

# Capstone Project

Ryosuke Honda

August 18, 2016

## 1 Definition

### 1.1 Project Overview

The improvement on the camera and the widespread of larger capacity hardwares lead us to store more digital pictures. However computers can't understand the meaning of the pictures, so it is essential for us humans to classify or search images. Image recognition(Figure.1) will help us to classify or search pictures without our intervention. These days, because of the development of the computational capacity, we can process large number of pictures with high level accuracy(nearly the human's recognition). In this project, I'll discuss the image recognition algorithm which will be useful in the classification of large number of images.



Figure 1: Image of the Algorithm

### 1.2 Problem Statement

The dataset I use in this task is CIFAR-10 dataset. The description about it is below

The CIFAR-10 dataset consists of 60,000  $32 \times 32 \times 3$  (width=height=32 and RGB=3) color images in 10 classes, with 6,000 images. 50,000 images are the training images and 10,000 images are the test images.

The classes are completely mutually exclusive. There is no overlap within each image class.

Many research have been done with the CIFAR-10 dataset. Some of the researches use very deep neural networks for training[1] and others have originality in pooling layer which is one of CNN components.[2]

In this task, I'll discuss the image recognition by Convolutional Neural Network(CNN) and its hyper-parameter tuning. When constructing CNN model,

hyper-parameter tuning is one of the essential and time-consuming tasks. Therefore, I propose a method, called "Bayesian Optimization", to tune hyper-parameters by not using "grid search" which is not appropriate for this task because of computationally expensive.

### 1.3 Metrics

The objective in this problem is multi class classification. Therefore, the metrics will be 'cross entropy'. The objective of the learning in this task is to minimize the cross entropy. Cross entropy is defined as the following formula[3].

$$C = - \sum_{j=1}^n d_j \log p_j \quad (1)$$

$d_1, \dots, d_n$  is the optimal output (Correct output).  $p_1, \dots, p_n$  is the probability of for the output class. This probability is calculated by softmax function which is defined as below.

$$p_j = \frac{e^{u_j}}{\sum_{k=1}^n e^{u_k}} \quad (2)$$

Classification Error is calculated by the difference between the optimal output  $d_1, \dots, d_n$  and predicted output  $p_1, \dots, p_n$ . The target output  $d_1, \dots, d_n$  takes the representation of 1-of-n. That means that only the correct label  $j$  becomes  $d_j = 1$  and the others  $k (\neq j)$  become  $d_k = 0$ .

## 2 Analysis

### 2.1 Data Exploration

CIFAR-10 dataset consists of 50,000 images of training data and 10,000 images of test data. For each data, it contains 10 different kind of classes (discussed in the next chapter) and the distribution of each class is the same in both training and test data. That means 5,000 images for each class in the training data and 1,000 images for each class in the test data. (Fig.3 and Fig.4).

### 2.2 Exploratory Visualization

There are 60,000 of images in total. Figure.2 (50,000 images for training data and 10,000 images for test data) shows the samples of the images. I plotted 10 images for each class. Figure.3 and Figure.4 shows the distribution of the data. For the training dataset, each class has 5,000 images and for the test dataset, each class has 1,000 images. The label 0 to 9 corresponds to 'Airplane', 'Automobile', 'Bird', 'Cat', 'Deer', 'Dog', 'Frog', 'Horse', 'Ship', and 'Truck'.

