectors, Parse column as type, Process numeric, and Sampling.

Data Wrangler: Data flow > prac3-16Mar24.flow > canvas-sample-housing.csv

Step 2. Data types

longitude (float) latitude (float) housing_median_age

-122.34 - -121.61	37.47 - 37.9	2 - 52
-122.23	37.88	41
-122.22	37.86	21
-122.24	37.85	52
-122.25	37.85	52
-122.25	37.85	52

Show in-column visualizations for the first 2,000 rows. Visualize the full dataset. Run Data quality and insights report

← Manage columns

Move, drop, duplicate or rename columns in the dataset. [Learn more](#).

Transform [i](#) Required

Drop column

Columns to drop Required

latitude x longitude x

Clear Preview Add

Ordinal Encoding for categorical column

Data Wrangler: Data flow > prac3-16Mar24.flow > canvas-sample-housing.csv

Step 3. Drop column

population (float) households (float) median_income (float) median_house_value (float) ocean_proximity (string)

18 - 12203	7 - 3701	0.4999 - 13.499	60000 - 5.0000e+5	3 Categories
322	126	8.3252	452600	NEAR BAY
2401	1138	8.3014	358500	NEAR BAY
498	177	7.2574	352100	NEAR BAY
558	219	5.6431	341300	NEAR BAY
665	259	3.8462	342200	NEAR BAY
413	193	4.0368	269700	NEAR BAY
1094	514	3.6591	299200	NEAR BAY

+ Add step

Steps

- 1. S3 Source
- 2. Data types
- 3. Drop column

Data Wrangler: Data flow > prac3-16Mar24.flow > canvas-sample-housing.csv

Step 3. Drop column

population (float) households (float) median_income (float) median_house_value (float) ocean_proximity (string)

18 - 12203	7 - 3701	0.4999 - 13.499	60000 - 5.0000e+5	3 Categories
322	126	8.3252	452600	NEAR BAY
2401	1138	8.3014	358500	NEAR BAY
498	177	7.2574	352100	NEAR BAY
558	219	5.6431	341300	NEAR BAY
665	259	3.8462	342200	NEAR BAY
413	193	4.0368	269700	NEAR BAY
1094	514	3.6591	299200	NEAR BAY
1157	647	3.12	241400	NEAR BAY
1206	595	2.0804	226700	NEAR BAY
1551	714	3.6912	261100	NEAR BAY
910	402	3.2031	281500	NEAR BAY
1504	734	3.2705	241800	NEAR BAY
1098	468	3.075	213500	NEAR BAY

← Add transform

Search transforms

Frequently used

Manage columns Move, drop, duplicate or rename columns in the dataset.

Custom

Custom formula Define a new column using a Spark SQL expression to query data in the current dataframe.

Custom transform Use PySpark, Pandas, or PySpark (SQL) to define custom transformations.

Standard

Balance data Balance the data for binary classification problems using random oversampling, random undersampling or SMOTE.

Dimensionality Reduction For the top K principal components, trains a model to project vectors to a lower dimensional space.

Encode categorical Convert categorical variables to numeric or vector representations.

Featureize date/time Extract features from dates, such as year, day, week, and more.

Featureize text Generate vector representations from natural language text.

Min-Max Scaling

Data Wrangler: Data flow > prac3-16Mar24.flow > canvas-sample-housing.csv

Step 4. Ordinal encode

total_rooms (float) total_bedrooms (float) population (float) households (float) median_income (float) median_house_value (float)

	12 - 28258	4 - 3864	18 - 12203	7 - 3701	0.4999 - 13.499	60000 - 5.0000
880	129	322	126	8.3252	452600	
7099	1106	2401	1138	8.3014	358500	
1467	190	496	177	7.2574	352100	
1274	235	558	219	5.6431	341300	
1627	280	565	259	3.8462	342200	
919	213	413	193	4.0368	269700	
2535	489	1094	514	3.6591	299200	
3104	687	1157	647	3.12	241400	
2555	665	1206	595	2.0804	228700	
3549	707	1551	714	3.6912	261100	
2202	434	910	402	3.2031	281500	
3503	752	1504	734	3.2705	241800	
2491	474	1098	468	3.075	213500	

Show in-column visualizations for the first 2,000 rows. Visualize the full dataset. Run Data quality and insights report

Process numeric

Transform: Scale values

Scaler: Min-max scaler

Input columns: total_bedrooms, total_rooms

Min: 0

Max: 1

Output column: (empty)

Preview Add

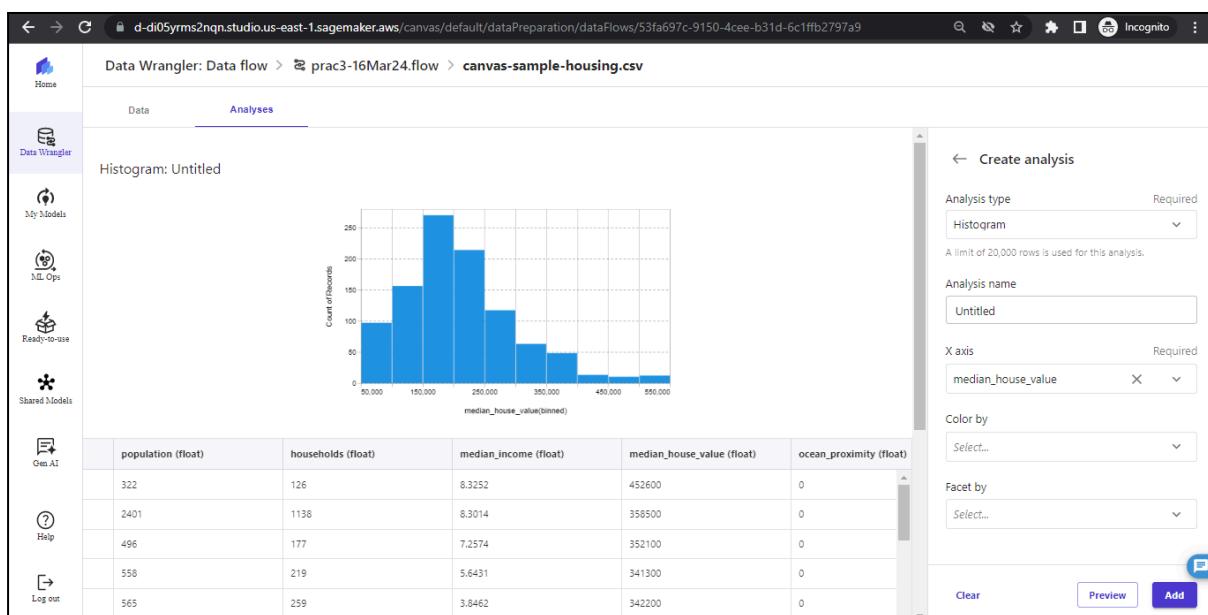
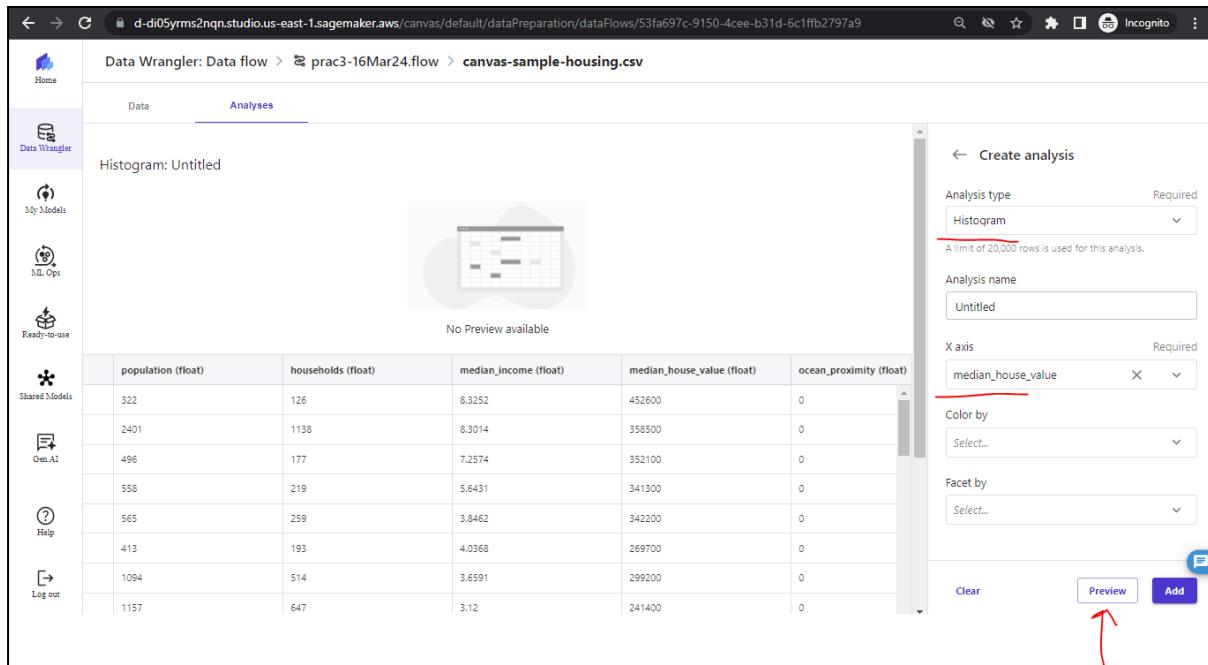
EDA using SageMaker DataWrangler

