

# Quick, Draw! Image Classification

Group 3

December 7, 2018

## 1 Introduction

Many recent industry innovations have included neural networks. An important research area within neural networks has been image classification. One of the best techniques to classify images is Convolutional Neural Networks (CNN). We have sought to better understand industry best practices by exploring different image classification techniques with neural networks. Our exploration starts by manually performing image classification and noting the level of human accuracy. We then use a fully connected neural network to perform image classification before implementing a CNN. We implemented options for tuning and debugging the neural network while implementing each method during our process. We were able to develop an intuition for image classification best practices using neural networks after our research.

Many supervised machine learning methods often differ from those employed when constructing neural networks. These supervised machine learning methods, such as logistic regression or decision tree classification, require feature engineering to successfully classify an image. Models requiring feature engineering rely on considerable effort by subject matter experts to achieve reasonable accuracy. Neural networks allow for an approach which does not require feature engineering to achieve similar or better accuracy. Industry engineering efforts regarding feature engineering can be shifted to neural networks in some cases such as image classification. Moreover, the availability of cheap computing resources sped up a shift towards the use of neural networks.

When considering modeling approaches, we examine both fully connected and convolutional neural networks. A fully connected neural network is known to be inefficient at classifying large images. These inefficiencies in a fully connected neural network partly arise because every neuron in the  $\ell - 1$  layer must be connected to the  $\ell$ th layer. This requires a series of matrix multiplications and additions across each layer  $\ell$ . Alternatively, a CNN allows for more efficient computations by: (1) efficient and automatic feature generation using local connectivity and parameter sharing of convolution operations, (2) dimensionality reduction using pooling layers. In this project, we attempt to understand these topics in depth from an empirical standpoint in relation to image classification using Google’s “Quick, Draw!” dataset.

## 2 Methods

Our methodology consists of reviewing statistical methodology, presenting theoretically informed hypotheses, and providing an overview of our data.

### 2.1 Overview of Statistical Modeling Techniques

We use human classification, fully connected neural networks, and a Convolutional Neural Network (CNN). We use more detail to describe techniques not covered during the course.

#### 2.1.1 Understanding the Loss Function

Most modern neural networks are trained using maximum likelihood. This means that the cost function is the negative log-likelihood, also described as the cross-entropy between the training data and the model distribution [1]. In our use case, we are interested in correctly classifying 10 categories of drawings. Our output unit use a Softmax function in our output layer because they are most often used for classifiers representing the probability distribution over  $n$  different classes. Given features  $h$ , weights  $W$ , and bias  $b$ , a linear layer first predicts unnormalized log probabilities:

$$z = W^T h + b \quad (1)$$

where  $z_i = \log \tilde{P}(y = i|x)$ . The softmax function can then exponentiate and normalize  $z$  to obtain the desired  $\hat{y}$ . The softmax function is given by

$$\text{softmax}(z)_i = \frac{\exp(z_i)}{\sum_j \exp(z_j)}. \quad (2)$$

As a cost function, we wish to maximize  $\log P(y = i; z) = \log \text{softmax}(z)_i$ . We can do this formally as

$$\log \text{softmax}(z)_i = z_i - \log \sum_j \exp(z_j). \quad (3)$$

which gives us the intuition for a single example. Overall, we want to learn parameters that drive the softmax to predict the fraction of counts of each outcome observed in the training set [1]:

$$\text{softmax}(z(x; \theta))_i \approx \frac{\sum_{j=1}^m 1_{y^{(j)}=i, x^{(j)}=x}}{\sum_{j=1}^m 1_{x^{(j)}=x}}. \quad (4)$$

This is how the `categorical_crossentropy(y_true, y_pred)` function is constructed in Keras [2] as well as the `softmax_cross_entropy_with_logits_v2(labels, logits)` in Tensorflow [3].