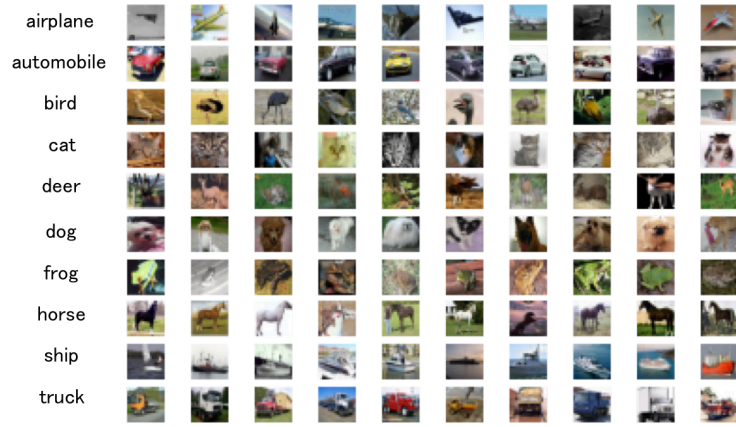


Figure 2: Sample of the Images

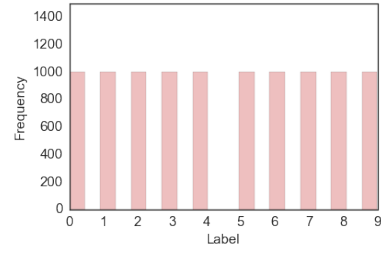
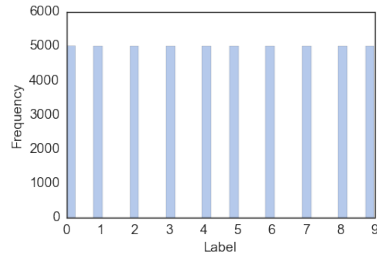


Figure 3: Distribution of Training Data    Figure 4: Distribution of Test Data

### 2.3 Algorithms and Techniques

In this task, I'll use Convolutional Neural Network(CNN). CNN has been successful in practical applications for image recognition. CNN consists of convolution layer, pooling layer and fully-connected layer and sometimes contains local contrast normalization(LCN). In this chapter, I'll discuss the convolution layer and pooling. As the name 'convolution' suggests, the network employs a mathematical operation called convolution. [3][4]

Fig.5 is the example of the architecture of the CNN. Fig.6 illustrates the image of the CNN.[6]

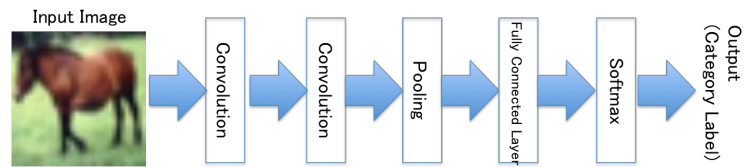


Figure 5: An example of CNN Architecture

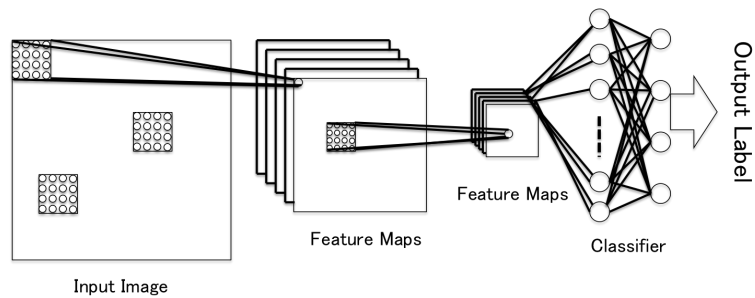


Figure 6: Overview of CNN

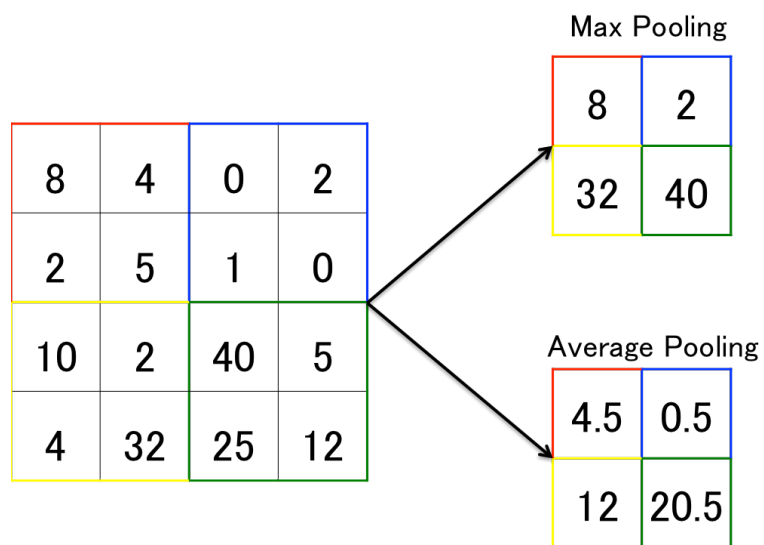


Figure 7: Example of Max Pooling and Average Pooling

Convolutional neural networks are biologically inspired variants of multilayer perceptron. They are emulated by the behavior of a visual cortex. CNN has three distinguishing features which traditional multilayer perceptron models don't have.[12]

- 3D volumes of neurons

The layers of a CNN have neurons arranged in 3 dimensions: width,height and weight. The neurons inside a layer are only connected to a small region of the layer before it which is so called 'receptive field'.