Table Summary

The screenshot shows the SageMaker Data Wrangler interface. The left sidebar has icons for Home, Data Wrangler (selected), My Models, ML Ops, Ready-to-use, Shared Models, Gen AI, Help, and a right arrow. The main area shows a data flow titled 'prac3-16Mar24.flow' with a file named 'canvas-sample-housing.csv'. The 'Analyses' tab is highlighted with a red circle. The 'Table Summary' section shows a preview of the data with columns: median_income (float), median_house_value (float), and ocean_proximity (float). Below this, a 'Create analysis' sidebar shows 'Analysis type: Table Summary' and 'Analysis name: Untitled'. At the bottom right of the sidebar, there are 'Clear', 'Preview' (highlighted with a red arrow), and 'Add' buttons.

This screenshot is similar to the one above, showing the 'Analyses' tab selected. The 'Table Summary' section now displays summary statistics for the columns: median_income, median_house_value, and ocean_proximity. The 'Create analysis' sidebar remains the same, with 'Analysis type: Table Summary' and 'Analysis name: Untitled'. The 'Preview' and 'Add' buttons are again highlighted with red arrows.

Checking if the target column is normally distributed (bell shaped graph)



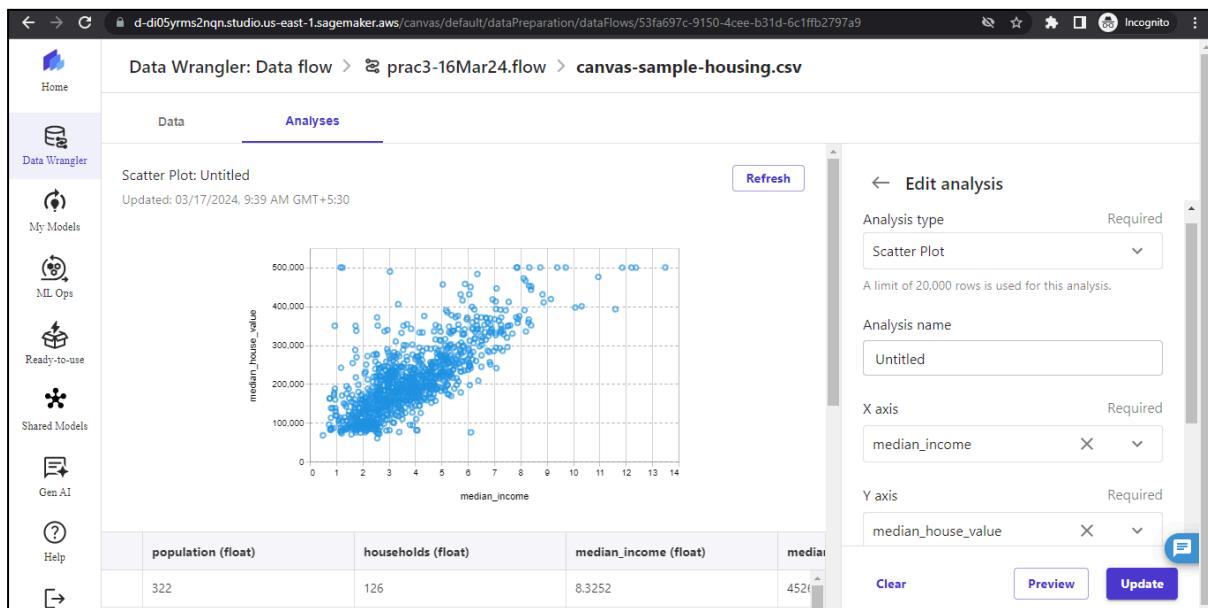
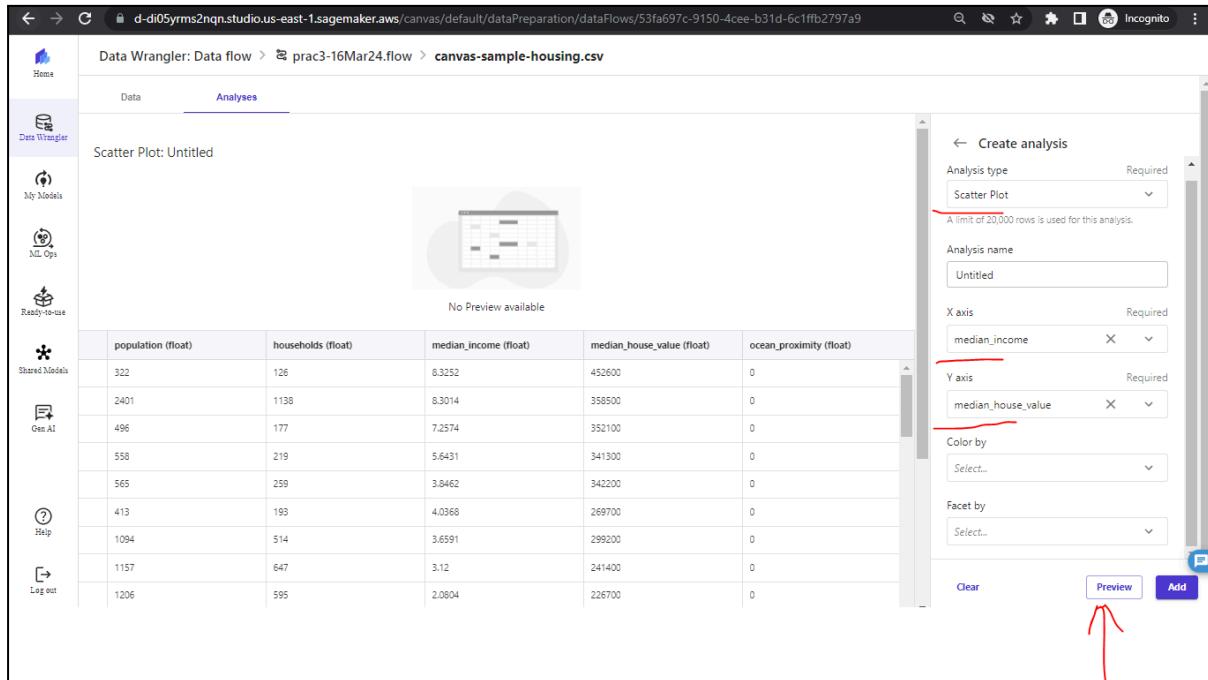
Data is normally distributed but slightly skewed towards the left

Since this is a sample dataset, we can proceed with a slightly skewed data also

Checking if the relationship between the predictor and target variable is linear

The **dependent variable** is called the **target** and the **independent variable** is called the **predictor**

