

Cloud Computing Amazon SageMaker GroundTruth

Farid Afzali, Ph.D., P.Eng.

AWS SageMaker GroundTruth



- AWS SageMaker GroundTruth is delineated as a service for data labeling, pivotal in the preparation of datasets for machine learning.
- The term 'ground truth' is defined within the machine learning lexicon as the 'gold standard'. This standard is used as the benchmark for evaluating the predictive accuracy of machine learning models.
- Ground Truth serves to identify the correct classification or outcome that a machine learning model is expected to predict, which is essential for supervised learning where the model's output is compared against known labeled data.
- The service encompasses various data labeling tasks, such as:
 - Image Classification: Categorizing images into predefined classes.
 - Bounding Box: Identifying and marking the location of objects within images using rectangular shapes.
 - Semantic Segmentation: Labeling specific pixels of an image to differentiate between objects and background, thus understanding the image at the pixel level.
 - Label Verification: Ensuring the accuracy of labels provided to the data, which is a quality control step to confirm whether the assigned labels are correct or incorrect.

AWS SageMaker GroundTruth Tasks



- <https://aws.amazon.com/sagemaker/data-labeling/features/>
- **Bounding Boxes:** This task involves drawing rectangles over objects within images. The coordinates of the box edges are used to train models to locate and identify objects in other images. It's a common approach used in computer vision for object detection and localization.
- **Image Classification:** In this task, images are categorized into predefined labels or classes. Each image is assigned a label that represents the content of the image. This is used to train models to recognize and classify images into these learned categories.
- **Semantic Segmentation:** This involves labeling each pixel in an image with a class that denotes what object the pixel belongs to. Unlike bounding boxes that locate objects, semantic segmentation differentiates between objects at the pixel level, enabling detailed understanding and analysis of images. It's particularly useful in applications like autonomous driving where the precise location of objects is crucial.
- **Text Classification:** This task is about assigning predefined labels to text data. For instance, classifying email messages into 'spam' or 'not spam'. This form of classification helps in training models to understand and categorize text data based on content.
- **Custom Tasks:** SageMaker GroundTruth also allows for the creation of custom labeling tasks. This flexibility means that if a project has specific requirements not covered by the standard tasks, a custom workflow can be designed. For example, labeling audio data, video frames, or complex document analysis that may require a combination of labeling techniques. Custom tasks can be tailored with specific instructions and interfaces to meet the unique needs of any machine learning project.

Tasks Available Data Labeling Workforces



Public Mechanical Turks: Amazon SageMaker connects customers with a large, on-demand workforce known as Mechanical Turk workers. This global pool consists of approximately 500,000 contractors who are available 24/7 to perform data labeling tasks.

Private: Customers have the option to use a private team of labelers for their data labeling needs. This can include the customer's own private labelers, and all management and coordination are conducted through SageMaker GroundTruth. It is noted that there is no requirement for these labelers to have an IAM (AWS Identity and Access Management) or an Amazon account.

Vendors: SageMaker GroundTruth also offers access to professional data labeling services through a network of third-party vendors. These vendors are curated and can provide specialized data labeling services. This option is for customers seeking professional and potentially more sophisticated data labeling services that may not be available from the public Mechanical Turk workforce or private labelers.

Four Primary Categories of Data Sources

Image/Video:

- This category includes visual data, both still and moving. Such data is rich in content and can be used for a variety of applications including facial recognition, object detection, video analytics, and more. Analyzing image and video data often requires advanced techniques like computer vision and deep learning.

Text (Corpus):

- Textual data sources, or corpora, consist of collections of written texts, which could be in the form of books, articles, social media posts, etc. Text data is used for natural language processing (NLP) tasks such as sentiment analysis, topic modeling, and language translation.

Audio/Sound:

- Audio data sources encompass sounds, speech, and other auditory signals. Analyzing audio data can involve speech recognition, music classification, and environmental sound analysis, among other tasks. Techniques such as signal processing and deep learning are commonly employed here.

Time Series/Signals:

- This refers to data that is collected over time, often at regular intervals. Examples include stock prices, weather data, and physiological signals like heart rate. Time series analysis is crucial for forecasting, anomaly detection, and understanding temporal dynamics.

