

# Machine Learning Engineer Nanodegree

## Capstone Project

Mofei Liu  
May 6th, 2019

### I. Definition

#### Project Overview

In recent years, sneaker culture and the industry itself has been grown a lot. With more and more people start get into it collecting sneakers and using them as an outlet to express themselves, brands like Nike & Adidas start going at a very fast paced when it comes to sneaker release. This made consumers are very overwhelmed by the new product they're seeing every day every week.

In this project, I created a web application that can identify the product for the customer giving an image. The application uses a classifier trained from images scraped from the internet.

#### Problem Statement

The goal is to create a web application with the task involved in the following:

1. Scraping image dataset from the internet, and preprocess the training data
2. Train a classifier that can identify sneaker model
3. Make the classifier run on a web app

The final application is expected to be useful identify sneakers giving an user upload picture

#### Metrics

To evaluation the model vs benchmark model, I'll be using the logarithmic loss method<sup>5</sup>, which works by penalizing the false classification. It should work well here since this is a multi-class classification problem.

Suppose, there are N samples belonging to M classes, then the Log Loss is calculated as below:

$$\text{LogarithmicLoss} = \frac{-1}{N} \sum_{i=1}^N \sum_{j=1}^M y_{ij} * \log(p_{ij})$$

Where,

y<sub>ij</sub>, indicates where sample i belongs to class j or not

$p_{ij}$ , indicates the probability of sample  $i$  belonging to class  $j$

Log Loss has no upper bound and it exists on the range  $[0, \infty)$ . Log Loss nearer to 0 indicates higher accuracy, whereas if the Log Loss is away from 0 then it indicates lower accuracy.

## II. Analysis

### Data Exploration



Fig.1 images from the dataset

The image from the dataset has dimension of 2200 x 1468 pixels, with all 0-360 degree of the product, this add up a very good image argumentation to the dataset. In total we have 52593 images in 179 categories within this dataset, which is sufficient to train the classifier.

```
df=pd.read_csv('output.csv')
```

```
df.head()
```

	brand	categories	colorway	deadstock_sold	gender	style_id	thumbnail_url	url	idx	label
0	adidas	['adidas', 'Yeezy']	Clay/Clay/Clay	1560.0	men	EG7490	https://stockx.imgix.net/adidas-Yeezy-Boost-35...	adidas-yeezy-boost-350-v2-clay	0	Adidas-Yeezy-350-V2
1	adidas	['adidas', 'Yeezy']	Grey/Grey/Grey	8542.0	men	EG7492	https://stockx.imgix.net/adidas-Yeezy-Boost-35...	adidas-yeezy-boost-350-v2-true-form	1	Adidas-Yeezy-350-V2
2	adidas	['adidas', 'Yeezy']	Geode/Geode/Geode	1718.0	men	EG6860	https://stockx.imgix.net/adidas-Yeezy-Boost-70...	adidas-yeezy-boost-700-v2-geode	2	Adidas-Yeezy-700-V2
3	adidas	['adidas', 'Yeezy']	Grey/Grey/Grey	1284.0	men	EG7491	https://stockx.imgix.net/adidas-Yeezy-Boost-35...	adidas-yeezy-boost-350-v2-hyperspace	3	Adidas-Yeezy-350-V2
4	adidas	['adidas', 'Yeezy']	Grey/Grey/Inertia	7665.0	men	EG7597	https://stockx.imgix.net/adidas-Yeezy-Boost-70...	adidas-yeezy-boost-700-inertia	4	Adidas-Yeezy-700

Fig.2 data annotation

The other annotation of the data is stored in a csv file, which have the following field:

- Brand
- Categories
- Colorway
- Deadstock\_sold
- Gender
- Style\_id
- Thumbnail\_url

- url
- index
- label

In this classifier we're only using label as the output of the prediction.

