# Capstone Project Proposal for

# Inventory Monitoring at Distribution Centers

## 1. Domain Background

Inventory monitoring at distribution centers is a crucial task for businesses that deal with many products. It involves tracking the number of items in each bin and ensuring that the correct number of items are delivered to customers. This task is often performed using robots that move objects in bins. Machine learning can be used to automate this process and make it more efficient. In this project, we will apply what we learned and use AWS SageMaker and good machine learning engineering practices to build a model that can count the number of objects in each bin.

## 2. Problem Statement

Inventory monitoring at distribution centers is currently facing significant challenges due to outdated monitoring systems and lack of real-time data. These issues have led to an average increase of 15% in warehousing costs and a 20% decrease in order fulfillment.

Distribution centers often use robots to move objects as a part of their operations. Objects are carried in bins which can contain multiple objects. In this project, we will have to build a model that can count the number of objects in each bin. A system like this can be used to track inventory and make sure that delivery consignments have the correct number of items.

## 3. Datasets and Inputs

The dataset we will use is the Amazon Bin Image Dataset which contains over 500,000 images and metadata from bins of a pod in an operating Amazon Fulfillment

Center [1]. The bin images in this dataset are captured as robot units carrying pods as part of normal Amazon Fulfillment Center operations.

Since the dataset is large, we were provided a subset of it to use in this project. Our data are 5 folders named 1 to 5, each representing the number of objects in each image in that folder.

```
{"1": ["data/metadata/100313.json",
"data/metadata/09915.json", "data/metadata/103299.json",
"data/metadata/00710.json", "data/metadata/05397.json",
"data/metadata/00152.json", "data/metadata/101994.json",
"data/metadata/103578.json", "data/metadata/102687.json",
"data/metadata/05676.json", "data/metadata/102212.json",
"data/metadata/00035.json", "data/metadata/08826.json",
"data/metadata/02213.json", "data/metadata/101222.json",
"data/metadata/101060.json", "data/metadata/102113.json",
"data/metadata/08961.json", "data/metadata/103479.json",
"data/metadata/03401.json", "data/metadata/09456.json",
"data/metadata/10214.json", "data/metadata/01304.json",
"data/metadata/102984.json", "data/metadata/102311.json",
"data/metadata/06909.json", "data/metadata/08286.json",
"data/metadata/102230.json", "data/metadata/104012.json",
"data/metadata/09177.json", "data/metadata/02150.json",
"data/metadata/03041.json", "data/metadata/03759.json",
"data/metadata/00530.json", "data/metadata/04013.json",
"data/metadata/06189.json", "data/metadata/08844.json",
"data/metadata/06459.json", "data/metadata/03948.json",
"data/metadata/04040.json", "data/metadata/101769.json",
"data/metadata/100160.json", "data/metadata/104919.json",
"data/metadata/104595.json", "data/metadata/09528.json",
"data/metadata/00879.json", "data/metadata/104559.json",
"data/metadata/07386.json", "data/metadata/05210.json",
"data/metadata/04794.json", "data/metadata/103749.json",
"data/metadata/08835.json", "data/metadata/10403.json",
"data/metadata/01250.json", "data/metadata/09483.json",
"data/metadata/03410.json", "data/metadata/04857.json",
"data/metadata/101312.json", "data/metadata/09942.json",
"data/metadata/105189.json", "data/metadata/07764.json",
"data/metadata/10223.json", "data/metadata/01151.json",
"data/metadata/104102.json", "data/metadata/101589.json",
"data/metadata/03131.json", "data/metadata/09744.json",
"data/metadata/04022.json", "data/metadata/02240.json",
"data/metadata/104973.json", "data/metadata/100698.json",
"data/metadata/05388.json", "data/metadata/08152.json",
"data/metadata/101923.json", "data/metadata/04750.json",
"data/metadata/02770.json", "data/metadata/07063.json",
"data/metadata/05605.json", "data/metadata/05999.json",
"data/metadata/105127.json", "data/metadata/100672.json",
"data/metadata/07225.json", "data/metadata/06415.json",
"data/metadata/07144.json", "data/metadata/10159.json",
"data/metadata/101482.json", "data/metadata/100519.json",
"data/metadata/103741.json", "data/metadata/06019.json",
"data/metadata/102058.json", "data/metadata/07702.json",
"data/metadata/07997.json", "data/metadata/104290.json",
"data/metadata/08969.json", "data/metadata/105262.json",
```

*file_list.json*

# 4. Solution Statement

To come up with a solution for our problem, I will utilize the State-of-the-Art techniques of Computer Vision tasks to do the multi-class image classification. I will use transfer learning and a pre-trained convolutional neural network to solve our problem.