- Local Connectivity

This is the most distinguishing feature of CNN. Typical multilayer perceptron are densely connected. Because of the connectivity, we can't utilize the architecture for image recognition because of the curse of dimensionality. As for CNN, the connectivity are inspired by the neuro-science(visual cortex). Therefore, the dimensionality is much smaller than the typical multilayer perceptron. The CNN architecture ensures that the learnt "filters" produce the strongest response to a spatially local input pattern. Stacking many such layers leads to non-linear "filters" that become increasingly "global". Because of this, the network can create good representations of small parts of input, then assemble representations of large areas from them.

- Shared weights

Each filter is replicated across the entire visual field in CNN. These replicated units share the same parameterization and form a feature map. Since the parameters are shared, the number of free parameters decrease. Therefore, when running CNNs model, the memory use will decrease, which means we can train more powerful(more layers) networks.

Pooling layer is put after the convolution layer. The function of the pooling layer is to progressively reduce the spatial size of the representation to reduce the amount of parameters and computation in the network, and hence to also control overfitting. By introducing pooling layer, not only the architecture will be more robust but also the dimensionality will be reduced. Max pooling and Average pooling are the typical pooling which are generally utilized. Figure.7 is the example of the pooling.The original map is the size of  $4 \times 4$ . The stride for the pooling is 2 and the pooling size is  $2 \times 2$ . "Max pooling" is to extract the maximum pixel from each region and "Average pooling" is to calculate the average value for each region.In this task, I utilized max pooling at the pooling layers.

## 2.4 Benchmark

For the CNNs architecture, it is quite important to decide the number of layers. Therefore, I first choose several convolutional layers and decide one of them as a benchmark.

When deciding the architecture, I set the mutual parameters as follows. The number of batch size is 32.The number of filters of each convolutional layer is 32. and finally the number of epochs is 20.

I tried 4 architecture of CNNs.

1. 1 Convolutional layer & 2 Fully connected layers

The simplest version in these model.The input and output dimension on each layer are below. The accuracy rate of training and validation data are also below.

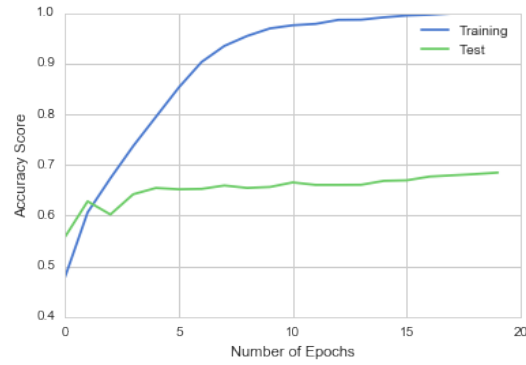


Figure 8: Accuracy rate of training and validation data

The accuracy rate on the training data is quite high, however the accuracy rate on the validation data is low. This means that this architecture falls into over-fitting.

## 2. 2 Convolutional layers & 2 Fully connected layers

The input and output dimension on each layer are below. The accuracy rate of training and validation data are also below.

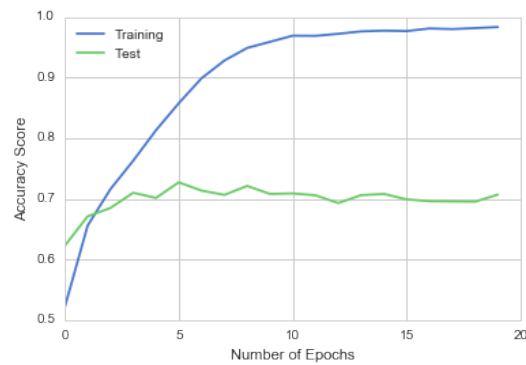


Figure 9: Accuracy rate of training and validation data

This architecture also caused over-fitting as mentioned above.