## Exploratory Visualization

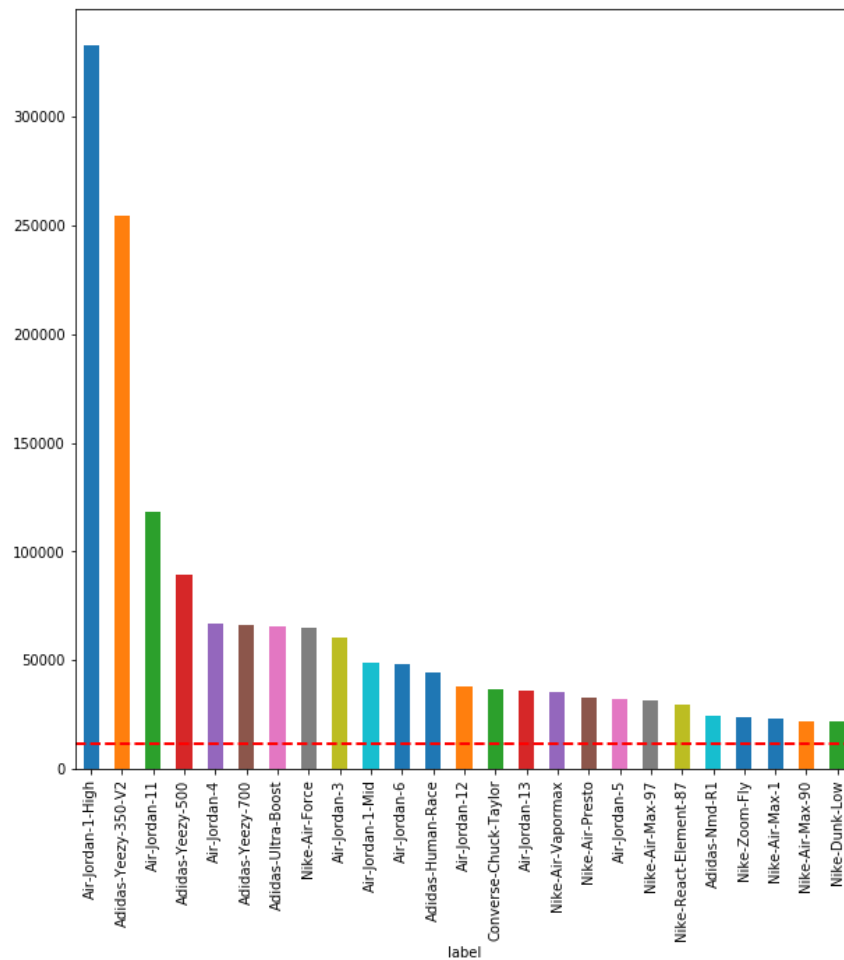


Fig.3 number of sold vs model

- Here's the top 25 most popular sneakers sold on StockX by #of pairs sold.
- The top 3 models are Air Jordan 1 High, Adidas Yeezy 350 V2 and Air Jordan 11.
- On average there's 11385 pair have been sold on the platform for each style
- The distribution is right tailed.

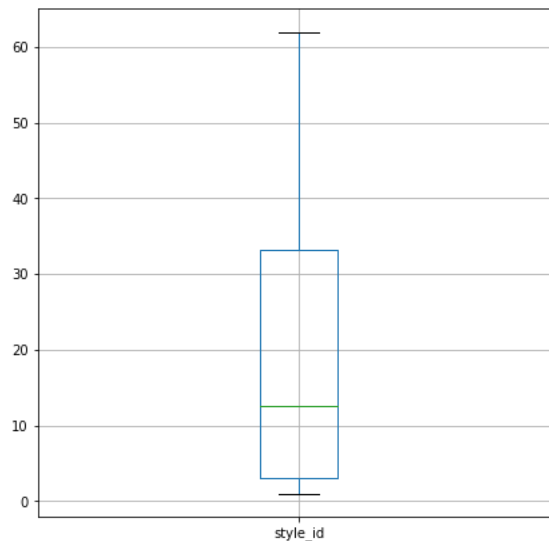


Fig.4 number of images per label

- Here's boxplot for how many data we have per each label, with max of 62 and min of 1 per label, median of 12.5
- I do think this data is not weighted properly, it could be affecting our result, I would either remove those labels with low data samples or gathering more data

## Algorithms and Techniques

The classifier is Convolutional Neural Network, which is the most advanced algorithm for image processing task such as classification.

These following parameters can be tuned to optimize the classifier:

- Training parameters
  - Training Length (Number of Epochs)
  - Batch Size (# of images look at once during a single step)
  - Optimizer type
  - Learning rate
  - Momentum
  - Weight decay
- Neural network architecture
  - Number of layers
  - Layer types
  - Layer parameters

During the training, both training & validation sets are loaded into RAM then random batches are selected to be loaded into GPU memory, the training is done using the Mini-batch gradient descent algorithm.

## Benchmark

1. Random Choice: A dummy classifier that random choice one of the labels.
2. A 4-layer CNN made from scratch is also used to measure the performance of transfer learning.