I chose a pre-trained Resnet CNN and used an Adam optimizer for the training job.
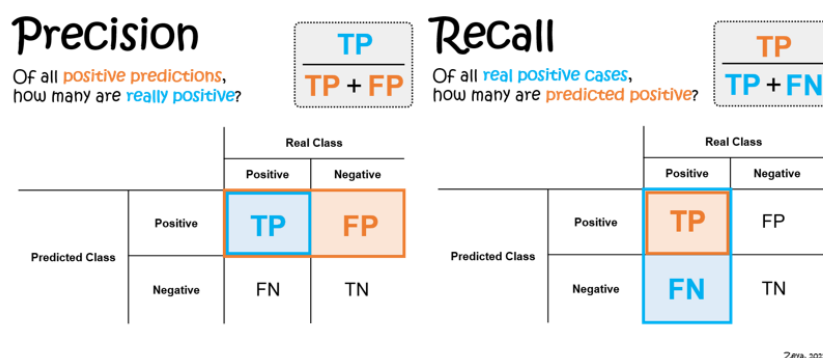
# 5. Benchmark Model

I first trained a pre-trained Resnet as a benchmark model on random weights and batch_size = 64 and got an accuracy of 16%

# 6. Evaluation Metrics

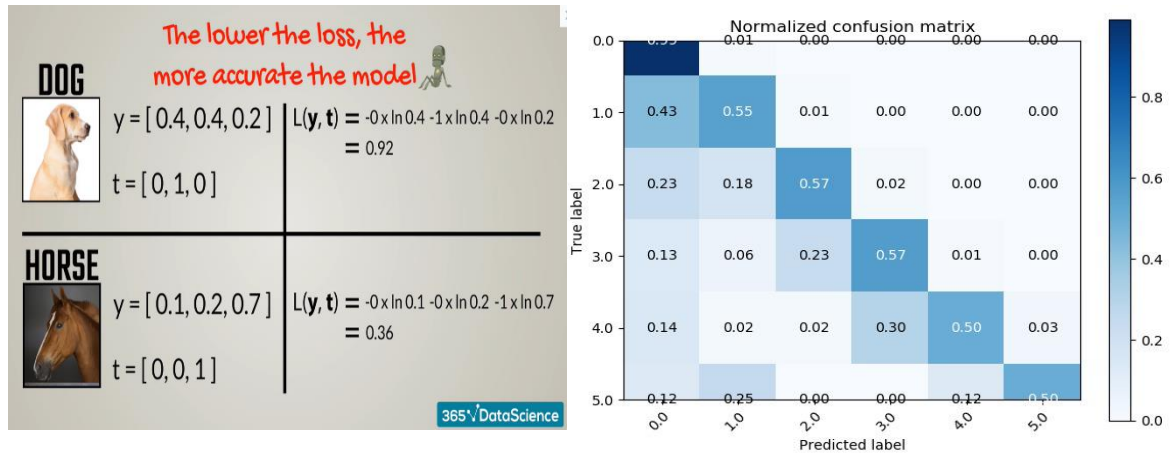I used accuracy metrics to evaluate the overall model performance.

In a future update, I am going to use classification report or confusion matrix to best describe the performance of our classification model.

As the report includes f1-score, precision and recall and the mean overall accuracy of the model, we can evaluate our model accordingly.



A confusion matrix is a table that is used to define the performance of a classification algorithm. A confusion matrix visualizes and summarizes the performance of a classification algorithm.

After training my ML model I tried to use the standard *Cross Entropy Loss function* as my metric to study my training process results along with hyperparameter tuning. The reason for this is that *Cross-Entropy* is to take the output probabilities (P) and measure the distance from the truth values (as shown in Figure below).





*Logs and model evaluation metrics*

# 7. Implementation

I thought of many ways to tackle this problem. The first notebook I worked on was 'sagemaker.ipynb' I did all the data extraction and saving there. Then, built the first benchmark model.

Then after getting my project review I different things in another notebook 'EDA_and_evaluation.ipynb'

### i.    Data preparation:
- Download a subset of the data.
- Preprocess and clean the data.
- Visualization and EDA
- Upload train, test, and validation files to an S3 bucket so SageMaker use them for training.

