

# **FINAL REPORT**

## **Retail Product Image Recognition**

### **Problem Statement**

With the recent times where digital transformation has been immense and there has been a significant shift towards online shopping. It still hasn't eradicated or has brought down the offline shopping experience.

People still like to go out to and buy the products, providing them with a sense of control and let them have the feel for their buy. And this is not changing even with the ease of online shopping.

But the problem still exists at the offline stores. The problem of time-consuming billing counters at place.

For the problem in hand, we are going to build an image recognition model that will help us redefine this situation and help reduce the time the customer has to spend at billing.

The model will automatically identify and count the items the user has selected and provide him with a bill of the respected items.

NOTE: At the initial stage of project we only look to identify products to their respective super-category and not to the item itself

## **Stage 1 – Data Wrangling & EDA**

The initial stage of the project was focused at collecting ample number of images for different products ranging to a variety of categories.

The data can be found at the following Kaggle repository.

<https://www.kaggle.com/diyer22/retail-product-checkout-dataset>

The images were divided into three directories, as follows-

Train = It contains 52,000 images of single items present within it

Test = It contains 24000 images of multiple items in a single image

Validation = This is to test the model performance in a more real-world examples where images contain number of items in a single image.

The data is provided to us in the json format having the following information.

```
|— info: dict 7
|— licenses: list 1
|— categories: list 200
|— __raw_Chinese_name_df: list 200
|— images: list 53739
└— annotations: list 53739
```

**Diagram1: Train json**

While doing some analysis it revealed that there are a total of 17 super-categories. These super –categories then in total comprised of 200 different products in them.

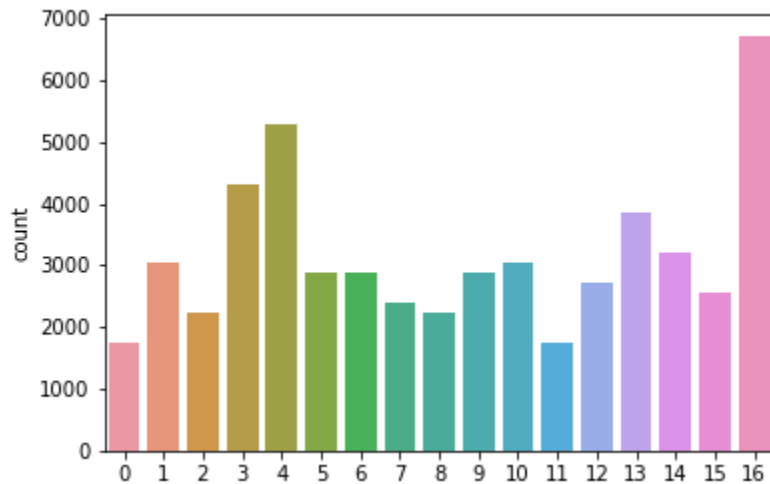
The 17 super-categories and their respective product count we have with us are-

|           | <b>supercategory</b> | <b>Count</b> |
|-----------|----------------------|--------------|
| <b>0</b>  | alcohol              | 11           |
| <b>1</b>  | candy                | 10           |
| <b>2</b>  | canned_food          | 14           |
| <b>3</b>  | chocolate            | 12           |
| <b>4</b>  | dessert              | 17           |
| <b>5</b>  | dried_food           | 9            |
| <b>6</b>  | dried_fruit          | 9            |
| <b>7</b>  | drink                | 15           |
| <b>8</b>  | gum                  | 8            |
| <b>9</b>  | instant_drink        | 11           |
| <b>10</b> | instant_noodles      | 12           |
| <b>11</b> | milk                 | 11           |
| <b>12</b> | personal_hygiene     | 10           |
| <b>13</b> | puffed_food          | 12           |
| <b>14</b> | seasoner             | 12           |
| <b>15</b> | stationery           | 7            |
| <b>16</b> | tissue               | 20           |

The images in the dataset are of these products themselves and are named by a unique code assigned to each product.

Using these image codes, I segregated the images from 1 directory to assigning them into respective super category folders.