Benchmark model architecture:

Layer (type)	Output Shape	Param #
conv2d_1 (Conv2D)	(None, 222, 222, 16)	448
max_pooling2d_1 (MaxPooling2D)	(None, 111, 111, 16)	0
conv2d_2 (Conv2D)	(None, 109, 109, 32)	4640
max_pooling2d_2 (MaxPooling2D)	(None, 54, 54, 32)	0
conv2d_3 (Conv2D)	(None, 52, 52, 32)	9248
max_pooling2d_3 (MaxPooling2D)	(None, 26, 26, 32)	0
conv2d_4 (Conv2D)	(None, 24, 24, 64)	18496
max_pooling2d_4 (MaxPooling2D)	(None, 12, 12, 64)	0
flatten_1 (Flatten)	(None, 9216)	0
dropout_1 (Dropout)	(None, 9216)	0
dense_1 (Dense)	(None, 179)	1649843
Total params: 1,682,675		
Trainable params: 1,682,675		
Non-trainable params: 0		

Fig.5 Benchmark CNN Architecture

- Conv2D: 4 Conv2D layer were used here with 16,32,32,64 filter sizes, relu activation function.
- MaxPooling2D: Add MaxPooling2D layer after each convolution layer.
- Flatten: Fully connected layer
- Dropout: One dropout layer with value of 0.3
- Dense: Fully connected layer with output corresponding number of labels.

### III. Methodology

#### Data Preprocessing

Here's the link of example data image:

[https://stockx-360.imgix.net/Nike-Air-Max-1-97-Sean-Wotherspoon-NA/Images/Nike-Air-Max-1-97-Sean-Wotherspoon-NA/Lv2/img01.jpg?auto=format,compress&q=90&updated\\_at=1538080256&w=1300](https://stockx-360.imgix.net/Nike-Air-Max-1-97-Sean-Wotherspoon-NA/Images/Nike-Air-Max-1-97-Sean-Wotherspoon-NA/Lv2/img01.jpg?auto=format,compress&q=90&updated_at=1538080256&w=1300)

From my observation, product link is constructed by the shoe name itself as the sub-link, and for the image link, in the folder it's in sequence for 0 to 360-degree image coded as img01.jpg to img35.jpg.

I created a simple python scraper using bs4 & urllib, downloading images into different folder named by the label, then separate entire folder with 80%/20% split for training & test data.

When using Keras CNN, it requires a 4D tensor as input, with shape of (nb\_samples, row, cols, channels), here I'm using the build in function **keras.preprocessing.image.ImageDataGenerator** to resize the image to 224 x 224 pixel and then convert the image to a 4D tensor, so after the processing, our input will have a dimension of (1,224,224,3)

#### Implementation

##### 1. Benchmark CNN

First, the initial benchmark CNN model is trained with pre-processed training data using validation split of 0.2, batch size of 256 and epoch of 50. As a result, this benchmark model achieved accuracy of 0.73

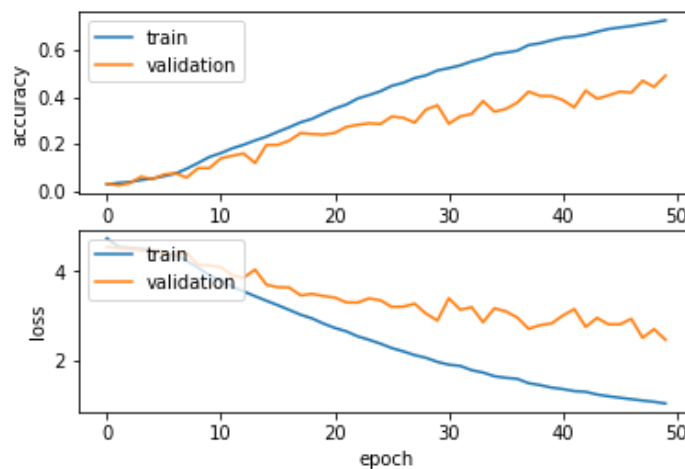


Fig.6 Benchmark CNN

## 2. VGG 16 Transfer Learning

For this model 2, transfer learning was used to create a CNN using the VGG-16 bottleneck features. First run through the images to the VGG-16 imagenet model to save the bottleneck feature, then train a fully connected top model on the bottleneck feature. This model achieved 0.83 of accuracy & 0.82 log-loss after 100 epochs.

