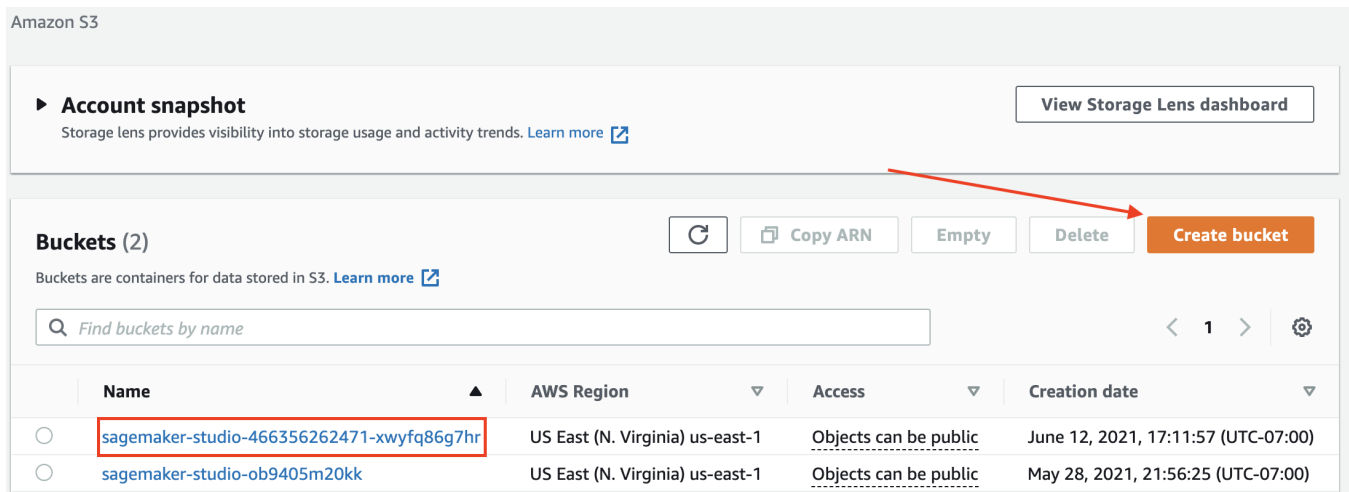# Solution: Data Wrangler

If you haven't done so already, you'll create an S3 bucket to upload files to and use throughout the course. Normally one will be created when you start using Sagemaker Studio, as pictured below.
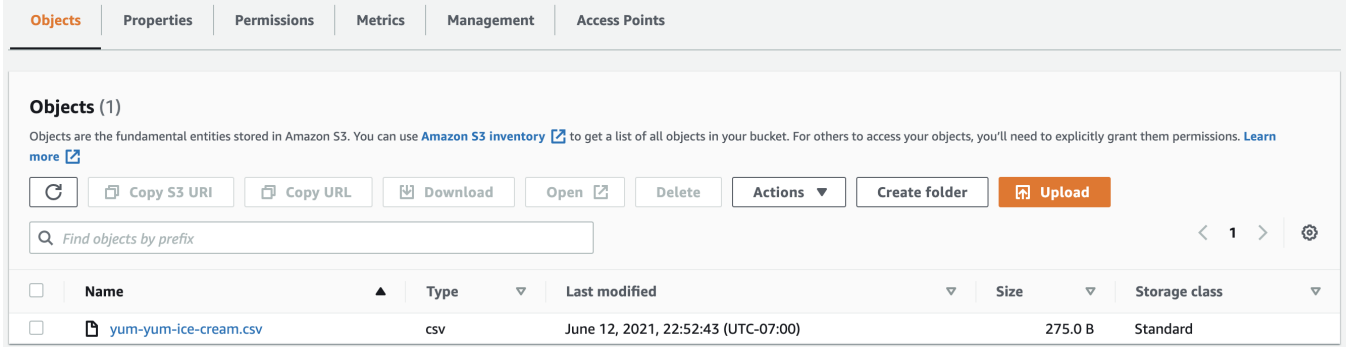


Choose or create a new bucket

Upload the `yum-yum-ice-cream.csv` file to your S3 bucket. The default settings that prompt you while uploading is acceptable.

Upload File to Bucket

When in the Sagemaker Studio IDE, you'll see a display for creating a `New data flow` , clicking on the `+` start your new Data Wrangler data flow.



Sagemaker Studio Data Wrangler

You already imported your file into S3, but in order to use it in Data Wrangler, you'll need to import it further from S3.

Data Wrangler S3



Import data file

After the file has been imported, the Data flow will be presented to you. Start transforming the data by clicking the `+` next to the `Data types` section. You'll be adding a transformation.

Data Flow

Transform the `date` field to have `year`, `month`, and `day` as their own features.



Extract Data Features



Extract Year, Month, Day

After transforming the date field, your data flow will be extended to visualize your changes.

Data Flow With Features

Create the two visualizations, `Table Summary` and `Histogram` .



Histogram and Table Summary

Now complete the exercise by exporting your data flow to S3. Save the results to S3.



Export to S3

When exporting your data flow to S3, it will create a Jupyter notebook called
`Save to S3 with a SageMaker Processing Job` . The notebook will have everything

configured for you to execute. To export your data, execute the cells up to the
`Optional Next Steps` . The optional steps are not required for this exercise.

