



Getting Started Resource Center ▾

Get Started ▾

Learn ▾

Get Connected ▾

Developer Tools

More Resources ▾

Explore by Role ▾

Getting Started / Hands-on / ...

Prepare Training Data for Machine Learning with Minimal Code

TUTORIAL



Overview

In this tutorial, you will learn how to prepare data for machine learning (ML) using [Amazon SageMaker Data Wrangler](#).

Amazon SageMaker Data Wrangler reduces the time it takes to aggregate and prepare data for ML from weeks to minutes. Using SageMaker Data Wrangler, you can simplify the process of data preparation and feature engineering and complete each step of the data preparation workflow, including data selection, cleansing, exploration, and visualization from a single visual interface.

In this tutorial, you will use Amazon SageMaker Data Wrangler to prepare data for a prediction model. You will use a version of the Brazil house rental dataset found in the Kaggle



data for bias, and lastly save the output to Amazon S3 to be used later for ML training.

What you will accomplish

In this guide, you will:

- Visualize and analyze data to understand key relationships
- Apply transformations to clean up the data and generate new features
- Automatically generate notebooks for repeatable data preparation workflows

Prerequisites

Before starting this tutorial, you will need:

- **An AWS account:** If you don't already have an account, follow the [Setting Up Your AWS Environment](#) getting started guide for a quick overview.

✓ **AWS experience**

Beginner



⌚ **Minimum time to complete**

30 minutes

฿ **Cost to complete**

See [Amazon SageMaker pricing](#) to estimate cost for this tutorial.

👤 **Requires**

You must be logged into an AWS account.

👤 **Services used**

Amazon SageMaker Data Wrangler

📝 **Last updated**

March 7, 2023



Step 1: Set up your Amazon SageMaker Studio domain

With Amazon SageMaker, you can deploy a model visually using the console or programmatically using either SageMaker Studio or SageMaker notebooks. In this tutorial, you deploy the model programmatically using a SageMaker Studio notebook, which requires a SageMaker Studio domain.

An AWS account can have only one SageMaker Studio domain per Region. If you already have a SageMaker Studio domain in the US East (N. Virginia) Region, follow the [SageMaker Studio setup guide](#) to attach the required AWS IAM policies to your SageMaker Studio account, then skip Step 1, and proceed directly to Step 2.

If you don't have an existing SageMaker Studio domain, continue with Step 1 to run an AWS CloudFormation template that creates a SageMaker Studio domain and adds the permissions required for the rest of this tutorial.



Choose the [AWS CloudFormation stack](#) link. This link opens the AWS CloudFormation console and creates your SageMaker Studio domain and a user named *studio-user*. It also adds the required permissions to your SageMaker Studio account. In the CloudFormation console, confirm that **US East (N. Virginia)** is the **Region** displayed in the upper right corner. **Stack name** should be **CFN-SM-IM-Lambda-catalog**, and should not be changed. This stack takes about 10 minutes to create all the resources.

This stack assumes that you already have a public VPC set up in your account. If you do not have a public VPC, see [VPC with a single public subnet](#) to learn how to create a public VPC.



This template is for an account without a pre-existing SageMaker Studio Domain & SageMaker User Profile - it creates these. This template also includes the additional permissions which are required to run this tutorial

Stack name

Stack name
CFN-SM-IM-Lambda-catalog
Stack name can include letters (A-Z and a-z), numbers (0-9), and dashes (-)

Parameters

Parameters are defined in your template and allow you to input custom values when you create or update a stack.

DomainName
The domain name of the Sagemaker studio instance
StudioDomain

UserProfileName
The user profile name for the SageMaker workshop
studio-user

Select **I acknowledge that AWS CloudFormation might create IAM resources**, and then choose **Create stack**.

Capabilities

The following resource(s) require capabilities: [AWS::IAM::Role]

This template contains Identity and Access Management (IAM) resources that might provide entities access to make changes to your AWS account. Check that you want to create each of these resources and that they have the minimum required permissions. [Learn more](#)

I acknowledge that AWS CloudFormation might create IAM resources.

Cancel Create change set Create stack

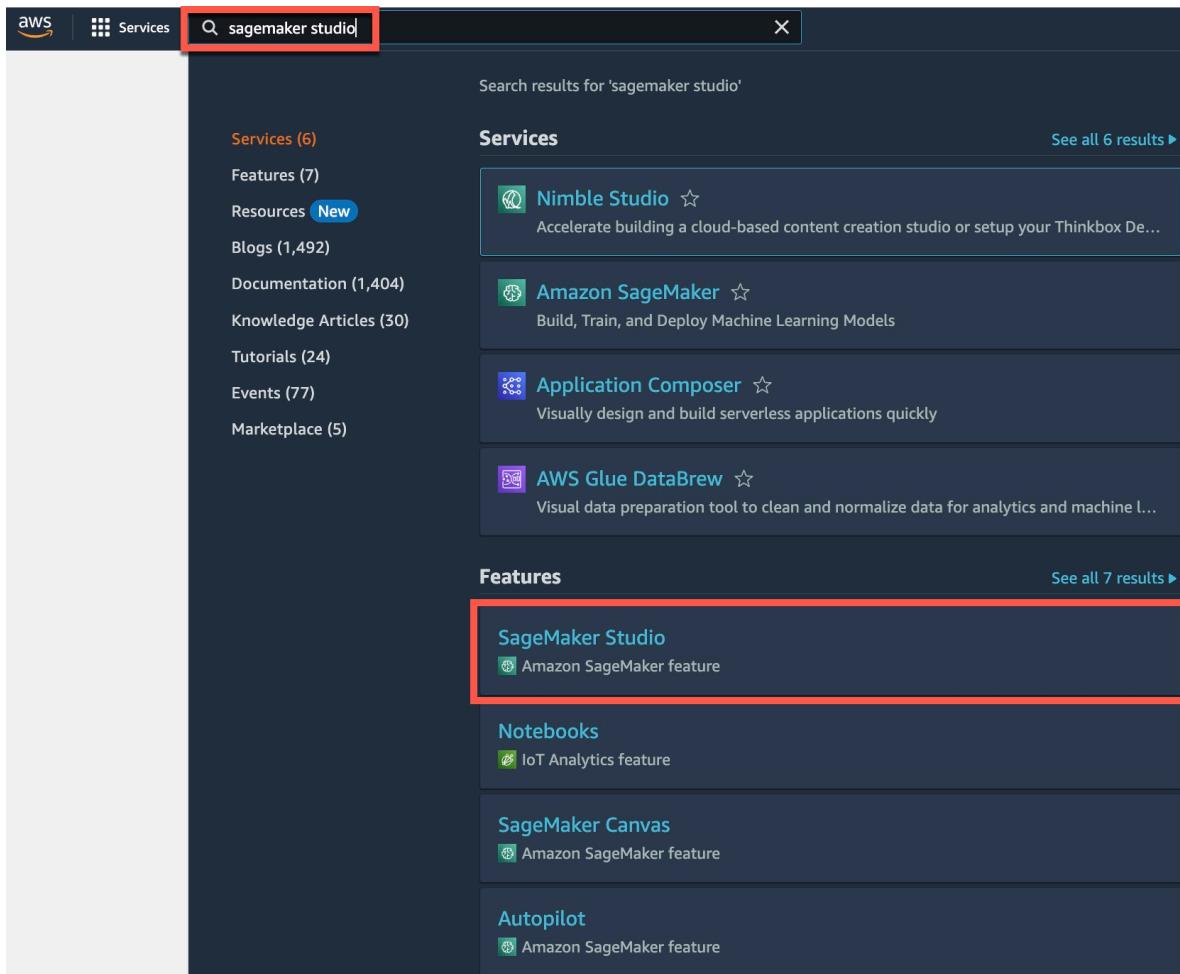
On the **CloudFormation** pane, choose **Stacks**. It takes about 10 minutes for the stack to be created. When the stack is created, the status of the stack changes from **CREATE_IN_PROGRESS** to **CREATE_COMPLETE**.

Stack name	Status	Created time	Description
CFN-SM-IM-Lambda-catalog	CREATE_COMPLETE	2022-06-14 18:45:17 UTC-0400	This template is for an account without a pre-existing SageMaker Studio Domain & SageMaker User Profile - it creates these. This template also includes the additional permissions which are required to run this tutorial

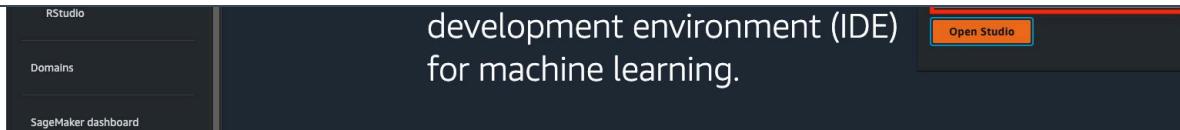
Step 2: Create a new SageMaker Data Wrangler flow

information about homes along with a target column indicating the rental amount of the property.