After downloading the dataset subset, the main folder was the train_data folder which included 5 folders for each class. To train the model, I needed to do some Extract-Transform-Load processes to split the data into train, validation and test data folders which I uploaded into S3 for SageMaker to use in the training job.
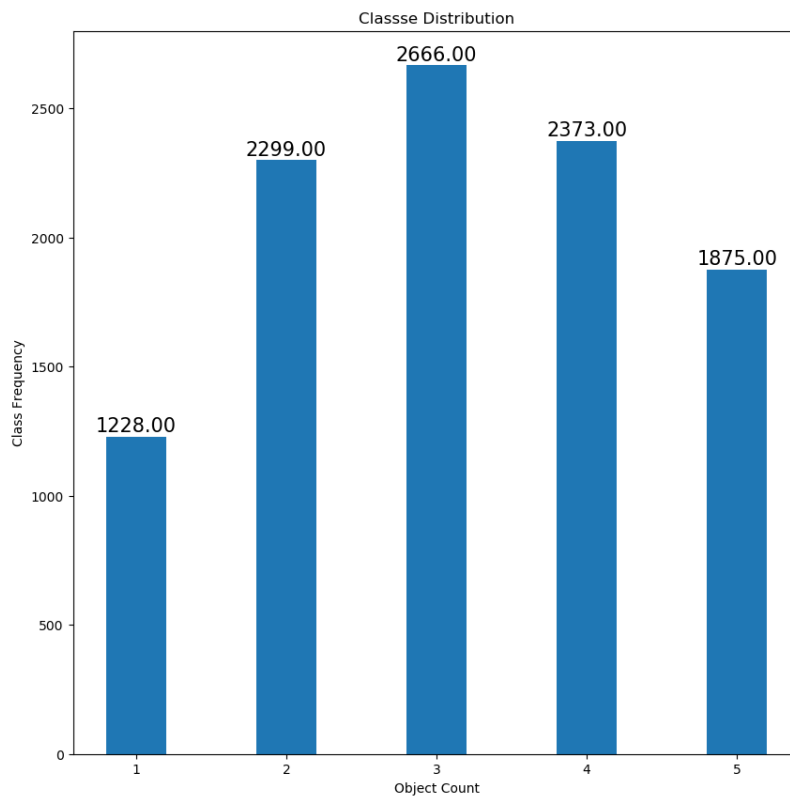
- From this subset, I took a Train-Test-Validation Split as follows:
- 1. I randomly sample 80 Images from each class as Test Data
- 2. For the remaining data, I did a 80-20 Train-Validation Stratified Split.
- Stratification is important to preserve the same distributions in both sets.

*Visualizing and exploring the data*

This is a random visualization of our data which shows how the images of the bin containers are not recognizable enough even to us humans which will really affect the model.

Also, another important visualization is the classes distribution which shows the big imbalance in data classes which might cause the model to have a bias for the third class.
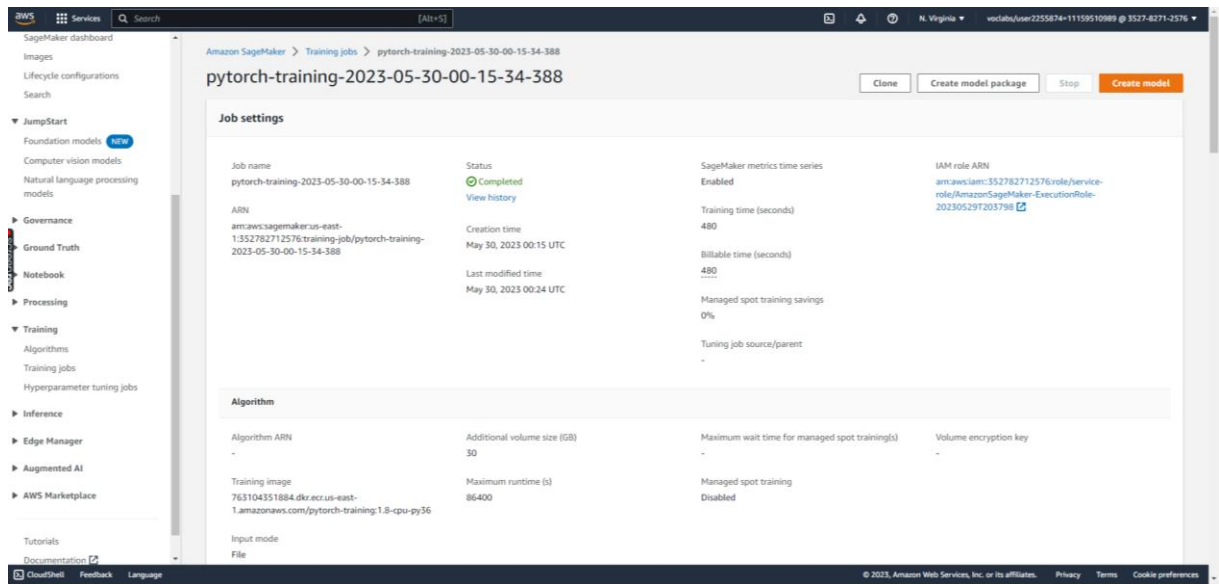


