

Sugar Cane Grading from photo using machine learning
2559:39

Miss Pham Thi Mai Phuong, Phuong, 56070503447, phuongmaipham@icloud.com

Advisor: Dr. Sally E. Goldin

Co-advisor: -

September, 29th, 2016

I have read this report and approve its content.

Abstract (English)

This project will work in co-operation with staff from Mitrphol Sugar Company to create a software test bed for improving cane quality control over a large area. The test bed will be able to analyze the sugar cane health from mobile phone photos. The project will use machine learning technique to train the software to discriminate photos based on cane quality. The test bed will first extract the sugar cane crucial features from mobile phone photos. This test bed will then classify the cane photos into different health categories based on extracted features using supervised machine learning. The results obtained from this project will be useful for developing a real world system to allow individual farmers to send photos of their fields, which can be analyzed and classified to get more detailed information about cane health over a wide area. This project is thus important because it will help sugar companies gain better information with a lower surveying cost.

Acknowledgements

This project cannot be completed without the support from

Table of Contents

<i>Chapter 1</i>	7
<i>Introduction</i>	7
<i>1.1 Problem Statement and Approach</i>	7
<i>This is a research - real world stakeholder project</i>	7
<i>1.2 Objectives</i>	7
<i>1.3 Scope</i>	7
<i>Deliverables for Term 1</i>	8
<i>Deliverables for Term 2</i>	8
<i>1.4 Tasks and Schedule</i>	8
<i>1.4.1 Task breakdown</i>	8
<i>1.4.2 Draft Schedule</i>	9
<i>Chapter 2</i>	11
<i>Background, Theory and Related Research</i>	11
<i>2.1 Digital Image Processing Concepts</i>	11
<i>2.1.1 Digital Image Representation</i>	11
<i>2.1.2 Image Processing Operation</i>	12
<i>2.2 Image Convolution</i>	13
<i>2.2.1 Concepts and Mathematics</i>	13
<i>2.2.2 Examples – Sobel, Blurring</i>	14
<i>2.3 Machine learning</i>	15
<i>2.3.1 Core concepts of machine learning</i>	15
<i>2.3.2 Training and testing strategies</i>	16
<i>2.3.3 Model evaluation</i>	16
<i>2.4 Multilayer feed-forward neural networks (MFNN)</i>	18
<i>2.4.1 Overviews</i>	18
<i>2.4.2 Back propagation</i>	19
<i>2.5 Convolutional neural networks (CNN)</i>	20
<i>2.5.1 Overviews</i>	20
<i>2.5.2 Convolutional layer</i>	20
<i>2.5.2.1 Local receptive field</i>	20
<i>2.5.2.2 Filter (Kernel)</i>	21
<i>2.5.2.3 Feature map and stride</i>	21
<i>2.5.2.4 The use of zero-padding</i>	22
<i>2.5.3 ReLU non-linearity</i>	22
<i>2.5.4 Pooling layer</i>	23
<i>2.5.4.1 Overlap pooling</i>	24
<i>2.5.4.2 Exploiting viewpoints in pooling</i>	24
<i>2.5.5 Fully connected layer</i>	25
<i>2.5.6 Softmax regression function</i>	25
<i>2.6 Learning in a Convolutional neural networks</i>	25

2.6.1	<i>Back propagation core concepts.....</i>	25
2.6.2	<i>Back propagation at the Softmax.....</i>	26
2.6.3	<i>Back propagation in the Fully connected layer (Neural networks).....</i>	26
2.6.4	<i>Back propagation in the ReLU layer</i>	28
2.6.5	<i>Back propagation in the pooling layer</i>	29
2.6.6	<i>Back propagation in the convolutional layer.....</i>	29
2.6.7	<i>What does a CNN learn?</i>	30
2.7	<i>Data Augmentation</i>	31
2.7.1	<i>Concept</i>	31
2.7.2	<i>Methods.....</i>	31
2.7.2.1	<i>Brightness adjustment.....</i>	31
2.7.2.2	<i>Contrast adjustment</i>	32
2.7.2.3	<i>Scaling.....</i>	32
2.8	<i>Data post-processing</i>	33
2.9	<i>GPU Concepts</i>	33
2.10	<i>CNN Software Frameworks</i>	34
2.11	<i>Related research</i>	36
2.11.1	<i>Researches using CNN to classify images</i>	36
2.11.2	<i>Variations in algorithms in CNN used for classifying fine-grained objects.....</i>	36

List of Figures

Figure 2.1-1 what happens as we reduce the resolution of an image while keeping its size the same. [9]	11
Figure 2.1-2 level grey scale. [9]	11
Figure 2.1-3 color space. [9]	12
Figure 2.1-4 color intensity. [11]	12
Figure 2.1-5 an image histogram[14]	13
Figure 2.2-1 a 3x3 neighborhood about point (u,v) in an image	14
Figure 2.2-2 the vertical Mask of Sobel operator.....	15
Figure 2.2-3 the horizontal Mask of Sobel operator	15
Figure 2.2-4 the 3x3 mean filter	15
Figure 2.4-1 fully connected layers. [8]	19
Figure 2.5-1 a convolutional neural networks. [16]	20
Figure 2.5-2 local receptive fields. Image from Deep learning - draft book in preparation, [7]	20
Figure 2.5-3 convolutional process - An input image of size 7x7x3 is filtered by 2 convolutional kernels which create 2 feature maps. [12].....	22
Figure 2.7-1 brightness adjustment. [14]	32
Figure 2.7-2 contrast adjustment. [14].....	32

Chapter 1

Introduction

1.1 Problem Statement and Approach

Sugarcane is an important crop in Thailand. In order to produce high quality sugar, sugar companies need detailed information on the cane conditions in different fields. Features such as color and size of the leaves are some indicators of cane conditions (cane health).

A field's condition, such as soil properties, seed quality or irrigation system, has some important consequences for those features. However, not every field is identical. Thus, the features and the resulting cane quality vary from one field to another. This will cause a problem in collecting data: there are too many fields and it is too complicated for sugar companies to do exhaustive surveys to get the information that they need.

To address this problem, this project will work in co-operation with staff from Mitrphol Sugar Company to create a software test bed for improving cane quality control over a large area. The test bed will be able to analyze the sugar cane health from mobile phone photos. The project will use machine learning to train the software to discriminate photos based on cane quality. The results obtained from this project will be useful for developing a real world system to allow individual farmers to send photos of their fields, which can be analyzed and classified to get more detailed information about cane health over a wide area. This project is thus important because it will help sugar companies gain better information with a lower surveying cost.

This is a research - real world stakeholder project.

1.2 Objectives

The goal of this project is to develop a software test bed for experimenting with mobile phone images of sugar cane using supervised machine learning technique. The test bed will first extract the sugar cane crucial features from mobile phone photos. This test bed will then classify the cane photos into different health categories based on extracted features using supervised machine learning.

1.3 Scope

This project does not attempt to deliver a final system, but rather a software testbed. The testbed is intended to experiment mobile phone images of sugarcane with different machine learning models and conclude which model can achieve the best result.

Deliverables for Term 1

- ✓ *Experimental data set*
- ✓ *Experimental design*
- ✓ *Some prototype using the selected framework*
- ✓ *Decision on what learning framework(s) to use, with justification*

Deliverables for Term 2

- ✓ *Complete experimental design of the test bed*
- ✓ *Software test bed with desirable results*
- ✓ *Results and data analysis*

1.4 Tasks and Schedule

1.4.1 Task breakdown

1. *Analyze and determine the requirements of the project*
2. *Plan the project schedule*
3. *Work on introduction chapter of the report (chapter 1)*
4. *Research emphasizing on the following topics:*
 - i. *Work by other researchers on discriminating between similar images using machine learning*
 - ii. *Machine learning methods for image classification and the available libraries*
 - iii. *Basic image processing concepts*
5. *Create the project proposal and get feedbacks*
6. *Test prototypes for various learning frameworks and make a decision on which learning framework to use*
7. *Collect and create dataset*
8. *Study and understand the dataset*
9. *Create experimental design*
10. *Complete progress report for the first semester*
11. *Work on theory and background chapter of the report (chapter 2)*
12. *Work on methodology chapter of the report (chapter 3)*
13. *Prepare for presentation for the first semester*
14. *Write pre-processing software to standardize images*
15. *Write scripts to control the experiments*
16. *Test the system and fix bugs*
17. *Train and test the system with different parameters*
18. *Analyze the results*
19. *Complete final report for the second semester (Result + conclusions, chapter 4 and 5)*
20. *Create poster and prepare for presentation for the second semester*

1.4.2 Draft Schedule

[illegible]

[illegible]

**** *H:holiday***

Chapter 2

Background, Theory and Related Research

2.1 Digital Image Processing Concepts

2.1.1 Digital Image Representation

Pixels are basic elements of an image. Each pixel depicts a light intensity value, which is represented in binary code. Digital image is a collection of pixels. There are two types of digital image: grayscale image and color image. Grayscale image is defined by a matrix of pixels, whereas color image is defined by a cube made of three matrixes of pixels.

The density of pixels in an image is called 'resolution'. With the same image size, the more pixels that we keep to describe an image, the more detailed the image.