# Save to S3 with a SageMaker Processing Job

💡 **Quick Start** To save your processed data to S3, select the Run menu above and click **Run all**
**the output S3 location.**

This notebook executes your Data Wrangler Flow `exercise2.flow` on the entire dataset using
processed data to S3.

This notebook saves data from the step `Featurize Date Time` from `Source: Yum-Yum-I`
go to Data Wrangler to select a new step to export.

## Contents

1. Inputs and Outputs
2. Run Processing Job
   A. Job Configurations
   B. Create Processing Job
   C. Job Status & S3 Output Location
3. Optional Next Steps
   A. Load Processed Data into Pandas
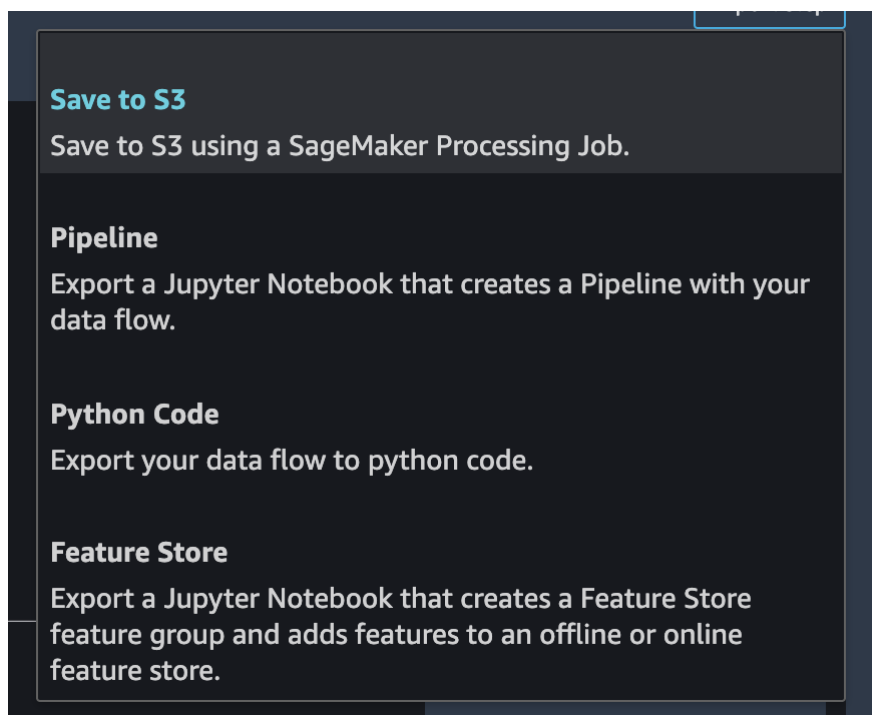   B. Train a model with SageMaker

Processing Job

You can view the status of your data flow export by navigating to `Processing Jobs` in the
Sagemaker console.

Amazon SageMaker  >  Processing jobs

**Processing jobs**                                          ↻    Actions ▼    **Create processing job**

🔍 Search processing jobs                                                    ‹  **1**  ›   ⚙

| | Name | | ARN | Creation time | Duration | Status | |
|---|---|---|---|---|---|---|---|
| ○ | data-wrangler-flow-processing-13-06-25-14-a5c433bb | ▽ | arn:aws:sagemaker:us-east-1:466356262471:processing-job/data-wrangler-flow-processing-13-06-25-14-a5c433bb | Jun 13, 2021 06:25 UTC | a minute | ⏱ InProgress | |

Processing Job in Progress

After the job is complete, the data will be available in S3 under a named `export-flow` directory. Below is an example of such a path.



Processed Job Results

The file that will be exported will look something like below. The naming convention is automatically created during the data job, but the data inside will reflect the date transformation including `date_year` , `date_month` , and `date_day` .

The exercise is a simple one, but the underlying principles are what make Data Wrangler so powerful. Because everything is managed through AWS, we can easily create and manage workflows creating complex processes that power our machine learning applications.

## part-00000-728b90e5-f33c-42fb-b833-6ca7970259c2-c000.csv

```
date,ice_cream_type,topping,location,date_year,date_month,date_day
2021-01-01,1,1,1,2021.0,0.0,0.0
2021-01-01,2,1,2,2021.0,0.0,0.0
2021-01-01,1,2,2,2021.0,0.0,0.0
2021-01-01,3,1,1,2021.0,0.0,0.0
2021-01-01,1,2,2,2021.0,0.0,0.0
2021-01-01,1,2,2,2021.0,0.0,0.0
2021-01-01,1,1,1,2021.0,0.0,0.0
2021-01-02,1,1,1,2021.0,0.0,1.0
2021-01-02,3,3,1,2021.0,0.0,1.0
2021-01-02,3,2,2,2021.0,0.0,1.0
2021-01-02,2,3,2,2021.0,0.0,1.0
2021-01-02,2,3,2,2021.0,0.0,1.0
2021-01-02,3,1,1,2021.0,0.0,1.0
2021-01-02,1,2,2,2021.0,0.0,1.0
```