### 2.1.2 Understanding the Optimization Technique

This course discussed optimization algorithms such as Stochastic Gradient Descent ( $m = 1$ ) and Batch Gradient Descent ( $m = n$  obs.). Current best practices for Neural Networks recommend Mini-Batch Gradient Descent ( $m = c$  batch size). In practice, we use the Adam optimization algorithm due to faster computation time. The Adam optimization algorithm is able to achieve a speedup due to its combination of two other techniques.

The Adam optimization algorithm takes an exponentially weighted average of the gradients, a technique called momentum, and then uses that gradient to update your weights. Secondly, we also want to slow down gradient descent in the horizontal direction with a technique called RMSprop to reduce oscillations from batch gradient descent. The Adam optimization method is combining momentum and RMSprop to perform gradient descent more quickly [4].

### 2.1.3 Convolutional Neural Network (CNN)

Convolutional networks, also known as convolutional neural networks(CNNs), are a specialized kind of neural network for processing data that has a known grid-like topology. Convolutional networks have been tremendously successful in practical applications. The name “convolutional neural network” indicates that the network employs a mathematical operation called convolution. Convolution is a specialized kind of linear operation. Convolutional networks are simply neural networks that use convolution in place of general matrix multiplication in at least one of their layers. See Figure 5 in the Appendix for a visual representation of a CNN.

CNN architectures are usually composed of the building blocks as described below:

1. **Convolution:** Convolution operation helps to extract features from the input image. The first argument to the convolution is often referred to as the input, the second argument as the feature detector (kernel) and the output as the feature map. The convolution operation takes the dot product of feature map and the input to detect features. Since each convolution operation only focuses on a part of an image enabling parameter sharing which helps to reduce the number of free parameters compared to a fully connected neural network, allowing the network to be deeper with fewer parameters.
  - **Strides:** Stride is the number of pixels jumps between each convolution operation performed on the input matrix. When the stride is 1 then we move the filters to 1 pixel at a time. Increasing the stride parameter reduces the output volume of the convolution since the jumps are larger.
  - **Padding:**Padding helps to detect features present at the border by adding a dummy values around the edge. Since convolution operations collect most information from the center of the image, adding padding allows the border pixels to be weighed more during the convolution operation.

See Figure 6 in the Appendix for an example of Stride and Padding.

2. **Pooling Layer:** Pooling layers section helps to reduce the number of parameters when the images are too large. Spatial pooling (subsampling or downsampling) reduces the dimensionality of each map, retaining important information. These can be of different types: (1) Max Pooling - takes the largest element (most important feature) from the rectified feature map (2) Average Pooling - averages all the pixel values (3) Sum Pooling. sum of all pixels values. See figure 7 for a visualization of a Pooling Layer.
3. **Fully Connected Layer:** The layer we call as FC layer, we flattened our matrix into vector and feed it into a fully connected layer like neural network.

## 2.2 Hypotheses

We have generated the following hypotheses based on theoretical expectations expressed in the literature [1, 4]. Performance is effected by both the model and data [4] so our hypotheses pertain to both effects.

### 2.2.1 Addressing model component of performance

*H1<sub>a</sub>* Fully Connected Neural Network will struggle to get reasonable accuracy.

*H2<sub>a</sub>* CNN should perform better than fully connected neural network.

### 2.2.2 Addressing data component of performance

*H3<sub>a</sub>* Small amounts of data will lead to overfitting.

*H4<sub>a</sub>* Imbalanced class prediction can be improved using data augmentation.

*H5<sub>a</sub>* Addition of a category to the trained model will reduce the performance of the model.

*H6<sub>a</sub>* Adding more training data will improve the performance of the model.

## 2.3 Data

Google originally collected the data from users by explicitly asking users to draw objects like forks, hotdogs, etc [5]. The original dataset is 70 GB in size but Google also provides a simplified version of the dataset- consisting only of the final images [6]. We used these '.npy' files to train all of our models. These numpy bitmap files consist of more than 100,000 rows each and 784 features, since they are basically  $28 \times 28$  images. Samples of these images are shown in Figure 1. The subset of this dataset sampled varies from experiment to experiment. More details can be found in the experiments mentioned below.