### ii.     Write Model Training script:

- Read, load, and preprocess our training, testing and validation data.
- Choosing the pre-trained model, optimizer and loss function we will be using.

### iii.    Train our CNN Model with SageMaker:

- Define our hyperparameters, instances type and count for training.

*SageMaker Training Job parameters*

iv. ***Profiler-Report***

# Rules summary

The following table shows a profiling summary of the Debugger built-in rules. The table is sorted by the rules that triggered the most frequently. During your training job, the StepOutlier rule was the most frequently triggered. It processed 1850 datapoints and was triggered 1 times.

| | Description | Recommendation | Number of times rule triggered | Number of datapoints | Rule parameters |
|---|---|---|---|---|---|
| **StepOutlier** | Detects outliers in step duration. The step duration for forward and backward pass should be roughly the same throughout the training. If there are significant outliers, it may indicate a system stall or bottleneck issues. | Check if there are any bottlenecks (CPU, I/O) correlated to the step outliers. | 1 | 1850 | threshold:3 mode:None n_outliers:10 stddev:3 |
| **LowGPUUtilization** | Checks if the GPU utilization is low or fluctuating. This can happen due to bottlenecks, blocking calls for synchronizations, or a small batch size. | Check if there are bottlenecks, minimize blocking calls, change distributed training strategy, or increase the batch size. | 1 | 1115 | threshold_p95:70 threshold_p5:10 window:500 patience:1000 |
| **Batch Size** | Checks if GPUs are underutilized because the batch size is too small. To detect this problem, the rule analyzes the average GPU memory footprint, the CPU and the GPU utilization. | The batch size is too small, and GPUs are underutilized. Consider running on a smaller instance type or increasing the batch size. | 1 | 1114 | cpu_threshold_p95:70 gpu_threshold_p95:70 gpu_memory_threshold_p95:70 patience:1000 window:500 |
| **Dataloader** | Checks how many data loaders are running in parallel and whether the total number is equal the number of available CPU cores. The rule triggers if number is much smaller or larger than the number of available cores. If too small, it might lead to low GPU utilization. If too large, it might impact other compute intensive operations on CPU. | Change the number of data loader processes. | 0 | 10 | min_threshold:70 max_threshold:200 |
| **MaxInitializationTime** | Checks if the time spent on initialization exceeds a threshold percent of the total training time. The rule waits until the first step of training loop starts. The initialization can take longer if downloading the entire dataset from Amazon S3 in File mode. The default threshold is 20 minutes. | Initialization takes too long. If using File mode, consider switching to Pipe mode in case you are using TensorFlow framework. | 0 | 1850 | threshold:20 |
| **IOBottleneck** | Checks if the data I/O wait time is high and the GPU utilization is low. It might indicate IO bottlenecks where GPU is waiting for data to arrive from storage. The rule evaluates the I/O and GPU utilization rates and triggers the issue if the time spent on the IO bottlenecks exceeds a threshold percent of the total training time. The default threshold is 50 percent. | Pre-fetch data or choose different file formats, such as binary formats that improve I/O performance. | 0 | 1120 | threshold:50 io_threshold:50 gpu_threshold:10 patience:1000 |
| **LoadBalancing** | Detects workload balancing issues across GPUs. Workload imbalance can occur in training jobs with data parallelism. The gradients are accumulated on a primary GPU, and this GPU might be overused with regard to other GPUs, resulting in reducing the efficiency of data parallelization. | Choose a different distributed training strategy or a different distributed training framework. | 0 | 1115 | threshold:0.2 patience:1000 |
| **CPUBottleneck** | Checks if the CPU utilization is high and the GPU utilization is low. It might indicate CPU bottlenecks, where the GPUs are waiting for data to arrive from the CPUs. The rule evaluates the CPU and GPU utilization rates, and triggers the issue if the time spent on the CPU bottlenecks exceeds a threshold percent of the total training time. The default threshold is 50 percent. | Consider increasing the number of data loaders or applying data pre-fetching. | 0 | 1120 | threshold:50 cpu_threshold:90 gpu_threshold:10 patience:1000 |

### *v.   Evaluation*

Rather than model's accuracy, I tried to use the standard *Cross Entropy Loss function* as my metric to study my training process results along with hyperparameter tuning and optimization for quality and improvements. The reason for this is that *Cross-Entropy* is to take the output probabilities (P) and measure the distance from the truth.