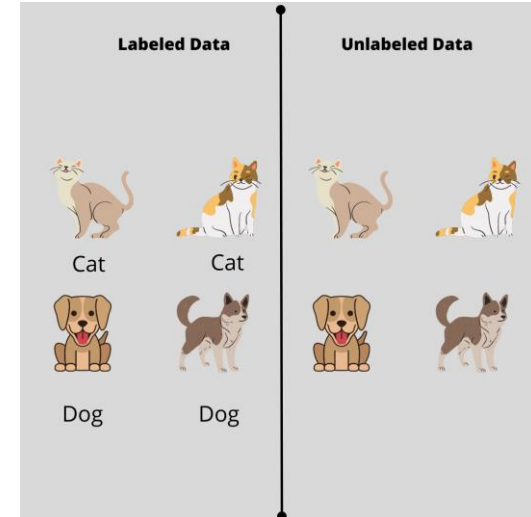
Labeled and Unlabeled Datasets

- **Unlabeled Dataset:**

- ❑ This type of dataset comprises data points without any associated explanations, classes, or tags. For instance, images of animals are provided without any information indicating the type of animal in each image. Unlabeled data is typically used in unsupervised learning where the goal is to identify patterns or structures within the data without prior knowledge.

- **Labeled Dataset:**

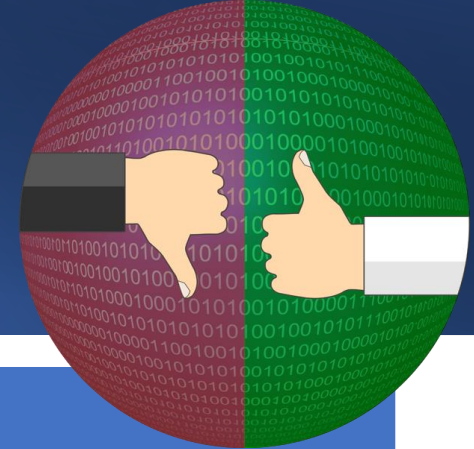
- ❑ In contrast, a labeled dataset contains data points that are each associated with a specific label, class, or tag that provides information about the data. Using the same example, images of animals would be labeled with tags such as "cat" or "dog" to indicate the content of the image. Labeled data is essential for supervised learning, where an AI model learns to predict the labels of new data based on the patterns it has learned from the labeled training data.



[Source](#)

Differentiate Good Data from Bad Data

[Source](#)



Good Data:

- **Abundant Samples:** Good data is characterized by having a large number of data points, which can provide a more comprehensive view of the phenomenon being studied.
- **Unbiased:** It is representative of the population or phenomena without any inherent biases that could skew results.
- **No Missing Data Points:** Good data sets are complete or have minimal missing values, which contribute to more accurate and reliable analyses.
- **Relevant Features:** It contains features (variables) that are important and relevant to the analysis or predictive modeling tasks at hand.
- **No Duplicate Samples:** Good data is devoid of duplicates, ensuring that each data point is unique and the analysis is not skewed by overrepresentation.

Bad Data:

- **Few Samples:** Bad data often suffers from a lack of enough data points, which can lead to overfitting in machine learning and unreliable conclusions in analysis.
- **Biased:** If data is biased, it does not accurately represent the population or phenomenon, leading to biased outcomes and decisions.
- **Missing Data Points:** Bad data often contains gaps or missing values, which can introduce significant uncertainty or bias into the results.
- **Irrelevant Features:** It includes features that are not useful for the analysis, which can complicate the modeling process and potentially reduce model performance.
- **Duplicate Samples:** Bad data may have duplicates, which can distort analytical results and affect the performance of predictive models.