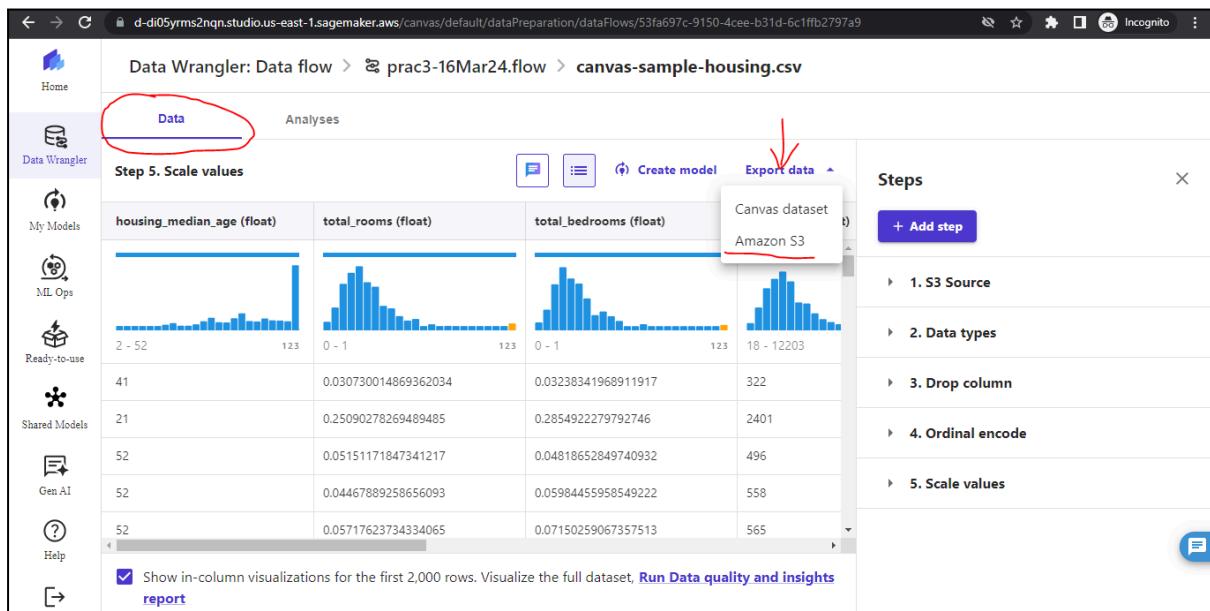
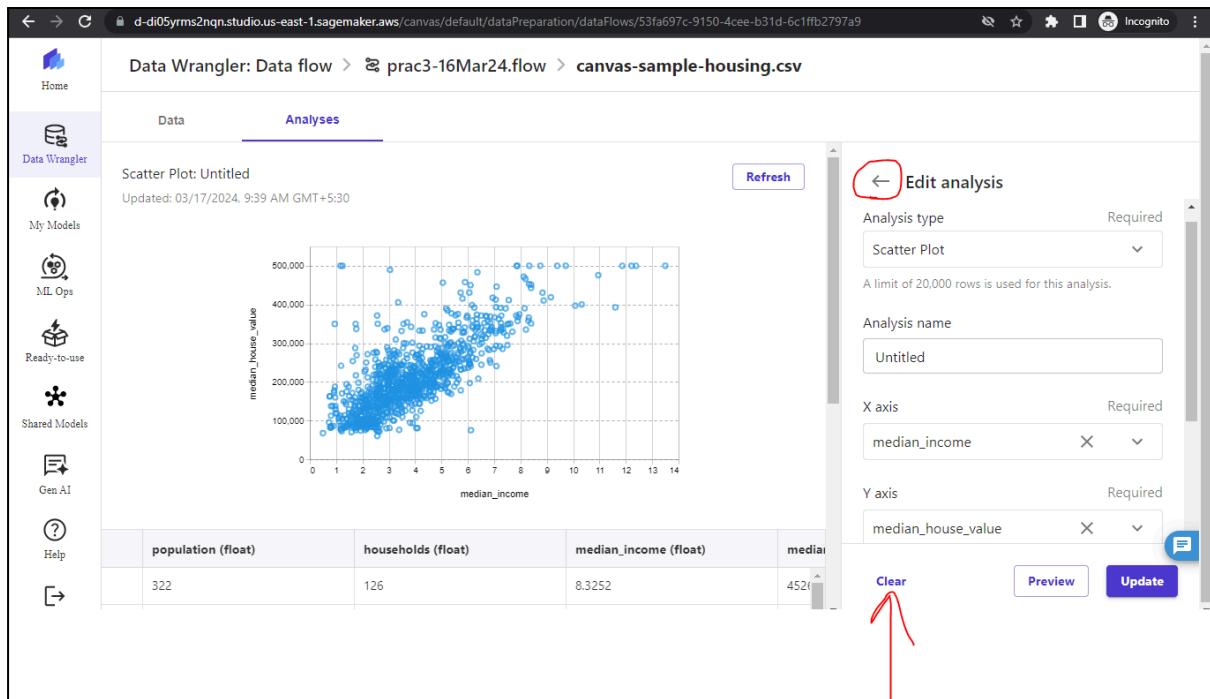
The dependent variable (target) depends on the independent variable (predictor)

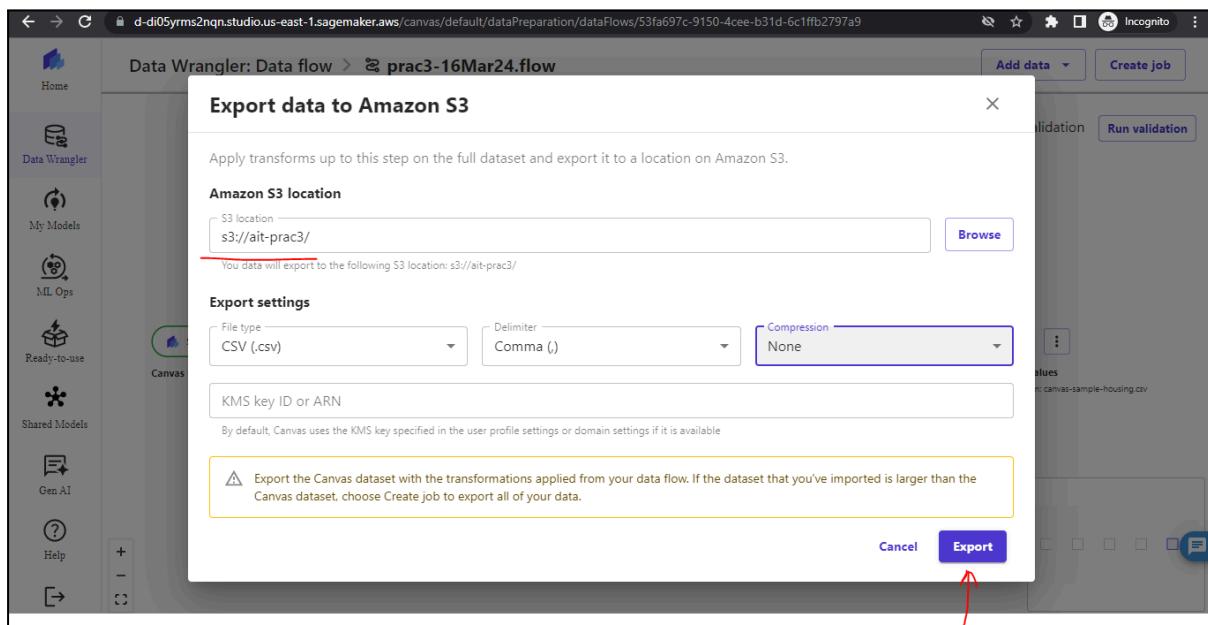


The columns' relation is roughly linear

The relationship looks roughly linear, so we can proceed with the linear regression.

Export prepared data





Logout of the Canvas after the data is exported



**Important: Delete the SageMaker-domain-user and the SageMaker-domain
Locate the exported CSV file in the bucket**

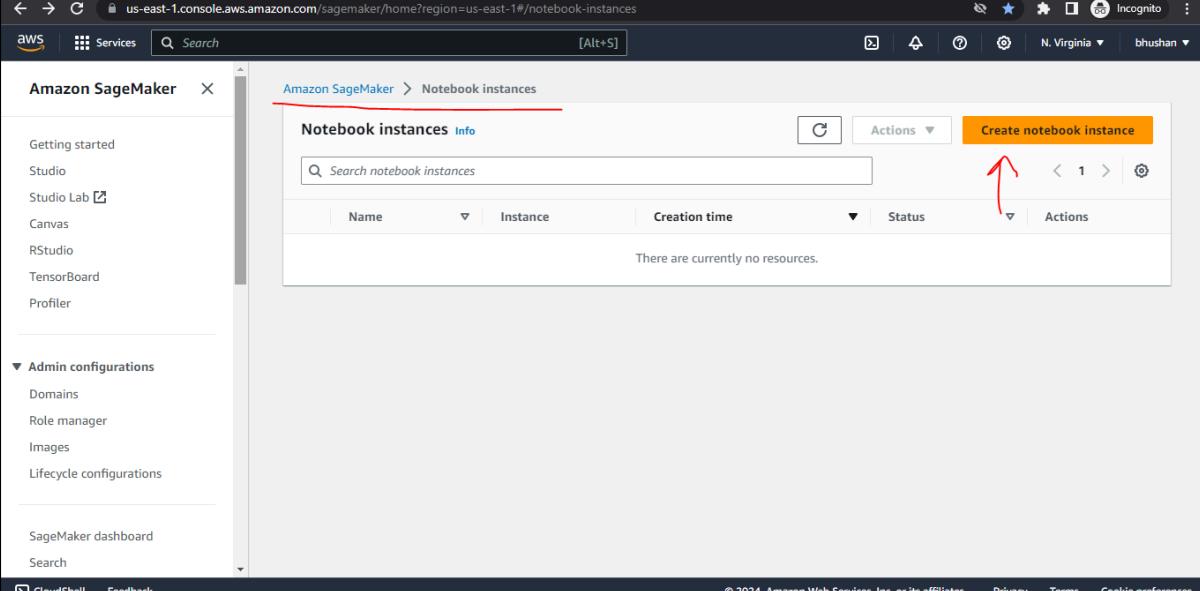
The screenshot shows the AWS S3 console interface. On the left, there's a sidebar with various AWS services like Buckets, Access Grants, and Storage Lens. The main area shows the 'ait-prac3' bucket. At the top, there are tabs for Objects, Properties, Permissions, Metrics, Management, and Access Points. Below that, there's a section for 'Objects (1) info' with a 'Create folder' and 'Upload' button. A red arrow points to the 'Upload' button. The object list shows one item: 'output_1710650355/' which is a Folder. Another red arrow points to this folder. The bottom of the screen shows copyright information and links for CloudShell, Feedback, Privacy, Terms, and Cookie preferences.