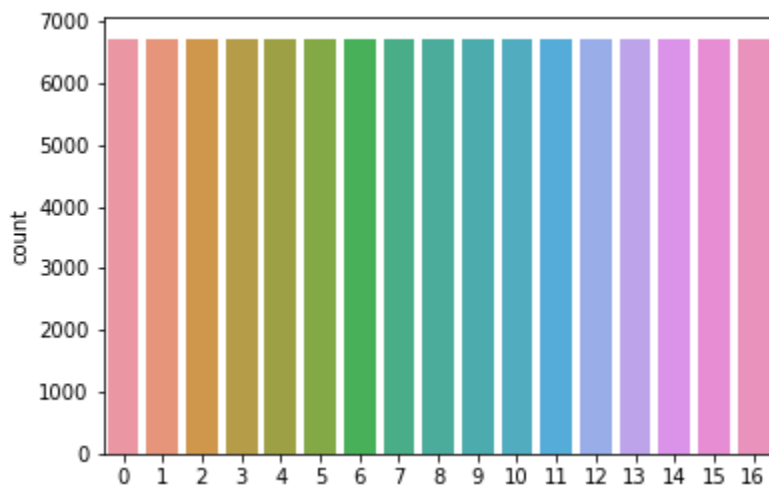
## Training Data Class Distribution



To manage the imbalance in the classes we used to smote to oversample the data.

Since the initial image array was 4 dimensional, we reshaped to a 2D array and applied SMOTE.

After oversampling the data our new dataset looks like,



The new dataset now has a total of 1,14,000 training image array.

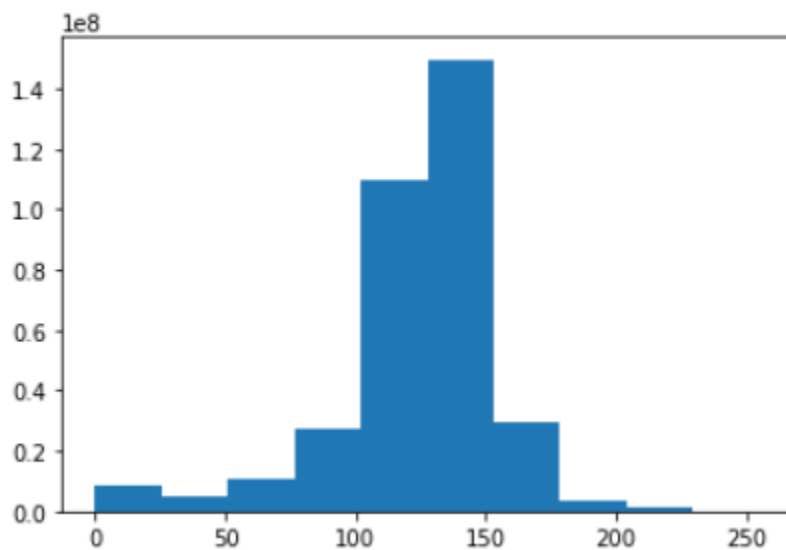
## **Stage 2 – Pre-Processing**

In the pre-processing stage we did the following steps to prepare our data for training our model.

1) Converting the images to their respective arrays.

Using the CV2 library we loaded the images as grayscale and then resizing each image into 80\*80 dimension

2) Having a look at histogram to check the threshold of the image.



3) Normalizing the data

Since the image threshold ranges from 0-255, we can normalize by dividing the dataset by 255.

4) Doing train/test split for training

Finally, we split our data into 80:20 ratio for our model training and evaluation.

### **Step 3 - Training**

In the pre-processing stage we converted each image to its respective array so that our model could interpret it.

We also integer encoded each of the super-category respective to its product.

The range of target variable is 0-16

We got our feature variable in array form as X and target variable in integer form as y.

Now once the data was normalized we began building a simple CNN architecture for our project.

