

Project report

Project: Amazon Distribution Center

By: Miguel Mayhuire

Definition

Project Overview

This project addresses a use case in Amazon logistics which is the number of items per bin. A solution to use case can leverage problems like orders assignments (is the number of items per order, correct?) or stock planning.

The proposed solution is an endpoint deployed in SageMaker. This endpoint is going to process images of the bins and it is going to respond with the number of items per bin. To do so, a computer vision model is trained and deployed. But before this step, the data is processed, analyzed, and finally ingested to a production machine learning model.

Problem Statement

The project and business need are very clear.

- Incorrect number of items per bin means incorrect order fulfillment and therefore low customer satisfaction.
- Incorrect number of items per bin means that items stock imbalance which means logistic costs and time.
- The strategy and means to deploy a solution are doable. Using Sagemaker capacities and photos coming from the warehouse machines, it is possible to deploy an end-to-end computer vision solution based on data.

Metrics

First, the problem is clearly a supervised machine learning problem with multi-labels. the model in its training loop uses log-loss metric and accuracy for evaluation. Accuracy is the right measure because what we are looking for is high degree of correct prediction for each class.

Analysis

Data Exploration

The data is provided by Amazon and can be found in the following link: <https://registry.opendata.aws/amazon-bin-imagery/>. The data is provided in a .zip file that has to be unzip. The unzip file have 5 files with names 1,2,3,4 and 5 that correspond the labels or number of items by bin. Each file has different number of images:

- 1 item 1K images
- 2 items 1.8K images
- 3 items 2.3K images
- 4 items 2K images
- 5 items 1.8K images

Exploratory Visualization

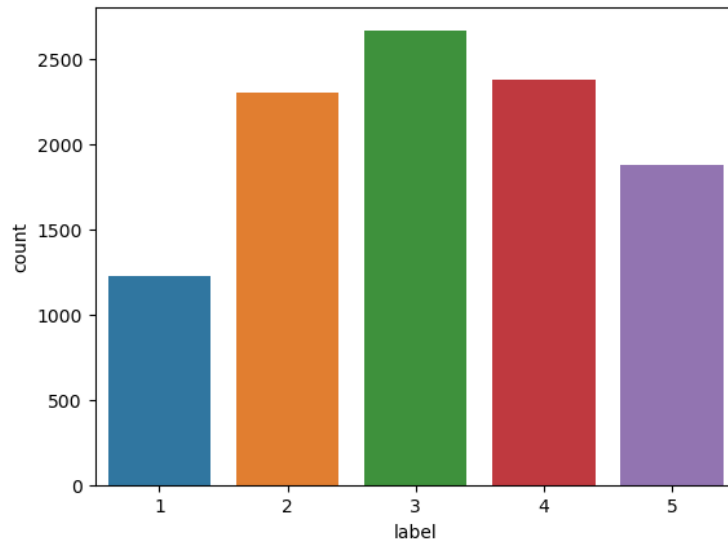
Here I share a sample of images with label 3 or 3 items in a bin:



the images are jpg and some characteristics are:

- images are somewhat dark
- some images seem to be difficult to classify (even using image inspection)
- some images seem to have incorrect labels
-

the distribution of the whole file images is:

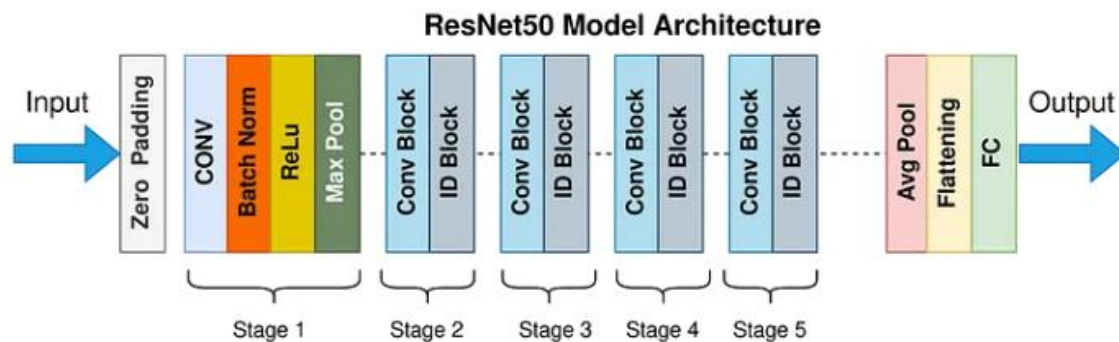


For the purpose of this use case, a stratified sample is taken 10% of validation data and 10% of testing data

Algorithms and Techniques

The chosen model for this task was ResNet50 (pretrained model) plus a custom layer for finetuning.