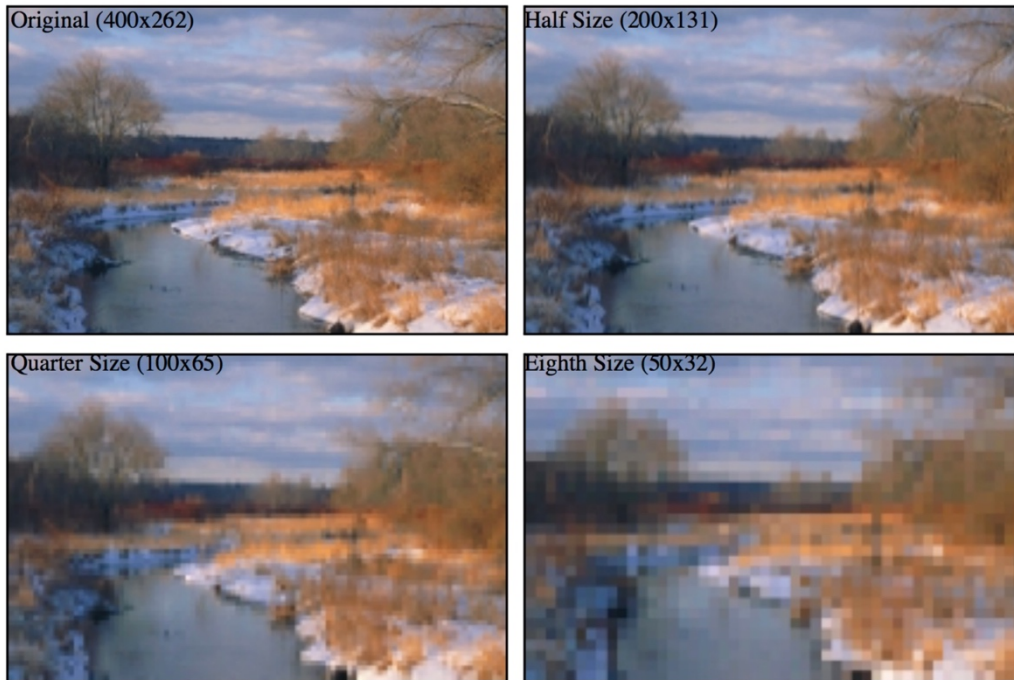


Figure 2.1-1 what happens as we reduce the resolution of an image while keeping its size the same. [9]

A monochrome image can be described in terms of a two-dimensional light intensity function $f(x,y)$, where the amplitude of $f(x,y)$ is the intensity (brightness) of the image at position (x,y) . The intensity of a monochrome image lies in the range

$$L_{\min} < f(x,y) < L_{\max}$$

Equation 2.1-1

where the interval $[L_{\min}, L_{\max}]$ is called the grey scale. There are two common grey scale storage methods: 8-bit storage and 1-bit storage.

The most common storage method is 8-bit storage. It can represent up to 2^8 colours for each pixel. The grey scale interval is $[0,255]$, with 0 being black and 255 being white.

The less common storage method is 1-bit storage. There are two grey levels, with 0 being black and 1 being white.



Figure 2.1-2 level grey scale. [9]

A colored image can be represented using multi-channel color models. The most widely used model is the RGB. Colored images are made up of three primary spectrums: red, green and blue (RGB). These three primary spectrums together create a three-dimensional color space where red defining one axis, green defining the second, and blue defining the third. Every existing color is described as a mixture of red, green, and blue light and located somewhere within the color space.

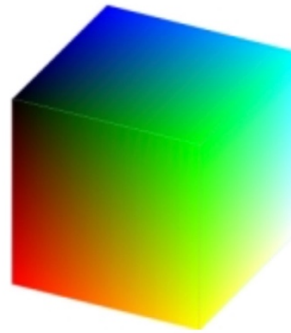


Figure 2.1-3 color space. [9]

Using RGB model, a colored image can be described using a three-channel intensity function

$$I_{RGB}(x, y) = (F_R(x, y), F_G(x, y), F_B(x, y))$$

Equation 2.1-2

where $F_R(x, y)$ is the intensity at position (x, y) in the red channel, $F_G(x, y)$ is the intensity at position (x, y) in the green channel, and $F_B(x, y)$ is the intensity at position (x, y) in the blue channel. Each channel usually uses an 8-bit storage which can describe up to 2^8 colours. Thus, computers commonly use a 24-bit storage to describe the intensity at position (x, y) , which can describe up to 2^{24} colours.

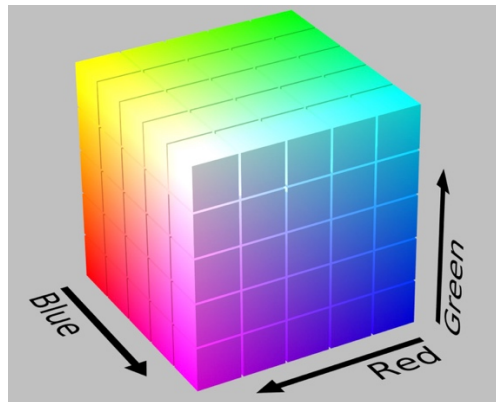


Figure 2.1-4 color intensity. [11]

2.1.2 Image Processing Operation

Image operations are calculations applied to each pixel (x, y) of an input image f to transform the input image $f(x, y)$ into an output image $g(x, y)$. Different operation types are applied to an image depending on what kind of results we want to achieve. In this section, we will focus on discussing three types of image operations: resampling, cropping and histogram modifications.

Resampling is a technique used to generate a new image with a different resolution. Increasing the number of pixels is called upsampling and decreasing the number of pixels is called downsampling. Upsampling softens the original image whereas downsampling sharpens the original image. There are several different resampling techniques. The main concept of those techniques is to interpolate the value of a new pixel from the existing pixels

of the original image. Nearest neighbor resampling is the simplest: The value of each pixel in the output image is the value of its nearest neighbor in the original image.

Cropping refers to the removal of the regions in the image that contains less useful information. A new pair of coordinate (x,y) is required to define the corners of the cropped image. The resulting image is a smaller image that contains useful data for further studies. An image's histogram is a graph shows the frequency of pixel intensities. For example, the histogram below shows the frequency of pixel values of a grayscale image. Where the x axis indicates the range of the pixel values and y axis indicates the frequency these values appear on the entire image.

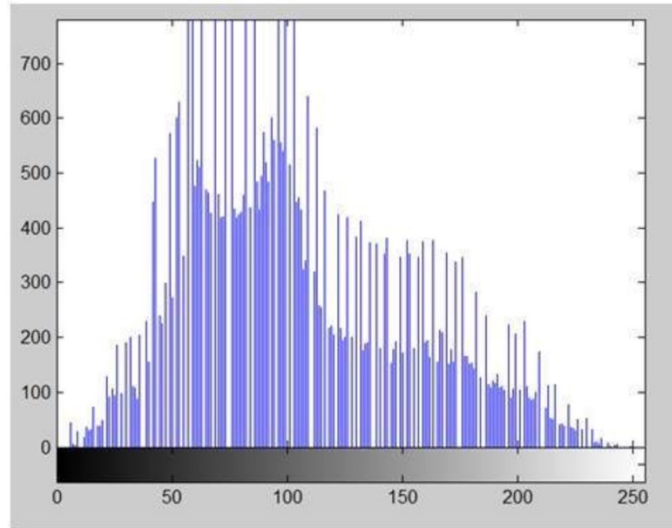


Figure 2.1-5 an image histogram[14]

There are various operations to modify an image's histogram. In this section we will focus on discussing about histogram sliding and histogram stretching.

Histogram sliding is a method used to change an image brightness. To increase the brightness, we shift an image's histogram to the right by adding a constant value to every pixel of an image. In contrast, to decrease the brightness, we shift an image's histogram to the left by subtracting a constant value from every pixel of an image.

Histogram stretching is a method used to increase the contrast. Each output pixel is calculated using this formula:

$$g(x, y) = \frac{f(x, y) - f_{\min}}{f_{\max} - f_{\min}} 2^{bpp}$$

Equation 2.1-3

Where $f(x, y)$ is the pixel intensities of the original image, f_{\min} is the minimum pixel intensity, f_{\max} is the maximum pixel intensity and bpp is the levels of gray. The purpose of this function is to 'spread' the old intensity values over a larger range.

2.2 Image Convolution

2.2.1 Concepts and Mathematics

Convolution is an image transformation technique using neighborhood (or area-based) operators. The objective of image convolution is to process an input image in order to enhance some important features that is suitable for a specific application. It has two main approaches: spatial-domain approach and frequency-domain approach.

Spatial-domain and frequency-domain approach can both be described in terms of convolutional function:

$$g(x,y) = h(x,y)*f(x,y)$$

Equation 2.2-1

given $f(x,y)$ is the input image, $h(x,y)$ is the kernel and $g(x,y)$ is the output image. In spatial-domain, we look at the value of each pixel varies with respect to scene whereas in frequency-domain, we look at the rate at which the pixel values change in spatial-domain. Thus, in frequency-domain approach, we will have to convert the pixel values into frequency-domain before applying convolution, then convert the result back into spatial-domain. We will mainly discuss about spatial-domain procedure in this section.

The procedure for convolution in spatial-domain is as follows: A filter (sometimes can be referred to as a mask, a kernel or a window) centered at point (u,v) is flipped in both dimensions and then slit around the input image. Each time the filter is placed at a new position, every pixel of the input image contained within the filter is multiplied by the corresponding filter coefficient and then summed together. The result from each multiplication and summation declares the next pixel of the output image. The described procedure is repeated this until all values of the image has been calculated. It is mathematically described as:

$$G(i,j) = \sum \sum H(u,v)F(i-u,j-v)$$

Equation 2.2-2

where F is the input image, H is the filter and G is the output image.

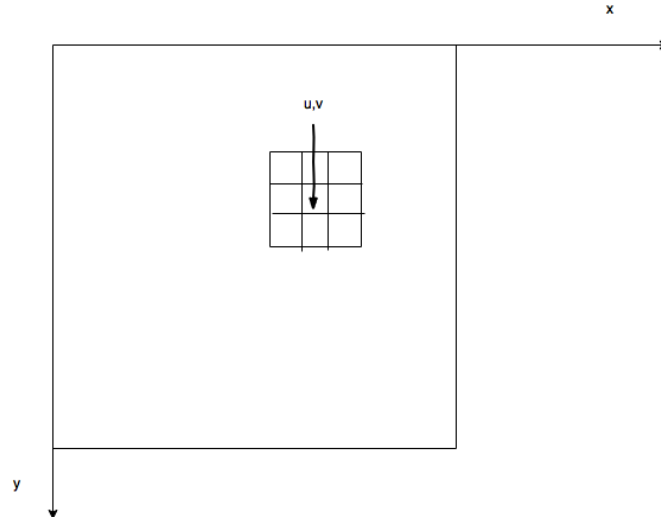


Figure 2.2-1 a 3x3 neighborhood about point (u,v) in an image

2.2.2 Examples – Sobel, Blurring

Below we will give two simple examples to illustrate the application of convolutional filtering.

The first application is to highlight edges in an image. Some known edge enhancement filters are Prewitt operator, Sobel operator, Robinson compass masks, Krisch compass Masks and Laplacian operator. The decision on which filter to use depends on our desired results. Here we will describe the use of Sobel edge detection. A Sobel filter is used to calculate edges in both horizontal and vertical direction.

