HOW TO RUN MACHINE LEARNING TRANSFORMS IN AWS GLUE



Here we show you how to do a machine learning transformation with Amazon Glue. Previous Glue tutorials include:

- How To Make a Crawler in Amazon Glue
- How To Join Tables in Amazon Glue
- How To Define and Run a Job in AWS Glue
- <u>AWS Glue ETL Transformations</u>

Now, let's get started.

Amazon's machine learning

A fully managed service from Amazon, AWS Glue handles data operations like ETL to get your data prepared and loaded for analytics activities. Glue can crawl S3, DynamoDB, and JDBC data sources.

Amazon called their offering **machine learning**, but they only have one ML-type function, **findMatches**. It uses an ML algorithm, but Amazon does not tell you which one. They even boast on their web page you don't need to know—but a <u>data scientist</u> would certainly want to know.

You can study their execution log to gain some insight into what their code is doing. Suffice it to say it is doing a type of <u>clustering algorithm and using Apache Spark</u> as a platform to execute that.

The process: Amazon Glue machine learning

Here is the general process for running machine learning transformations:

- 1. Upload a csv file to an S3 bucket. Then you <u>set up a crawler</u> to crawl all the files in the designated S3 bucket. For each file it finds, it will create a metadata (i.e., schema) file in Glue that contains the column names.
- 2. Set up a **FindMatches** machine learning task in Glue. It's an iterative process. It takes your input date, created in the crawler process, and makes a label file. These labels are like a k-means clustering algorithm. It looks at the input data and all of the columns in the data set. Then it put the data into groups, each labeled with a **labeling_set_id**.
- 3. Download the label file. There will be an empty column called **label**. You are invited to add your own label to classify data however you see fit. For example, it could be borrower risk rating, whether or not a patient has diabetes, or whatever. Labels should be a single value, like A, B, C or 1, 2, 3. A data scientist would say they must be **categorical**.
- 4. Upload the labelled file to a different S3 bucket. Do not use the same bucket where you put the original input data, as the crawler will attempt to crawl that and create another metadata file.
- 5. Rerun Step 2, above, and it creates another labelled file. Do this iteratively until it supplies the most accurate result. In this example, there was no improvement from one run to the next. Repeating machine learning runs is standard practice for improving accuracy. However, at some point, the gain in accuracy will level off.
- 6. Generate and then inspect the **Quality Metrics**. Perhaps change some of the parameters and run the **Tune** operation, which means to run the algorithm again.

Tutorial: Amazon Glue machine learning

Now, let's run an example to show you how it works.

I have copied the Pima Native American database <u>from Kaggle</u> and put it on GitHub, <u>here</u>. You have to add a primary key column to that data, which Glue requires. Download the data <u>here</u>. I have also copied the input data and the first and second label files <u>here</u>, in a Google Sheet, so that you can see the before and after process.

The data looks like this:

recordID,Pregnancies,Glucose,BloodPressure,SkinThickness,Insulin,BMI,Diabetes
PedigreeFunction,Age,Outcome

2,6,148,72,35,0,33.6,0.627,50,1 3,1,85,66,29,0,26.6,0.351,31,0 4,8,183,64,0,0,23.3,0.672,32,1 5,1,89,66,23,94,28.1,0.167,21,0 6,0,137,40,35,168,43.1,2.288,33,1

Then copy it to an Amazon S3 bucket as shown below. You need to have installed the Amazon CLI (command line interface) and run **aws configure** to configure your credentials. Importantly, the data must be in the same <u>Amazon zone</u> as the instance you are logged into.

aws s3 cp diabetes.csv s3://sagemakerwalkerml

Add a label

The diabetes data is already labelled in the column **outcome**. So, I used Google Sheets to copy that value into the label **column**. You do this after the first run, like this:

- 1. Upload the original data.
- 2. Run a training model
- 3. Download the resulting labels file.

At that point you can populate the label with some kind of categorical data. You might put the outcome of logistic regression on your input data set into this label, but that's optional. You don't need a label at all.

The algorithm does not require a label the first time it runs. Glue says:

As you can see, the scope of the labels is limited to the labeling_set_id. So, labels do not cross labeling_set_id boundaries.

In other words, when there is no label, it groups records by **labeling_set_id** without regards to the **label** value. When there is a label then the **labeling_set_id** is within the label.

In other words, given this:

labeling_set_id	labeling_set_id	other columns
123	blank	
123	blank	
456	blank	
The first two rows are grou	oed together. But if we add a	a label:
labeling_set_id	label	other columns
labeling_set_id 123	label A	other columns
0		other columns
123	А	other columns

А

Then the first two rows are matched together while rows 456, even if they had matched without the label, are groups separately. Remember Amazon said they "don't cross boundaries."

Of course, that does not mean only the label determines what is considered to be a match. (That would be of little use.) It's the other columns that determine what matches. The label just confines that matching to records with that label. So, it's matching within a subset of records. It's like having n number of files with no label instead of one file with n labels, so you can run the process one time and not n times.