ResNet is a good model because it is a powerful model for computer vision. In short, the architecture relies on a very long chain of convolutional cells, each cell contains convolutional layers from the previous cell that is used as a residual layer that helps to add noise to the CNN and to address overfit



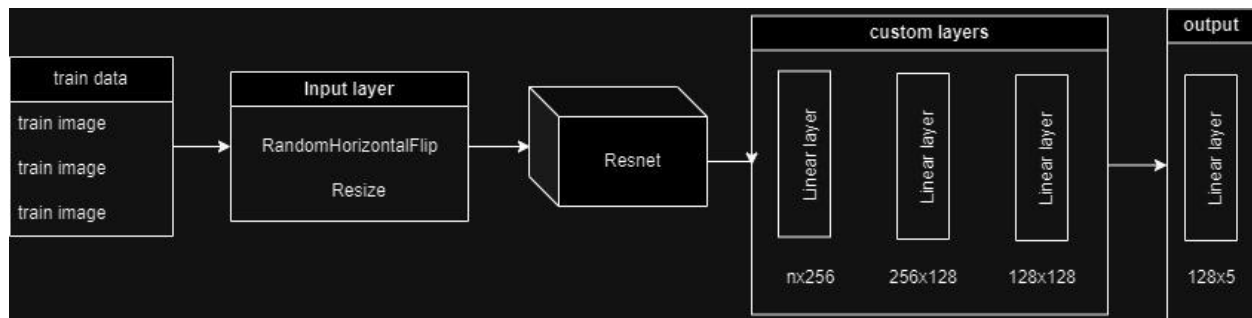
for hyperparameter tuning I used:

- learning rate tuning between 0.001 and 0.1
- batch size options: 16, 32 and 64
- epochs between 2 and 6

Benchmark

In the previous section I described the in general terms what kind of algorithm was suitable for this use case. In this section I am going to explain in depth how the base line model is going to be set.

The base architecture of the base line is:



Description:

- first an input layer where there is a simple data augmentation operation called random horizontal flip. After the data augmentation there is a resizing layer that helps the input fit in the next step
- a pretrained Resnet model with fixed parameters or weights
- then the custom layers are composed of 3 linear layers with activation function Relu. The initial layer is $n \times 256$ where n is the number of end notes from the Resnet. The next layers are 256×128 and 128×128
- finally an end layer that takes 128 input from the previous layer and outputs 5 because of the number of labels

now that the base architecture is set, best parameters can be found and thanks to hyper parameter tuning job it is possible to train many models and to evaluate them completely:

For this use case I set 8 jobs or models

Training job status counter						
Completed	In Progress	Stopped	Failed	(Retractable: 0, Non-retractable: 0)		

Training jobs						
Sorting by objective metric value will display only jobs that have metric values.						
<input type="text" value="Search training jobs"/> View logs View instance metrics Stop Create model						
Name	Status	Final objective metric value	Creation time	Training Duration		
pytorch-training-230904-0923-008-7ee44eec	Completed	0.10148215293884277	9/4/2023, 11:55:41 AM	3 minute(s)		
pytorch-training-230904-0923-007-0a747749	Completed	0.05024878680706024	9/4/2023, 11:52:28 AM	3 minute(s)		
pytorch-training-230904-0923-006-5b5698e4	Completed	0.025828825309872627	9/4/2023, 11:48:02 AM	4 minute(s)		
pytorch-training-230904-0923-005-2d6c6454	Completed	0.025746991857886314	9/4/2023, 11:43:35 AM	4 minute(s)		
pytorch-training-230904-0923-004-6457485a	Completed	0.025875814259052277	9/4/2023, 11:39:11 AM	4 minute(s)		
pytorch-training-230904-0923-003-0b18583d	Completed	0.02587529644370079	9/4/2023, 11:34:48 AM	4 minute(s)		
pytorch-training-230904-0923-002-c81fc01d	Completed	0.025889519602060318	9/4/2023, 11:30:23 AM	4 minute(s)		
pytorch-training-230904-0923-001-8196581e	Completed	0.02585635334253311	9/4/2023, 11:23:06 AM	5 minute(s)		