The vertical Mask of Sobel operator is as follows:

-1	0	1
-2	0	2
-1	0	1

Figure 2.2-2 the vertical Mask of Sobel operator

The pixels at the corresponding position to the area declared by this filter of the input image are respectively multiply by -1,0,1,-2,0,2,-1,0 and 1. The results of these nine multiplications are then summed. This process is repeated until all values of the image has been calculated. The horizontal Mask of Sobel operator is as follows:

-1	-2	-1
0	0	0
1	2	1

Figure 2.2-3 the horizontal Mask of Sobel operator

Similarly, the pixels at the corresponding position to the area declared by this filter of the input image are respectively multiply by -1,-2,-1,0,0,0,1,2 and 1. The results of these nine multiplications are then summed. This process is repeated until all values of the image has been calculated.

This give more weight age to the pixel values around the edge region. It thus increases the edges intensities. As a result, the output image edges become enhanced comparatively to the original image.

The first application is to blur an image. There are three common type of filters that are used to perform blurring: Mean filter, weighted average filter and Gaussian filter. We are going to discuss about mean filter. In a mean filter, there are an odd number of filter elements, all of which are the same and can be summed to one. For example, a 3x3 can be declared as followed:

1/9	1/9	1/9
1/9	1/9	1/9
1/9	1/9	1/9

Figure 2.2-4 the 3x3 mean filter

The pixels at the corresponding position to the area declared by this filter of the input image are respectively multiply by 1/9. The results of these nine multiplications are then summed. This operation can be used to discard false effects that may be present in a digital image as a result of poor sampling system or transmitting channel. This process is repeated until all values of the image has been calculated.

2.3 Machine learning

2.3.1 Core concepts of machine learning

Machine learning studies computers' abilities to automatically recognize patterns and generate intelligent conclusion from the given data. It involves two learning types: supervised and unsupervised.

Supervised learning (SVM) is equivalent to data classification. Computers learn patterns from the labeled examples then use them to make intelligent classification on the unknown data.

Unsupervised learning (USVM) is equivalent to clustering. Initially, there is no label associated with the data. Computers try to divide the dataset into clusters to discover classes

within the data. USVM cannot narrate semantic meaning of the clusters in their learnt models.

Data classification consists of learning step and classification step. In the learning phrase, computers build a classification model by studying the **training set** made up of pre-labeled tuples. A tuple X is described by a set of features. Each tuple is represented by a vector of n -dimension $X = (x_1, x_2, \dots, x_n)$ and categorized with a class label attribute x_c . The class label attribute has a discrete-valued that indicates which class tuple X tuple belongs.

In the classification phrase, the learned model is used for predicting the class label for given data. To evaluate this model, a **testing set** is used. Each testing set is also made up of pre-labeled tuples and is independent of training sets.

When the model is tailored to fit the random noise in one specific sample rather than reflecting the overall population, we have an **overfitting** problem. There are various methods to reduce **overfitting**. One is to put some random noise in the data. Another is to stop training before 100% accuracy is achieved.

2.3.2 Training and testing strategies

The given dataset is divided into training and testing sets using one of these three following methods:

Hold out: The given data set is divided into two independent sets: training set and testing set. Two-thirds of the data are in the training set and one-third of the data is in the testing set.

k-fold cross-validation: The data set is divided into k subsets, and the holdout method is repeated k times. Each time, one of the k subsets is used as the test set and the other $k-1$ subsets are put together to form a training set. Then the average error across all k trials is computed. The advantage of this method is that it matters less how the data gets divided.

Bootstrap: A common bootstrap method is the .632 bootstrap. Given a set of d tuples. This dataset will be sample d times with replacement. Each time a tuple is selected, it is re-added into the data pool and likely to be selected again. The data that do not make it into the training set will eventually be added into the testing set. The training set and the testing set are not independent. In .632 bootstrap, 63.2% of the dataset will end up in the bootstrap sample, and the 36.8% that did not make it to the bootstrap sample will together form the test set. Bootstrapping is especially useful when we want to replicate the distributional properties of a random variable in the real world.

2.3.3 Model evaluation

Confusion matrix is often used to measure the performance of our classifier. It is represented a table of size $j \times j$, where j is the number of output classes. A table entry a_{ij} is in row i^{th} and column j^{th} , indicates the number of tuples that are actually in class i but was classified by the classifier to be in class j .

For calculating convention, for every output class j , we will categorize the the tuples classified by the classifier into four groups, as suggested by Jiawei et al. [8]

True positive (TP): For each given class j , tuples that are pre-labeled as class j and correctly recognized by the classifier to be in class j are true positive. TP represents the number of true positive tuples.

True negative (TN): For each given class j , tuples that are pre-labeled as class i and correctly recognized by the classifier to be in class i are true negative. TN represents the number of true negative tuples.

False positive (FP): For each given class j , tuples that are pre-labeled as class i and mistakenly recognized by the classifier to be in class j are false positive. FP represents the number of false positive tuples.

False negative (FN): For each given class j , tuples that are pre-labeled as class i and

mistakenly recognized by the classifier to be in any class other than j are false positive. FP represents the number of false negative tuples.

Table 2.3-1 and 2.3-2 shows two examples of TP , TN , FP , FN for class 0 and class 1 respectively.

Actual class	Predicted class				
	0	1	2	...	j
0	TP	FN	FN	FN	FN
1	FP	TN	FN	FN	FN
2	FP	FN	TN	FN	FN
...	FP	FN	FN	TN	FN
j	FP	FN	FN	FN	TN

Table 2.3-1 Confusion matrix for class 0

Actual class	Predicted class				
	0	1	2	...	j
0	TN	FP	FN	FN	FN
1	FN	TP	FN	FN	FN
2	FN	FP	TN	FN	FN
...	FN	FP	FN	TN	FN
j	FN	FP	FN	FN	TN

Table 2.3-2 Confusion matrix for class 1

By dividing classified tuples into for groups, we can now evaluate our classifiers. The classifier accuracy or recognition rate on the given test set is calculated as follow:

$$\text{accuracy} = \frac{TP + TN}{P + N}$$

Equation 2.3-1

where

$$P = TP + FP$$

Equation 2.3-2

and

$$N = TN + FN$$

Equation 2.3-3

Similarly, the classifier error rate or misclassification rate is:

$$\text{Error rate} = 1 - \text{accuracy} = \frac{FP + FN}{P + N}$$

Equation 2.3-4

2.4 Multilayer feed-forward neural networks (MFNN)

2.4.1 Overviews

MFNN is a classic machine learning algorithm that simulates the way human brain works. MFNN is very useful for studying large datasets due to its ability to tolerate noise data as well as to recognize patterns on which they have little knowledge. It is especially efficient for real world data, where we do not know much about the relationships between attributes and classes. Generally, MFNN requires a large amount of data in order to well perform.

Each MFNN consists of an input layer, one more hidden layers and an output layer. Layers are connected in acyclic graph. Each layer is made up of computational elements called neurons. Neurons between two adjacent layers are pairwise connected, but neurons in one layer share no connection. No direct connection exists between input and output layer.

Inputs are fed into the neurons making up the input layer. The outputs produced by this layer are weighted and passed simultaneously to the first hidden layer. This hidden layer outputs are again weighted and input to an another hidden layer and so on. It is arbitrary how many hidden layers there should be. The weighted outputs of the last hidden layer are input to the output layer, where the prediction for the given tuples will be produced.

Neurons in the input layer are 'input units'. Neurons in the hidden layers and the output layers are called 'neurodes' or sometimes referred to as 'output units'. The number of input units are not necessarily equal number of input units. There can be more or less number of hidden units than number of input or output units. Each output unit applies a nonlinear (activation) function to its input. The activation function will be described in section 2.5.6.

The output is suggested as in the following function:

$$a_k^l = f\left(\sum_{j=1}^m w_{jk} a_j^{l-1} + b_k^l\right)$$

Equation 2.4-1

Where f is the neural activation function, a_k^l is output of neuron k at layer l , a_j^{l-1} is output of neuron j at layer $(l-1)$, w_{jk} is the weight of the connection between node j and k , b_k^l is the bias of node k .

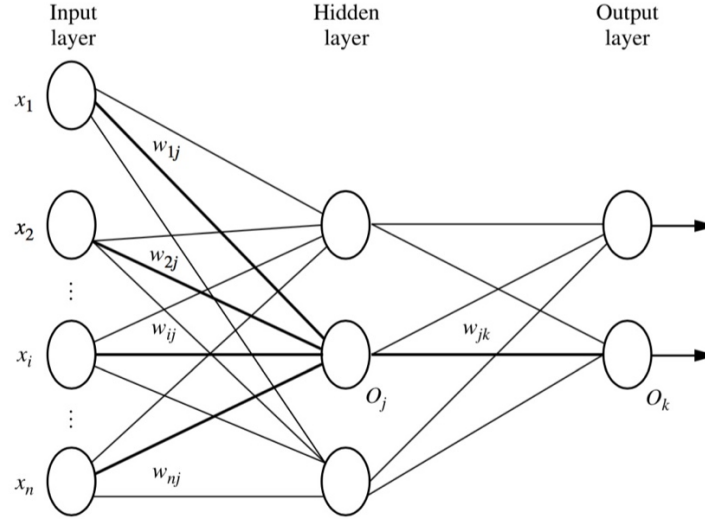


Figure 2.4-1 fully connected layers. [8]

2.4.2 Back propagation

Back propagation is the learning algorithm of neural networks. Total error is the summation of error at each output unit, which can be calculated using the squared error function:

$$E = \sum_N \frac{1}{2} (y_{o_N}^o - a_{o_N}^o)^2$$

Equation 2.4-2

Where

$$a_{o_N}^o = \frac{1}{1 + e^{-z_{o_N}^o}}$$

Equation 2.4-3

The error rate at each output unit can thus be calculated from this total error:

$$\delta_{o_N}^o = \frac{\partial E}{\partial z_{o_N}^o} = \frac{\partial E}{\partial a_{o_N}^o} * \frac{\partial a_{o_N}^o}{\partial z_{o_N}^o} = -(y_{o_N}^o - a_{o_N}^o) a_{o_N}^o (1 - a_{o_N}^o)$$

Equation 2.4-4

This error is then propagated back to the networks:

$$\delta_{z_j}^l = \frac{\partial E}{\partial a_j^l} * \frac{\partial a_j^l}{\partial z_j^l} = \frac{\partial E}{\partial z_i^{l+1}} * \frac{\partial z_i^{l+1}}{\partial a_j^l} * \frac{\partial a_j^l}{\partial z_j^l} = \left(\sum_{0 \leq i' \leq m} \delta_{z_{i'}}^{l+1} * w_{ji'}^{l+1} \right) * \frac{\partial a_j^l}{\partial z_j^l}$$

Equation 2.4-5

With learning rate η , this error will affect the corresponding weight by an amount of:

$$w_{ij}^l = w_{ij}^l \pm \eta \Delta w_{ij}^l$$

Equation 2.4-6

where

$$\Delta w_{ij}^l = \frac{\partial E}{\partial w_{ij}^l} = \frac{\partial E}{\partial a_j^o} * \frac{\partial a_j^o}{\partial z_j^o} * \frac{\partial z_j^o}{\partial w_{ij}^l} = \delta_{z_j}^l a_i^{l-1}$$

Equation 2.4-7

2.5 Convolutional neural networks (CNN)

2.5.1 Overviews

CNN is a “state-of-the-art technique for image recognition” [5]. It is a multilayer neural networks (as illustrated in figure 2.5-1). It consists of one or more convolutional layers, followed by pooling layers (sometimes called subsampling layers), ReLU layers, and one or more fully connected layers. Convolutional layers are used for extracting important features from the input images. The features learnt by a convolutional layer are often summarized by a pooling layer. ReLU layer eliminates negative outputs produced by the pooling layer preceding it. The learnt features are eventually passed into fully connected layers, where each input image is mapped with a suitable output class.