## 3. 3 Convolutional layers & 2 Fully connected layers

The accuracy rate of training and validation data are also below.

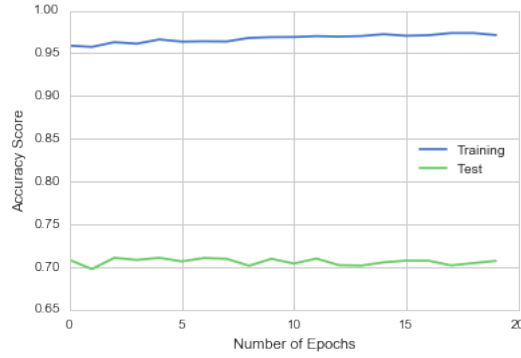


Figure 10: Accuracy rate of training and validation data

This architecture also falls into over-fitting.

#### 4. 4 Convolutional layers & 2 Fully connected layers

The accuracy rate of training and validation data are also below.

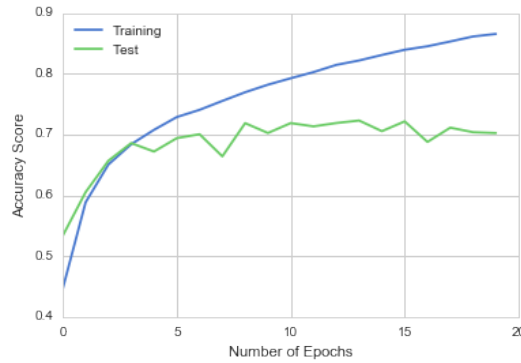


Figure 11: Accuracy rate of training and validation data

This architecture is far better than the others. The accuracy rate of both training and validation data is high as the number of epochs increases. The accuracy rate of this architecture is 70.2%. Actually the highest accuracy rate in CIFAR-10 is more than 95.0% and the benchmark is quite lower than the highest one. My target in this task is to find the optimal parameters for the architecture. Therefore, I'll utilize this 4 layers convolution & 2 fully connected layers architecture as a benchmark.

From these results, I choose the fourth architecture(4 convolution layers & 2 fully connected layers) as the benchmark for this task.

## 3 Methodology

### 3.1 Data Preprocessing

Apart from the models of image processing or computer vision, CNN doesn't need complex preprocessing. However, when analyzing the data, data preprocessing plays a crucial role. One of the first steps is the normalization of the data. This step is essential when dealing with parameters of different units and scales. In this dataset, pixel values range from 0 to 255. I process normalization to this dataset. Normalization scales all numeric variables in the range of [0,1]. The formula is given below.

$$x_{new} = \frac{x - x_{min}}{x_{max} - x_{min}} \quad (3)$$

The minimum pixel value is 0 and maximum pixel value is 255. Therefore, I normalized the data by following.

$$x_{new} = \frac{x}{255} \quad (4)$$

### 3.2 Implementation

As I introduced in the section 2.4 (Benchmark), I utilized the 4 convolutional layers and 2 fully connected layers. In this section, I explain the architecture of the model in more depth.

The objective of the learning is to minimize the cross entropy in the learning process. The architecture of the model is the following.

Layer	Batch	Stride	Map Size	Function
data	—	—	$3 \times 32 \times 32$	—
conv1	$3 \times 3$	1	$32 \times 32 \times 32$	ReLu
pool1	$2 \times 2$	2	$32 \times 16 \times 16$	—
conv2	$3 \times 3$	1	$32 \times 16 \times 16$	ReLu
pool2	$2 \times 2$	2	$32 \times 8 \times 8$	—
conv3	$3 \times 3$	1	$32 \times 8 \times 8$	ReLu
pool3	$2 \times 2$	2	$32 \times 4 \times 4$	—
conv4	$3 \times 3$	1	$32 \times 4 \times 4$	ReLu
pool4	$2 \times 2$	2	$32 \times 2 \times 2$	—
fc5	—	—	512	ReLu
fc6	—	—	10	Softmax

Figure 12: Architecture of the model

As for the optimizer, I utilized 'Adam'. With this setting, I got the accuracy rate of 70.2%

To train the CNN model, I used Keras which is one of the neural network libraries like Theano or Tensorflow.