Enter **SageMaker Studio** into the console search bar, and then choose **SageMaker Studio**.



Choose **US East (N. Virginia)** from the Region dropdown list on the upper right corner of the SageMaker console. Browse to the **Getting Started** section in the left-hand navigation and then choose **Studio**. Then select the studio you created and then choose the **Open Studio** button.

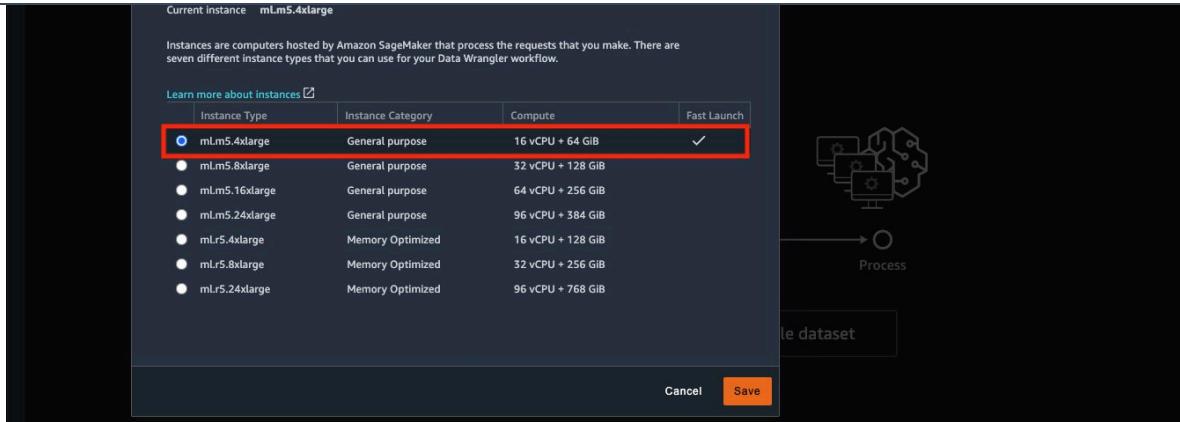


Open the **SageMaker Studio** interface. On the navigation bar, choose **Data Wrangler** on the left-hand side, and then choose the **Import Data** button.



Note that you can change the Flow's compute instance type using the upper right button showing the current Compute instance. You may decide to change the compute instance type based on your scenario's dataset size and can scale it up or down when your requirements change. For the purposes of this tutorial, you can use the default **ml.m5.4xlarge**.





In the **Data Import** tab, under **Import** data, choose **Amazon S3**.



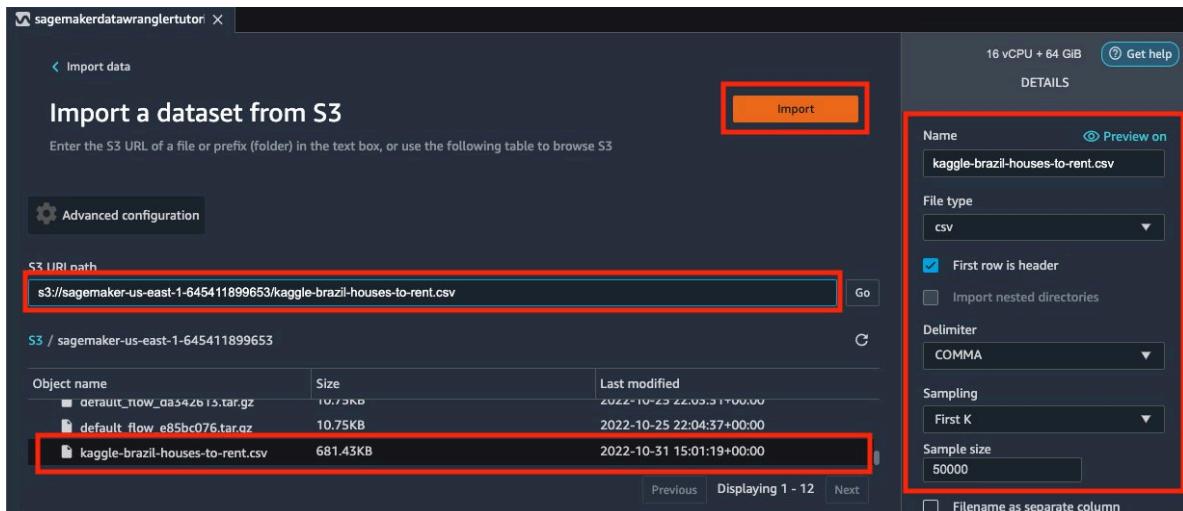
The screenshot shows the 'Data sources' section of the AWS Glue Data Catalog. At the top, there is a search bar and a 'Available 6' button. Below this, there is a row of icons for different data sources: Amazon Athena, Amazon Redshift, Snowflake, Amazon EMR, and Databricks. The first icon, 'Amazon S3', is highlighted with a red box. Below this row, there is a heading 'Set up new data sources 42' followed by a list of 42 additional data sources arranged in a grid. The data sources listed in the grid are: Amplitude, CircleCI, Datadog, DocuSign Monitor, Domo, Dynatrace, Facebook Ads, Facebook Page Insights, Freshdesk, GitHub, GitLab, Google Ads, Google Analytics, Google Analytics, v4, Google Search Console, Infor Nexus, Instagram Ads, Jira Cloud, LinkedIn Ads, Mailchimp, and Marketo.

In the S3 URI Path field, enter `s3://sagemaker-sample-files/datasets/tabular/brazil_houses/kaggle_brazil_houses_rental_data.csv`, and then choose **Go**. Under **Object name**, select `kaggle_brazil_houses_rental_data.csv`.





In the S3 import details panel, note that you can change the default delimiter and the sampling method when necessary. For the purposes of this tutorial, you can use the default **comma delimiter** and **First K sampling method**. Then choose **Import**.