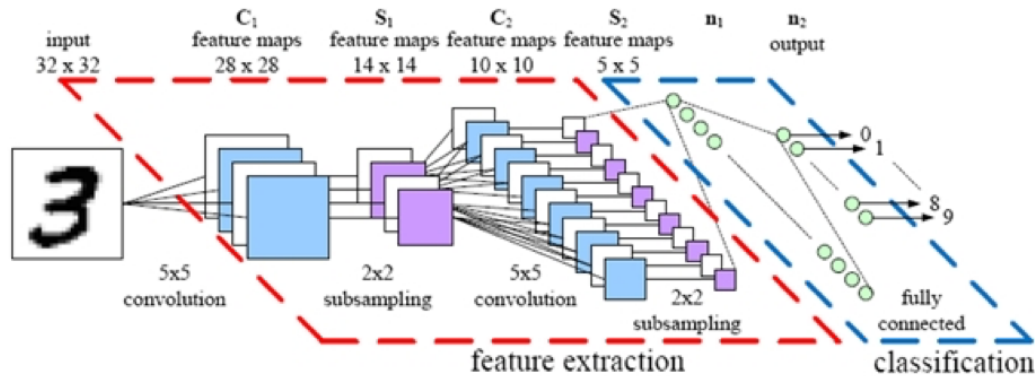


Figure 2.5-1 a convolutional neural networks. [16]

The networks build a sufficient model for the input dataset by imposing a set of forward and back propagations. At each iteration, the CNN model tries to reduce the number of wrong predictions by updating the weights that are associated with the location of the errors. The updating process iterates until all the weights in the network converges.

2.5.2 Convolutional layer

Convolutional layer extracts the important features of images via image convolution. Four parameters are required in convolutional layer: the number of layer K , receptive field size F , the stride S and the amount of zero padding P . This layer accepts input of volume size $Width_{in}Height_{in}Dimension_{in}$ and produces an output of volume size $[(Width_{in} - F + 2P)/S + 1][(Height_{in} - F + 2P)/S + 1][K]$.

2.5.2.1 Local receptive field

Think of the input image as a square of $n \times n$ neurons. The value of each neuron is the corresponding pixel intensity of the input image. We will map a localized region of the input neurons to a neuron in the hidden layer. These localized regions of input neuron are called local receptive fields.

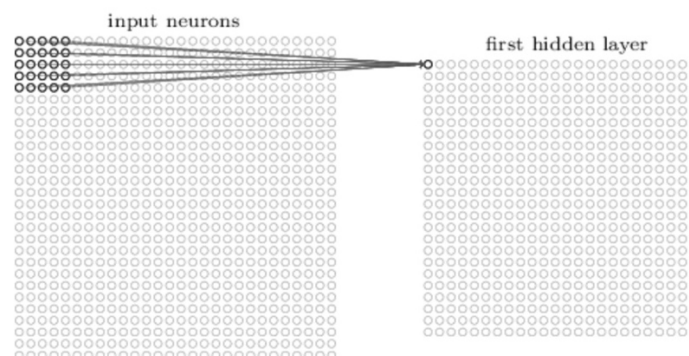


Figure 2.5-2 local receptive fields. Image from Deep learning - draft book in preparation, [7]

2.5.2.2 Filter (Kernel)

Each mapping from input layer to a hidden neuron at the hidden layer next to it learns a weight $w_{a,b}^l$ and each corresponding hidden neuron learns a bias $b_{x,y}^l$. This weight and bias together make up a kernel (sometimes called a filter). The kernel value is shared among all neurons in the same hidden layer. The output generated by the x^{th} , y^{th} hidden neuron is a convolutional function between an input neuron at position $(x-a, y-b)^{th}$'s value and its weight:

$$a_{x,y}^l = \sigma \left(\sum_a \sum_b w_{a,b}^l a_{x-a,y-b}^{l-1} + b_{x,y}^l \right)$$

Equation 2.5-1

Where σ is neural activation function, $a_{x,y}^l$ is the real output value at position (x,y) , $a_{x-a,y-b}^{l-1}$ is the input pixel value at position $(x-a, y-b)$, $b_{x,y}^l$ is the shared value for the bias, $w_{a,b}^l$ the weight of the connection at the a^{th} row and the b^{th} column of the receptive field. Notice that the weights used are inversed weights (as discussed in 2.2).

2.5.2.3 Feature map and stride

The combination of outputs generated by a hidden layer is called a feature map. To calculate the value of the next pixel in the feature map, we will move our kernel k pixel along both width and height concurrently (in which case we would say a stride length of k is being used) and re-apply the convolution function discussed above, until we run out of input neuron. Figure 2.5-3 shows an example of the convolutional process using $2 \times 3 \times 3 \times 3$ convolutional kernels with stride 2. To calculate the first pixel in dimension one of the output, we will perform the following. We will first concurrently convolute each dimension of filter $W0$ with the first 3×3 receptive field in the corresponding dimension of the input. We will then sum up the results from the convolutional process and the bias $b0$ together. The next pixels in dimension one of the output are calculated similarly with one variation: the receptive field is slit to the right or to the bottom by 2 pixels each time. Likewise, the pixels in dimension two of the output are calculated with the same fashion but the filter $W0$ is replaced by filter $W1$.

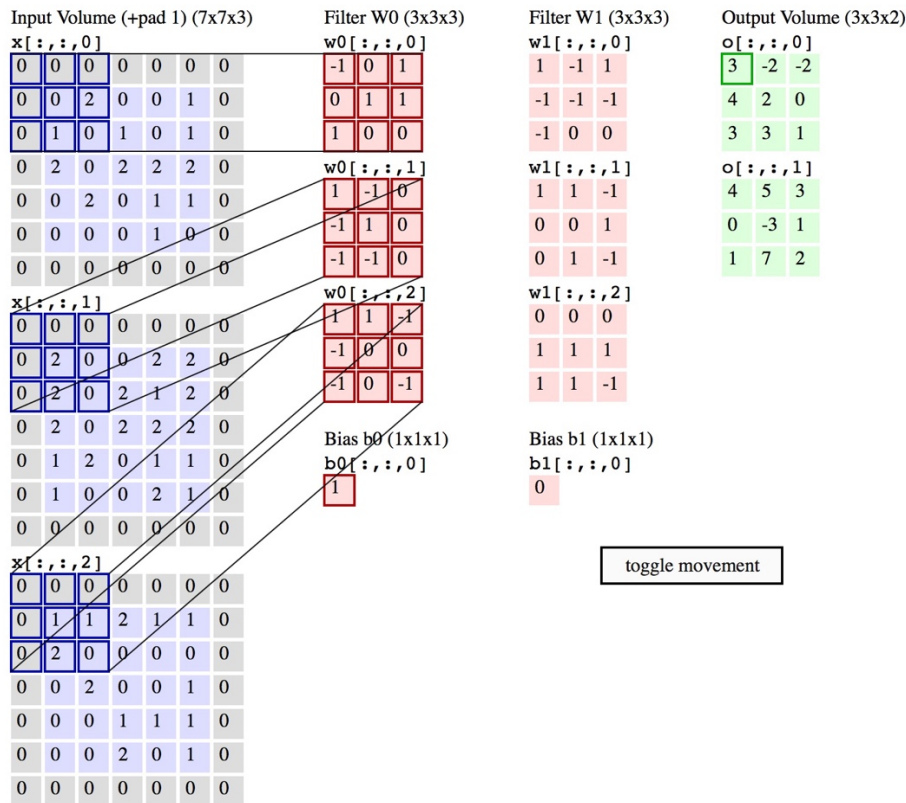


Figure 2.5-3 convolutional process - An input image of size 7x7x3 is filtered by 2 convolutional kernels which create 2 feature maps. [12]

Each feature map can detect one feature at different locations of the same input image. To detect many features from one input image, we need several features maps. To generate several feature maps, we will need several kernels. A complete convolutional layer consists of several different features maps.

2.5.2.4 The use of zero-padding

Zero padding is a simple method of padding the borders of the input image in order to control the dimensionality of the output image. Specifically, the relationship among padding size P , input height/length W , output height/length O and filter size K can be described as follow:

$$O = \frac{W - K + 2P}{2} + 1$$

Equation 2.5-2

2.5.3 ReLU non-linearity

A standard way to model a neuron output is with tanh function:

$$f(x) = (1 + e^{-x})^{-1}$$

Equation 2.5-3

Rectified Linear Unit is a modern way to represent an output f as a function of input x . It computes the function:

$$f(x) = \max(0, x)$$

Equation 2.5-4

Where

$$x = w * x + b$$

Equation 2.5-5

which means the output is 0 when the input is less than 0 and the output is a linear with slope 1 when x the input is greater than 0.

ReLU has two advantages over tanh function. Deep learning neural networks with ReLU was found to train significantly faster than networks with tanh unit. This is because saturating nonlinearities are generally slower than non-saturating nonlinearities in terms of training time with gradient descent. Secondly, while tanh function involves expensive operations, ReLU can be implemented easily by thresholding a matrix of activations at zero. A convolutional layer is often followed by ReLU.