Figure 1: Sample images from Quick, Draw! dataset

Our analysis required us to use a subset of this data. Our dataset sample selected for this experiment had 10 categories: fish, fork, hotdog, flamingo, airplane, alarm clock, baseball bat, bicycle, dolphin, elephant. We had 100,000 examples in total, with 10,000 example for each of the aforementioned categories.

## 2.4 Modeling Strategy

We provide our modeling strategy as it relates to the exploration of our stated hypotheses.

### 2.4.1 Baseline image classification using humans

Two group members completed a Human Intelligence Task by manually classifying images into one of our ten defined categories. This approach helped us get a benchmark which we could use to define “reasonable accuracy”.

### 2.4.2 $H1_a$ : Fully connected neural network will struggle to get reasonable accuracy

In order to perform this experiment, we built a fully connected neural network using Tensorflow. Our train:test split was chosen to be 80:20, an industry standard. Our initial objective was to build a model that overfit our given dataset, with the idea that once our model is complex enough to overfit the training data, we could regularize it in order to achieve a more optimal fit. The fully connected neural network had 20,067 parameters.

After experimenting with various architectures, we narrowed down to a 3 layer model as specified in Figure 3 in the Appendix. The model was run on a machine with 13 GB of RAM, Intel Xeon CPU 2.3 GHz (1 core, 2 threads), 1xTesla K80 GPU having 2496 CUDA cores with 12 GB GDDR5 VRAM.

### 2.4.3 $H2_a$ CNN should perform better than fully connected neural network

As for the CNN, we started off with the same approach with attempting to overfit the model using convolution operations and a fully connected layer with an objective to categorise 10

categories, softmax was chosen as the activation function for the final layer, while every other neuron had ReLU as the activation function owing to ReLU’s faster convergence. See Figure 4 for a representation of this first CNN model.

The number of parameters is much greater than the number of observations used to train the model resulting in overfitting. In order to reduce the number of parameters, we incorporated Pooling (Max Pooling) post each convolution operation extracting only the most significant features resulting in reduction of the volume of the convolutional output, giving the model a smaller memory footprint as well as a smaller number of trainable parameters. Figure 2 has a representation of this second CNN model. We employ strategies to reduce overfitting related to dropouts and stride in addition to adding more training data.

## 3 Results

We report our results by examining the validity of each hypothesis.

### 3.1 Addressing model component of performance

In order to perform this experiment we built a fully connected neural network using tensorflow. Our train:test split was chosen to be 80:20, an industry standard.

Our initial objective was to build a model that overfit our given dataset, with the idea that once our model is complex enough to overfit the training data, we could regularize it in order to achieve a more optimal fit.

After experimenting with various architectures, we narrowed down to a 3 layer model as specified in Figure 3.

The training set accuracy observed was 93%, while the test accuracy was at 73%, which clearly indicated an overfit in the model, as designed.

The model took 1500 epochs to converge taking around 45 minutes on a machine with 13 GB of RAM, Intel Xeon CPU 2.3 GHz (1 core, 2 threads), 1xTesla K80 GPU having 2496 CUDA cores with 12 GB GDDR5 VRAM. In order to manage variance, we attempted regularizing using dropouts. Unfortunately, that effort was futile. We then doubled the amount of data sampled to train the model, eventually helping us achieve training set accuracy of 86% and test set accuracy of 82%. Our learning rate was specified as 0.001.

While our initial hypothesis assumed that a regular fully connected network would struggle to achieve a reasonable accuracy, our experiment proved our hypothesis as invalid by achieving a commendable accuracy. Having said that, the model did take a longer to converge than expected.