Data Labeling

- **Requirement of Large Labeled Datasets:** Contemporary deep learning models, especially those considered state-of-the-art, rely heavily on extensive volumes of labeled data. The size of the dataset can significantly impact the performance of these models, as they need substantial amounts of data to learn the intricate patterns necessary for tasks such as image and speech recognition, natural language processing, and predictive analytics.
- **Need for a Large Human Labeling Team:** To create these large labeled datasets, a considerable workforce is often necessary. Human labelers are employed to manually annotate and label the data, which can be a time-consuming and labor-intensive process. This is particularly true for tasks that require a high level of precision and cannot be easily automated.
- **Inherent Biases:** Human involvement in data labeling can introduce biases, whether from individual perceptions, cultural contexts, or misunderstanding of the labeling task. Such biases can adversely affect the model's performance, as it may learn and perpetuate these biases. Controls and standardization are necessary to minimize this risk, ensuring that the labels are as objective and accurate as possible.
- **Time Consumption for Data Scientists:** Data curation and labeling can consume a significant portion of data scientists' time. Before even beginning to develop and train models, data scientists often find themselves spending a great deal of time preparing the data—cleaning, labeling, and ensuring it is suitable for use in model training. This highlights the need for efficient data labeling processes and tools that can help streamline this aspect of the data scientist's role.

JSON Lines format

```
{"id": 1, "name": "John Doe"}  
{"id": 2, "name": "Jane Doe"}  
{"id": 3, "name": "John Smith"}
```

[Source](#)

- The JSON Lines format, often abbreviated as JSONL, is a convenient format for storing structured data, allowing it to be processed one record at a time. This can be particularly useful when dealing with large datasets or streaming data, where it's beneficial to read, process, and write data incrementally rather than loading the entire dataset into memory.
- In the context of the JSONL provided:
 - Each line is a complete JSON object.
 - Every line represents a separate data record.
 - This format is line-delimited, which means each new record is placed on a new line.
- JSON Lines is particularly suited for data labeling as it allows each image or data point to be paired with its label in a format that's easy to read and write programmatically.

```
{"Image 1": "Cat"}  
{"Image 2": "Dog"}  
{"Image 3": "Lion"}
```

Manifest Files

[Source](#)

```
<?xml version="1.0" encoding="utf-8"?>
<omexManifest xmlns="http://identifiers.org/combine.specifications/omex-manifest">
  <content location="."
    format="http://identifiers.org/combine.specifications/omex"/>
  <content location="models/model.xml"
    format="http://identifiers.org/combine.specifications/sbml"/>
  <content location="simulation.xml" master="true"
    format="http://identifiers.org/combine.specifications/sed-ml"/>
  <content location="doc/article.pdf"
    format="http://purl.org/NET/mediatypes/application/pdf"/>
  <content location="metadata.rdf"
    format="http://identifiers.org/combine.specifications/omex-metadata"/>
</omexManifest>
```

- "Manifest Files," which are used in the context of Amazon SageMaker GroundTruth jobs for labeling data in machine learning projects, particularly in supervised learning tasks such as image classification:
- **Purpose of Manifest Files:** They are used in computer vision classification problems to store data about both the inputs (like images) and their corresponding outputs (labels). This is crucial for supervised machine learning models that require labeled data to learn from.
- **Structure of Manifest Files:** The files are structured in JSON Lines (JSONL) format, meaning each line in the file is a valid JSON object. This format is advantageous for processing large datasets because it allows for reading and writing one record at a time, which can be more efficient and require less memory than processing an entire JSON array.
- **Content of a Manifest File:** Each line in a manifest file corresponds to an individual data record, such as an image, and includes labeling information. This might include the source reference (such as the image's location on S3), metadata about the label (like the class name and confidence level), and additional information like whether the label was human-annotated.
- **Manifest File Data Points:** Each manifest file contains 'N' JSON objects, where 'N' is the total number of images or data points that have been used in the dataset for the GroundTruth job.

Labeling Project



020.jpg



019.jpg



018.jpg



017.jpg



016.jpg



015.jpg



014.jpg



013.jpg



012.jpg



011.jpg



010.jpg



009.jpg



008.jpg



007.jpg



006.jpg



005.jpg



004.jpg



003.jpg



002.jpg



001.jpg

aws

Services

Search

[Alt+S]

Step 1

Specify job details

Step 2

Select workers and configure tool

Specify job details

All fields are required unless otherwise specified

Job overview

Job name

Maximum of 63 alphanumeric characters. Can include hyphens (-), but not spaces. Must be unique within your account in an AWS Region.

☐ I want to specify a label attribute name different from the labeling job name.

Label attribute name is the key where your labels are stored in the augmented manifest. Ground Truth uses the labeling job name as the default label attribute name.

Input data setup [Info](#)

Use the automated setup to have Ground Truth automatically identify your dataset in S3. Use the manual setup if you have an input manifest file.