Copy the S3 URI of the csv file (to be used later in the python program)

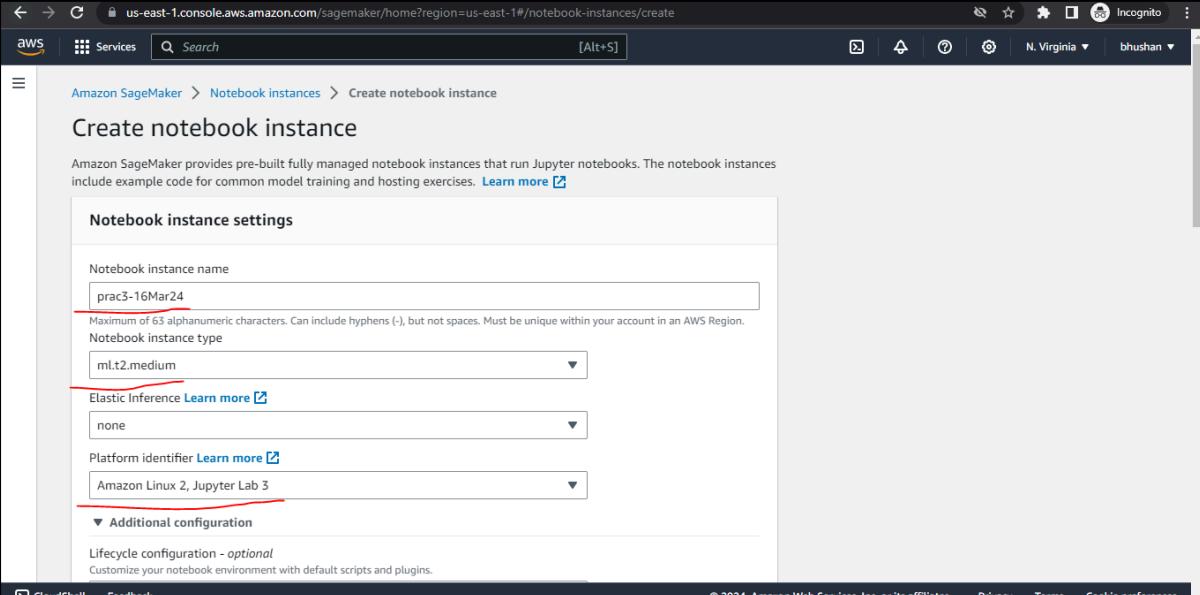
This screenshot shows the contents of the 'output_1710650355/' folder within the 'ait-prac3' bucket. The top right of the object list has a 'Copy S3 URI' button, which is highlighted with a red arrow. The object list itself shows one item: 'c000.csv'. This file is a CSV file, as indicated by the 'csv' extension in the table. The table includes columns for Name, Type, Last modified, Size, and Storage class. The 'Last modified' column shows 'March 17, 2024, 10:09:17 (UTC+05:30)'. The 'Size' column shows '75.3 KB'. The 'Storage class' column shows 'Standard'. The bottom of the screen shows copyright information and links for CloudShell, Feedback, Privacy, Terms, and Cookie preferences.

Simple Linear Regression using SageMaker

Create a Notebook



The screenshot shows the AWS SageMaker console with the 'Notebook instances' page open. On the left, there's a sidebar with various options like 'Getting started', 'Studio', 'Studio Lab', etc. The main area shows a table with columns 'Name', 'Instance', 'Creation time', and 'Status'. A message at the bottom says 'There are currently no resources.' An orange button labeled 'Create notebook instance' is visible at the top right of the main area. A red arrow points to this button.



The screenshot shows the 'Create notebook instance' configuration page. It has a header 'Create notebook instance' and a sub-header 'Notebook instance settings'. It includes fields for 'Notebook instance name' (containing 'prac3-16Mar24'), 'Notebook instance type' (set to 'ml.t2.medium'), 'Elastic Inference' (set to 'none'), and 'Platform identifier' (set to 'Amazon Linux 2, Jupyter Lab 3'). There's also a section for 'Additional configuration' and a note about 'Lifecycle configuration - optional'. A red arrow points to the 'Notebook instance name' field.

Click on "create" to create the notebook

Open in Jupyter to add the code

The screenshot shows the 'Notebook instances' page in the Amazon SageMaker console. On the left, there's a sidebar with links like 'Getting started', 'Studio', 'Canvas', etc. The main area lists notebook instances with columns for Name, Instance, Creation time, Status, and Actions. One instance, 'prac3-16Mar24', is selected and has a green 'InService' status. A red arrow points to the 'Open Jupyter' button in the Actions column for this instance.

The screenshot shows the Jupyter interface within the SageMaker environment. The top navigation bar includes 'Open JupyterLab', 'Quit', and 'Logout'. Below it, a sidebar has tabs for 'Files', 'Running', 'Clusters', 'SageMaker Examples', and 'Conda'. A message says 'The notebook list is empty.' To the right, a 'New' button is highlighted with a red arrow. A dropdown menu is open, listing various notebook types: R, Sparkmagic (PySpark), Sparkmagic (Spark), Sparkmagic (SparkR), conda_python3 (which is underlined in red), conda_pytorch_p310, conda_tensorflow2_p310, Other: Text File, Folder, and Terminal.

Add the following code

```
# Using linear regression algorithm from sklearn library
import pandas as pd
import boto3
from io import StringIO
from sklearn.model_selection import train_test_split
from sklearn.linear_model import LinearRegression
from sklearn.metrics import r2_score
# Load the dataset
s3_bucket = 'ait-prac3'
```

```
s3_key =
'output_1710650355/part-00000-d8518c19-068e-43b2-a81d-d6dfe090c63f-c000.csv'
s3_client = boto3.client('s3')
response = s3_client.get_object(Bucket=s3_bucket, Key=s3_key)
data = pd.read_csv(response['Body'])
# data = pd.read_csv(response['Body'], na_values=['NA', 'N/A', 'NaN', ""])
# Split the data into features (X) and target variable (y)
X = data[['median_income']]
y = data[['median_house_value']] # Notice the double square brackets to keep it as a
DataFrame
# Split the data into training and testing sets
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=42)
# Initialize and train the linear regression model
model = LinearRegression()
model.fit(X_train, y_train)
# Make predictions on the test set
predictions = model.predict(X_test)
#print(predictions)
# Extract the coefficients
# slope = model.coef_[0][0]
# intercept = model.intercept_[0]
# Print the equation of the linear regression line
# print(f"Linear Regression Equation: y = {slope:.2f}x + {intercept:.2f}")
## Evaluate the model's accuracy
test_predictions = model.predict(X_test)
r_squared = r2_score(y_test, test_predictions)
print("R-squared:", r_squared)
# using SageMaker's built-in Linear Regression algorithm
# import sagemaker
# from sagemaker import get_execution_role
# from sagemaker.amazon.amazon_estimator import get_image_uri
# from sagemaker.session import s3_input
# Define your S3 bucket and prefix
# bucket = 'your-s3-bucket'
# prefix = 'your-prefix'
```

```
# Set up SageMaker session and role
# sagemaker_session = sagemaker.Session()
# role = get_execution_role()
# Specify the location of your dataset in S3
# train_data = 's3://{}//{}//{}'.format(bucket, prefix, 'canvas-sample-housing.csv')
# Create a LinearLearner estimator
# linear_container = get_image_uri(sagemaker_session.boto_region_name, 'linear-learner')
# linear = sagemaker.estimator.Estimator(linear_container,
#                                         role,
#                                         train_instance_count=1,
#                                         # train_instance_type='ml.m4.xlarge',
#                                         train_instance_type="",
#                                         output_path='s3://{}//{}//output'.format(bucket, prefix),
#                                         sagemaker_session=sagemaker_session)

# Set hyperparameters
# linear.set_hyperparameters(predictor_type='regressor',
#                            mini_batch_size=50,
#                            epochs=3)

# Train the model
# linear.fit({'train': s3_input(train_data, content_type='text/csv')})
```

Output

The screenshot shows a Jupyter Notebook interface with the following sections:

- Initialize and train the linear regression model**

```
In [36]: model = LinearRegression()
model.fit(X_train, y_train)
```

Out[36]:

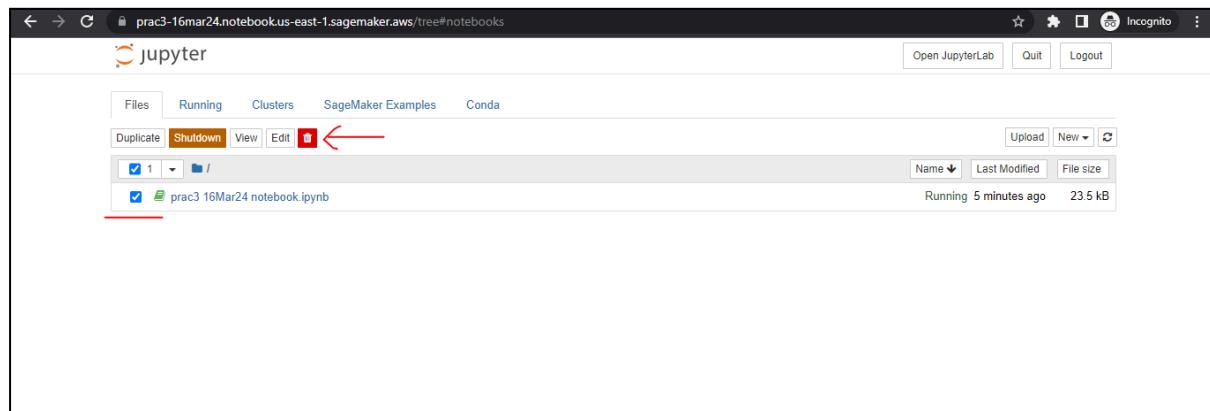
```
+ LinearRegression
LinearRegression()
```
- Make predictions on the test set**

```
In [37]: predictions = model.predict(X_test)
#print(predictions)
```
- Evaluate the model's accuracy**

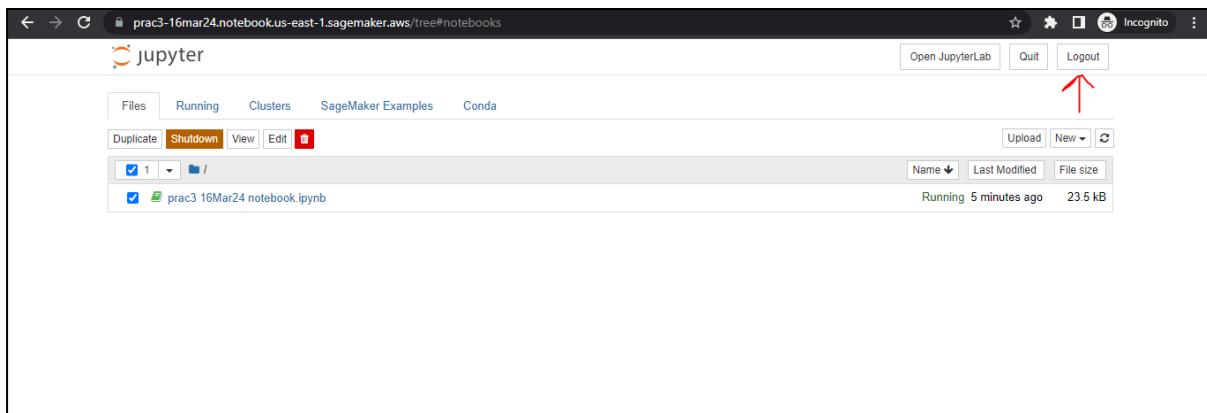
```
In [38]: test_predictions = model.predict(X_test)
r_squared = r2_score(y_test, test_predictions)
print("R-squared:", r_squared)
R-squared: 0.668378885837168
```
- in [39]:** ##### using SageMaker's built-in linear Regression algorithm
import sagemaker
from sagemaker import get_execution_role
from sagemaker.amazon.amazon_estimator import get_image_uri
from sagemaker.session import s3_input

Since the R^2 value is 0.66, we can say that the linear regression model is 66% accurate

Delete the Jupyter Notebook



Logout from Jupyter



CLEAN UP

1. Stop all the services that are running.
2. Delete all the resources that were created for this practical (including the buckets and the notebooks, etc.)

Conclusion: We saw data Preparation and EDA using AWS SageMaker DataWrangler, linear regression using SageMaker.

Name of Student: Pushkar Sane

Roll Number: 45

Lab Assignment Number: 7, 8, 9

Title of Lab Assignment: XGBoost Algorithm, Data Split for Sagemaker, Hyperparameter Search Optimization, Hyperparameters Optimization Strategies, Bias Variance Tradeoff, L2 Regularization (Ridge Regression), L1 Regularization (Lasso Regression), Hyperparameters Optimization Using GridSearchCV.

DOP: 21-02-2024

DOS: 28-03-2024

CO Mapped:

CO4

PO Mapped:

**PO1, PO2, PO3, PO4, PO5,
PO6, PO7, PO8, PO9, PO11,
PO12, PSO1, PSO2**

Signature:

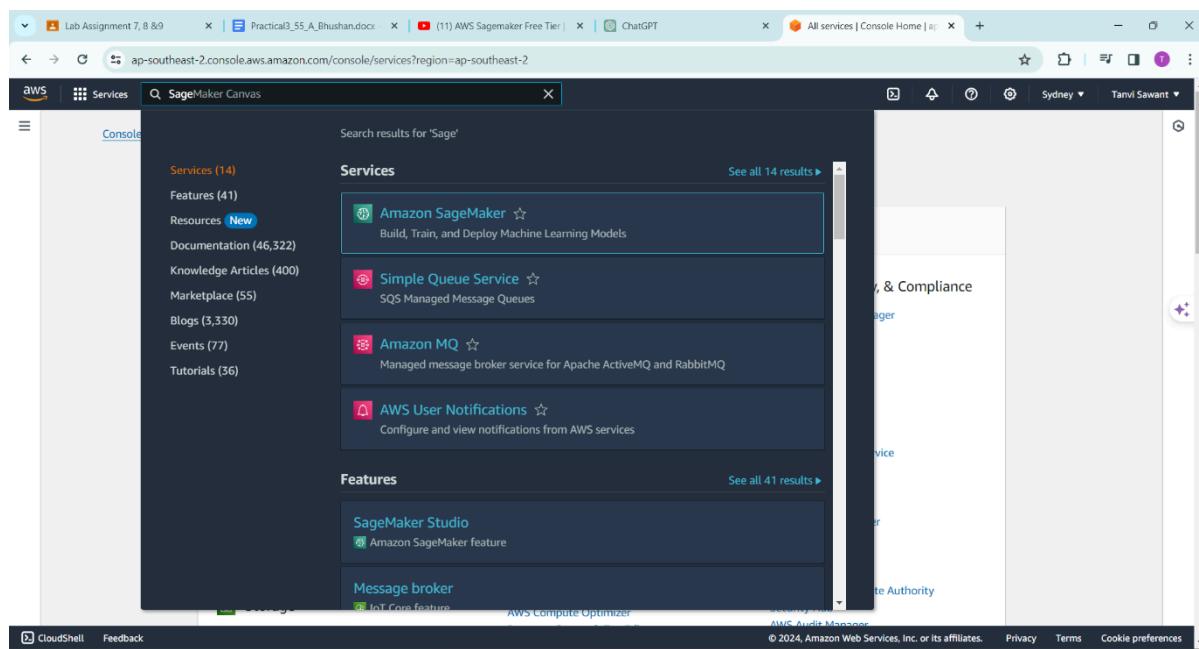
Practical No. 7, 8, 9

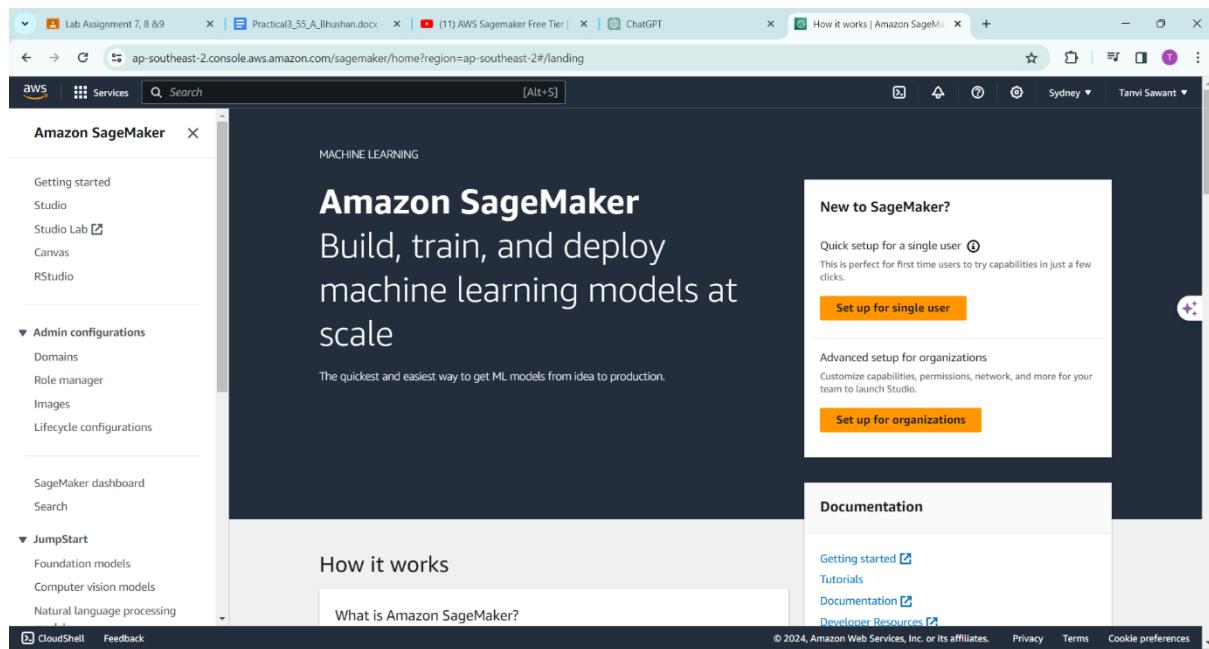
Sign Up for AWS Account:

- Go to the AWS website and sign up for a free tier account if you haven't already.
- Follow the steps to create your account, providing necessary information and setting up your billing preferences.