As for the CNN, we started off with the same approach with attempting to overfit the model using convolution operations and a fully connected layer. Since our objective is to categorise

10 categories, softmax was chosen as the activation function for the final layer, while every other neuron had ReLU as the activation function owing to ReLU's faster convergence.

Number of Trainable Parameters: 320,890

The number of clearly being much greater than the number of observations used to train the model, we need to reduce the number of parameters vis-a-vis number of observations. From the available strategies to do so, we incorporated MaxPooling [TODO: why?] post each convolution operation, in order to extract only the most significant features. As a consequence, this reduced the volume of the output, giving the model a smaller memory footprint as well as a smaller number of trainable parameters.

The observable training accuracy is 0.90, whereas the validation accuracy is 0.89. The model converged after 10 epochs, which took 3 minutes.

Some of the other unsuccessful strategies used to decrease the number of parameters and to regularize the network were strides and dropouts which sacrificed way too much in the way of accuracy and thus had to be dropped. Since most of our drawings had most of their information in the center of the image, padding, in theory wouldn't be of much help. However, actually implementing padding improved our accuracy by 1% at a cost of 6000 parameters.

Conclusion

A slight (5%) improvement in the accuracy of a CNN when compared to a fully connected network, keeping the number of parameters constant, was observed. However, the CNN happened to be around 13 times faster in practice. This delta in accuracy and speed is likely to be amplified as we scale the model up with respect to model complexity or data.

## 3.2 Addressing data component of performance

ok

## 4 Discussion

Talk about our final slide and tie the hypotheses together

## References

- [1] Ian Goodfellow, Yoshua Bengio, and Aaron Courville. *Deep Learning*. MIT Press, 2016. <http://www.deeplearningbook.org>.
- [2] François Chollet et al. Keras. <https://keras.io>, 2015.

- [3] Martín Abadi, Ashish Agarwal, Paul Barham, Eugene Brevdo, Zhifeng Chen, Craig Citro, Greg S. Corrado, Andy Davis, Jeffrey Dean, Matthieu Devin, Sanjay Ghemawat, Ian Goodfellow, Andrew Harp, Geoffrey Irving, Michael Isard, Yangqing Jia, Rafal Jozefowicz, Lukasz Kaiser, Manjunath Kudlur, Josh Levenberg, Dan Mané, Rajat Monga, Sherry Moore, Derek Murray, Chris Olah, Mike Schuster, Jonathon Shlens, Benoit Steiner, Ilya Sutskever, Kunal Talwar, Paul Tucker, Vincent Vanhoucke, Vijay Vasudevan, Fernanda Viégas, Oriol Vinyals, Pete Warden, Martin Wattenberg, Martin Wicke, Yuan Yu, and Xiaoqiang Zheng. TensorFlow: Large-scale machine learning on heterogeneous systems, 2015. Software available from tensorflow.org.
- [4] Andrew Ng. Deep learning — coursera. <https://www.coursera.org/specializations/deep-learning>, 2018. (Accessed on 12/07/2018).
- [5] Google. Quick, draw! doodle recognition challenge — kaggle. <https://www.kaggle.com/c/quickdraw-doodle-recognition>, 2018. (Accessed on 12/07/2018).
- [6] Google. googlecreativelab/quickdraw-dataset: Documentation on how to access and use the quick, draw! dataset. <https://github.com/googlecreativelab/quickdraw-dataset>, 2016. (Accessed on 12/07/2018).

## 5 Appendix

$$Z_1 = \underset{25 \times 784}{W_1} X + \underset{25 \times 1}{b_1} \Rightarrow A_1 = \text{ReLU}(Z_1)$$

$$Z_2 = \underset{12 \times 25}{W_2} A_1 + \underset{12 \times 1}{b_2} \Rightarrow A_2 = \text{ReLU}(Z_2)$$

$$Z_3 = \underset{10 \times 12}{W_3} A_2 + \underset{10 \times 1}{b_3} \Rightarrow \text{output layer} = \text{SoftMax}(Z_3)$$

$$(25 \times 784 + 25 \times 1)_{Z_1} + (12 \times 25 + 12 \times 1)_{Z_2} + (10 \times 12 + 10 \times 1)_{Z_3} = 20,067 \text{ parameters}$$