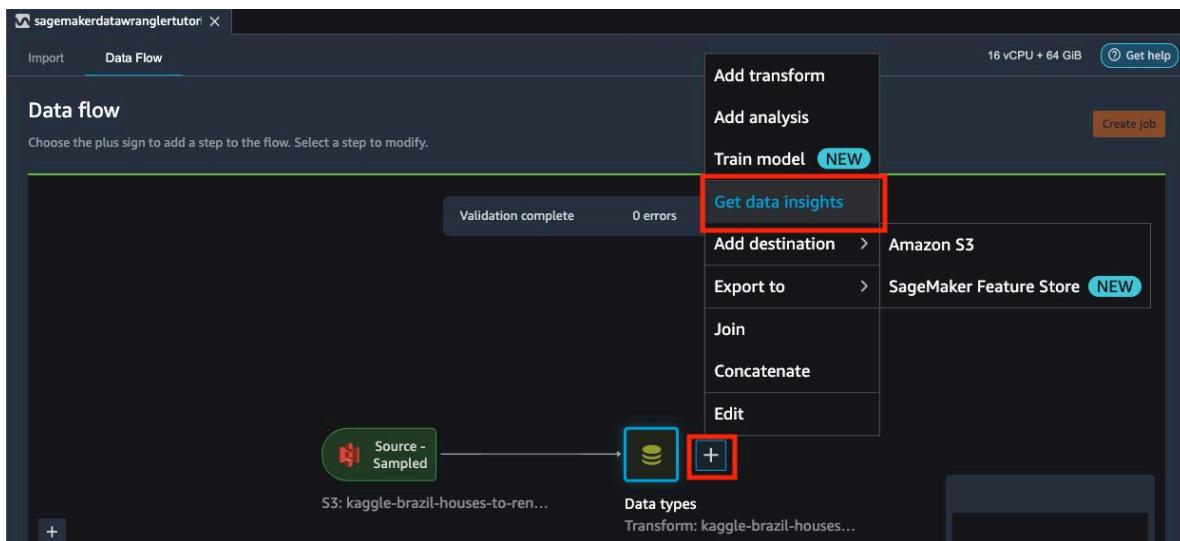


Step 3: Explore the data

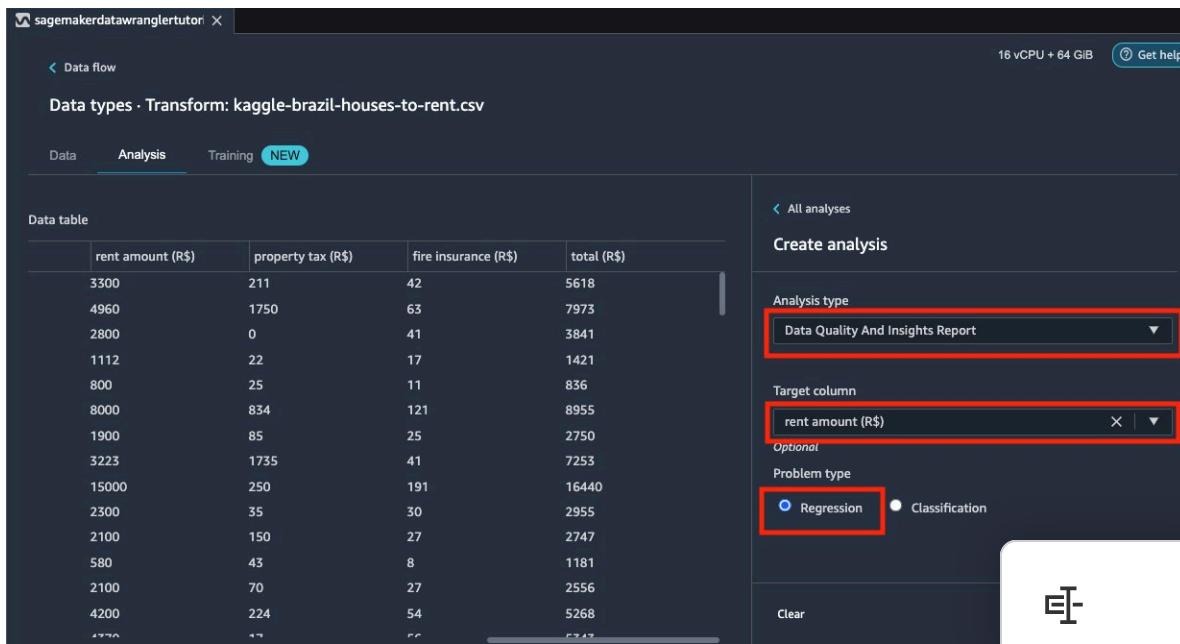
In this step, you use SageMaker Data Wrangler to assess and explore the quality of the training dataset for building machine learning models. Use the Data Quality and Insights report feature to understand your dataset quality, and then use the Quick Model feature to estimate the expected prediction quality and the predictive power of the features in your dataset.



your attention on the most important areas to improve the data. On the **Data flow** tab, in the data flow diagram, choose the **+** icon, then choose **Add analysis**. Then choose **Get data insights**.



From the **Data Insights** pane, choose **rent amount** as the **Target column**. Then choose **Regression** as the **Problem type**. Then choose **Create**.



Further before building the ML model. For this specific dataset, the Data Insights report has highlighted two possible issues: the first is related to **duplicate rows** in the dataset and the second is related to possible **target leakage** such that one feature is highly correlated with the output and may indicate a duplicate of the target **rent** column. The report can also be downloaded to a PDF file and shared with colleagues on your team.

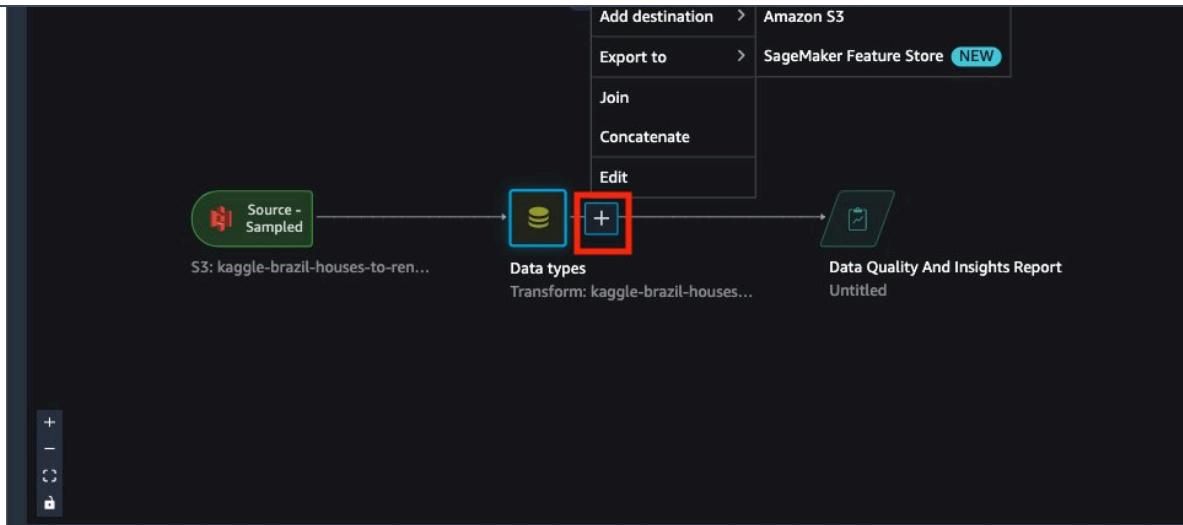
The screenshot shows the AWS SageMaker Data Wrangler interface with a Data Insights report for the 'kaggle-brazil-houses-to-rent.csv' dataset. The report is titled 'Transform: kaggle-brazil-houses-to-rent.csv'. At the top right, it shows '16 vCPU + 64 GiB', 'Get help', 'TARGET COLUMN: rent amount (R\$)', 'TYPE: Regression', 'DATASET: kaggle-brazil-houses-to-rent.csv', and 'DATE: October 31, 2022 at 10:45 AM CDT'. A red box highlights the 'SUMMARY' section, which contains 'Dataset statistics' and a table:

Key	Value	Feature type	Count
Number of features	13	numeric	8
Number of rows	10692	categorical	2
Missing	0%	text	0
Valid	100%	datetime	0
Duplicate rows	5.65%	binary	2
		unknown	0

Below the summary is a 'High Priority Warnings' section with two items, each highlighted by a red box:

- Duplicate rows** (High): We found that 5.65% of the data are duplicate. Some data sources could include valid duplicates and in other cases these duplicates could point to problems in data collection. Duplicate samples resulting from faulty data collection, could derail machine learning processes that rely on splitting to independent training and validation folds. For example quick model scores, prediction power estimation and automatic hyper parameter tuning. Duplicate samples could be removed from the dataset using the **Drop duplicates** transform under **Manage rows**.
- Target leakage** (High): The feature **fire insurance (R\$)** predicts the target extremely well on its own. A feature that predictive often indicates an error called target leakage. The cause is typically data that is not available at time of prediction. For example, a duplicate of the target column in the dataset can result in target leakage. Alternatively, if the machine learning task is "easy", then a single feature can have legitimately high prediction power. If you think that a single feature is very highly predictive, you don't need to do anything further. However, if you think there's target leakage, we recommend that you remove the highly predictive column from the dataset using the **Drop column** transform under **Manage columns**.

For further data analysis and exploration, you can create additional analytical artifacts including correlation matrices, histograms, scatter plots, and summary statistics as well as custom visualizations. For example, choose the **+** icon, then choose **Add analysis**.



Under the **Create analysis** panel, for **Analysis type**, select **Histogram** and name it **RentHistogramByRooms**. For **X axis**, select **rooms**.

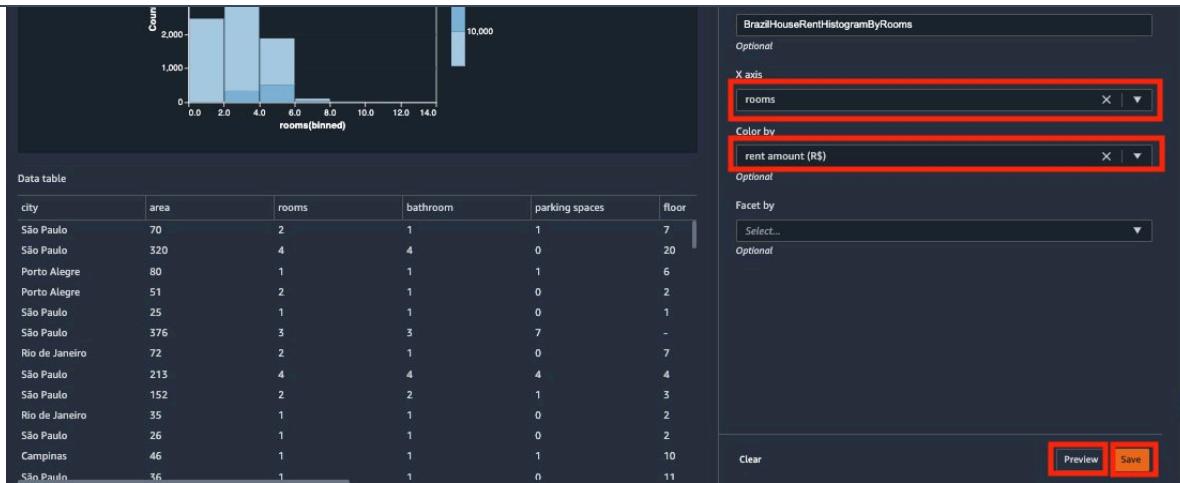
For **Color by**, select **Rent amount**.