☒ Automated data setup

Provide the S3 location of the dataset you want labeled and let Ground Truth automatically connect to and use this dataset for your job.

☐ Manual data setup

Provide the S3 location of a file (an input manifest file) that identifies the data objects you want labeled

Data setup

S3 location for input datasets [Info](#)

This is the location in S3 where your dataset objects are stored. Ground Truth will use all data objects in this location for your labeling job.

View

Browse S3

S3 location for output datasets [Info](#)

This is the location in S3 where your labeling job output data is stored.

☒ Same location as input dataset

All fields are required unless otherwise specified

Job name

fashiontestlabeljob

Maximum of 63 alphanumeric characters. Can include hyphens (-), but not spaces. Must be unique within your account in an AWS Region.

☐ I want to specify a label attribute name different from the labeling job name.

Label attribute name is the key where your labels are stored in the augmented manifest. Ground Truth uses the labeling job name as the default label attribute name.

Input data setup [Info](#)

Use the automated setup to have Ground Truth automatically identify your dataset in S3. Use the manual setup if you have an input manifest file.

- Automated data setup

Provide the S3 location of the dataset you want labeled and let Ground Truth automatically connect to and use this dataset for your job.

☐ Manual data setup

Provide the S3 location of a file (an input manifest file) that identifies the data objects you want labeled

Data setup

S3 location for input datasets [Info](#)

This is the location in S3 where your dataset objects are stored. Ground Truth will use all data objects in this location for your labeling job.

🔍 s3://fashiondatatestgroundtruth/dataset ✕

View

Browse S3

S3 location for output datasets [Info](#)

This is the location in S3 where your labeling job output data is stored.

- Same location as input dataset

AWS SageMaker GroundTruth

aws

Services

Search [Alt+S]

🔍

🔔

?

⚙️

Central ▾

Farid_winter2024 ▾

☰

Data setup

S3 location for input datasets [Info](#)

This is the location in S3 where your dataset objects are stored. Ground Truth will use all data objects in this location for your labeling job.

🔍 s3://fashiondatatestgroundtruth/dataset ✕

View [↗](#)

Browse S3

S3 location for output datasets [Info](#)

This is the location in S3 where your labeling job output data is stored.

☒ Same location as input dataset

☐ Specify a new location

Data type

Image

Supported formats are .jpg, .jpeg, and .png.

IAM Role [Info](#)

Provide the ID or ARN for your own AWS KMS encryption key for Amazon SageMaker to access your S3 bucket. Choose a role or let us create a role with the [AmazonSageMakerFullAccess](#) IAM policy attached.

AmazonSageMaker-ExecutionRole-20240228T143506

Use this button to process and complete your input data setup.

Complete data setup

✔️ Input data connection successful. [View more details](#)

▶ Additional configuration - optional

Dataset object selection, encryption

Task type [...](#)

AWS SageMaker GroundTruth

aws

Services

Search

[Alt+S]

Global

Farid_winter2024

Amazon S3

Buckets

Access Grants

Access Points

Object Lambda Access Points

Multi-Region Access Points

Batch Operations

IAM Access Analyzer for S3

Storage Lens

Dashboards

Storage Lens groups

AWS Organizations settings

Amazon S3

Buckets

fashiondatatestgroundtruth

dataset-20240228T143631.manifest

dataset-20240228T143631.manifest

Info

Copy S3 URI

Download

Open

Object actions

Properties

Permissions

Versions

Object overview

Owner

8e839f9356326f80911d3bb208b23891075c1b0d05b6f25740574720f9c95a66

AWS Region

Canada (Central) ca-central-1

Last modified

February 28, 2024, 14:36:34 (UTC-05:00)

S3 URI

s3://fashiondatatestgroundtruth/dataset-20240228T143631.manifest

Amazon Resource Name (ARN)

arn:aws:s3:::fashiondatatestgroundtruth/dataset-20240228T143631.manifest

Entity tag (Etag)