2.5.4 Pooling layer

A convolutional layer is usually followed by a pooling layer. Pooling layer's function is to simplify every feature map in the previous layer at every depth slide spatially. Thus it reduces the amount of parameters, computation in the network and therefore overfitting. A successful pooling layer should be able to preserve critical information while being "invariant to troublesome deformations" [9].

Pooling layer requires two parameters: filter size F and stride S . It accepts input with volume size $Width_{in} Height_{in} Dimension_{in}$ and produces an output of volume size $[(Width_{in} - F) / S + 1][(Height_{in} - F) / S + 1][Dimension_{in}]$.

Pooling is done via a summary statistic over a region of neurons in the preceding layer. The summary statistic in l_i is defined by the l_i norm of inputs in the pools. If node $a_{x+p,y+q}^{l-1}$ (where $0 \leq p, q \leq k$) from the preceding layer are in the pool, the output of the pooling process for each pixel at position (x,y) can be mathematically represented as:

$$a_{x,y}^l = \left(\sum_{0 \leq p, q \leq k} |a_{x+p,y+q}^{l-1}|^i \right)^{1/i}$$

Equation 2.5-6

Several widely used pooling techniques can be derived from function (2.5-6). The most known one is **max pooling**. This is when $p \rightarrow \infty$

$$a_{x,y}^l = \max_{0 \leq p, q \leq k} (a_{x+p,y+q}^{l-1})$$

Equation 2.5-7

Figure 2.5-4 shows an example of max pooling operation on 4 4 x 4 slices. Each slice is filtered by a 2x2 kernel with stride 2 which results in 4 pooling feature maps, each of size 2 x 2. [12]

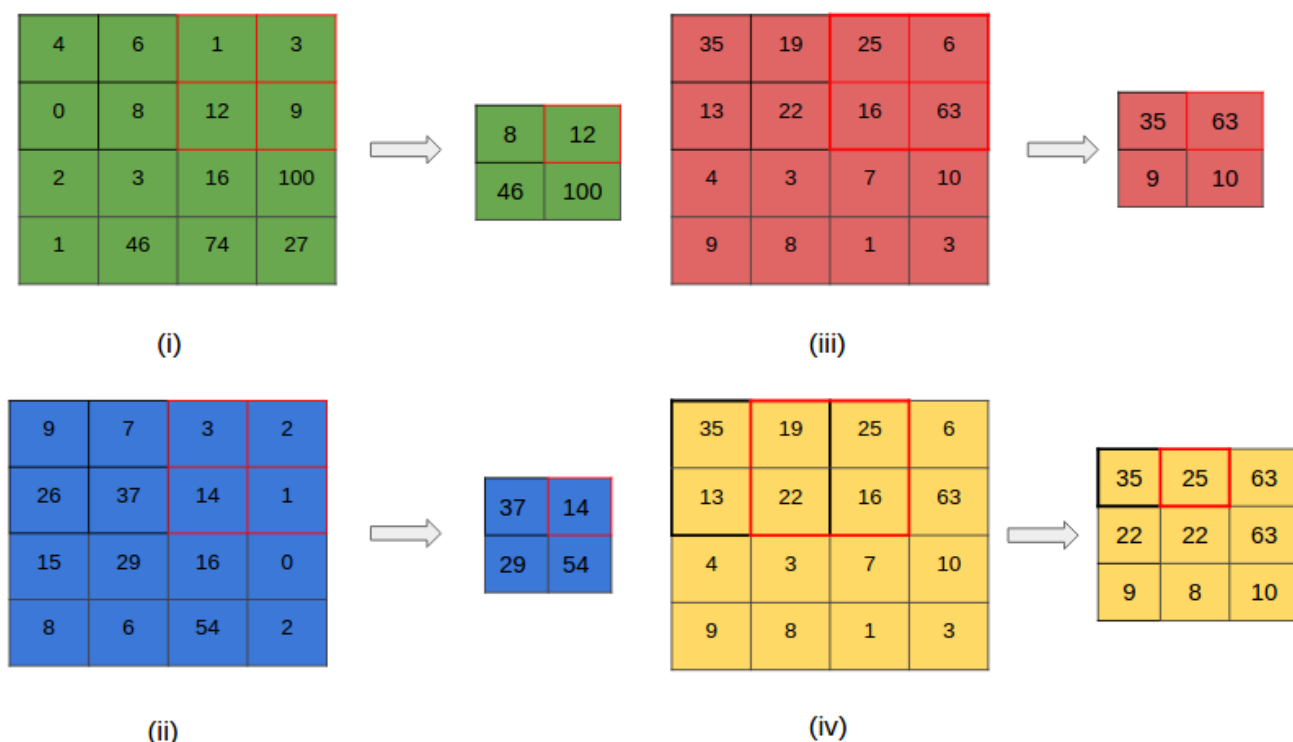


Figure 2.5-4 a max pooling layer

Other popular pooling techniques are average pooling ($p=1$), l_2 pooling ($p=2$) and subsampling. Details about the application of subsampling can be found in Dominik et al. paper [18]. Average pooling was often used historically but has recently fallen out of favor because it gives less competitive results than the others.

A filter size of 2×2 and a stride of length 2 are often applied with these pooling techniques. According to some research [4,9], there is no best pooling technique. We may need to try applying several pooling techniques to our problem to see which one yields the best result.

2.5.4.1 Overlap pooling

A. Krizhevsky et al. has taken a step further to the traditional pooling [1]. Suppose a pooling layer consists of a grid of pooling units, each summarizes a neighborhood of size $z \times z$ and stride s . Traditionally, we set $z = s$, where we obtained non overlap pooling. Now we will set $s < z$ to obtain overlap pooling. A. Krizhevsky et al. has pointed out this is a more effective pooling technique as it reduces the top-1 error rate by 0.4% and top-5 error rate by 0.3%.

2.5.4.2 Exploiting viewpoints in pooling

This approach was introduced by Sander et al. in 2014. After preprocessing and data augmentation, viewpoints are extracted by the combination of flipping, rotating, cropping to the input image. Notice that the combination of those operations is necessary because it reduces redundancy. All viewpoints are presented to the same convolutional architecture. The resulting feature maps from each viewpoint are first concatenated, then processed by a set of fully connected layer to obtain predictions. The benefit of exploiting viewpoints in pooling is that it allows the network to “look at” the image at different angles.

Figure 2.5-5 shows how viewpoints can be exploited in pooling. Different viewpoints are extracted from the original image (step 2,3). Each viewpoint is fed to a separated convolutional architecture (step 4). Results generated by each convolutional architecture are concatenated and feed into dense layers (step 5) to obtain prediction (step 6). [2]

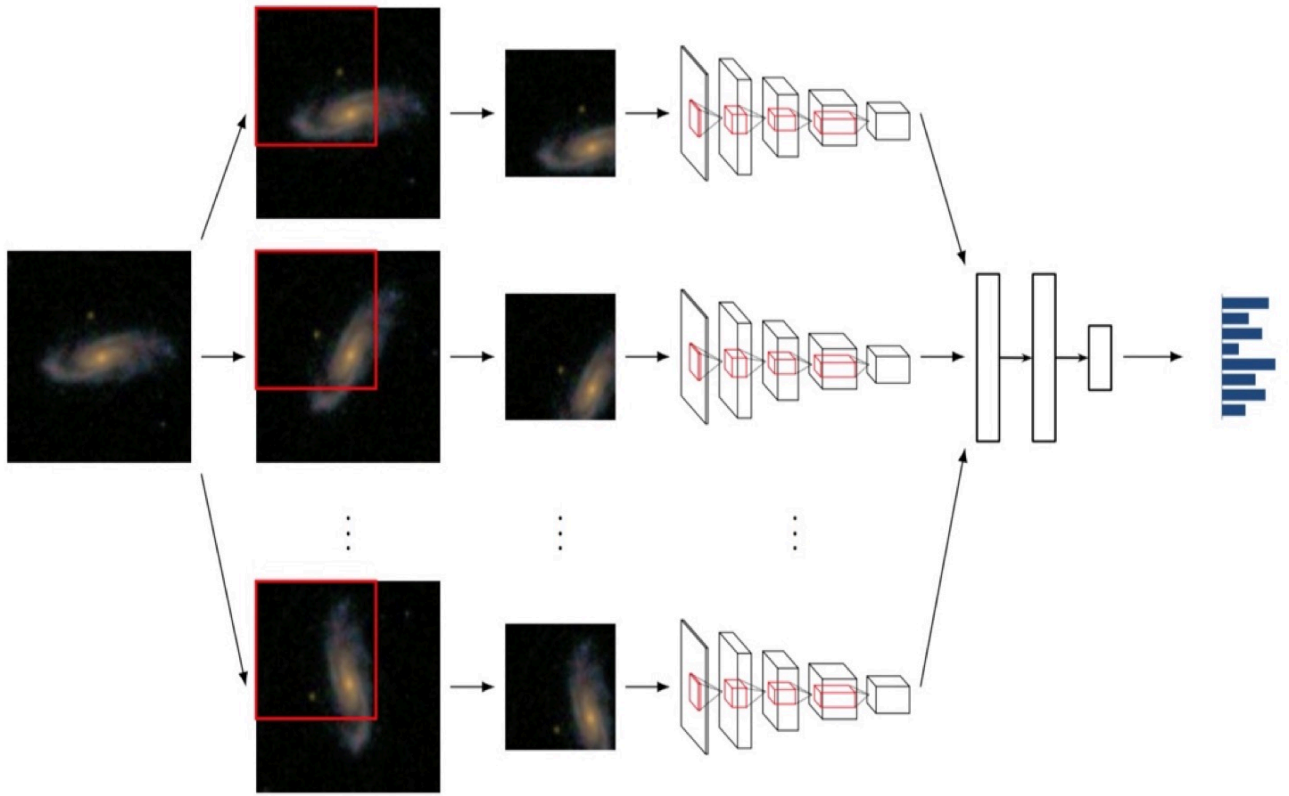


Figure 2.5-5 overview viewpoints exploitation in pooling.

2.5.5 Fully connected layer

Flatten output from the last pooling layer is proceeded by a set of fully connected layers. Fully connected layers are originated from multi-layers feed forward neural networks. More details on MFNN has been discussed in section 2.4.