Choose **Preview** to generate a **histogram** of the **rent amount** field, color-coded by the **rooms** variable.

Choose **Save** to save this analysis to the data flow.



E-T-

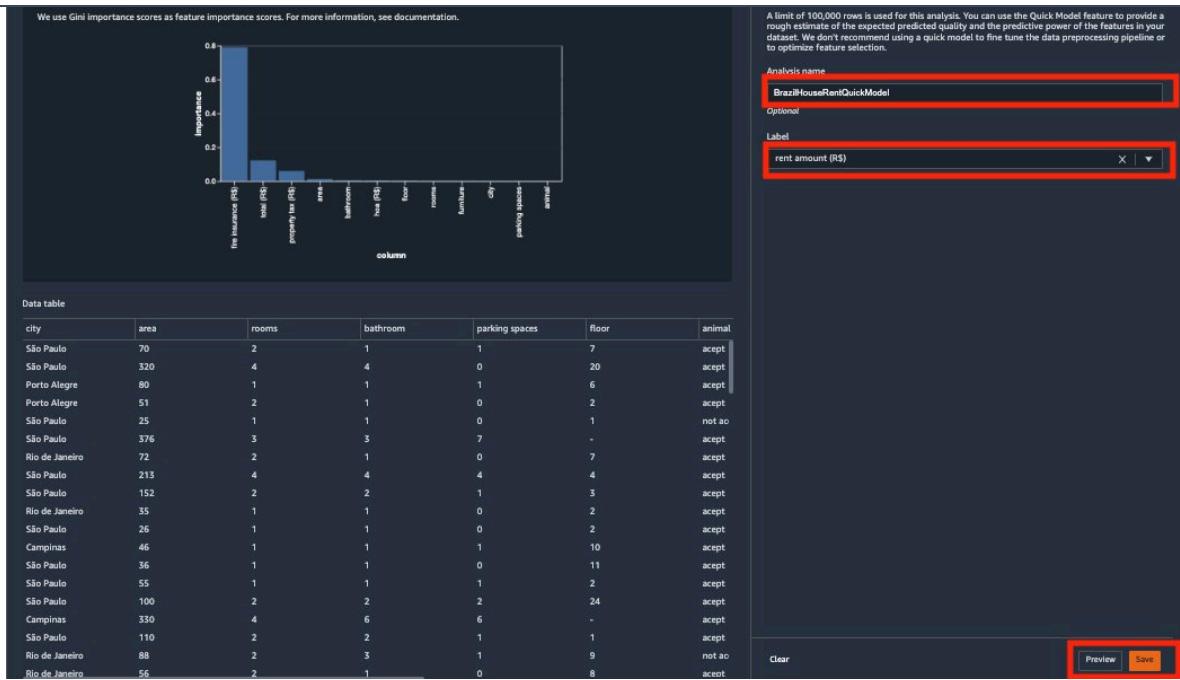


Next, to gain higher confidence that the underlying data has some predictive power, we are going to create a Quick Model. Under the **Create analysis** pane, for **Analysis type**, choose **Quick Model** and name it **RentQuickModel**.

Then for **Label**, select **rental amount** and then choose **Preview**.

The **Quick Model** may take several minutes to complete, then the pane shows a brief overview of the Random Cut Forest model built and trained with default hyperparameters. The model generated also displays some statistics, including the Mean Square Error (MSE) score and feature importance to help you evaluate the quality of the dataset. Choose **Save**.





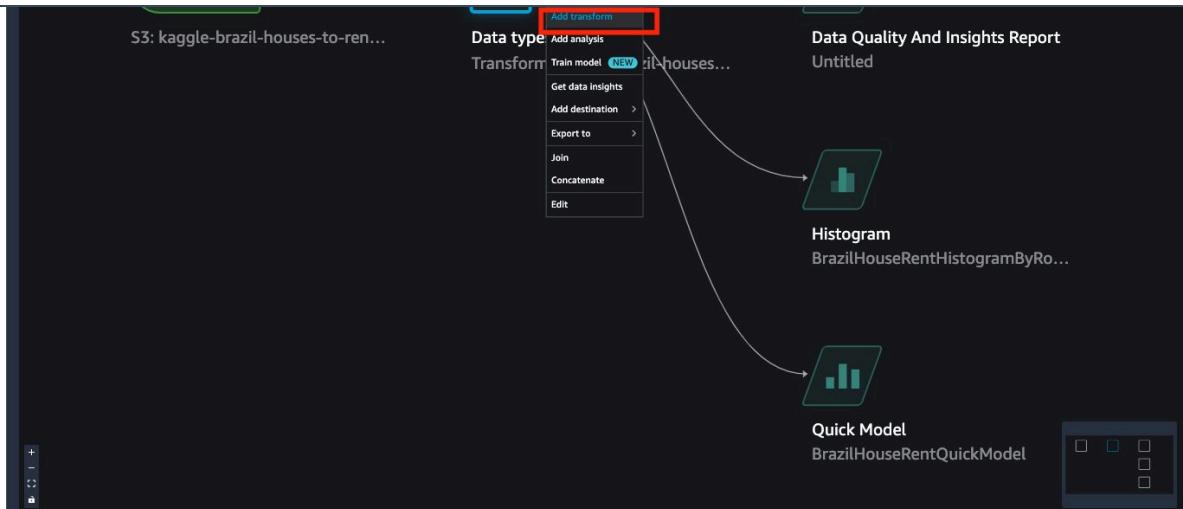
Step 4: Add transformations to the data flow

SageMaker Data Wrangler simplifies data processing by providing a visual interface with which you can add a wide variety of pre-built transformations. You can also write your custom transformations when necessary using SageMaker Data Wrangler. In this step, you change the type of a string column, rename columns, and drop unnecessary columns using the visual editor.

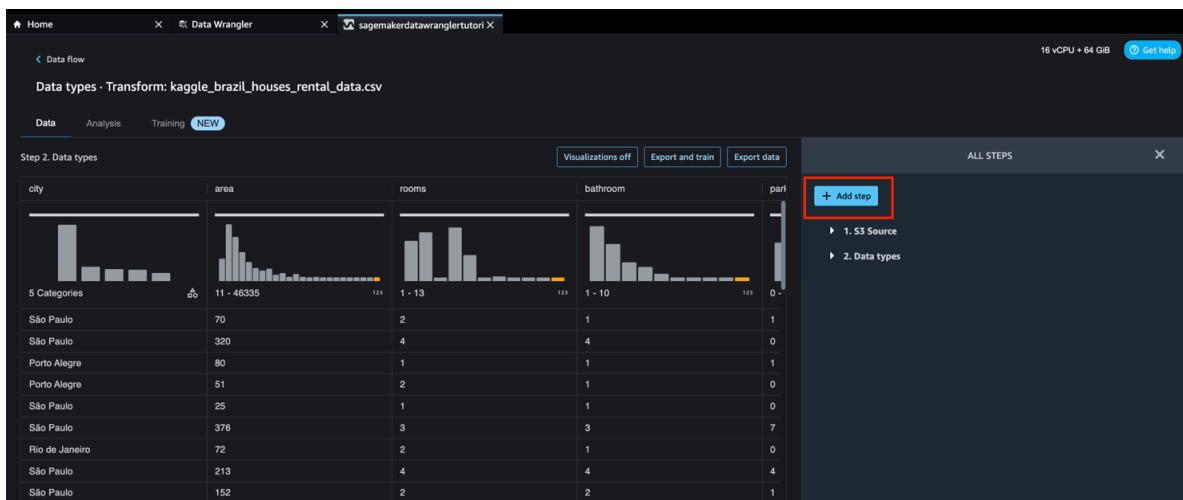


To navigate to the data flow diagram, choose **Data flow**. On the data flow diagram, choose the **+** icon, then **Add transform**

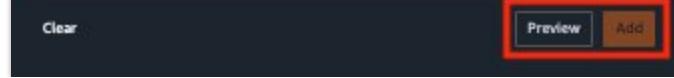




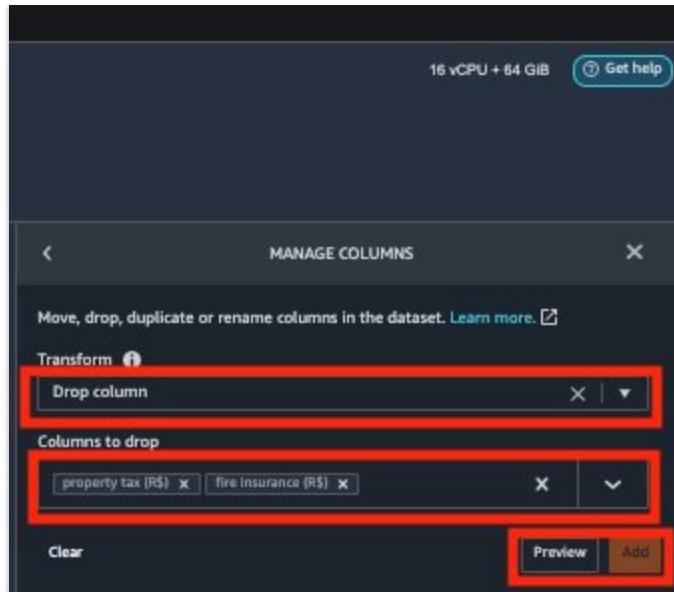
Under the **ALL STEPS** pane, choose **Add step**.