Navigate to Amazon SageMaker:

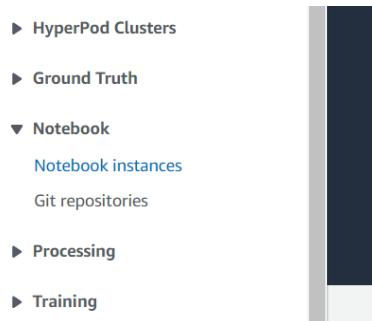
- Once logged in to your AWS Management Console, navigate to the Amazon SageMaker service.
- Click on "Services" in the top-left corner, then type "SageMaker" in the search bar, and click on the SageMaker service.





Create a Notebook Instance:

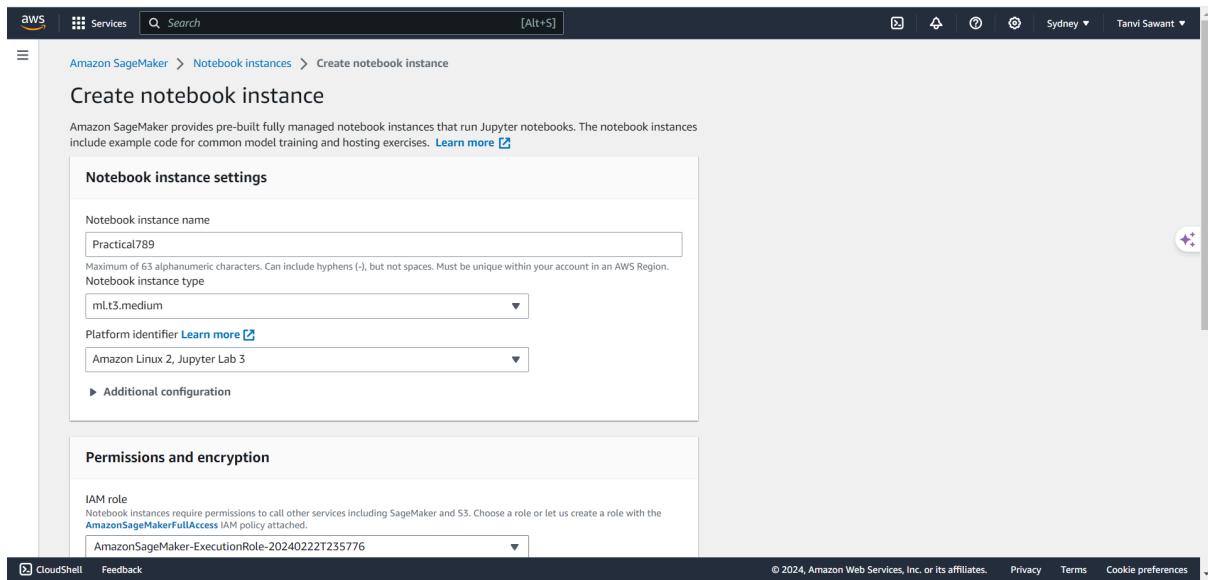
- In the SageMaker dashboard, click on "Notebook instances" in the left sidebar.
- Click on the "Create notebook instance" button.
- Enter a name for your notebook instance, choose an instance type eligible for the free tier (like **ml.t2.medium**), and leave other settings as default.



Configure Notebook Instance:

- In the "Notebook instance settings" section, give your notebook instance a name. Something descriptive like "MyMLNotebook" would work.
- For "Notebook instance type," choose an instance type that is eligible for the free tier, such as **ml.t2.medium**. This ensures you don't incur any charges beyond the free tier limits.

- Leave other settings as default for now. You can explore advanced settings later as needed.



Create IAM Role (Optional):

If you haven't created an IAM role before, you might be prompted to do so. This role grants necessary permissions to SageMaker to access other AWS services on your behalf.

Click on "Create a new role" if prompted, and follow the instructions to create a new IAM role. You can give it a name like "SageMakerRole".

Encryption (Optional):

If you want to encrypt the data stored on the notebook instance's EBS volume, you can choose an AWS Key Management Service (KMS) key. This step is optional for the purposes of this tutorial.

Networking Configuration (Optional):

You can choose a VPC (Virtual Private Cloud) and Subnet for your notebook instance if you have specific networking requirements. For simplicity, you can leave this as default.

Create Notebook Instance:

Once you've configured the settings according to your preferences, scroll down and click on the "Create notebook instance" button at the bottom.

Amazon SageMaker will now provision your notebook instance, which may take a few minutes.

This screenshot shows the first step of the 'Create Notebook Instance' wizard. It is titled 'Permissions and encryption'. Under the 'IAM role' section, it says 'Notebook instances require permissions to call other services including SageMaker and S3. Choose a role or let us create a role with the **AmazonSageMakerFullAccess** IAM policy attached.' A dropdown menu shows 'AmazonSageMaker-ExecutionRole-20240222T235776'. Below this is a 'Create role using the role creation wizard' button. Under 'Root access - optional', the 'Enable' radio button is selected. Under 'Encryption key - optional', it says 'Encrypt your notebook data. Choose an existing KMS key or enter a key's ARN.' and shows 'No Custom Encryption' selected. There are also sections for 'Network - optional' and 'Git repositories - optional'.

This screenshot continues the 'Create Notebook Instance' wizard. It shows the same configuration as the previous screenshot, including the IAM role selection and optional root access and encryption settings. The sections for 'Network - optional' and 'Git repositories - optional' are also present. At the bottom right, there is a large orange 'Create notebook instance' button.

This screenshot shows the 'Notebook instances' page in the AWS SageMaker console. A green success message at the top says 'Success! Your notebook instance is being created. Open the notebook instance when status is InService and open a template notebook to get started.' Below this, the 'Notebook instances' table lists one instance: 'Practical789' (ml.t3.medium), created on 3/28/2024, 9:16:38 PM, with a status of 'Pending'. There is a 'View details' button next to the instance name.

Notebook instance settings

Name	Status	Notebook instance type	Platform identifier
Practical789	Pending	mLt3.medium	Amazon Linux 2, Jupyter Lab 3 (notebook-al2-v2)

ARN	Creation time	Volume Size	Minimum IMDS Version
arn:aws:sagemaker:ap-southeast-2:058264124894:notebook-instance/Practical789	Mar 28, 2024 15:46 UTC	5GB EBS	2

Lifecycle configuration	Last updated
-	Mar 28, 2024 15:46 UTC

Git repositories

Name	Repository URL	Type
There are currently no resources.		

Notebook instances

Name	Instance	Creation time	Status	Actions
Practical789	mLt3.medium	3/28/2024, 9:16:38 PM	InService	Open Jupyter Open JupyterLab

Open Jupyter Notebook:

- Once the notebook instance is created (it may take a few minutes), click on "Open Jupyter" next to your notebook instance name.
- This will open a Jupyter Notebook interface in a new tab.

practical789.notebook.ap-southeast-2.sagemaker.aws/tree

jupyter

Files Running Clusters SageMaker Examples Conda

Select items to perform actions on them.

The notebook list is empty.

Dataset: Titanic: Machine Learning from Disaster

This dataset contains information about passengers aboard the Titanic, including details like age, gender, passenger class, ticket fare, and whether they survived the disaster. It's commonly used for classification tasks and is suitable for practicing data preprocessing, model building, and evaluation.