CloudShell

Feedback

© 2024 Amazon Web Services, Inc. or its affiliates

Privacy

Terms

Cookie preferences

AWS SageMaker GroundTruth

```
{"source-ref": "s3://fashiondatatestgroundtruth/001.jpg"}  
{"source-ref": "s3://fashiondatatestgroundtruth/002.jpg"}  
{"source-ref": "s3://fashiondatatestgroundtruth/003.jpg"}  
{"source-ref": "s3://fashiondatatestgroundtruth/004.jpg"}  
{"source-ref": "s3://fashiondatatestgroundtruth/005.jpg"}  
{"source-ref": "s3://fashiondatatestgroundtruth/006.jpg"}  
{"source-ref": "s3://fashiondatatestgroundtruth/007.jpg"}  
{"source-ref": "s3://fashiondatatestgroundtruth/008.jpg"}  
{"source-ref": "s3://fashiondatatestgroundtruth/009.jpg"}  
{"source-ref": "s3://fashiondatatestgroundtruth/010.jpg"}  
{"source-ref": "s3://fashiondatatestgroundtruth/011.jpg"}  
{"source-ref": "s3://fashiondatatestgroundtruth/012.jpg"}  
{"source-ref": "s3://fashiondatatestgroundtruth/013.jpg"}  
{"source-ref": "s3://fashiondatatestgroundtruth/014.jpg"}  
{"source-ref": "s3://fashiondatatestgroundtruth/015.jpg"}  
{"source-ref": "s3://fashiondatatestgroundtruth/016.jpg"}  
{"source-ref": "s3://fashiondatatestgroundtruth/017.jpg"}  
{"source-ref": "s3://fashiondatatestgroundtruth/018.jpg"}  
{"source-ref": "s3://fashiondatatestgroundtruth/019.jpg"}  
{"source-ref": "s3://fashiondatatestgroundtruth/020.jpg"}
```

AWS SageMaker GroundTruth

aws

Services

Search

[Alt+S]

Central

Farid_winter2024

Select the type of data being labeled to view available task templates for it or select 'Custom' to create your own.

Image

Task selection


Select the task that a human worker will perform to label objects in your dataset.

☒ Image Classification (Single Label)

Get workers to categorize images into individual classes. [Info](#)

☒ Basketball

☐ Soccer




☐ Image Classification (Multi-label)

Get workers to categorize images into one or more classes. [Info](#)

☒ Human


☒ Vehicle

☐ Animal




☐ Bounding box

Get workers to draw bounding boxes around specified objects in your images. [Info](#)



☐ Semantic segmentation

Get workers to draw pixel level labels around specific objects and segments in your images. [Info](#)



AWS SageMaker GroundTruth

aws

Services

Search

[Alt+S]

Central

Farid_winter2024

Amazon SageMaker

Labeling jobs

Create labeling job

Step 1

Specify job details

Step 2

Select workers and configure tool

Select workers and configure tool

All fields are required unless otherwise specified

Workers

Info

Worker types

☐

Amazon Mechanical Turk

An on-demand 24/7 workforce of over 500,000 independent contractors worldwide powered by Amazon Mechanical Turk.

☒

Private

A team of workers that you have sourced yourself, including your own employees or contractors for handling data that needs to stay within your organization.

☐

Vendor managed

A curated list of third party vendors that specialize in providing data labeling services, available via the AWS Marketplace.

Team name

Maximum of 63 alphanumeric characters. Can include hyphens, but not spaces. Must be unique within your account in an AWS Region. The name can't be changed later.

Invite private annotators

Enter email addresses of workers that will work on this job.

Add email addresses, separated by commas to add annotators.

Enter up to 20 addresses and use a comma between each one.

Task timeout

The maximum time a worker can work in a single task. Please see [here](#) for information on default and maximum values.

0

hours

5

mins

0

secs

CloudShell

Feedback

© 2024, Amazon Web Services, Inc. or its affiliates.

Privacy

Terms

Cookie preferences

AWS SageMaker GroundTruth

aws

Services

Search

[Alt+S]

Central

Farid_winter2024

H1H2BI^A

Good example

Enter description to explain the correct label to the workers

Add image here

Bad example


Enter description of an incorrect label

Add image here

Enter a brief description of the task

Enter a brief description of the task

PREV



NEXT

Select an option

Add up to 30 labels

Add new label

You can add 28 more labels.