The first step is following the Data Insights Report recommendations regarding high risk items and removing the duplicate rows. So as the first transformation, choose **Manage Rows**, and then select the **Drop duplicates** operation, choose **Preview and Save**.



Second, we are going to remove the dataset features highlighted as possible sources of target leakage and not appropriate for a machine learning model predicting the rental amount. From the **ADD TRANSFORM** list, choose **Manage columns**. Then choose **Drop column** and choose **property tax** and **fire insurance**. Choose **Preview** then **Save**.



Next, change the data type of the **floor** column from **string** to **long**. Machine learning models can benefit from using numerically typed columns and this step will allow us to perform further processing later on.



2. Data types	
Column name	Type
city	String
area	Long
rooms	Long
bathroom	Long
parking spaces	Long
floor	Long
animal	String
furniture	String
hoa (R\$)	Long
rent amount (R\$)	Long
property tax (R\$)	Long
fire insurance (R\$)	Long
total (R\$)	Long

Clear

Preview

Update



Then rename several columns to improve the readability of the input data set and later analysis.

From the **ADD TRANSFORM** list, choose **Manage columns**. Then choose **Rename column**. Then choose **bathroom** as the input column and **bathrooms** as the output column. Choose **Preview** then **Save**.

Repeat this renaming column process for **hoa** [originally from **hoa (R\$)**], **rent** [originally from **rent amount (R\$)**], and **total** [originally from **total (R\$)**].

E-T-



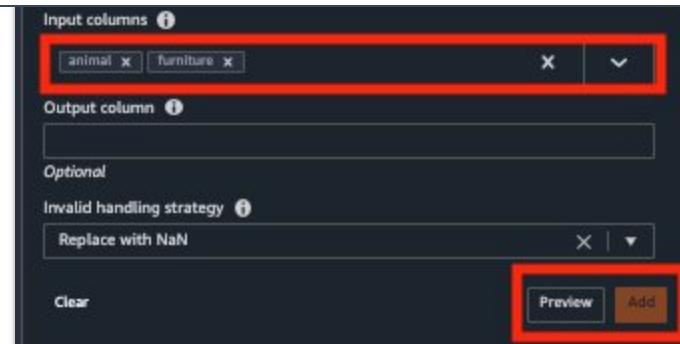
Step 5: Add categorical encoding and numeric scaling transformations to data flow

In this step, you encode categorical variables and scale numerical variables. Categorical encoding transforms string data type categories into numerical features. It's a common preprocessing task because the numerical features can be used in a wide variety of machine learning model types.

In the dataset, the rental property's **animal** and **furniture** classification is represented by various strings. In this step, you convert these string values to a binary representation, 0 or 1.

Under the **ALL STEPS** pane, choose **+ Add step**. From the **ADD TRANSFORM** list, choose **Encode categorical**. SageMaker Data Wrangler provides three transformation types: Ordinal encode, One hot encode, and Similarity encode. Under the **ENCODE CATEGORICAL** pane, for **Transform**, use the default **Ordinal encode**. For **Input columns**, select **animal** and **furniture**. Ignore the **Invalid handling strategy** box for this tutorial. Choose **Preview**, then **Add**.





To scale the numerical columns **area** and **floor**, apply a scaler transformation to normalize the distribution of the data in these columns: Under the **ALL STEPS** pane, Choose **+ Add step**. From the **ADD TRANSFORM** list, choose **Process numeric**. For **Scaler**, select the default option **Standard scaler**. For **Input** columns, select **area** and **floor**. Choose **Preview**, and then **Add**.



E-T-



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.

Featurize date/time
Encode date/time values to numeric and vector representations.

Featurize text
Generate vector representations from natural language text.

Format string
Clean and prepare strings using standard string formatting operations.

Group by
Add an aggregated column after group by as a new column.

Handle missing
Replace, drop, or add indicators for missing values.

Handle outliers
Remove or replace outlier numeric and categorical values.

Handle structured column
Flatten JSON and perform other operations on structured data

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

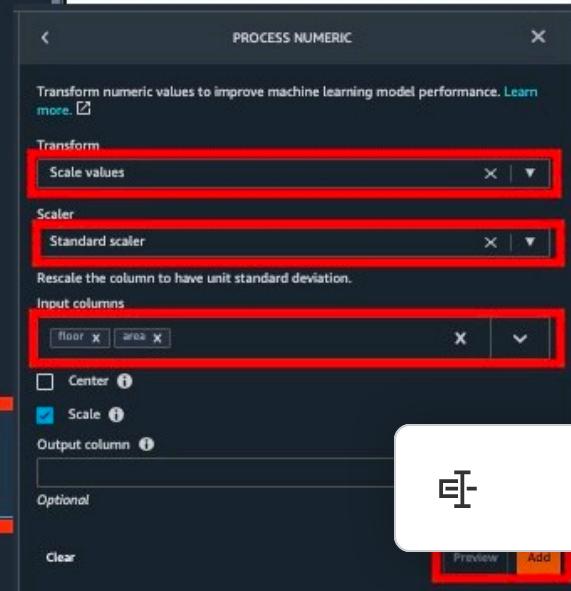
Manage rows
Sort, shuffle or drop duplicate rows.

Manage vectors
Expand or create vector columns.

Parse column as type
Cast a column to a new data type.

Process numeric
Transform numeric values to improve machine learning model performance.

Sampling
Create a sample of your data set.



Like
Dislike
Comment

Finally, we will follow another recommendation from the Data Insight report and replace the 0s in the Home Owner Association (*hoa*) feature with **NaN** because they indicate missing data and should not be treated as valid inputs that might skew the model. Under the **ALL STEPS** pane, choose **+ Add step**. From the **ADD TRANSFORM** list, choose **Search and edit** and then **Convert regex to missing**. Choose **hoa** as the **Input column**, specify **0** as the **Pattern**. Click **Preview**, and then choose **Add**.



Featurize date/time
Encode date/time values to numeric and vector representations.
Featurize text
Generate vector representations from natural language text.
Format string
Clean and prepare strings using standard string formatting operations.
Group by
Add an aggregated column after group by as a new column.
Handle missing
Replace, drop, or add indicators for missing values.
Handle outliers
Remove or replace outlier numeric and categorical values.
Handle structured column
Flatten JSON and perform other operations on structured data
Manage columns
Move, drop, duplicate or rename columns in the dataset.
Manage rows
Sort, shuffle or drop duplicate rows.
Manage vectors
Expand or create vector columns.
Parse column as type
Cast a column to a new data type.
Process numeric
Transform numeric values to improve machine learning model performance.
Sampling
Create a sample of your data set.
Search and edit
Find, replace, split, and otherwise transform input string values using search and edit functions.
Split data
Split an input dataframe into new dataframes.
Time Series
Transformers to preprocess and manipulate time series.
Validate string
Validate the format of string values using standard string functions.



A small white rectangular box with a black border and a black edit icon inside, located at the bottom right of the main content area.



Step 6: Check for data bias

In this step, check your data for bias using Amazon SageMaker Clarify, which provides you with greater visibility into your training data and models so you can identify and limit bias and better explain predictions.



Choose **Data flow** in the upper left to return to the data flow diagram. Choose the + icon, **Add analysis**. In the **Create analysis** pane, for **Analysis type**, select **Bias Report**. For **Analysis name**, enter **RentalDataBiasReport**. For **Select the column your model predicts (target)**, select **rent**. Then select **Threshold** as the predicted column type since this is a regression problem. Specify **3000** as the **predicted threshold** which corresponds to the average of the **rent** column in the dataset. Then select **city** as the column to analyze for bias because we are interested in whether the dataset is imbalanced and over-represents some cities instead of others. Then for **Choose bias metrics**, keep the default selections. Then choose **Check for bias** and then **Save**.

