

# Capstone 2: Yelp Dataset Image Classification

## Milestone Report

Nils Madsen

---

### Introduction

#### **Goal and Utility**

When photos are uploaded to Yelp, not all information will necessarily be filled in by the user. Certain fields of information for the photo, such as the label, will likely be left empty for many pictures. Yet, having missing data in the database can make it difficult to run analyses, so there is a motivation to fill in data holes wherever possible. Having labels for all images also helps visitors to the website find information they want about a business.

The sheer number of pictures uploaded to Yelp on a daily basis would make it uneconomical to hire people to classify photos manually. However, modern machine learning algorithms, especially convolutional neural networks (CNN), are powerful, accurate, high-throughput tools for image classification, and therefore are well suited to solving this problem.

The goal of this project will be to develop a CNN model that can accurately classify the photographs in the Yelp Dataset into broad categories.

#### **Data**

Yelp has published the Yelp Dataset for academic uses. The larger dataset includes a variety of fields for each rated business, including location, name, type, and star rating, as well as full-text reviews. For each user, there is information on join date, number of reviews written, and number of compliments received. Each photo entry has a corresponding label, caption, and business it is tied to.

This project will focus entirely on the 'label' field for each photograph. This is a broad classification, with only five possible entries ('food', 'drink', 'inside', 'outside', and 'menu'). There are a total of 280,992 photographs.

Round 12 of the Yelp Dataset was collected from the Yelp website:

<https://www.yelp.com/dataset/challenge>

## **Data Cleaning, Inspection and Management**

### **Data Format**

The Yelp Dataset comes in the form of a set of json files, accompanied by a set of images in .jpg format. To be more usable in Python, the json files were loaded into pandas dataframes before manipulation.

### **Missing Values**

There are no missing values in the label field, either in the form of NaN or as empty strings. All values are members of the set of five categories noted above: 'food', 'drink', 'inside', 'outside', and 'menu'.

### **Class Imbalance**

There is a dramatic degree of class imbalance in the data, with the majority of images being classified as food, and very few being classified as menus. This will likely lead to the model learning to classify high-representation classes well, but performing poorly on the lower-representation classes.

<b>Class</b>	<b>Image Count</b>	<b>Abundance (%)</b>
food	184456	65.64
inside	61620	21.93
outside	23214	8.26
drink	10350	3.68
menu	1352	0.48

Table 1: Class imbalance in the Yelp image set

### **Duplicate Images**

There are a small number of duplicated images in this dataset, which might hinder the evaluation of model performance by contaminating the validation and test sets with images that are also in the training set. However, almost all the duplicate images appear to be stock images, for example a promo image for McDonald's McCaffe. Since these images will likely show up again in future unseen Yelp images, it might be more relevant to this task to keep the duplicates in the dataset, as long as this dataset is, in fact, a random sample of Yelp images. Therefore I have elected not to remove any duplicate images.

### **Image Sorting**

The ImageDataGenerator used to feed images to Keras models expects image files to be separated into different directories based on their classification. Since all the Yelp Dataset images start in one directory when downloaded, Python's system utilities were used to automatically sort them into appropriate directories based on their label.

### **Mislabeling**

A cursory inspection of the images and labels indicates that there is a small degree of mislabeling in the dataset. For example, some images which are labeled as outside are clearly

of indoor areas, and a small amount of images that are labeled as 'inside' or 'drink' are actually of food. In the case of food being misclassified as drinks, it is usually soft food such as yogurt or ice cream, served in drink-shaped containers such as mugs or deep bowls. In the case of food being misclassified as an indoor area, there are instances where the picture of the food has been taken at such an angle that the indoor area can be seen in the background.

Without manual human inspection of all 280,000 photo/label pairs it is not clear how this issue could be addressed, and for the most part the labels appear to be accurate enough to be useful for this task. The mislabeling present will, however, lower the model's max possible performance to some degree.

## **Approach to Model Design and Evaluation**

### **Overview of Approach**

CNN architectures are very time consuming to design, and I had limited computing hardware. For this reason, the approach taken for this project was to fine-tune a premade CNN model for transfer learning rather than attempt to design and train an architecture from scratch. The VGG neural network architecture designed by the Visual Geometry Group at Oxford scored second in the 2014 ImageNet classification challenge, and is available as a template in the Keras API for Tensorflow, with or without pretrained weights. This made it a convenient model to use for this task.