Additional instructions - optional

AWS SageMaker GroundTruth

aws

Services

Search

[Alt+S]

Central

Farid_winter2024

Provide labeling instructions with examples below for workers. Workers will be viewing these instructions when they perform your task. Workers can choose up to 30 labels. See guidelines for [See guidelines for creating high-quality instructions](#)

H1 H2 B I A

Good example

Enter description to explain the correct label to the workers

Add image here

Bad example

Enter description of an incorrect label

Add image here

Enter a brief description of the task

labeling performance

PREV

NEXT

Select an option

Add up to 30 labels

BAG

X

SUNGLASS

X

SHOES

X

WATCH

X

Add new label

You can add 26 more labels.

Additional instructions - optional

Cancel

Previous

Create

AWS SageMaker GroundTruth



Previewing Answers Submitted by Workers

This message is only visible to you and will not be shown to Workers.

You can test completing the task below and click "Submit" in order to preview the data and format of the submitted results.



Instructions

Shortcuts

labeling performance



Select an option

BAG	1
SUNGLASS	2
SHOES	3
WATCH	4

AWS SageMaker GroundTruth

aws

Services

Search

[Alt+S]

Central

Farid_winter2024

Amazon SageMaker

Getting started

Control panel

Studio

Studio Lab

RStudio

SageMaker dashboard

Images

Lifecycle configurations

Search

Ground Truth

Labeling jobs

Labeling datasets

Labeling workforces

Notebook

Processing

Training

Inference

ca-central-1_MIYHbG3MP

App client

3ig1n27fimnsvs3ao9kijn04lm

Labeling portal sign-in URL

https://ejbxkzdboh.labeling.ca-central-1.sagemaker.aws

Workforce status

Active

VPC Add

Private teams (1)

A team of workers from your private workforce. Only one team can work on a labeling job or review task (jobs). Each team can be assigned to multiple jobs.

Search private teams by name

	Name	ARN	Creation time
	cloudcourselabelers	arn:aws:sagemaker:ca-central-1:730335673928:workteam/private-crowd/cloudcourselabelers	Feb 28, 2024, 7:57 PM UTC