Bias Report

A limit of 100,000 rows is used for this analysis.

Analysis name

RentalDataBiasReport

Optional

Select the column your model predicts (target)

rent

Is your predicted column a value or threshold?

Value Threshold

Predicted threshold

3000

Select the column to analyze for bias

city

Is your column a value or threshold?

Value Threshold

Column value(s) to analyze for bias

Enter column value(s)

Optional

Choose bias metrics

Class imbalance (CI) ⓘ

Difference in Positive Proportions in Labels (DPL) ⓘ

JS divergence (JS) ⓘ

Conditional Demographic Disparity in Labels (CDDL) ⓘ

To measure CDDL, select a column in the dataset to be used as the group variable.

Select...

Optional

Would you like to analyze additional metrics?

Yes No

Clear

Check for bias

Save

After several seconds, SageMaker Clarify generates a report, which shows the target and feature columns score on a number of bias-related metrics

consider a bias remediation method, such as using SageMaker Data Wrangler's built-in SMOTE transformation. For the purpose of this tutorial, skip the remediation step. Choose **Save** to save the bias report to the data flow.

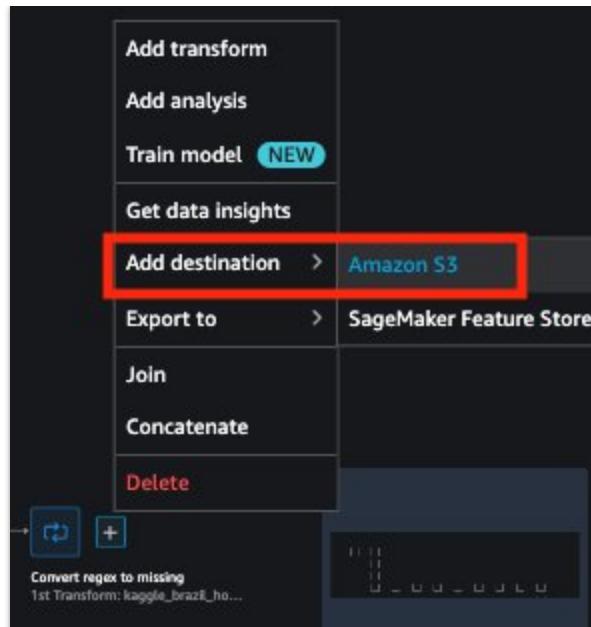
The screenshot shows the SageMaker Data Wrangler interface with the 'Analysis' tab selected. A bias report for the 'rent' column is displayed. The 'Predicted value or threshold' is set to 3000.0. The 'Column analyzed for bias' is 'city' and the 'Column value or threshold analyzed for bias' is 'São Paulo'. A red box highlights the 'Class imbalance (CI)' section, which shows a value of -0.11. Below it, other metrics like 'Difference in Positive Proportions in Labels (DPL)' and 'Jensen-Shannon Divergence (JS)' are listed with their respective values (-0.26 and 0.055).

Step 7: Review, integrate, and export your data flow

From the **Data Flow** tab, review your end-to-end data flow graph including the data source, analytical artifacts, and data transformations. You can easily navigate, view, modify, and delete data flow steps iteratively.



Data Wrangler further streamlines the automation process of exporting the output of the data flow to a persistent destination and can orchestrate the schedule of the flow's execution. First, set the storage destination to Amazon S3.

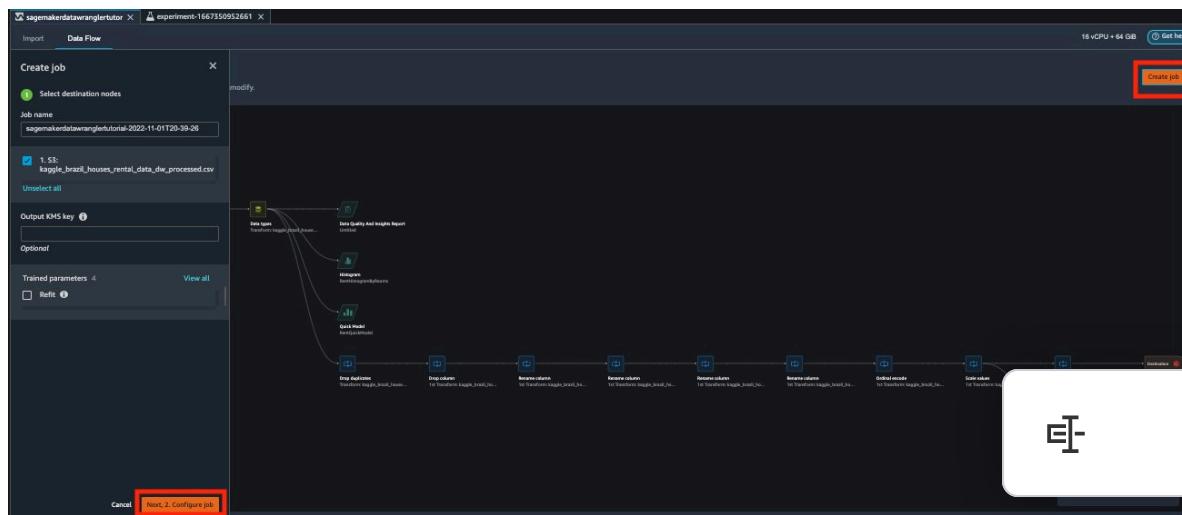


Then specify the output dataset name (**(kaggle_brazil_houses_rental_data_dw_processed.csv)**) and the Am location as your preferred S3 bucket. Then choose **Add destination**.

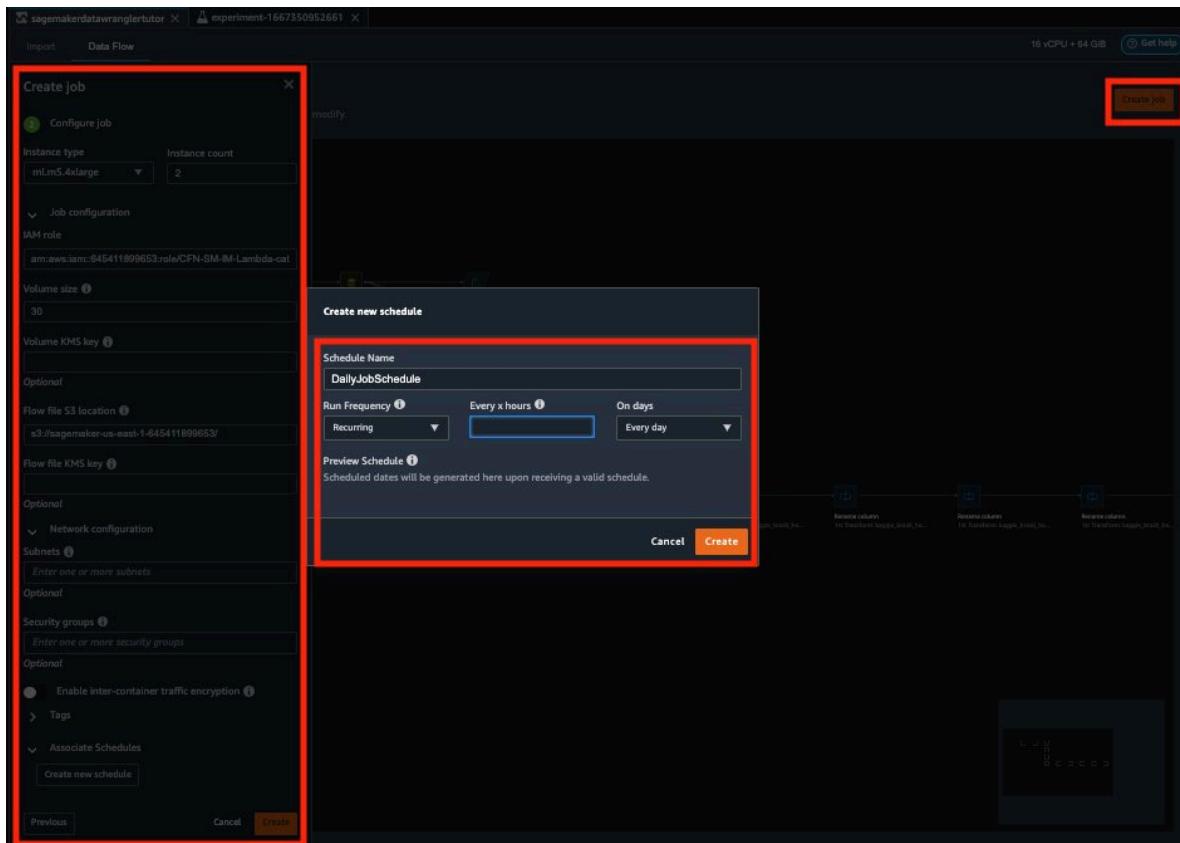
E-T



Lastly, create the scheduled job that will export the data flow output to Amazon S3 by choosing the **Create job** button from the **Data Flow** diagram pane, and then choosing **Configure job**.