2.5.6 Softmax regression function

Softmax regression is a generalized logistic regression model which is used to turn a single neuron into a linear classifier so that the neuron can handle multiple classes. For each input a_i^{l-1} , the real output a_j^l can be computed. This real output a_j^l is then used to interpret in which of the K classes does the input a_i^{l-1} belongs

$$a_j^l = P(y_i^{l-1} = k | a_i^{l-1}; w) = \frac{e^{(w_{ij}a_i^{l-1} + b_j^l)}}{\sum_{k=1}^K e^{(w_{ik}a_i^{l-1} + b_k^l)}} = \frac{z_j^l}{\sum_{k=1}^K e^{(w_{ik}a_i^{l-1} + b_k^l)}}$$

Equation 2.5-8

where y_i^{l-1} represents the predicted class k for input a_i^{l-1} , K is the number of classes at output layer, $(w_{ij}a_i^{l-1} + b_j^l)$ is the output produced by activated equation given the input a_i^{l-1} and the weight of its connection to class k , $(w_{ik}a_i^{l-1} + b_k^l)$ is the output produced by activated equation given the input a_i^{l-1} and the weight of its connection to each output class j , where $j = 1, \dots, K$

This function can map an input to an output in a range between 0 and 1, and all calculated probabilities of different classes must sum to 1.

2.6 Learning in a Convolutional neural networks

2.6.1 Back propagation core concepts

Back propagation is an iteration process to adjust the learnt weight and bias by comparing

the network's prediction with the tuple's known target value. It takes place after each training example is presented to the network. The aim is to minimize the error between the network's prediction and the known target value. The target value can either be a known class label or a continuous value. The weight and bias modification process is done in a backward direction, from softmax through fully connected, pooling and ReLU layers to the convolutional layer. The forward and backward process is repeated until all the weights in the network converges. The steps are as follows:

1. Initialize all weights and bias in the network with some small random values and choose a learning rate.
2. Propagate the inputs forward the networks using our initialized weights and bias until we reach the output layer.
3. Calculate error rate and consequently update the weight at each layer.
4. Repeat the process until when one of the following conditions is reached:
 - All Δw in the previous epoch are below some pre-specified threshold
 - A pre-specified number of epochs has reached
 - The percentage of misclassified tuples in the previous run is lower than some predefined threshold

2.6.2 Back propagation at the Softmax

The total error E_{total} (will be refer to later on as E) is the sum of errors for each output neuron and can be calculated using log-loss (sometimes called cross entropy) function:

$$E_{total} = \frac{-1}{N} \left(\sum_{k=1}^N y_k (\log a_k) + (1 - y_k) \log(1 - a_k) \right)$$

Equation 2.6-1

where N is the number of output neurons, y_k is the expected output of neuron k and z_k is the real output of neuron k .

We will look at how the networks adjust its weight using E_{total} by looking at E_{total} effect on each layer.

2.6.3 Back propagation in the Fully connected layer (Neural networks)

The real output a_j^l at fully connected layer can be expressed in terms of net output z_j^l by function () in section

$$a_j^l = \frac{e^{(w_{ij}a_i^{l-1} + b_j^l)}}{\sum_{k=1}^K e^{(w_{ik}a_i^{l-1} + b_k^l)}} = \frac{z_j^l}{\sum_{k=1}^K e^{(w_{ik}a_i^{l-1} + b_k^l)}}$$

Equation 2.6-2

Therefore, we can calculate

$$\frac{\partial a_j^l}{\partial z_j^l} = \frac{\sum_{k=1}^K e^{(w_{ij}a_i^{l-1} + b_j^l)} - z_j^{l^2}}{e^{(w_{ij}a_i^{l-1} + b_j^l)^2}}$$

Equation 2.6-3

Each neuron at the output layer has a net value defined by the following function

$$z_k^l = \sum_{j=1}^m w_{jk} a_j^{l-1} + b_k^l$$

Equation 2.6-4

where z_k^l is the net output of neuron k at layer l , a_j^{l-1} is the real output of neuron j at layer $(l-1)$ where $j = \{1, 2, 3, \dots, m\}$, w_{jk} is the weight of the connection between neuron j and k , b_k^l is the bias value of neuron k .

The amount of error that we want to adjust at each neuron is equal to:

$$\delta_{o_N}^o = \frac{\partial E}{\partial z_{o_N}^o} = \frac{\partial E}{\partial a_{o_N}^o} * \frac{\partial a_{o_N}^o}{\partial z_{o_N}^o}$$

Equation 2.6-5

(as a result of chain rule)

where N is the number of output neurons, o stands for output layer, $\delta_{o_N}^o$ is the error rate at neuron N in the output layer, $z_{o_N}^o$ is the net output of neuron N in the input layer, $a_{o_N}^o$ is the real output of neuron N in the input layer.

We will propagate this error forward. Accordingly, this rate will affect the weight at neuron N by an amount of

$$\Delta w_{jo_N}^o = \frac{\partial E}{\partial w_{jo_N}^o} = \frac{\partial E}{\partial a_{o_N}^o} * \frac{\partial a_{o_N}^o}{\partial z_{o_N}^o} * \frac{\partial z_{o_N}^o}{\partial w_{jo_N}^o}$$

Equation 2.6-6

(as a result of chain rule)

where $\Delta w_{jo_N}^o$ is the amount of weight to adjust between neuron j and neuron o_N , $a_{o_N}^o$ is the real output of neuron o_N at output layer, $z_{o_N}^o$ is the net output of neuron o_N at output layer, $w_{jo_N}^o$ is the current weight between neuron j and neuron o_N .

From function (2.6-5) and (2.6-6), function (2.6-6) can be rewritten as:

$$\Delta w_{jo_N}^o = \frac{\partial E}{\partial w_{jo_N}^o} = \delta_{o_N}^o a_j^{h1}$$

Equation 2.6-7

where $\Delta w_{jo_N}^o$ is the amount of weight to adjust between neuron j and neuron o_N , $w_{jo_N}^o$ is the current weight between neuron j and neuron o_N , $\delta_{o_N}^o$ is the error rate at neuron N in the output layer, a_j^{h1} is the real output of neuron a_j at hidden layer $h1$ where $h1$ precedes o .

To decrease the error, we then subtract this value from the current weight

$$w_{jo_N}^o = w_{jo_N}^o - \eta * \Delta w_{jo_N}^o$$

Equation 2.6-8

where η denotes the learning rate, which is normally set to 0.5

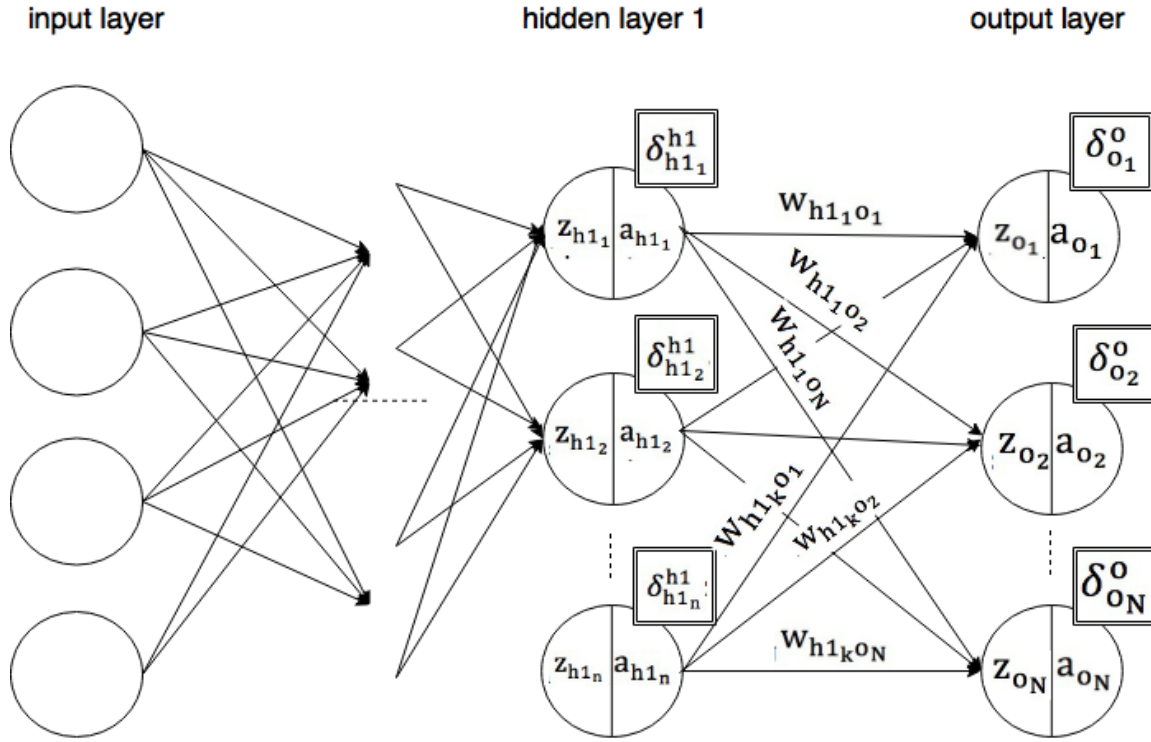


Figure 2.6-1 fully connected layer

$$\delta_{o_N}^o z_{o_N} w_{h1_k o_1}$$

For every neuron at each hidden and input layer of the neural networks, we can calculate δ and Δw by applying the same technique. Thus our calculation can be generalized as follows:

$$\delta_{z_j}^l = \frac{\partial E}{\partial a_j^l} * \frac{\partial a_j^l}{\partial z_j^l} = \frac{\partial E}{\partial z_i^{l+1}} * \frac{\partial z_i^{l+1}}{\partial a_j^l} * \frac{\partial a_j^l}{\partial z_j^l} = \left(\sum_{0 \leq i' \leq m} \delta_{z_{i'}}^{l+1} * w_{ji'}^{l+1} \right) * \frac{\partial a_j^l}{\partial z_j^l}$$

Equation 2.6-9

where $\delta_{z_j}^l$ is the error rate at neuron j in layer l , z_j^l is the net output of neuron j in layer l , a_j^l is the real output of neuron j in layer l , z_i^{l+1} represents the total net output of all neurons in layer $(l+1)$ that are pairwise connected with neuron j , m is the total number of neurons in layer $(l+1)$, $\delta_{z_{i'}}^{l+1}$ is the error rate at neuron i' in layer $(l+1)$ and $w_{ji'}^{l+1}$ is the weight between neuron j and neuron i' where $i' = \{0, 1, \dots, m\}$, as illustrated by figure 2.6-1.