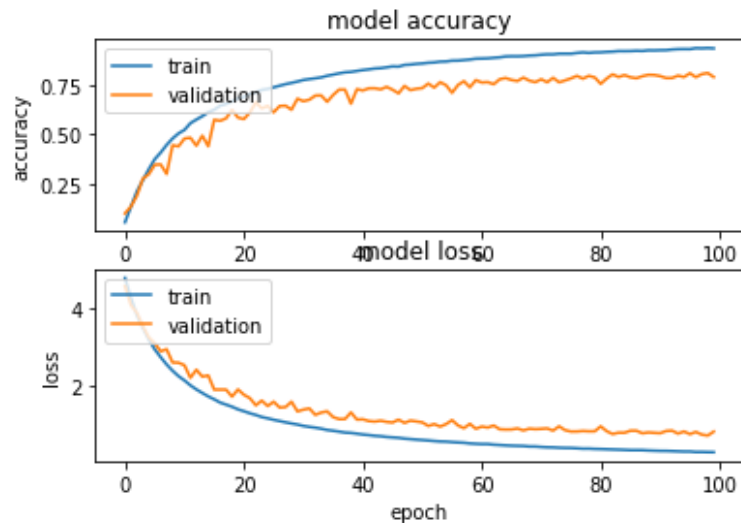


Fig.7 VGG-16 Transfer Learning

## 3. VGG 16 Fine Tuning

For this model, we freeze the first 15 convolutional layer in the VGG-16, only train last few to prevent overfitting, also using the top model weight we got from model 2. As a result, this model achieved 0.92 accuracy in 50 epochs.

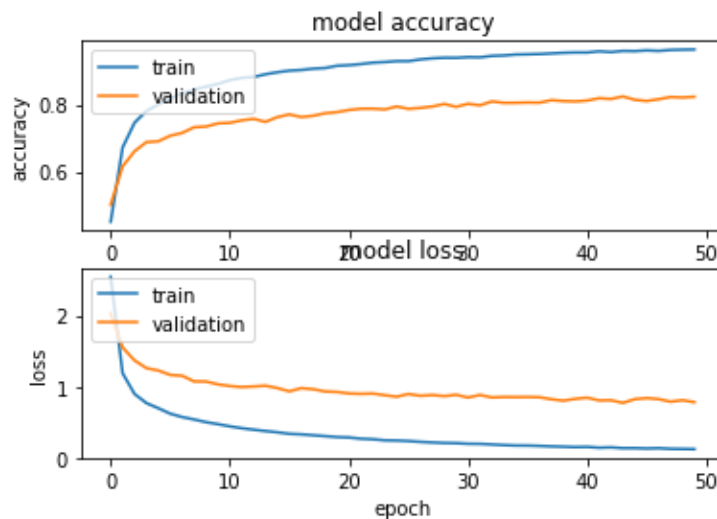


Fig.8 VGG-16 Fine Tuning

## Refinement

After changing the batch size and learning rate, I was able to get a better accuracy on the fine-tuning model with 0.94 accuracy, which meet my expectations considering tuning gain vs time.

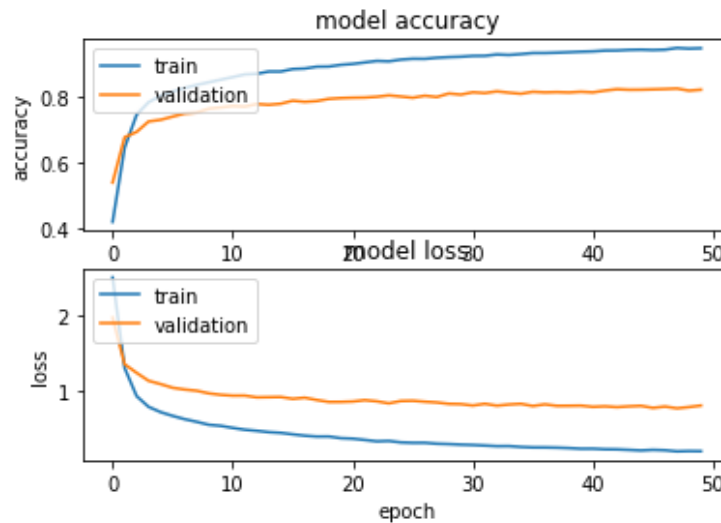


Fig.9 VGG-16 Fine Tuning Refinement

Since this model is trained on stock picture, which all have white background, pictures taken from user phone with varies background will have a lot of noise in the input so the model can't reach desired performance, to solve this issue, I added a simple cv grubcut function to remove the background of the input image for prediction.