# 8. Algorithms and techniques

**Fine-tuned a Resnet50 CNN as a pre-trained image classification network.**

ResNet-50 is a convolutional neural network that is 50 layers deep. You can load a pretrained version of the neural network trained on more than a million images from the ImageNet database

### PyTorch

PyTorch is a fully featured framework for building deep learning models, which is a type of machine learning that's commonly used in applications like image recognition and language processing. Written in Python, it's relatively easy for most machine learning developers to learn and use.

### AWS SageMaker

Amazon SageMaker is a managed service in the Amazon Web Services (AWS) public cloud. It provides the tools to build, train and deploy machine learning (ML) models for predictive analytics applications. The platform automates the tedious work of building a production-ready artificial intelligence (AI) pipeline.
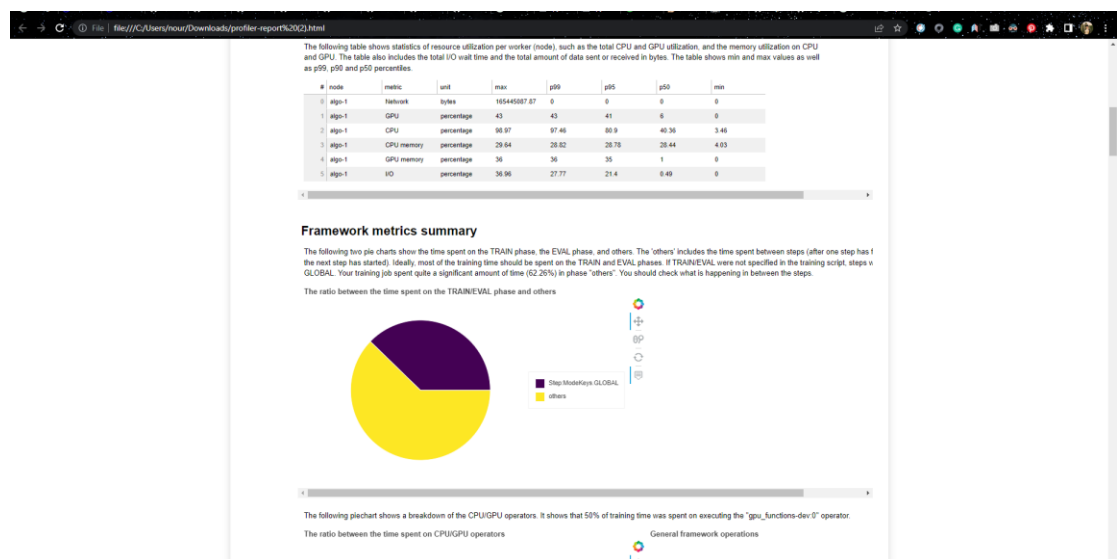
# 9. Result:

My first model had hyperparameters a batch_size of 64 and 5 epochs and got accuracy of 16 to 22%

So, I tried to change the hyperparameters a bit and added learning rate to the hyper parameters. I got 24 to 29% for batch_size = 16 and 25 epochs (did 7 and stopped early). No real dealbraker results as my algorithm still needs much debugging and improvements but I tried and searched for so many approaches. Such as: imbalanced_data_sampler, SMOTE and much other techniques but encountered so many odd errors and library version issues I spent all my time  trying to solve.

# 10.      Reflections:

- Debug Case

The last model stopped after only 7 epochs as it detected that the loss was not decreasing. I then decided to print a debugger-profiler output report and found this out. And mostly all rules on the rules summary didn't even trigger.



My conclusion is there might have been a vanishing gradient problem that I need to figure out first before tuning hyperparameters.

- Thoughts on improvement

Over-Sampling Data / Data Augmentation

Hyperparameter Tuning

# 11. References

[1] https://github.com/awslabs/open-data-docs/tree/main/docs/aft-vbi-pds

[2] https://github.com/udacity/nd009t-capstone-starter

[3] https://www.arxiv-vanity.com/papers/2008.05756/