CNNs are also notorious for taking enormous computational resources to train. Using the entire set of 280,992 images for fine-tuning of the model design would be impractical, as training would take a very long time and lead to large gaps between design iterations. Given the class imbalance in the image set, downsampling was used to both equalize the representation of classes, and create a smaller image set on which rapid design iteration could be achieved.

### **Splitting into training, validation, and test image sets**

The 280,992 images were randomly divided into training (230,992), validation (20,000), and test (30,000) image sets. The sampling was not stratified by class; however, inspection of class counts shows nearly identical relative abundance of each class in each set.

<b>Class</b>	<b>Train</b>	<b>Validation</b>	<b>Test</b>
food	65.62	66.21	65.42
inside	21.94	21.32	22.24
outside	8.29	8.20	8.12
drink	3.67	3.76	3.71
menu	0.47	0.53	0.50

Table 2: Class abundances (%) in each image set

## **Downsampling**

Downsampling is an established and effective method for addressing class imbalance. This technique involves randomly sampling the higher-representation classes without replacement to create a dataset where all classes have the same number of instances as the least-represented class. However, this solution comes at the cost of producing an image set that is very small in size. Since there are only 1096 menu images in the train set, and 105 menu images in the validation set, the sizes of these sets after downsampling are 5480 and 525, respectively. This is two orders of magnitude smaller than the 230,992 and 20,000 images in the full training and validation sets, and will likely confer a penalty on overall performance.

## **Tools**

Pandas was used as the main data management and manipulation tool, while Keras API for Tensorflow was used for building and training the CNN. Python's built-in system utilities were used for automated image file management. Sci-kit learn was utilized for various model evaluation metrics.

The model was trained on an Nvidia GTX 1060 graphics card with 6GB VRAM. The graphics card was fed by an Intel i7-3770k CPU with access to 16GB system memory.

## **Choice of performance metric**

Due to the large degree of class imbalance in this dataset, overall accuracy is not the optimal performance metric. The loss metric will also suffer from a tendency to favor performance on the high-representation classes over performance on the low-representation classes.

Therefore, three additional metrics have been included. The first of these is the macro-averaged recall (MacR, macro-recall), which is the arithmetic mean of the per-class recalls (true positives divided by class count). Macro-recall thus weighs all classes equally, and is completely insensitive to class imbalance. Conveniently, in the case where all classes have equal representation, such as in the downsampled image set, macro-recall is equivalent to accuracy.

Macro-averaged precision is a less valuable metric in this case due to its sensitivity to class imbalance. In a class-imbalanced dataset, improving the recall of the major classes will improve macro-precision more than improving the recall of the minor classes, so macro-precision has a tendency to favor performance on the major classes over performance on the minor classes.

That being said, macro-averaged F1 score (MacF1, macro-F1) and AUC score (MacAUC, macro-AUC) have been included to help integrate the precision performance into the model evaluation.

## **Model architecture**

The VGG-16 architecture was imported from the Keras database, and the convolutional layers were kept, since they are the aspect of the model that extracts features from the images. The fully connected layers at the end of the VGG model were discarded, as they are specific to the ImageNet task. The output of the convolutional layers was flattened, and two newly-initialized,

fully connected (FC) layers were added after the flatten layer. Finally, a new output layer with 5 nodes was added to the end of the model.

ReLu was used as the activation function for the FC layers, while the output layer was transformed with softmax. Categorical cross-entropy was used as the loss function, and Adam was used as the optimizer.

Keras layer	Shape	Param #
InputLayer	(224, 224, 3)	0
Conv2D	(224, 224, 64)	1,792
Conv2D	(224, 224, 64)	36,928
MaxPooling2D	(112, 112, 64)	0
Conv2D	(112, 112, 128)	73,856
Conv2D	(112, 112, 128)	147,584
MaxPooling2D	(56, 56, 128)	0
Conv2D	(56, 56, 256)	295,168
Conv2D	(56, 56, 256)	590,080
Conv2D	(56, 56, 256)	590,080
MaxPooling2D	(28, 28, 256)	0
Conv2D	(28, 28, 512)	1,180,160
Conv2D	(28, 28, 512)	2,359,808
Conv2D	(28, 28, 512)	2,359,808
MaxPooling2D	(14, 14, 512)	0
Conv2D	(14, 14, 512)	2,359,808
Conv2D	(14, 14, 512)	2,359,808
Conv2D	(14, 14, 512)	2,359,808
MaxPooling2D	(7, 7, 512)	0
Flatten	(25088)	0
Dense	(500)	12,544,500
Dropout	(500)	0
Dense	(500)	250,500
Dropout	(500)	0
Output	(5)	2,505