The best hyper parameter for the baseline model is:

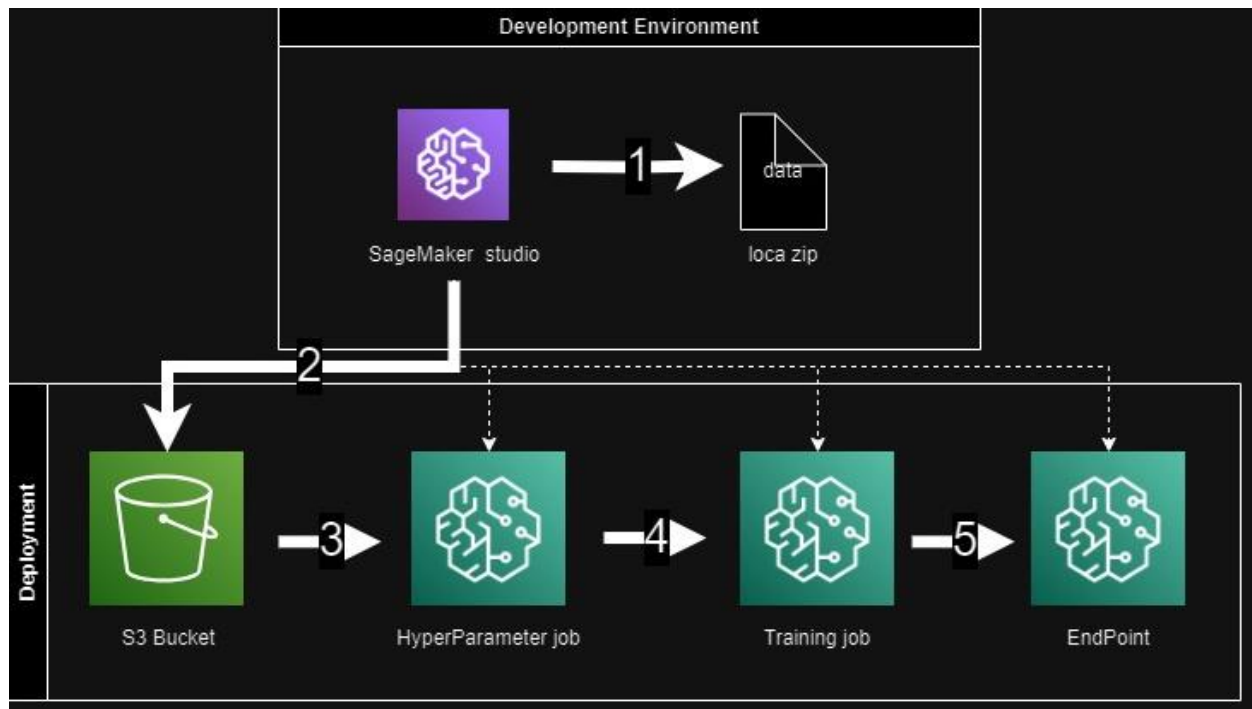
```
'batch_size': 64,
'epochs':2,
'lr': 0.06524108473213124
```

The criteria to choose the best model was the one with the least log loss error and which log loss error was lower than 0.003 using the evaluation data (part of the training phase). Using 500 random image sample from the test set the accuracy is 0.297



Methodology

Here I share a general view of the full pipeline for this project



Data Preprocessing

the images are jpg and some characteristics are:

images are somewhat dark

some images seems to be difficult to classify (even using image inspection)

some images seems to have incorrect labels

export the data to s3 such that execute training jobs – In order to run training jobs, data must be in an S3 bucket. The bucket structure has train, validation and test files

Implementation

Search for best hyper parameters – This stage ingests the data from s3 and trains different models using different parameter settings. Note that to improve execution time, multi-instance was set and it helped to improve execution time

Best parameters:

'batch_size': 64,

'epochs':2,

'learning rate': 0.06524108473213124

Model training – using the best parameter training jobs where launched:

Training jobs <small>Info</small>						
<input type="text" value="Search training jobs"/>						
<div><div>Actions</div><div>Create training job</div></div>						
<div>< 1 ></div>						
Name	Creation time	Duration	Job status	Warm pool status	Time left	
<input type="radio"/> pytorch-training-2023-09-04-12-22-25-264	9/4/2023, 2:22:25 PM	an hour	Completed	-	-	
<input type="radio"/> pytorch-training-2023-09-04-12-17-37-335	9/4/2023, 2:17:37 PM	4 minutes	Failed	-	-	

The first training job failed because there were some problems with code for validation. Once the problem was resolved, a next training job was successfully executed. Here I share some results:

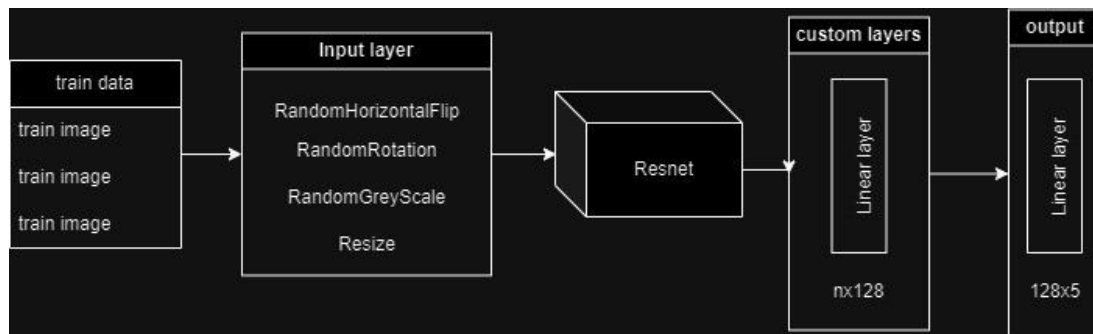
```
PoorWeightInitialization: IssuesFound
Validation set: Average accuracy: 33.42911877394636%
Validation set: Average accuracy: 34.2911877394636%
Validation set: Average accuracy: 30.07662835249042%
Validation set: Average accuracy: 32.662835249042146%
Validation set: Average accuracy: 29.214559386973182%
Validation set: Average accuracy: 32.56704980842912%
Test set: Average loss: 0.02433362703003952, Accuracy: 29.856459330143544%
Saving the model
```

- profiler results
 - 1 training hour
 - just 30% of the execution time is spent in the validation and training then the rest of the time is spent in other tasks
 - of the training time, the model is spending mainly in executing the convolutional layers.
- Debugger:
 - Poor weight initialization was found
- Validation
 - Validation accuracy is low even though log loss is low

Refinement

- Model performance can improve if more training data is provided and further architectures are tested
- Model training time can be improved if GPU resources are used
- Maybe vanishing gradient hooks are making the task to take longer, (because the machine is assessing if the model weight optimization is doing well) I would suggest to remove this hook and see if the model trains faster

A new architecture was set to compete with the baseline model:



The changes:

- Further data augmentation layers in the input layer. Random rotation (goal generalize better the item distribution inside a bin), random grey scale (increase the samples with low light contrast)
- The custom layer is shortened and just a single layer is used Nx128 meaning that just a layer of 128 nodes are trained to make this block of layers less complex
- Epochs for hyper parameter tuning is set 4 and 20 epochs goal is to train the model many more times compared to the baseline model

After hyperparameter tuning using the new architecture, the best parameters for this model are:

```

'batch_size': 64,
'epochs': 20,
'lr': 0.09980513226333607