Here are the steps to download the dataset from Kaggle:

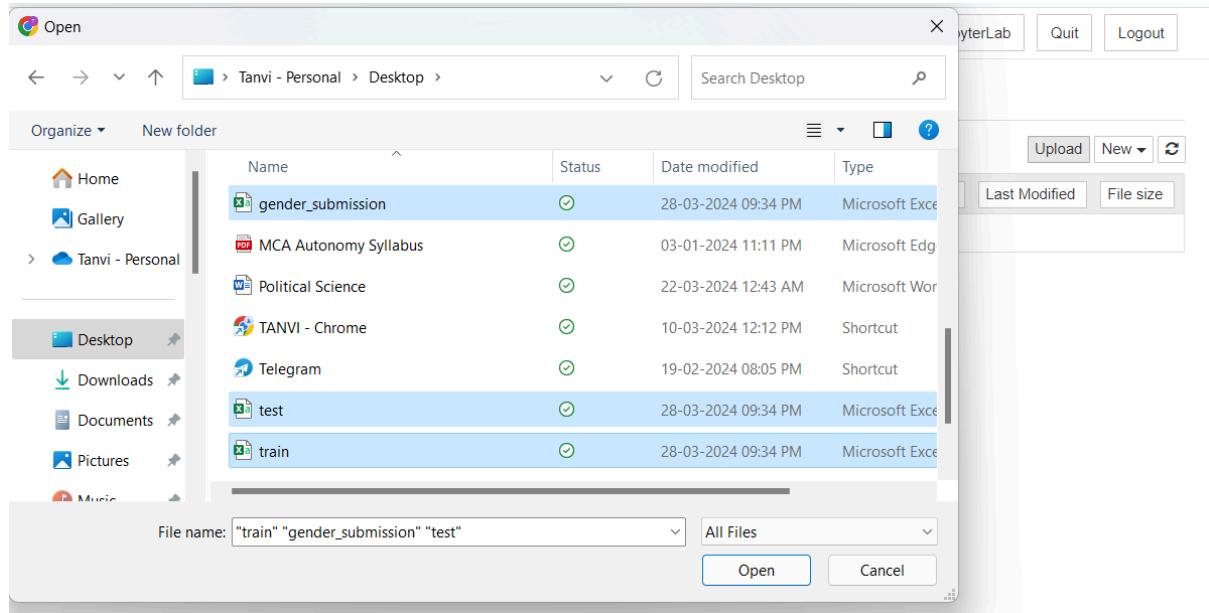
1. Sign in to your Kaggle account or create a new one.
2. Navigate to the Titanic competition page.
3. Scroll down to the "Data" section.
4. Click on the "Download All" button to download the entire dataset as a zip file.
5. Extract the zip file to access the dataset files (e.g., train.csv, test.csv).

Once you have the dataset files, you can upload them to your Jupyter Notebook environment on AWS SageMaker as explained earlier. This dataset is well-suited for practicing data preprocessing, model training, and evaluation tasks.

The screenshot shows the Kaggle competition page for 'Titanic - Machine Learning from Disaster'. At the top, there is a search bar and a user profile icon. Below the title, there are tabs for Overview, Data (which is selected), Code, Models, Discussion, Leaderboard, and Rules. The Data tab displays a summary of the dataset with four bars: 892, 1309, 0, and 1. Below this, there is a table with 12 rows, each containing a number from 892 to 901. To the right of the table, there is a sidebar titled 'Summary' showing 3 files and 25 columns. At the bottom of the page, there is a terminal window displaying the command: >_ kaggle competitions download -c titanic. There are also navigation icons at the bottom right.

Prepare Your Data:

- Upload the dataset files (**train.csv**, **test.csv**, **gender_submission.csv**) to your Jupyter Notebook environment.



practical789.notebook.ap-southeast-2.sagemaker.aws/tree

jupyter

Files Running Clusters SageMaker Examples Conda

Select items to perform actions on them.

The notebook list is empty

Name	Last Modified	File size
train.csv		
test.csv		
gender_submission.csv		

Upload Cancel

Upload Cancel

Upload Cancel

Files Running Clusters SageMaker Examples Conda

Select items to perform actions on them.

The notebook list is empty

Name	Last Modified	File size
gender_submission.csv	seconds ago	3.26 kB
test.csv	seconds ago	28.6 kB
train.csv	seconds ago	61.2 kB

Upload New

- Use pandas to load the datasets into separate DataFrames.
- Perform any necessary data cleaning and preprocessing steps.

Create a New Notebook:

- In the Jupyter Notebook interface, click on the "New" button on the top right corner.
- Select "Python 3" to create a new Python notebook.

Prepare Your Data:

- In this step, you'll upload your datasets (test.csv, train.csv, gender_submission.csv) to your Jupyter Notebook environment and start working with them.

```
In [16]: import pandas as pd

# Upload the datasets
train_data = pd.read_csv('train.csv')
test_data = pd.read_csv('test.csv')
submission_data = pd.read_csv('gender_submission.csv')

# Now you can start working with your datasets
```

What is XGBoost Algorithm?

XGBoost is a robust machine-learning algorithm that can help you understand your data and make better decisions.

XGBoost is an implementation of gradient-boosting decision trees. It has been used by data scientists and researchers worldwide to optimize their machine-learning models.

What is XGBoost in Machine Learning?

XGBoost is designed for speed, ease of use, and performance on large datasets. It does not require optimization of the parameters or tuning, which means that it can be used immediately after installation without any further configuration.

XGBoost Features

XGBoost is a widespread implementation of gradient boosting.

- ❖ XGBoost offers regularization, which allows you to control overfitting by introducing L1/L2 penalties on the weights and biases of each tree.
- ❖ Another feature of XGBoost is its ability to handle sparse data sets using the weighted quantile sketch algorithm. This algorithm allows us to deal with non-zero entries in the feature matrix while retaining the same computational complexity as other algorithms.

- ❖ XGBoost also has a block structure for parallel learning. It makes it easy to scale up on multicore machines or clusters. It also uses cache awareness, which helps reduce memory usage when training models with large datasets.
- ❖ XGBoost offers out-of-core computing capabilities using disk-based data structures instead of in-memory ones during the computation phase.

XgBoost Formula

The XGBoost algorithm builds an ensemble of decision trees in a sequential manner, where each subsequent tree corrects the errors made by the previous ones. The final prediction is a weighted sum of the predictions from all the trees in the ensemble. The general formula for predicting the output of a new data point using XGBoost can be expressed as follows

$$\hat{y}_i = \phi(\mathbf{x}_i) = \sum_{k=1}^K f_k(\mathbf{x}_i)$$

where:

- \hat{y}_i is the predicted output for the i -th data point.
- \mathbf{x}_i represents the feature vector for the i -th data point.
- K is the total number of trees in the ensemble.
- $f_k(\mathbf{x}_i)$ is the prediction made by the k -th tree for the i -th data point.

Conclusion: Overall, XGBoost is widely used in both academia and industry due to its excellent performance, versatility, and efficiency in handling various machine learning tasks. Performing data splitting, hyperparameter search optimization, and hyperparameter optimization strategies in AWS SageMaker involves several steps. Below is a detailed guide on how to perform these tasks:

1. Data Split for SageMaker.

Before training a machine learning model, it is essential to split the dataset into training, validation, and test sets. SageMaker provides utilities to perform this data split.

```
import sagemaker
from sagemaker import get_execution_role
from sklearn.model_selection import train_test_split
import pandas as pd
import numpy as np
```

```
# Assume you have a dataframe df with features and target
# Split data into train, validation, and test sets
train_data, temp_data = train_test_split(df, test_size=0.3, random_state=42)
val_data, test_data = train_test_split(temp_data, test_size=0.5, random_state=42)

# Save the data to CSV files
train_data.to_csv('train.csv', index=False, header=False)
val_data.to_csv('validation.csv', index=False, header=False)
test_data.to_csv('test.csv', index=False, header=False)

# Upload the data to S3
role = get_execution_role()
sess = sagemaker.Session()
bucket = sess.default_bucket()

train_data_s3_path = sess.upload_data(path='train.csv', bucket=bucket, key_prefix='data')
val_data_s3_path      =      sess.upload_data(path='validation.csv',      bucket=bucket,
key_prefix='data')
test_data_s3_path = sess.upload_data(path='test.csv', bucket=bucket, key_prefix='data')
```

2. Hyperparameter Search Optimization

SageMaker provides built-in hyperparameter optimization (HPO) to automate the process of finding the best hyperparameters for your model.

```
from sagemaker.tuner import HyperparameterTuner, IntegerParameter,
ContinuousParameter

# Define hyperparameter ranges
hyperparameter_ranges = {
    'learning_rate': ContinuousParameter(0.001, 0.1),
    'num_trees': IntegerParameter(50, 200),
    'max_depth': IntegerParameter(5, 15)
}
```

```
# Create an estimator
estimator = sagemaker.estimator.Estimator(image_uri='image-uri', # specify your image uri
                                            role=role,
                                            instance_count=1,
                                            instance_type='ml.m5.large',
                                            hyperparameters={'objective': 'regression'},
                                            output_path='s3://{}{}/output'.format(bucket))
```

```
# Create a HyperparameterTuner object
tuner = HyperparameterTuner(estimator=estimator,
                             objective_metric_name='validation:rmse', # specify the metric to optimize
                             objective_type='Minimize',
                             hyperparameter_ranges=hyperparameter_ranges,
                             max_jobs=20,
                             max_parallel_jobs=4)
```

```
# Fit the tuner
tuner.fit({'train': train_data_s3_path, 'validation': val_data_s3_path})
```

3. Hyperparameter Optimization Strategies

AWS SageMaker's HPO supports multiple optimization strategies, such as random search, Bayesian optimization, and more. The strategy parameter in the HyperparameterTuner can be used to specify the optimization strategy.