The important parameters training the model is as follows. Optimization parameters: Learning rate, beta1, beta2,

The number of stride or map size are shown in Fig.12. As for the initialization of the random weights in convolutional layers, I used uniform distribution.



I used normalized uniform distribution which is proposed by Glorot[10] for the initialization of the weights in fully connected layers I used ReLU(Rectified Linear Unit) as the activation except for the final layer and I used softmax as the final layer of activation.

$$a_{ReLU} = \log(1 + \exp(x)) \simeq \max(0, x) \quad (5)$$

### 3.3 Refinement

Finding the optimal parameters for deep learning is quite difficult though it is important. When it comes to typical machine learning algorithm(Decision tree, Support Vector Machine etc.),grid search is taken to search the optimal parameters. However, it's almost impossible to apply grid search in deep learning because of the computational time.It has reported that the grid search strategies are inferior to random search.[11] Therefore, other methods are indispensable. A good choice is Bayesian optimization, which has been shown to outperform other state of the art global optimization algorithms on a number of challenging optimization benchmark functions.

Bayesian Optimization[7]

For continuous functions,Bayesian optimization typically works by assuming the unknown function was sampled from a Gaussian process and maintains a posterior distribution for this function as observations are made or, in our case, as the results of running learning algorithm experiments with different hyper-parameters are observed.

There are two major choices that must be made when performing Bayesian Optimization. At first, we must choose a prior over functions that will express assumptions about the function being optimized. For this we select the Gaussian process prior, because of its flexibility and tractability. Second, we must select acquisition function, which is used to build a utility function from the model posterior, allowing us to determine us to decide the next point to evaluate.

Gaussian Process is a probabilistic model of regressor which can predict the new output  $y^{(n+1)}$  for  $x^{(n+1)}$  when given the data  $D = \{(x^{(n)}, y^{(n)})\}_{n=1}^N$  For more overview of Gaussian process, please refer to Rasmussen and Williams[8]

As for the acquisition function for Bayesian Optimization, I used GP Upper Confidence Bound

$$a_{UCB}(x; \{x_n, y_n\}, \theta) = \mu(x; \{x_n, y_n\}, \theta) + \kappa \sigma(x; \{x_n, y_n\}, \theta) \quad (6)$$

Here,  $\kappa$  is tunable parameter to balance exploitation against exploration.I set  $\kappa$  as 2.576.

The power of the Gaussian process to express a rich distribution on functions rests solely on the covariance function.The automatic relevance determination(ARD) squared exponential kernel is often used for Gaussian process regression.

$$K_{SE}(x, x') = \theta_0 \exp\left(-\frac{1}{2}r^2(x, x')\right) \quad (7)$$

$$r^2(x, x') = \sum_{d=1}^D (x_d - x'_d)^2 / \theta_d^2 \quad (8)$$

However, sample functions with this covariance function are unrealistically smooth for practical use. Therefore, I use ARD Matern 5/2 kernel.

$$K_{M5/2}(x, x') = \theta_0 (1 + \sqrt{5r^2(x, x')} + \frac{5}{3}r^2(x, x')) \exp(-\sqrt{5r^2(x, x')}) \quad (9)$$

I optimized the number of the dimensions of the first fully-connected layer and the learning rate in the optimizer(Adam) in this task.

## 4 Result

### 4.1 Model Evaluation, Validation and Justification

By utilizing the bayesian optimization, I optimized the number of dimensions of first layer of fully connected layer and learning rate. I got the optimal number of first layer of fully connected layer and the learning rate of the optimizer. Those numbers are 204 and 0.000597, respectively. By using those numbers, I got the test accuracy which is 72.1%

From the result of Bayesian Optimization, I chose 204 of first fully connected layer and 0.000597 of learning rate as the final model parameters.

The training and test accuracy in this model is as follows(Fig.13).