### **MODEL ARCHITECTURE**

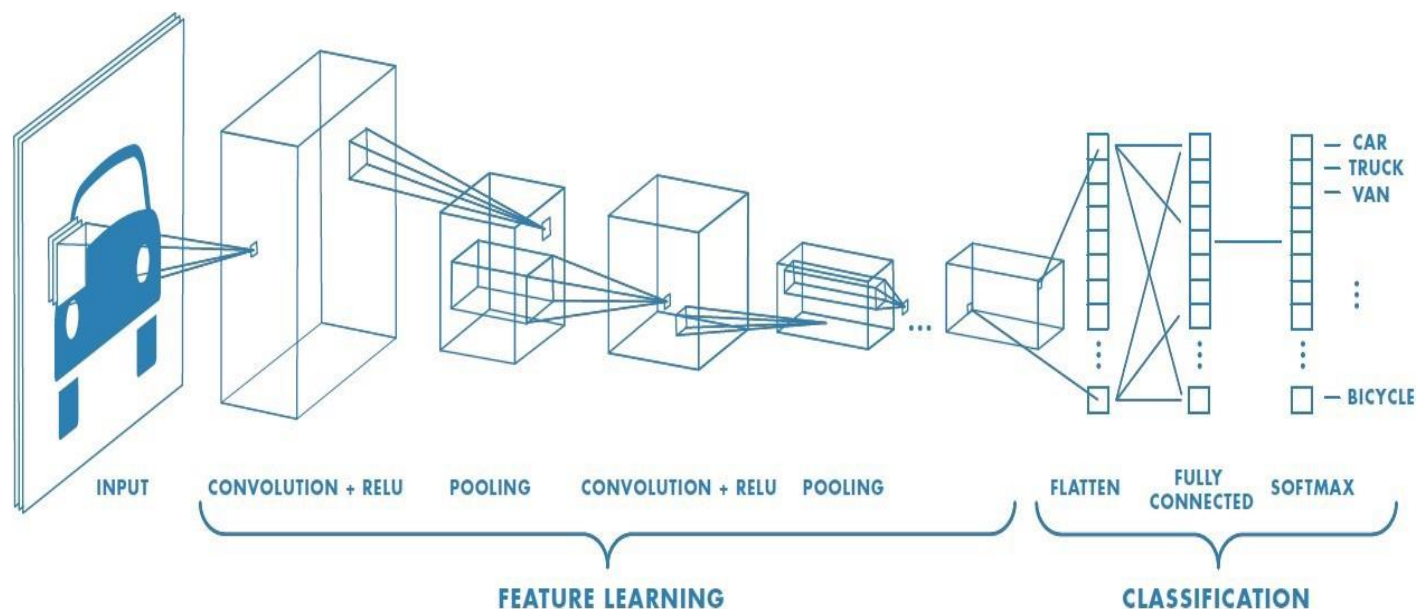
We chose Sequential model for this project.

In the model we added two 64 feature densely connected convolution layer with relu activation function.

Each of the densely connected conv layer was followed by a max pooling layer to reduce the spatial dimensionality, using a 2\*2 step.

After that we connected to a dense layer with 64 features to it again having a relu activation function.

Finally, since we had 17 target features, we added our final dense layer with 17 nodes and having the activation function as SoftMax.