Figure 2: Fully Connected Neural Network Architecture

28 x 28 x 1 [Input Layer]	3 x 3 x 4 [ReLU] [Convolution]	2 x 2 x 8 [ReLU] [Convolution]	64 [Relu] [Dense]	10 [Softmax] [Output Layer]
------------------------------	--------------------------------------	--------------------------------------	-------------------------	-----------------------------------

Figure 3: (CNN Model 1) Loss Function: Categorical Cross Entropy, Optimizer: Adam, Trainable Parameters: 320,890



28 x 28 x 1 [Input Layer]	3 x 3 x 4 [ReLU] [Convolution] [MaxPool - 2x2]	2 x 2 x 8 [ReLU] [Convolution] [MaxPool - 2x2]	64 [Relu] [Dense]	10 [Softmax] [Output Layer]
------------------------------	---	---	-------------------------	-----------------------------------

Figure 4: (CNN Model 2) Loss Function: Categorical Cross Entropy, Optimizer: Adam, Trainable Parameters: 19,322

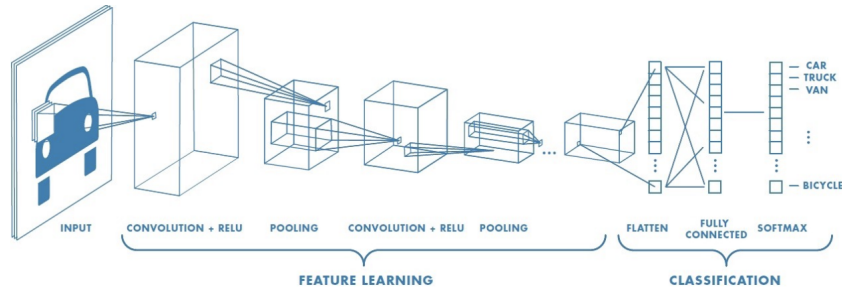


Figure 5: CNN Architecture Example

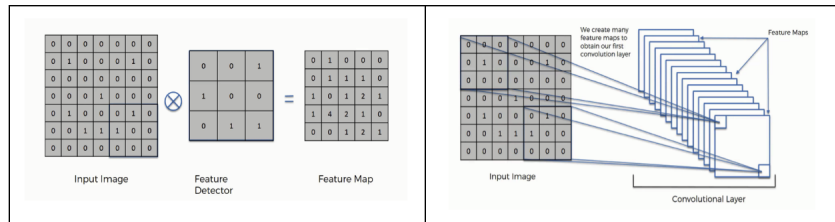


Figure 6: Convolution operation with stride=1, padding=None

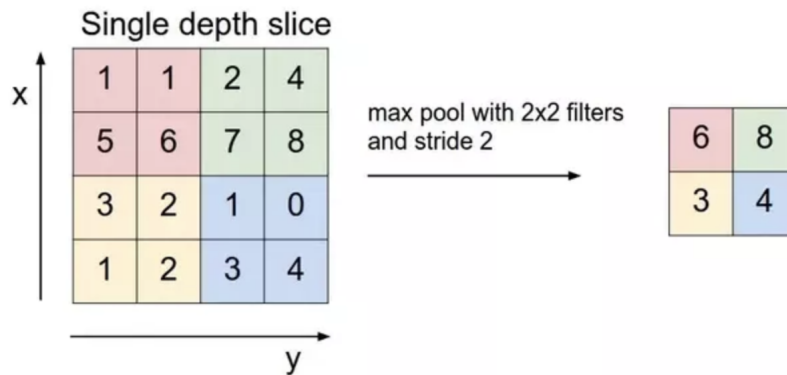


Figure 7: Example of Max Pooling