Table 3: Layers of the adapted VGG16 architecture

## **Model Hyperparameters**

### **Approach to tuning**

Tuning hyperparameters in neural networks is a very random process, as the ultimate performance of a model depends on the initialization state of the weights, which is different for each new model instance trained. For this reason, each condition in the tests below was run in triplicate. Any model that did not converge to >70% accuracy was discarded. The mean performance among the models that did converge within each condition is reported below. Any

condition which is reported as ‘did not converge’ (DNC) had no models converge out of the three replicates. Macro-recall is not reported in the performance tables in this section due to its being equivalent to accuracy in the downsampled image set.

While having more replicates within each condition is preferred to get a more accurate picture of the expected performance across different hyperparameter values, this requires more computing time.

### **ImageNet weights vs new weights**

The first aspect of the model to investigate is whether the convolutional weights pretrained on the ImageNet dataset are useful for this dataset. To evaluate this, three model conditions were compared. In the first condition, the ImageNet weights were discarded and all weights were trained from scratch. In the second condition, the ImageNet weights were kept, but the model was not able to alter the weights of the convolutional layers during training. In the final condition, the ImageNet weights were kept, and the model was able to further train these weights to fine-tune the feature extraction for this dataset.

Keeping and further training the ImageNet weights performed the best, while the baseline ImageNet weights outperformed new weights trained from scratch, according to validation loss (Table 4). Therefore, the ImageNet weights are useful for this task, and can be made to perform even better by fine-tuning them to this dataset.

### **Image augmentation**

The next aspect of the model that needs to be tested is whether image augmentation aids this task. Image augmentation is the act of randomly altering images to increase, in a way, the amount of data the model has available to learn from. Since image augmentation has a regularizing effect (i.e. it makes the model less likely to overfit) it is important to establish image augmentation before tuning parameters such as model capacity.

Image augmentation, however, needs to be executed in a way that is sensible for the task. Flipping images upside down, for example, is not a sensible transformation because people do not usually upload upside-down images to Yelp. Flipping horizontally was included because food is food regardless of whether the meat is to the left or right of the potatoes, and an indoor space is an indoor space regardless of whether the door is to the left or right of the bar. One aspect of this task that could be affected by a horizontal flip is the text on menus, however there are other features of the images that the model can use to identify this class.

Random horizontal and vertical shifting were also performed (up to 20% of total width/height), as users may choose to upload off-center images for stylistic reasons. Finally, a random amount of zoom (up to 20%) was applied because users may take images at differing distances from the object of interest. Any gaps in the image caused by shifts were filled by reflection of the image across the border.

The model trained using image augmentation had similar performance on the validation set as the model without image augmentation. However, the performance on the training set

decreased to come more into line with the performance on the validation set. This suggests that the model trained without image augmentation was overfitting to some degree. The lack of improvement in validation performance could indicate that the performance is limited by some other factor, such as a model capacity that is too small.

Convolution Weights	Mean Acc	Mean Loss	Mean F1	Mean AUC
ImageNet Baseline	0.877	0.904	0.878	0.979
ImageNet Trainable	0.872	0.375	0.872	0.981
New Weights	0.707	0.763	0.708	0.924
ImageNet with Augmentation	0.854	0.381	0.853	0.980

Table 4: Performance on downsampled validation set across different weight conditions

### Model capacity

The number of nodes in the two FC layers is an important hyperparameter to tune, as it affects the tendency of the model to overfit vs generalize. Having FC layers with too many nodes will lead to a model which will overfit to the training data, and have worse performance on the validation and test data. Having FC layers with too few nodes will lead to poor performance on the training and validation images, as there will not be enough model capacity to capture all the important patterns in the dataset.