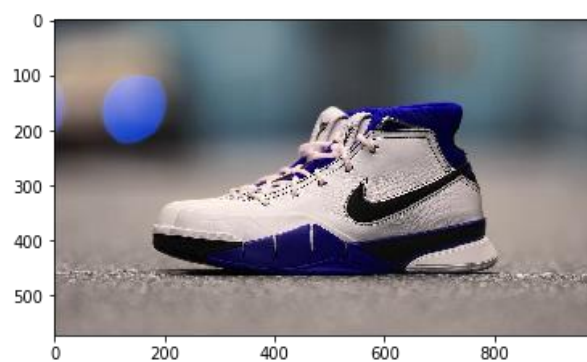


Fig.10-1 Original Image

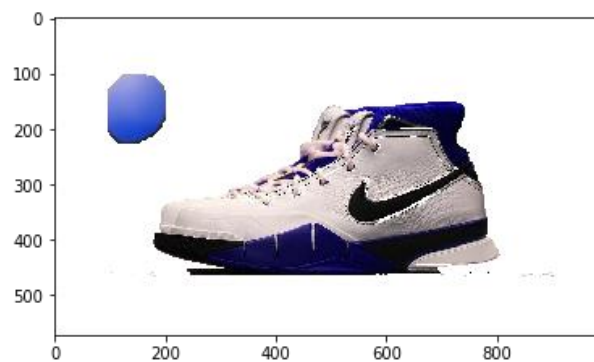


Fig.10-2 Original Image



## IV. Result

### Model Evaluation and Validation

The final fine-tuned model outperformed the benchmark models and other models above. As seen from fig.9, both train/validation accuracy increased through each epoch and after 50 epoch the model achieved accuracy of 0.938 and log-loss of 0.254.

The test data have 10404 images, after making predictions it has an accuracy of 0.946 on the test set.

### Justification

Model	Accuracy	Log-loss
Benchmark CNN	<b>0.73</b>	<b>1.65</b>
VGG-16 Bottleneck Feature	<b>0.83</b>	<b>0.82</b>
VGG-16 Fine-tuning	<b>0.94</b>	<b>0.25</b>

The final model had a significant improvement from the benchmark model accuracy from 0.73 to 0.94 and log-loss from 1.65 to 0.25

## V. Conclusion

### Free-Form Visualization

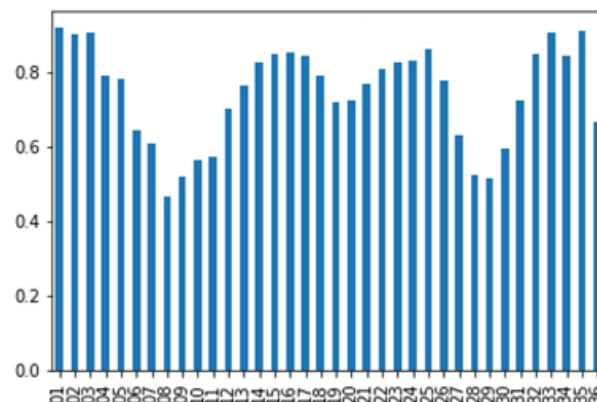


Fig.11 Accuracy vs Angle

From this bar chart we can see this model performance best on the side view (0-50 degree), have more error on front and back pictures (80-90 degree & 280-290 degree) which have fewer feature in the picture for model to predict also there are certain sneakers does have a same heel or very similar looks.

## Reflection

The process of this project can be summarized as follows:

1. Propose the problem, looking for data
2. No public dataset available, build scraper
3. Pre-process the dataset
4. Build benchmark model for the classifier & initial training
5. Transfer learning training
6. Fine-tuning
7. Flask Web App
8. Adding features to improve result (removing background)

Through out the process, I feel like the fine-tuning part is most difficult in my experience, because it is computationally very expensive to train each time on GPUs when I try to tune the hyperparameters, this quite took me a few days, especially when the optimizer doesn't have suitable momentum I can only see the accuracy/loss value ticking between epochs.

Still it is very satisfying to see the final model does perform exceptionally well than the benchmark model built from scratch.

## Improvement

As mentioned above, this dataset I used still is not the image what we would see people taking from their phone, so there's still background noise will affect the result. For improvement of the future work I would like to try to use other machine learning technique to help pre-process the image before the prediction, removing/replacing background with white so there's minimal noise here.

Another improvement can be done on the data augmentation, as this dataset still unbalanced, few labels have way less images available than others, by doing this or gathering new images this dataset will be well-balanced. Training the model on a better-balanced dataset may give us a boost in the performance.