Anyway, that's the conclusion I draw from this design. Perhaps yours will differ.

Crawl S3

789

We start with the crawlers. Here is the metadata extracted from the diabetes.csv file in S3:

Schema	ı		
	Column name	Data type	Partition k
1	recordid	bigint	
2	pregnancies	bigint	
3	glucose	bigint	
4	bloodpressure	bigint	
5	skinthickness	bigint	
6	insulin	bigint	
7	bmi	double	
8	diabetespedigreefunction	double	
9	age	bigint	
10	outcome	bigint	

_____It created these tables in the

database.

				Choose a d	ata source	
				AWS Glue d	ata catalog	
Q,	Filter by attributes or search	h by keyword				
	Name	÷	Database	÷	Location -	Classi
۲	diabetes_csv		diabetes		s3://sagemakerwalkerml/diabetes.csv	CSV
0	s3_		diabetes		s3://sagemakerwalkerml/s3:/	Unkno
0	sagemaker		diabetes		s3://sagemakerwalkerml/sagemaker/	Unkno

that has access to S3 and give the transformation a name.

Pick an IAM role

Configure transform properties]
Transform name	
diabetes	
Description (optional)	
Enter description	
IAM role ()	
AWSGlueServiceRole-S3IAMRole ~ 2	
Ensure that this role has permission to your Amazon S3 sources, temporary directory, and labeling files. Create	
 Task run properties (optional) 	
 Tags (optional) 	
Next	The data must have a prin

key. The matching algorithm requires that to do its matching logic.

	Choose a primary key		
hoose a	a primary key column. This column typically contains a unique identifier for every record in the data source.		
Q, Filter	r by attributes or search by keyword		
		Showing: 1 -	
	Column name	Data type	
۲	recordid	bigint	
0	pregnancies	bigint	
0	glucose	bigint	Then it

you to tune the transformation. These are tradeoffs between cost and accuracy:

- Cost is financial.
- Cost function is data science and computing.

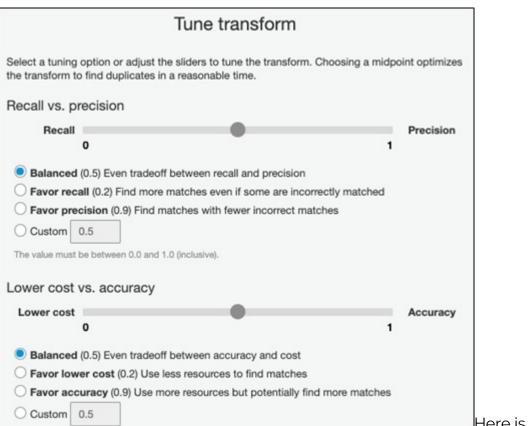
(Pricing is based on resources (DPUs) you consume, which *Lover below*.)

The data science-related tuning parameters are between **recall** and **precision**.

Recall is:

 $\frac{true \ positives}{predicted \ results} = \frac{true \ positives}{true \ positives + false \ positives}$ Precision is:

true positives _____ true positives actual results true positive + false negative



```
Here is a summary of the
```

parameters:

	Review
	Machine learning transforms
Transform name	diabetes
IAM role	arn:aws:iam::782976337272:role/service-role/AWSGlueServiceRole-S3IAMRole
Data source	diabetes_csv
Primary key	recordid
Precision-recall tradeoff	0.5
Accuracy-cost tradeoff	0.5
Force output to match labels	false
Tags	-
	Task run properties
Worker type	G.2X
Number of workers	10
Task timeout (minutes)	2880
Number of retries	0
Spark Version	2.4
	Back

- **The first time** let it generate a label file for you. It will match records based on all of the data points taken together.
- **The second time** it will incorporate labels in its matching algorithm should you choose to add one.

Teach the transform using la	abels
Teach your machine learning algorithms by providing examples, calle provide examples of matching and non-matching records.	ed labels. For your transform,
Labeling file	
 I do not have labels I have labels 	
Generate the labeling file AWS Glue extracts records from your source data and suggests potential matching records. The file will contain approximately 100 data samples for you to work with.	Generate labeling file
You can download the file once it has been generated.	Download labeling file
Next	

Edit the file in Excel or

Google Sheets to both review it and optionally add a label. Copy it back to S3, putting it in a different bucket than the original upload file. Then run transformation again (called **train**). It will produce yet another label file which is the results of the matching aka grouping process.

Upload labels		
Labeling file		
◯ I do not have labels		
I have labels		
Upload labels from S3		
The completed labeling file must be in the correct format and in Amazon S3.	Upload labeling file from S3	
Back Next		It asks for the bucket name:
]
Generate labeling	file	
S3 path to store the generated label file		
s3://bucket/prefix/object	b	
Generate		
		You download the labels from

this screen.