1. Random Search

Random search is a simple and effective method for hyperparameter optimization. SageMaker's HyperparameterTuner supports random search as one of its optimization strategies.

```
from sagemaker.tuner import HyperparameterTuner, ContinuousParameter,
IntegerParameter
```

```
# Define hyperparameter ranges
hyperparameter_ranges = {
    'learning_rate': ContinuousParameter(0.001, 0.1),
    'num_trees': IntegerParameter(50, 200),
    'max_depth': IntegerParameter(5, 15)
}
```

```
# Create an estimator
estimator = sagemaker.estimator.Estimator(image_uri='image-uri', # specify your image uri
                                            role=role,
                                            instance_count=1,
                                            instance_type='ml.m5.large',
                                            hyperparameters={'objective': 'regression'},
                                            output_path='s3://{}{}/output'.format(bucket))
```

```
# Create a HyperparameterTuner object with random search strategy
tuner_random = HyperparameterTuner(estimator=estimator,
                                     objective_metric_name='validation:rmse',
                                     objective_type='Minimize',
                                     hyperparameter_ranges=hyperparameter_ranges,
                                     max_jobs=20,
                                     max_parallel_jobs=4,
                                     strategy='Random')
```

```
# Fit the tuner
tuner_random.fit({'train': train_data_s3_path, 'validation': val_data_s3_path})
```

2. Bayesian Optimization

Bayesian optimization is an advanced method for hyperparameter optimization, which builds a probabilistic model of the objective function and uses it to select the most promising hyperparameters

```
# Create a HyperparameterTuner object with Bayesian optimization strategy
tuner_bayesian = HyperparameterTuner(estimator=estimator,
                                       objective_metric_name='validation:rmse',
                                       objective_type='Minimize',
                                       hyperparameter_ranges=hyperparameter_ranges,
                                       max_jobs=20,
                                       max_parallel_jobs=4,
                                       strategy='Bayesian')
```

Fit the tuner

```
tuner_bayesian.fit({'train': train_data_s3_path, 'validation': val_data_s3_path})
```

3. Customized Optimization Strategy

You can also implement a customized optimization strategy by subclassing the HyperparameterTunerStrategy class and implementing the required methods.

```
from sagemaker.tuner import HyperparameterTunerStrategy

class CustomStrategy(HyperparameterTunerStrategy):
    def hyperparameter_tuning_job_config(self):
        config = super(CustomStrategy, self).hyperparameter_tuning_job_config()
        config['HyperParameterTuningJobConfig']['Strategy'] = 'Custom'
        return config

# Create a HyperparameterTuner object with custom strategy
tuner_custom = HyperparameterTuner(estimator=estimator,
                                    objective_metric_name='validation:rmse',
                                    objective_type='Minimize',
                                    hyperparameter_ranges=hyperparameter_ranges,
                                    max_jobs=20,
                                    max_parallel_jobs=4,
                                    strategy=CustomStrategy)

# Fit the tuner
tuner_custom.fit({'train': train_data_s3_path, 'validation': val_data_s3_path})
```

4. Bias-Variance Tradeoff

The bias-variance tradeoff is a fundamental concept in machine learning that relates to the model's ability to generalize to unseen data. A high-bias model is too simplistic and may underfit the training data (high training error), whereas a high-variance model is too complex and may overfit the training data (low training error but high validation error).

To manage the bias-variance tradeoff, you can adjust hyperparameters and model complexity. SageMaker's Hyperparameter Optimization (HPO) can help find the optimal hyperparameters that balance bias and variance by searching for the hyperparameter values that result in the best validation performance.

1. Define Hyperparameter Ranges

```
from sagemaker.tuner import HyperparameterTuner, ContinuousParameter,  
IntegerParameter
```

```
# Define hyperparameter ranges  
hyperparameter_ranges = {  
    'learning_rate': ContinuousParameter(0.001, 0.1),  
    'num_trees': IntegerParameter(50, 200),  
    'max_depth': IntegerParameter(5, 15)  
}
```

2. Create an Estimator

```
import sagemaker
```

```
# Create an estimator  
estimator = sagemaker.estimator.Estimator(image_uri='image-uri', # specify your image uri  
                                            role=role,  
                                            instance_count=1,  
                                            instance_type='ml.m5.large',  
                                            hyperparameters={'objective': 'regression'},  
                                            output_path='s3://{}{}/output'.format(bucket))
```

3. Create a HyperparameterTuner Object

```
# Create a HyperparameterTuner object  
tuner = HyperparameterTuner(estimator=estimator,  
                           objective_metric_name='validation:rmse',  
                           objective_type='Minimize',  
                           hyperparameter_ranges=hyperparameter_ranges,  
                           max_jobs=20,  
                           max_parallel_jobs=4,  
                           strategy='Bayesian') # Use Bayesian optimization as the optimization  
strategy
```

4. Fit the HyperparameterTuner

```
# Fit the tuner
tuner.fit({'train': train_data_s3_path, 'validation': val_data_s3_path})
```

5. Analyze the Results

After the hyperparameter tuning job is completed, you can analyze the results to understand how different hyperparameters affect the bias-variance tradeoff. You can use the Amazon SageMaker console or SDK to analyze the results and identify the optimal hyperparameters that balance bias and variance.

```
# Get the best hyperparameters
best_hyperparameters = tuner.best_estimator().hyperparameters()
```

```
print("Best hyperparameters:")
print(best_hyperparameters)
```

5. L2 Regularization (Ridge Regression) and L1 Regularization (Lasso Regression)

```
hyperparameters={
    'alpha': 1.0, # L2 regularization parameter
    'objective': 'reg:linear'
}
```

```
hyperparameters={
    'alpha': 1.0, # L1 regularization parameter
    'objective': 'reg:linear',
    'lambda': 1.0 # L2 regularization parameter (for Lasso in XGBoost)
}
```

6. Hyperparameters Optimization Using GridSearchCV

```
from sklearn.model_selection import GridSearchCV
from sagemaker.sklearn.estimator import SKLearn
```

```
# Define hyperparameter grid
param_grid = {
    'alpha': [0.1, 1.0, 10.0],
    'learning_rate': [0.001, 0.01, 0.1],
    'max_depth': [3, 5, 7],
```

```
'n_estimators': [50, 100, 200]
}

# Create a SKLearn estimator
sklearn_estimator = SKLearn(entry_point='script.py', # your script
                           role=role,
                           instance_count=1,
                           instance_type='ml.m5.large',
                           framework_version='0.23-1',
                           sagemaker_session=sess)

# Create GridSearchCV object
grid_search = GridSearchCV(estimator=sklearn_estimator,
                           param_grid=param_grid,
                           cv=3,
                           scoring='neg_mean_squared_error',
                           n_jobs=-1)

# Fit GridSearchCV
grid_search.fit({'train': train_data_s3_path})
```

Conclusion

In this guide, we explored the essential aspects of hyperparameter optimization (HPO) in AWS SageMaker, focusing on managing the bias-variance tradeoff, which is a fundamental concept in machine learning.

We began by defining the hyperparameter ranges that we wanted to optimize, such as the learning rate, number of trees, and maximum depth. These hyperparameters play a crucial role in determining the complexity of the model and, consequently, its ability to generalize to unseen data.

Next, we created an estimator using SageMaker's Estimator class, specifying the necessary configurations such as the Docker image URI, IAM role, instance type, and hyperparameters.

We then utilized SageMaker's HyperparameterTuner to perform the hyperparameter optimization. We selected Bayesian optimization as the optimization strategy, which is an advanced method that builds a probabilistic model of the objective function to select the most promising hyperparameters. The HyperparameterTuner automatically searched for the

optimal hyperparameters within the defined ranges, aiming to minimize the validation root mean squared error (RMSE).

After fitting the HyperparameterTuner, we analyzed the results to retrieve the best hyperparameters that balance the bias-variance tradeoff effectively. By understanding how different hyperparameters affect the model's generalization performance, we can make informed decisions to adjust the model's complexity and improve its ability to generalize to unseen data.

In conclusion, AWS SageMaker's Hyperparameter Optimization provides a powerful and efficient way to find the optimal hyperparameters that balance bias and variance, thereby enhancing the model's performance and generalization capabilities. By carefully tuning the hyperparameters and analyzing the results, machine learning practitioners can develop models that deliver superior performance on real-world tasks.