To find the optimal amount of nodes, the model was trained with a FC width ranging from 400 to 1600, in 200 node increments. The model improved in accuracy and loss up to a width of 800-1000 nodes. At around 800 nodes and greater, the model's performance was essentially the same. This indicates that 800-1000 nodes is the minimum width required for optimal performance on the downsampled image set. The model might be able to make use of more capacity on the full image set, but to optimize this parameter on the full set would require too much computation time to be feasible. 1000 nodes was chosen as the model width because it is the safe choice for having at least the minimum capacity needed.

Width	Mean Acc	Mean Loss	Mean F1	Mean AUC
400	0.857	0.403	0.857	0.977
600	0.878	0.327	0.879	0.984
800	0.884	0.306	0.884	0.985
1000	0.886	0.305	0.885	0.985
1200	0.888	0.311	0.888	0.983
1400	0.882	0.330	0.882	0.983
1600	0.887	0.284	0.887	0.987

Table 5: Performance on downsampled validation set across different widths

### Learning rate

The learning rate of the Adam optimizer is also an important hyperparameter to tune. A learning rate that is too low or too high could lead to the model finding a less optimal solution, or failing to learn at all.



The model was tested on learning rates of 0.0005, 0.0002, 0.0001, 0.00005, and 0.00002. A learning rate of 0.0001 was optimal, with 0.00005 and 0.0002 performing slightly worse and 0.00002 producing a much worse solution, according to the validation loss. A learning rate of 0.0005 did not converge to a solution.

Learning Rate	Mean Acc	Mean Loss	Mean F1	Mean AUC
0.0005	DNC	DNC	DNC	DNC
0.0002	0.815	0.608	0.814	0.959
0.0001	0.876	0.325	0.875	0.984
0.00005	0.879	0.330	0.880	0.984
0.00002	0.863	0.370	0.861	0.982

Table 6: Performance on downsampled validation set across different learning rates

### Dropout

Dropout is a powerful regularization technique used to boost validation set performance. In this technique, a proportion of the nodes in a given layer are randomly selected to be dropped for each training batch. The remaining nodes are then trained on that batch. This forces the model to store a more robust, redundant understanding of the dataset across its nodes. Dropout thus allows for a larger model capacity without overfitting the training data.

Dropout was tested for each FC layer individually. FC layer width was increased in proportion to the degree of dropout. For example, at a dropout of 0.2, width was increased to 1250 so that 1000 nodes were active after 20% were randomly dropped. Dropout on the first layer was tested first, as earlier dropout usually has a greater impact on the model.

Dropout proportions ranging from 0.1 to 0.4 in the first FC layer did not improve the performance of the model over a baseline of 0 dropout. Testing dropout in the second FC layer revealed that a dropout of 0.2 performed best.



Dropout 1	Mean Acc	Mean Loss	Mean F1	Mean AUC
0	0.887	0.302	0.886	0.986
0.1	0.885	0.320	0.885	0.985
0.2	0.868	0.349	0.867	0.982
0.3	0.848	0.382	0.847	0.979
0.4	0.863	0.372	0.862	0.979

  

Dropout 2	Mean Acc	Mean Loss	Mean F1	Mean AUC
0	0.872	0.353	0.872	0.982
0.1	0.874	0.330	0.874	0.982
0.2	0.896	0.294	0.896	0.986
0.3	0.865	0.343	0.865	0.982
0.4	0.872	0.353	0.872	0.982

Table 7: Performance on downsampled validation set across different dropouts

### Final model design

In summary, the final model design that was adopted for this task includes:

1. VGG-16 convolutional layers with pretrained weights from the ImageNet competition, with the capability of further tuning these weights for this dataset
2. Two fully connected layers composed of 1200 nodes each, with ReLU activations
3. An output layer composed of 5 nodes, with softmax activation and categorical cross-entropy as the loss function
4. Image augmentation in the form of horizontal flip, horizontal and vertical shift, and zoom, with fill by reflection
5. A learning rate of 0.0001 for the Adam optimizer
6. No dropout regularization in the first FC layer, and a dropout of 0.2 in the second FC layer

### Next Steps

The next steps for this project will include training the tuned model on the full image set, developing strategies for handling the class imbalanced if necessary, and concluding with a discussion of the results and future directions.