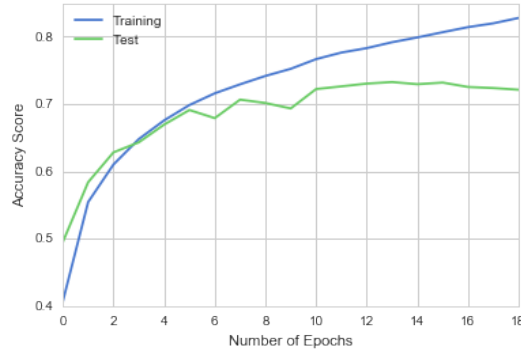


Figure 13: Accuracy rate of training and validation data

The loss of the model is shown in Fig.14

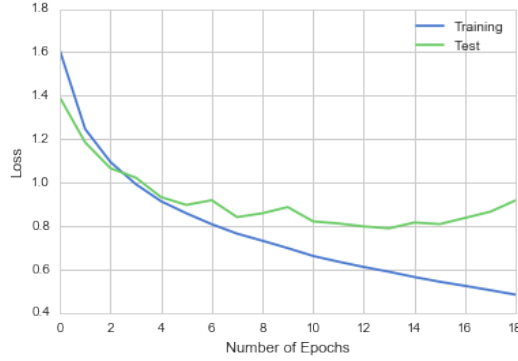


Figure 14: Loss of training and validation data

The model improves the accuracy rate by 2.0%. The final model doesn't seem to fall into the overfitting according to the Fig.13 and when compared to the Fig.11 which is the benchmark result, the difference between the training and test accuracy rate is much smaller. That means the final model is more robust than the benchmark model. Therefore, the model improved compared to the benchmark result.

I only optimize two parameters, learning rate and the number of first fully connected layer's neurons. The number of random sampling points of finding the optimal score in bayesian network is 25. When using grid search, we can only search 5 values for each parameters( $5 \times 5 = 25$ ). Apparently 5 different values for each parameter is not enough to search optimal parameters for grid search. Therefore, much more different values are necessary, resulting much computationally expensive. The difference will be much distinguishable when I try to optimize more parameters. In this case, grid search won't be a choice to find optimal parameter because of the scalability. As for bayesian optimization, it is highly possible that I can't find the best parameters within 25 sample points. However, the parameters found within 25 sample points are better result than the benchmark model parameters. When I choose more sample points, it will probably get better results.

## 5 Conclusion

### 5.1 Free-Form Visualization

Some data are miss-classified in this task. I'll look into which label is the poorest result and I'll put some of the images which are miss-classified.

The distribution of the correctly classified images are Fig.15

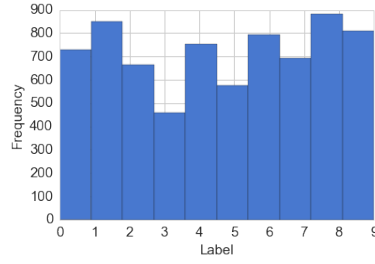


Figure 15: Distribution of the correct prediction

As you can see, the number of label 3 which is "cat" is apparently small. That means "cat" is overly miss-classified in this dataset.

The Fig.16 shows the distribution of the predicted output. The label 4("deer") and 8("ship") are conspicuously higher than the others. This shows that some of the label are miss classified.

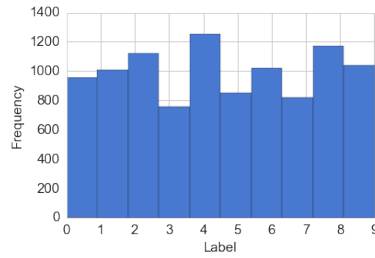


Figure 16: Distribution of the predicted label

The heatmap below clearly explains the miss-classified data. Label 3("cat") is miss-classified as label 5("dog"). Label 2("bird") and label 7("horse") are miss-classified as label 4("deer"). Label 1("automobile") is miss-classified as label 8("ship"). Judging from the heatmap, the miss-classified data are somewhat understandable since some cats are similar to dogs and some automobiles may be similar to ship because of the color and shape (Low resolution may be another reason).