```

Using 100 random images from the test data, the accuracy score is 33.8%

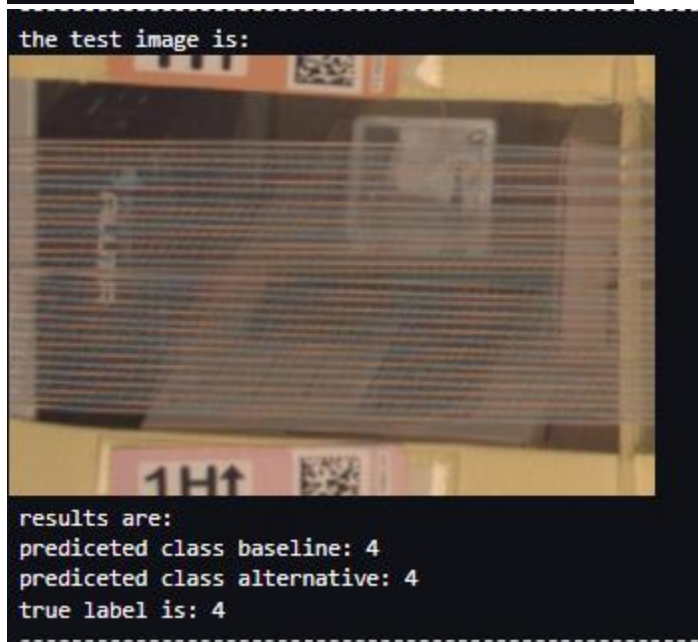
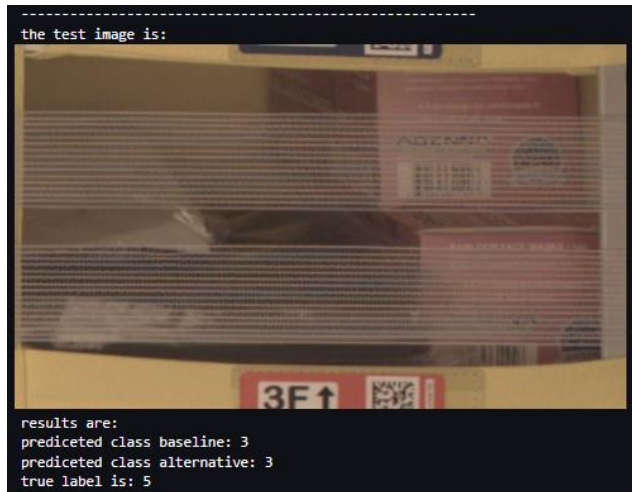
Results

Model Evaluation and Validation

Comparative table:

Items	Baseline	Alternative
Data augmentation	Horizontal Flip	Horizontal Flip, rotation and grey color sclaing
Epochs	2	20
Learning rate	0.065	0.0998
Custom Layers	3	1
Training time	1 hrs	6.3 hrs
Sample Accuracy	0.297	0.338

It seems that the alternative model is better than the base line model, lets do some inspections:



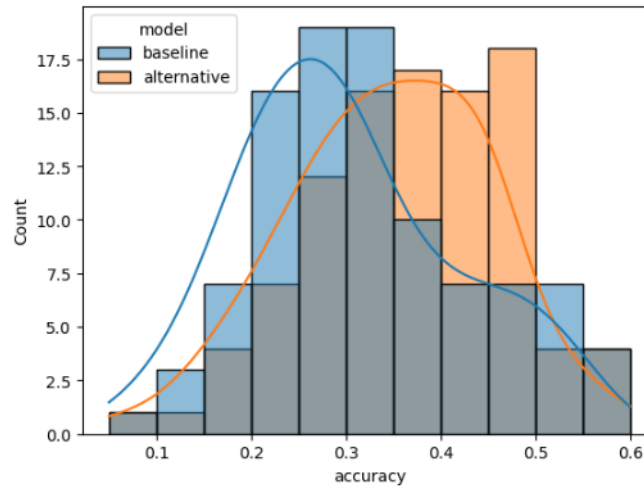
What can be spotted in both models:

- Still general accuracy is low
- Dissimilar predictions
- Endpoint response time is fairly fast

To be more rigorous with the evaluation of both models I will run some tests in the following section.

Justification

Here I share some statistical analysis of both models, baseline and alternative model. I took 100 random samples with replacement of 10 images (bootstrapping) for both models and the results are:



Metric	baseline	alternative
Mean accuracy	0.30	0.345
Standard deviation	0.115	0.106

T test p value: 0.00516 (alpha 0.05) therefore means are different

Thanks to bootstrapping and t test for two population over accuracy metrics on different sample test data we see that:

- Alternative model has better accuracy
- And t test demonstrate that this value is statistically different than the baseline accuracy 0.30

Therefore we can conclude that the alternative model is a best model than the baseline.