Diag 2. Model Architecture

Features of model compilation -

Optimizer – Adam

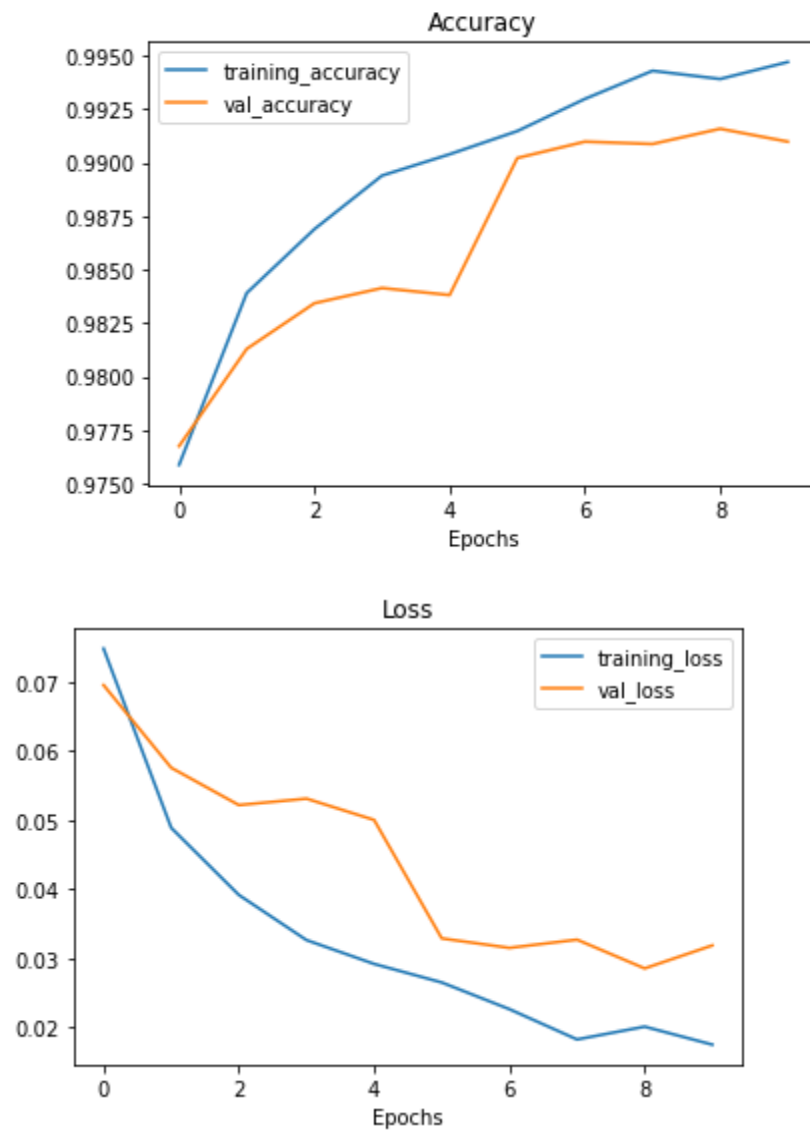
Loss function – sparse-categorical cross entropy

Metrics – Accuracy

## Performance

Accuracy Score – 90.95%

### Loss and Accuracy curve





## Classification Report -

|              | precision | recall | f1-score | support |
|--------------|-----------|--------|----------|---------|
| 0            | 0.95      | 0.90   | 0.92     | 363     |
| 1            | 0.89      | 0.84   | 0.86     | 622     |
| 2            | 0.68      | 0.98   | 0.81     | 485     |
| 3            | 0.90      | 0.96   | 0.93     | 868     |
| 4            | 0.96      | 0.89   | 0.92     | 1024    |
| 5            | 0.94      | 0.85   | 0.89     | 569     |
| 6            | 0.82      | 0.94   | 0.88     | 594     |
| 7            | 0.91      | 0.90   | 0.90     | 471     |
| 8            | 0.89      | 0.98   | 0.93     | 425     |
| 9            | 0.92      | 0.92   | 0.92     | 563     |
| 10           | 0.98      | 0.93   | 0.95     | 660     |
| 11           | 0.96      | 0.69   | 0.80     | 373     |
| 12           | 0.96      | 0.89   | 0.92     | 519     |
| 13           | 0.97      | 0.89   | 0.93     | 774     |
| 14           | 0.89      | 0.94   | 0.91     | 636     |
| 16           | 0.95      | 0.97   | 0.96     | 1302    |
| micro avg    | 0.91      | 0.91   | 0.91     | 10248   |
| macro avg    | 0.91      | 0.90   | 0.90     | 10248   |
| weighted avg | 0.92      | 0.91   | 0.91     | 10248   |

In the list below we can match the class number to its label.

```
['alcohol', 'candy', 'canned food', 'chocolate', 'dessert',  
'dried food', 'dried fruit', 'drink', 'gum', 'instant  
drink', 'instant noodles', 'milk', 'personal hygiene',  
'puffed food', 'seasoner', 'stationery', 'tissue']
```

In the classification report we could notice we had a poor precision for class canned food and for personal hygiene we had a poor recall.

### **Future Scope-**

- 1) We aim to detect more than one object in the image and successfully able to identify the category of it.
- 2) After detecting multiple super-category in one image we aim to improve the model to identify individual products themselves
- 3) We aim to improve the model on classes that are not performing well.

## **Conclusion**

We were able to build a simple CNN model to detect an image with single item and successfully identify the category it belonged to.