We will propagate this error forward

$$\Delta w_{ij}^l = \delta_{z_j}^l a_i^{l-1}$$

Equation 2.6-10

where Δw_{ij}^l is the amount of weight to adjust between neuron i and neuron j , δ_j^l is the error rate at neuron j in layer l , a_i^{l-1} is the real output of neuron i in layer $(l-1)$.

2.6.4 Back propagation in the ReLU layer

The output a_j^l for each neuron can be written as

$$a_j^l = \max(0, a_i^{l-1})$$

Equation 2.6-11

where a_i^{l-1} is the real output value of neuron i of ReLU layer $(l-1)$. Thus for each output at layer ReLU layer, we want to adjust an error amount equivalent to

$$\delta_{a_i}^{l-1} = \frac{\partial E}{\partial a_i^l} = \begin{cases} = 0, \text{ if } (a_i^{l-1} \leq 0) \\ \frac{\partial E}{\partial a_j^l} * \frac{\partial a_j^l}{\partial a_i^{l-1}} = \frac{\partial E}{\partial a_j^l} = \delta_{a_i}^l \text{ otherwise} \\ \text{(since } a_j^l = a_i^{l-1} \text{ in this case)} \end{cases}$$

Equation 2.6-12

We need not to adjust weight at this layer because it has no weight.

2.6.5 Back propagation in the pooling layer

The output a_j^l for each neuron can be represented by

$$a_{x,y}^l = \max_{0 \leq p,q \leq k} (a_{x+p,y+q}^{l-1})$$

Equation 2.6-13

where $a_{x,y}^l$ is the the real output value of neuron at position (x,y) of max pooling layer l, $a_{x+p,y+q}^{l-1}$ is the real output value of neuron at position (x+p,y+q) of convolutional layer (l-1) where (p,q) is the position in a kernel of size kxk.

For each output pixel at max pooling layer, we want to adjust an amount of

$$\delta_{a_{x+p,y+q}}^{l-1} = \frac{\partial E}{\partial a_{x+p,y+q}^{l-1}} = \begin{cases} = 0, \text{ if } (a_{xy}^l \neq a_{x+p,y+q}^{l-1}) \\ \frac{\partial E}{\partial a_{xy}^l} = \delta_{a_{xy}}^l \text{ otherwise} \end{cases}$$

Equation 2.6-14

We need not to adjust weight at this layer because it also has no weight.

2.6.6 Back propagation in the convolutional layer

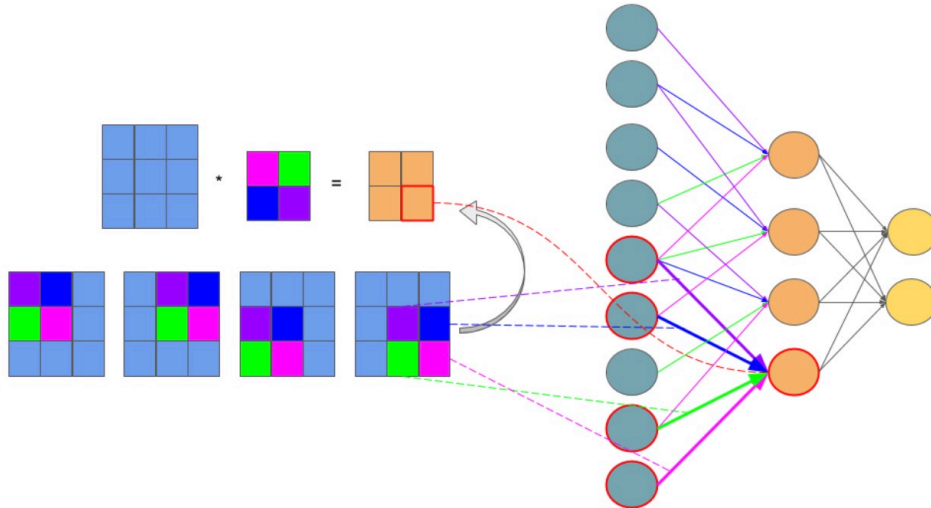


Figure 2.6-2 feedforward in CNN is identical with convolution operation. [15]
The net output $z_{x,y}^{l+1}$ of each neuron at convolutional layer can be represented by

$$z_{x,y}^l = \sum_a \sum_b w_{a,b}^l a_{x-a,y-b}^{l-1} + b_{x,y}^l$$

Equation 2.6-15

as discussed in section 2.5.2

The amount of error that we want to adjust at each neuron is equal to:

$$\delta_{z_{x,y}}^l = \frac{\partial E}{\partial z_{x,y}^l} = \frac{\partial E}{\partial a_{x,y}^l} * \frac{\partial a_{x,y}^l}{\partial z_{x,y}^l} = \frac{\partial E}{\partial z_{i,j}^{l+1}} * \frac{\partial z_{i,j}^{l+1}}{\partial a_{x,y}^l} * \frac{\partial a_{x,y}^l}{\partial z_{x,y}^l}$$

Equation 2.6-16

(as discussed in fully connected layer)

function (2.6-16) is thus equivalent to

$$\delta_{z_{x,y}}^l = \sum_{i'} \sum_{j'} \left(\frac{\partial E}{\partial z_{i',j'}^{l+1}} * \frac{\partial z_{i',j'}^{l+1}}{\partial z_{i',j'}^l} * \frac{\partial a_{i',j'}^l}{\partial z_{i',j'}^l} \right) = \sum_{i'} \sum_{j'} \left(\delta_{i',j'}^{l+1} * w_{(x,y),(i',j')}^{l+1} * \frac{\partial a_{i',j'}^l}{\partial z_{i',j'}^l} \right)$$

Equation 2.6-17

where $z_{x,y}^l$ is net output of pixel at position (x,y) at layer l, $a_{x,y}^l$ is real output of pixel at position (x,y) at layer l, $z_{i,j}^{l+1}$ represents the total net output of all pixels at layer (l+1) that are pairwise connected with pixel at position (x,y), $\delta_{z_{x,y}}^l$ is the error rate of pixel at position (x,y) at layer l, $\delta_{i,j}^{l+1}$ is the error rate of pixel at position (x,y) at layer (l+1), $w_{i,j}^{l+1}$ is the weight between (i',j') and (x,y), $z_{i',j'}^l$ is the net output of pixel (i',j'), $a_{i',j'}^l$ is the real output of pixel (i',j').

This error is contributed into the current weight as follows:

$$\frac{\partial E}{\partial w_{(a,b),(x,y)}^l} = \sum_x \sum_y \frac{\partial E}{\partial a_{x,y}^l} * \frac{\partial a_{x,y}^l}{\partial z_{x,y}^l} * \frac{\partial z_{x,y}^l}{\partial w_{(a,b),(x,y)}^l} = \sum_x \sum_y \delta_{z_{x,y}}^l * \sigma'(z_{x-a,y-b}^{l-1})$$

Equation 2.6-18

where $w_{(a,b),(x,y)}^l$ is the weight between pixel (a,b) and pixel (x,y), $z_{x,y}^l$ is net output of pixel (x,y) at layer l, $a_{x,y}^l$ is real output of pixel (x,y), $\delta_{z_{x,y}}^l$ is the error rate of pixel (x,y) at layer l.

2.6.7 What does a CNN learn?

In this section we will examine the final CNN and how it represents learning: as adapted filters at convolutional layers and weights at fully connected layers. We will then discuss about the use of learnt filters and weights.

Visualization of the adapted weights at a convolutional layer once training is complete can be illustrated as blocks of multiple nxn filters (figure 2.6-1). More weights are represented by more color intensive blocks.

Each filter can detect a feature at different locations of the same input image, which is a useful thing since a sugarcane image is translationally invariant. By learning a variety of selective filters, the networks improve its ability to detect robust features. Vice versa, examining the adjusted filters can help us visualize what kind of features are being extracted. For example, in figure 2.6-2, the filter at row 1, cols 1 extracts the discriminative parts of the image, as shown in the the corresponding image patch. The visualization in figure also shows that the extracted features are strongly grouped within each kernel. Interestingly, by studying these visualizations, we can debug the problems with the model (for instance: Are the kernels informative enough?) and improve our results accordingly.

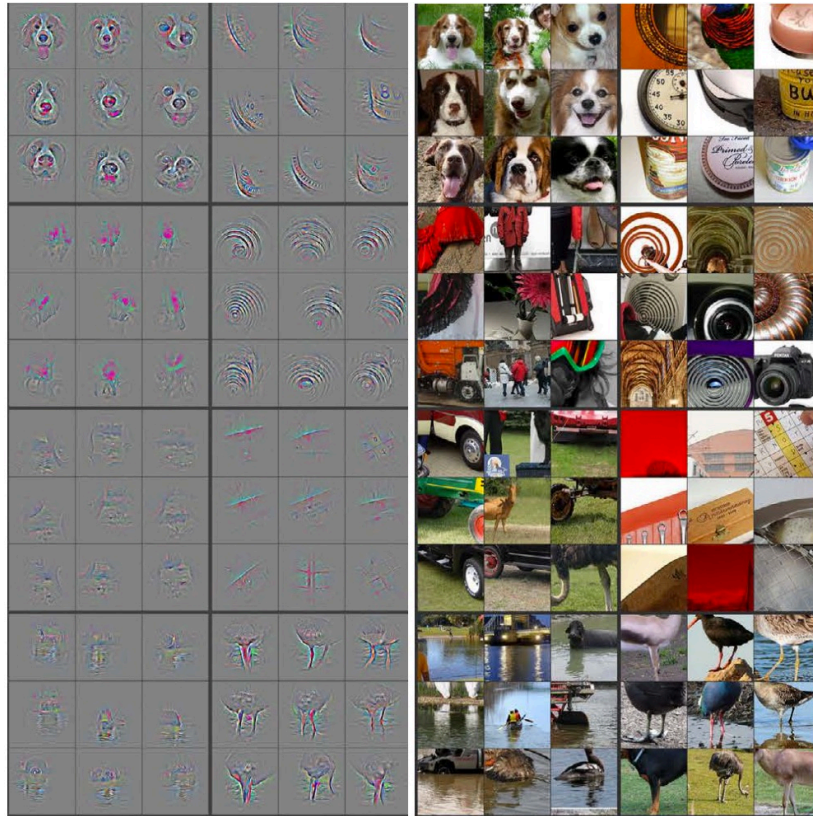


Figure 2.6-3 Visualization of features in a fully trained model.[17]

Learning at fully connected layers on the other hand is represented with adapted weights at each connection. Those weights are the networks' reference when mapping features representing each object with the output class. Therefore, the more precise the weight the more precise the model prediction.