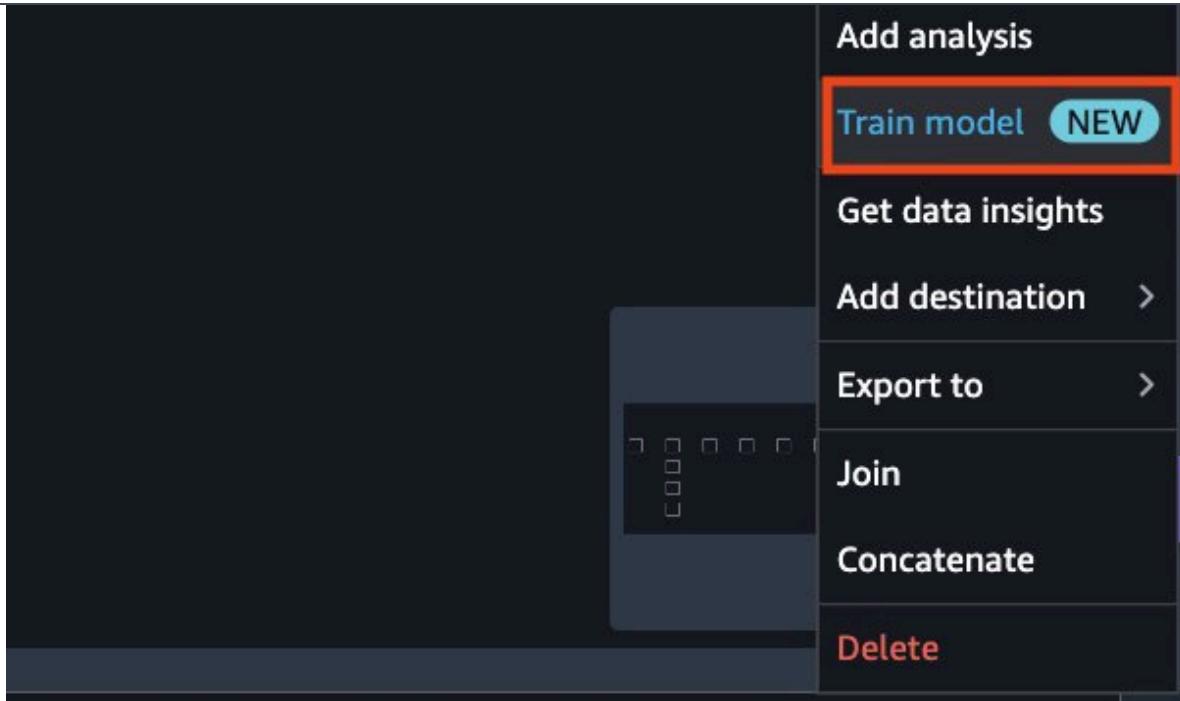
Then you can decide on the job instance type, instance count, the job's IAM security role, and the job schedule.



SageMaker Autopilot Integration

You can also integrate your data flow with [SageMaker Autopilot](#) which automates key tasks of training and deploying a machine learning model. From the **Data Flow** tab, choose the + icon and then choose **Train model**.





Choose **Export and Train** to export the Data Wrangler flow and associate its output with the Autopilot Experiment input.

Choose the **S3 location** where the Data Wrangler flow saved the processed input dataset and specify the **target** column as **rent** for the Autopilot model.

Specify the Autopilot **Training method**. You can choose Ensembling, Hyperparameter Optimization, or Auto. For the purposes of this tutorial, choose **Auto**.

For **Deployment**, select the machine learning problem type as **Regression** with the object metric as **MSE**.

Confirm the Autopilot Experiment deployment settings and then choose **Create experiment**. This action launches a SageMaker Autopilot job that ins input data, generates and evaluates multiple ML models, and then se E-T best model for subsequent deployment according to the desired perf

attribution and explainability statistics using SageMaker Clarify.

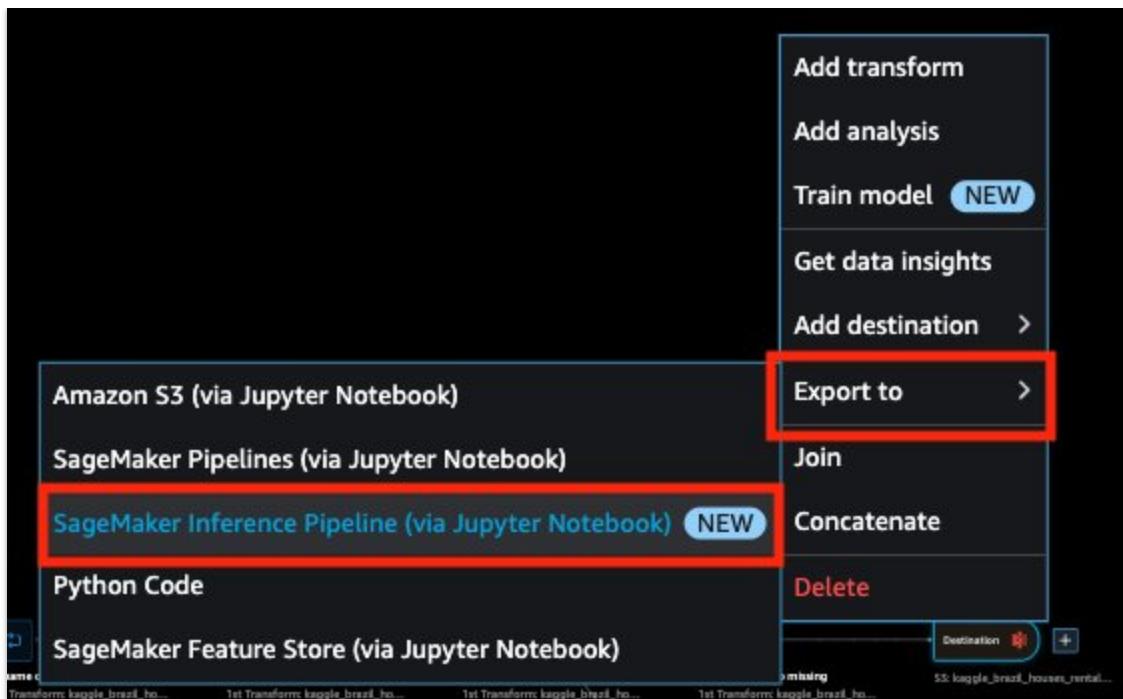
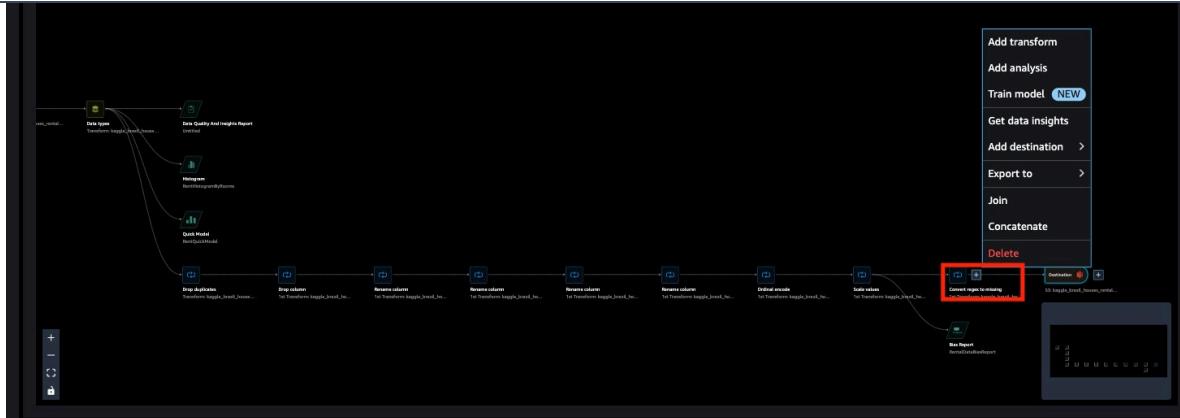
Model name	Objective: Mse	Status	Start time
Experiment-16673509526619b15sg8l-010-d83aa984	328196.75	Completed	35 minutes ago
Experiment-16673509526619b15sg8l-022-92...	366718048	Completed	32 minutes ago
Experiment-16673509526619b15sg8l-044-3f...	340098080	Completed	27 minutes ago
Experiment-16673509526619b15sg8l-031-85...	325408864	Completed	30 minutes ago
Experiment-16673509526619b15sg8l-026-f5...	320702976	Completed	32 minutes ago
Experiment-16673509526619b15sg8l-024-57...	310899936	Completed	32 minutes ago
Experiment-16673509526619b15sg8l-037-a8...	164719440	Completed	29 minutes ago

SageMaker Pipeline Integration

Data Wrangler can also be integrated with SageMaker Inference Pipelines to process data at the time of inference, thereby streamlining the steps between data processing and model inference. When you export one or more steps from the data flow to an inference endpoint, Data Wrangler creates a Jupyter notebook that you can use to define, instantiate, customize, run, and manage the inference pipeline. To create the inference endpoint, choose the **+** next to the final transformation step (Convert regex to missing) and choose **Export to**, and then choose **SageMaker Inference Pipeline (via Jupyter Notebook)**. Then inspect and run that Jupyter notebook.