Here is the first label

file it created. You can't see all of the columns because it's too wide. But you can see the labeling_set_id, thus how it grouped the data:

	A	8	с		D	E	- F		G	н	1	
1	labeling_set_id	label	pregnancies	bmi		outcome	age		diabetespedigre	bloodpressure	skinthickness	ree
2	244fcf87-d3a6-36cb-81be-43a29dece209		0 0)	40.5	(0	44	1.781	82		0
3	244fcf87-d3a6-36cb-81be-43a29dece209		0 1	1	22.4	(0	27	0.207	122		0
-4	244fcf87-d3a6-36cb-81be-43a29dece209		1 10)	32.4	1	1	42	0.272	66		0
5	244fcf87-d3a6-36cb-81be-43a29dece209		0 1	1	24.3	(0	21	0.187	58		0
6	244fcf87-d3a6-36cb-81be-43a29dece209		0 0)	18.4	(0	27	0.582	76		0
7	244fcf87-d3a6-36cb-81be-43a29dece209		0 2	2	25.3	(0	22	0.881	62	1	10
	244fcf87-d3a6-36cb-81be-43a29dece209		0 6	3	30.8	(0	37	0.122	70	3	32
9	244fcf87-d3a6-36cb-81be-43a29dece209		0 1	1	32	(0	22	0.389	68	3	35
10	244fcf87-d3a6-36cb-81be-43a29dece209		0 0)	20	(5	22	0.236	68	2	22
11	244fcf87-d3a6-36cb-81be-43a29dece209		1 13	3	42.3		1	44	0.257	114		0
12	28603e1d-6ae9-38af-9801-c5bb2d7ff7c3		0 7	7	27.4	(5	40	0.294	64		0
13	28603e1d-6ae9-38af-9801-c5bb2d7ff7c3		1 7	7	30.4	1	1	36	0.383	66		0
14	28603e1d-6ae9-38af-9801-c5bb2d7ff7c3		0 2	2	28.7	0	5	25	0.092	68	2	22
15	28603e1d-6ae9-38af-9801-c5bb2d7ff7c3		1 1	1	43.3	1	1	41	0.282	0		0
16	28603e1d-6ae9-38af-9801-c5bb2d7ff7c3		0 0)	22.1	(0	21	0.207	62	1	17
17	28603e1d-6ae9-38af-9801-c5bb2d7ff7c3		1 3	3	31.6	1	1	28	0.268	70	() (I	18
18	28603e1d-6ae9-38af-9801-c5bb2d7ff7c3		0 3	3	26.3	(0	24	0.107	74	1	15
19	28603e14.6ae9.38af.9801.c5bb24787c3		0 0	1	36.3	(1	23	0.804	80		37

Evaluation metrics

This screen lets you calculate accuracy. I have yet to figure out where you can see the results as the screen mentioned in the documentation does not exist. (I will update this tutorial once I get a response on the user forum.)

edictions using a subset of	Estimate quality metrics (option bility to find matches. Estimates are calculated by of your labeled data against the labels you have pr	comparing t	
proximate. To improve yo	our transform quality, provide more labels.		
	that you can monitor in the History pane of the tra ed in the Quality metrics pane of the transform.	nsform. Whe	en the task complete
w estimates can be view	ed in the quality metrics pare of the transform.		
Estimate transform qualit	ty		
You have not vet es			
Vou baue pot ust or			
• • • • • • • • • • • • • • • • • • • •	stimated the quality of the transform. Click "Estima transform.	te transform	quality to get
the quality of your t		te transform	quality to get
• • • • • • • • • • • • • • • • • • • •		Result	Last modified
the quality of your t	ransform.		
the quality of your t Quality metric Area under the	Definition Single number summarizing the performance		
the quality of your t Quality metric Area under the Precision-Recall curve	Definition Single number summarizing the performance of the transform When your transform predicts a match, how		
the quality of your t Quality metric Area under the Precision-Recall curve Precision	Definition Single number summarizing the performance of the transform When your transform predicts a match, how often is it correct? For an actual match, how often does your		Last modified
the quality of your t Quality metric Area under the Precision-Recall curve Precision Recall upper limit	Definition Single number summarizing the performance of the transform When your transform predicts a match, how often is it correct? For an actual match, how often does your transform predict a match? Indicates transform's accuracy. Harmonic mean of Precision and Recall.		Last modified

Pricing

Price is by DPU. I used 10 DPUs for about 30 minutes. It's \$0.44 for each multiple or fraction of an hour. So presumably I spent \$0.44*10=\$4.40.

These ETL jobs run on Amazon's Spark and Yarn infrastructure. If you want to write code to do transformations you need to set up a **Development Endpoint**. Basically, the development endpoint is a <u>virtual machine</u> configured to run Spark and Glue. We explained how to use a Development Endpoint <u>here</u>. Then you can run Python or Scala and optionally <u>use a Jupyter Notebook</u>.

Important note: When you don't need your development endpoint, be sure to delete it—it gets expensive quickly! (I spent \$1,200 on that in a month.)

Additional resources

For more tutorials like this, explore these resources:

- BMC Machine Learning & Big Data Blog
- Apache Spark Guide, with 15 articles and tutorials
- AWS Guide
- <u>Amazon Braket Quantum Computing: How To Get Started</u>