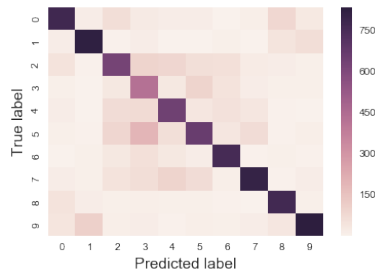


Figure 17: Heatmap of the data

## 5.2 Reflection

In this project, I solved the image recognition problem. These days, deep neural networks surpass the other typical image recognition algorithms and even surpass human recognition. However, one of the problems of deep learning is tuning hyper-parameters. As is often said, the architecture of the neural network is black art. Putting more layers may tend to grasp the feature of the input data. However, because of the back-propagation process, the gradient will vanish or explode while calculating. Therefore, the deep neural networks don't always reach the good result. What's more, the deep networks tend to take much more time than the shallower ones. Personally, I don't have any GPU environment so that trying deep network was quite tough. Therefore, I tried training the shallow networks with hyper-parameter tuning.

First, I decided the number of layers of convolutional layers and then I optimized the parameters with bayesian optimization which is one of the methods to tune hyper-parameters. I chose 4 layers of convolutional layers and 2 layers of fully connected layers for the benchmark, then I optimized the number of dimension of the first layer of the fully connected layer with bayesian optimization.

One of the difficulties of this task was parameter tuning. At first, I tried to tune by grid search which is quite popular in the machine learning field. However, I soon realized the method isn't reasonable because of the number of the parameters and the computational time. Thus I chose other hyper-parameter tuning method which is bayesian optimization.

## 5.3 Improvement

In this task, I chose 4 layers of convolutions and 2 layers of fully connected layers. It is reported that the number of layers is the crucial factor in the high accuracy rate. For example, GoogLeNet, VGG, SPP are around 20 convolution layers. In 2015, ResNet which surpassed the human recognition had maximum 152 layers. In fact, deep networks can represent certain function classes far more efficiently than shallow ones. That means basically deeper networks can get the higher accuracy rate.

However, the stacking of several non-linear transformations in convolutional feed-forward network architectures typically result in poor propagation of activations and gradients. Therefore, training deep network is difficult. ResNet takes the strategy of skipping connections between layers. When the number of layers increases, finding optimal parameters becomes more and more important. Bayesian optimization can be one of the methods to find the optimal parameter.

## References

- [1] Rupesh Kumar Srivastava, Klaus Greff, Jurgen Schmidhuber, Training Very Deep Networks, arXiv:1507.06228v2[cs.LG] 23 Nov 2015, 2015
- [2] Benjamin Graham. Fractional Max-Pooling, arXiv:1412.6071v4[cs.CV] 12 May 2015, 2015
- [3] 岡谷貴之.”機械学習プロフェッショナルシリーズ 深層学習”, 講談社, 2015

- [4] 神島敏弘, 麻生英樹, 安田宗樹, 前田新一, 岡野原大輔, 岡谷貴之, 久保陽太郎, ボレガラ ダヌシカ, 深層学習, 近代科学社, 2015
- [5] 浅川伸一, Python で体験する深層学習-Caffe, Theano, Chainer, TensorFlow-, コロナ社, 2016
- [6] Sebastian Raschka, Python Machine Learning, Packt Publishing, 2015
- [7] Jasper Snoek, Hugo Larochelle, Ryan P. Adams, Practical Bayesian Optimization of Machine Learning Algorithms
- [8] Carl E. Rasmussen and Christopher Williams. Gaussian Process for Machine Learning, MIT Press, 2006.
- [9] Diederik Kingma and Jimmy Ba. Adam: a method for stochastic optimization, In the 3rd International Conference for Learning Representations ICLR 2015, 2015
- [10] Xavier Glorot, Yoshua Bengio, Understanding the difficulty of training deep feedforward neural networks, JMLR W&CP 9:249-256, 2010
- [11] James Bergstra and Yoshua Bengio. Random search for hyper-parameter optimization. Journal of Machine Learning Research, 13:281-305, 2012
- [12] [https://en.wikipedia.org/wiki/Convolutional\\_neural\\_network](https://en.wikipedia.org/wiki/Convolutional_neural_network)