You can optionally export your Data Wrangler data flow to a Jupyter notebook to run the flow steps as a SageMaker Processing job.





Step 8: Clean up resources

It is a best practice to delete resources that you are no longer using so that you don't incur unintended charges.

To delete the S3 bucket, do the following:



object to delete and enter **permanently delete** into the **Permanently delete objects** confirmation box.

- Once this is complete and the bucket is empty, you can delete the **sagemaker-<your-Region>-<your-account-id>** bucket by following the same procedure again.

The screenshot shows the Amazon S3 console interface. The top navigation bar shows 'Amazon S3 > Buckets > sagemaker-us-east-1-'. Below it, the bucket name 'sagemaker-us-east-1-' is displayed with an 'Info' link. The main area has tabs for 'Objects' (which is selected), 'Properties', 'Permissions', 'Metrics', 'Management', and 'Access Points'. Under the 'Objects' tab, there's a sub-section titled 'Objects (9)'. It includes a toolbar with buttons for 'C' (Create), 'Copy S3 URI', 'Copy URL', 'Download', 'Open', 'Delete' (which is highlighted with a red box), 'Actions', 'Create folder', and 'Upload'. Below the toolbar is a search bar with the placeholder 'Find objects by prefix'. A table lists the objects, with columns for 'Name', 'Type', and 'Last modified'. The first object listed is 'data_wrangler_flows/' (highlighted with a red box). The 'Actions' column for this object contains a 'Delete' link.

The Data Science kernel used for running the notebook image in this tutorial will accumulate charges until you either stop the kernel or perform the following steps to delete the apps. For more information, see [Shut Down Resources](#) in the *Amazon SageMaker Developer Guide*.

To delete the SageMaker Studio apps, do the following: On the SageMaker Studio console, choose **studio-user**, and then delete all the apps listed under **Apps** by choosing **Delete app**. Wait until the **Status** changes to **Deleted**.



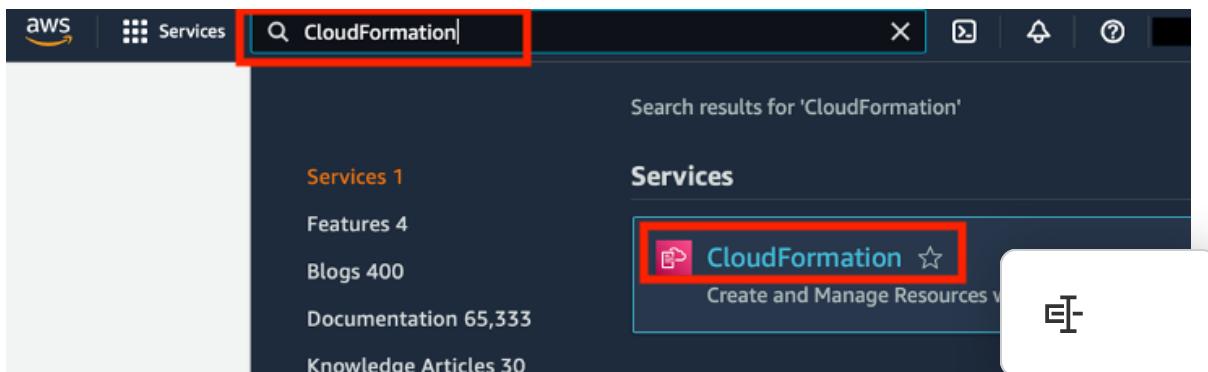
Apps				
App name	Status	App type	Created	Action
datascience-1-0-ml-t3-medium-1abf3407f667f989be9d86559395	⌚ Ready	KernelGateway	Sat Apr 09 2022 15:25:16 GMT-0400 (Eastern Daylight Time)	<button>Delete app</button>
default	⌚ Ready	JupyterServer	Sat Apr 09 2022 15:22:55 GMT-0400 (Eastern Daylight Time)	<button>Delete app</button>

If you used an existing SageMaker Studio domain in Step 1, skip the rest of Step 8 and proceed directly to the conclusion section.

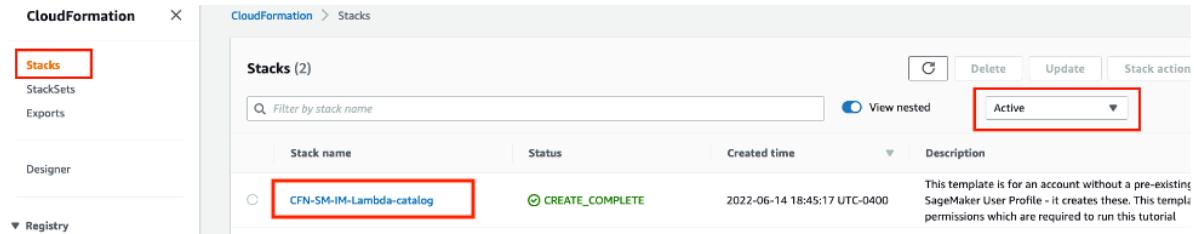
If you ran the CloudFormation template in Step 1 to create a new SageMaker Studio domain, continue with the following steps to delete the domain, user, and the resources created by the CloudFormation template.



To open the CloudFormation console, enter **CloudFormation** into the AWS console search bar, and choose **CloudFormation** from the search results.



In the **CloudFormation** pane, choose **Stacks**. From the status dropdown list, select **Active**. Under **Stack name**, choose **CFN-SM-IM-Lambda-catalog** to open the stack details page.



On the **CFN-SM-IM-Lambda-catalog** stack details page, choose **Delete** to delete the stack along with the resources it created in Step 1.



Conclusion

Congratulations! You have completed the **Prepare Training Data for Machine Learning with Minimal Code** tutorial.

You have successfully used Amazon SageMaker Data Wrangler to prepare data for a machine learning model. SageMaker Data Wrangler offers 300+ preconfigured transformations, such as convert column type, one-hot encoding, impute missing values, and mean



Train a deep learning model

Learn how to build, train, and tune a TensorFlow deep learning model.

[Next »](#)

Create an ML model automatically

Learn how to use AutoML to develop ML models without writing code.

[Next »](#)

Find more hands-on tutorials

Explore other machine learning tutorials to dive deeper.

[Next »](#)



Learn About AWS

[What Is AWS?](#)

[What Is Cloud Computing?](#)

[AWS Accessibility](#)

[What Is DevOps?](#)

Resources for AWS

[Getting Started](#)

[Training and Certification](#)

[AWS Solutions Library](#)

[Architecture Center](#)

[Product and Technical FAQs](#)

[Analyst Reports](#)

Developers on AWS

[Developer Center](#)

[SDKs & Tools](#)

[.NET on AWS](#)

[Python on AWS](#)

[Java on AWS](#)

[PHP on AWS](#)



What is Machine Learning (ML)?

[AWS Cloud Security](#)

[What's New](#)

[Blogs](#)

[Press Releases](#)

Help

[Contact Us](#)

[Get Expert Help](#)

[File a Support Ticket](#)

[AWS re:Post](#)

[Knowledge Center](#)

[AWS Support Overview](#)

[Legal](#)

[AWS Careers](#)



Amazon is an Equal Opportunity Employer: *Minority / Women / Disability / Veteran / Gender Identity / Sexual Orientation / Age.*

Language

| [Arabic](#)

| [Bahasa Indonesia](#)

| [Deutsch](#)

| [English](#)

| [Español](#)

| [Français](#)

| [Italiano](#)



日本語 |

한국어 |

中文 (简体) |

中文 (繁體)

Privacy

|

Accessibility

|

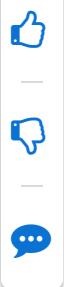
Site Terms

|

Cookie Preferences

|

© 2024, Amazon Web Services, Inc. or its affiliates. All rights reserved.



≡