Workers

All workers in your private workforce

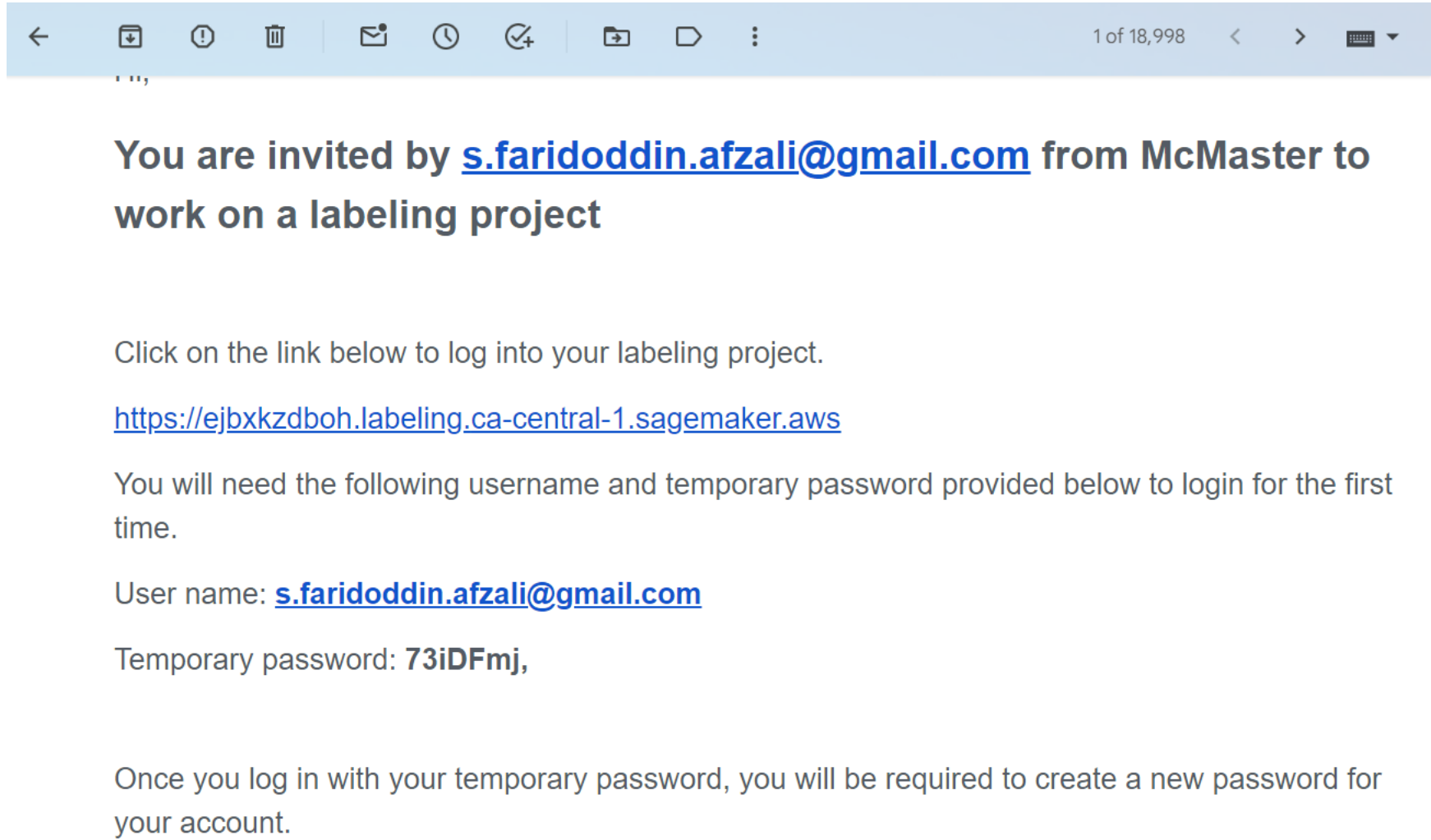
Search workers by email

	Email	Status	Cognito status	Enabled	Username
	s.faridoddin.afzali@gmail.com	Invitation sent	Force_change_password	Yes	s.faridoddin.afzali@gmail.com

CloudShell Feedback

© 2024, Amazon Web Services, Inc. or its affiliates. Privacy Terms Cookie preferences

AWS SageMaker GroundTruth



AWS SageMaker GroundTruth

Hello, s.faridoddin.afzali@gmail.com

Log out

☐ Show instructions

Jobs (1)

Start working

< 1 >

	Task title	Customer ID	Status	Creation time
<input checked="" type="radio"/>	Image Classification (Single Label): labeling performance	730335673928	Available	February 28, 2024 19:57:15 UTC

AWS SageMaker GroundTruth

Hello, s.faridoddin.afzali@gmail.com

Customer ...

Task description: Cate...

Task time: 0:23 of 4 Min 59 Sec

Decline task

Release task

Stop and resume later

Instructions

Shortcuts

labeling performance



Select an option

BAG	1
SUNGLASS	2
SHOES	3
WATCH	4



z

o

m

f

Submit

AWS SageMaker GroundTruth

aws

Services

Search

[Alt+S]

Global

Farid_winter2024

Amazon S3

Buckets

Access Grants

Access Points

Object Lambda Access Points

Multi-Region Access Points

Batch Operations

IAM Access Analyzer for S3

Block Public Access settings for this account

Storage Lens

Dashboards

Storage Lens groups

AWS Organizations settings

Feature spotlight

AWS Marketplace for S3

Amazon S3 > Buckets > fashiondatatestgroundtruth > fashiontestlabeljob/

fashiontestlabeljob/

Copy S3 URI

Objects

Properties

Objects (3) Info

Refresh

Copy S3 URI

Copy URL

Download

Open

Delete

Actions

Create folder

Upload

Objects are the fundamental entities stored in Amazon S3. You can use [Amazon S3 inventory](#) to get a list of all objects in your bucket. For others to access your objects, you'll need to explicitly grant them permissions. [Learn more](#)

Find objects by prefix

< 1 > ⚙

<input type="checkbox"/>	Name	Type	Last modified	Size	Storage class
<input type="checkbox"/>	annotation-tool/	Folder	-	-	-
<input type="checkbox"/>	annotations/	Folder	-	-	-
<input type="checkbox"/>	manifests/	Folder	-	-	-

AWS SageMaker GroundTruth

[illegible]