2.7 Data Augmentation

2.7.1 Concept

Artificially enlarge the dataset with label-preserving transformations is the simplest and the most known strategy to reduce overfitting on image data. To each randomly chose sample image, we will apply n random transformations. Each of these random transformations is a combination of several elementary forms transformation, which we will describe shortly. The benefit of using little computation is that we will not have to store the pre-processed image on the disk. Image transformation contains two major approaches: point operators (sometimes called 1-to-1 pixel transforms) and neighborhood (or area-based) operators. In point operators, each output pixel value is strictly a function of the corresponding input pixel value. Brightness and contrast adjustments are two examples of such transformation. Techniques like convolution, which we have discussed in part 2.2, are not 1-to-1 transform.

2.7.2 Methods

2.7.2.1 Brightness adjustment

We add or subtract a constant amount of light to all input pixel in order to change an image brightness, as suggested in the following function:

$$g(i,j) = f(i,j) + \beta$$

Equation 2.7-1

Where $f(i,j)$ is the pixel located in the i -th row and j -th column of the input image, $g(i,j)$ is the

pixel located in the i -th row and j -th column of the output image and β is a bias parameter which is used to control the image brightness.

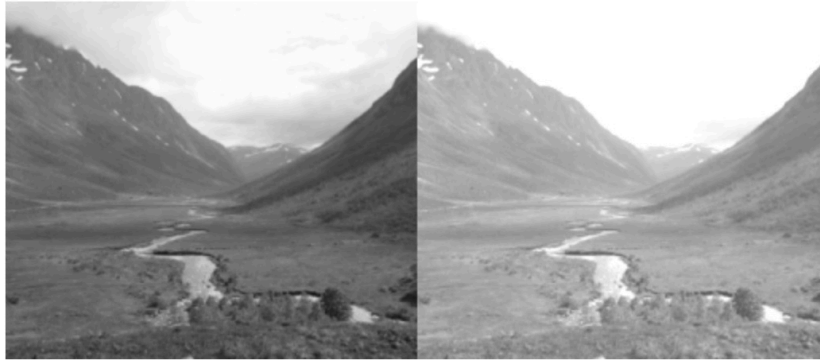


Figure 2.7-1 brightness adjustment. [14]

2.7.2.2 Contrast adjustment

Contrast is the difference between the maximum and the minimum pixel intensity of an image. To change the contrast of an image, we change the range of the luminance value presented in the input pixels. It is mathematically suggested by the following function:

$$g(i,j) = \alpha f(i,j)$$

Equation 2.7-2

Where $f(i,j)$ is the pixel located in the i -th row and j -th column of the input image, $g(i,j)$ is the pixel located in the i -th row and j -th column of the output image and α is a weight parameter which is used to control the image contrast.



Figure 2.7-2 contrast adjustment. [14]

We can combine the brightness and contrast control in a single operation. It is mathematically represented as:

$$g(i,j) = \alpha f(i,j) + \beta$$

Equation 2.7-3

2.7.2.3 Scaling

Scaling can be done by multiplying the scale of the patch by a factor x . Artificially enlarge the dataset using the combination of these three forms of image translation in image pre-processing has two extra benefits. Firstly, it allows the network to “look at” the seed image at different perspectives. Secondly, it allows the CNN model to

work properly, because CNN often require to be fed with a considerably large dataset. Notice that data augmentation is different from image pre-processing.

2.8 Data post-processing

Combination of models is a very efficient way to reduce testing errors. At the same time, it significantly increases the number of networks parameters and thus computers processing speed. Dropout is a very efficient method for model combination. Dropout sets the output of each hidden layer to zero with the probability of 0.5. Propagation and back propagation processes will not consider the 'dropout neurons' into their calculations. Dropout reduces the existence of co-adaptations neurons: A neuron is forced to learn more ruggedized features to make predictions because it cannot rely on the existence of other neurons.

2.9 GPU Concepts

The enormous amount of input data of deep learning may considerably reduce computation time. This problem can be overcome by spreading and training the networks across processing units. There are two types of processing units: central processing unit (CPU) and graphics processing unit (GPUs).

In compare with CPU, GPUs are far more powerful and efficient in parallel computing. They can be used to train far larger training sets in considerably less time. Current GPUs allow cross-GPU parallelization to be read and write into one another's memory directly. Even though each GPU has limited memory, which may restrict the size of networks to be trained on one, it can still distribute the networks size cross one another.

Recent deep learning toolkits are mostly develop based on CUDA library. This library only supports NVIDIA GPU card with compute capability ≥ 3.0 .

2.10 CNN Software Frameworks

<i>Framework</i>	<i>Platform</i>	<i>Base language</i>	<i>API</i>	<i>Parallel support</i>	<i>Computational graph and automatic differentiation</i>	<i>Single GPU execution speed</i>	<i>Ready-to-use low-level operators for writing new models</i>	<i>GPU memory for training large models</i>	<i>Ease to use</i>
<i>Tensorflow</i>	<i>Linux, Mac OS X</i>	<i>Python and C++</i>	<i>Python and C/C++ Mathematical operations are supported with numpy</i>	<i>Use CUDA library to support multi-GPU</i>	<i>Yes. However, the use of pure Python slows down the computational graph speed</i>	<i>Slower than other frameworks</i>	<i>Fairly good</i>	<i>Not so good</i>	<i>Sample code and tutorials available Error messages are difficult to understand, therefore can be unhelpful</i>
<i>Caffe</i>	<i>Ubuntu, OS X, AWS, unofficial Android port, Windows support by Microsoft Research, unofficial Windows port</i>	<i>C++, Python</i>	<i>C++, command line, Python, MATLAB</i>	<i>Use CUDA library to support multi-GPU Needs to write C++ / CUDA for implementing new GPU layers</i>	<i>Yes</i>	<i>Slower than other frameworks</i>	<i>Fairly good</i>	<i>Better than Torch</i>	<i>Sample code and to tutorials are somewhat confusing because different versions are developed by different people Code might not need to be written to train models</i>

<i>Theano</i>	<i>Cross-platform</i>	<i>Python</i>	<i>Python Mathematical operations are supported with numpy</i>	<i>Use CUDA library to support experimental multi-GPU</i>	<i>Yes</i>	<i>Takes considerably long time to train a large model</i>	<i>Many basic operations</i>	<i>Great</i>	<i>Sample code and tutorials available</i>
<i>MXNET</i>	<i>Ubuntu, OS X, Windows, AWS, Android, iOS, JavaScript</i>	<i>C++, Python, Julia, Matlab, R, Scala</i>	<i>C++, Python, Julia, Matlab, JavaScript, R, Scala</i>	<i>Distributed computing</i>	<i>Yes</i>	<i>Takes considerably long time to train a large model</i>	<i>Very few</i>	<i>Excellent</i>	<i>Sample code and to tutorials are somewhat confusing because different versions are developed by different people</i>
<i>Torch</i>	<i>Linux, Android, Mac OS X, iOS, Windows</i>	<i>Lua</i>	<i>Lua</i>	<i>Yes</i>	<i>Yes</i>	<i>Competitive with Theano</i>	<i>Many basic operations</i>	<i>Fair</i>	<i>Easy to set up Error messages are helpful</i>

2.11 Related research

2.11.1 Researches using CNN to classify images

CNN has proved to be very efficient in object recognition. A. Krizhesy et al. (2012) successfully modeled a deep CNN “to classify the 1.2 million high-resolution images in the ImageNet [...] into the 1000 different classes” with “top-1 and top-5 error rates of 37.5% and 17.0%”. The networks contained five convolutional layers, three max pooling layers, three fully connected layers and a 1000-way softmax. Data argumentation, overlap pooling and dropout were used to reduce data overfitting. To improve the processing speed, the networks were split into different parts, which were trained on multiple GPUs [1]. Hokuto et al. (2014) developed a food detection and recognition deep CNN trained from 20,000 samples of food items. The networks consisted of two convolutional layers, one ReLU layer and was able to detect up to 93.8% and recognize up to 72.39% of testing food items [5].

2.11.2 Variations in algorithms in CNN used for classifying fine-grained objects

Jonathan et al. (2014) exploited deep CNN in fine-grained object recognition. The objective was to identify different car models. The networks were built by adapting the model from A. Krizhesy et al. with little variation. A deep CNN model consists of two convolutional layers, three fully connected layers and a softmax loss were used to extract useful features from the seed image. The seed image then is used to retrieve its nearest neighbors (those that has the same pose with it). Only parts with highest energy detected from this seed image and its neighbors are chosen because they are likely to be important when describing the seed image. Learnt features obtained earlier from CNN are then pooled in the regions of each selected part. Those regions are said to describe critical parts for the class where seed image belongs. This networks were able to categorize testing data with accuracy up to 73.9%. However, the described technique is only useful for recognizing pictures with the same pose [4]. Aäron et al (2015) attempted to build a CNN to classify grayscale 30000 images of plankton into one of 121 classes. Input images were preprocessed in by downsampling. 300 CNN models were trained and several of them were combined to improve the final accuracy. A significant variation made when training the model was to exploit the viewpoints in pooling. One example of the models that works well consisted of 10 convolutional layers, 4 pooling layers and 3 fully connected layers. Models were trained on the NVIDA GPUs using various overfitting-reducing techniques: data argumentation, ReLU, dropout and overlap pooling. The best models had an accuracy of 82% on the testing set and top 5% accuracy of 98%[3]. Sander et al. (2015) built a deep CNN to measure approximately 900,000 galaxy morphology by exploiting viewpoints in pooling. The model was claimed to “reproduce their consensus with near-perfect accuracy (> 99%) for most questions” in Galaxy Challenge Contest 2015. Techniques discussed in A. Krizhesy et al. were also applied to reduce data overfitting [2].

Reference

- [1] A. Krizhevsky, I. Sutskever, G.E. Hinton. *ImageNet classification with deep convolutional neural networks*. In *Advances in Neural Information Processing Systems 25 (NIPS'2012)*, 2012
- [2] Sander Dieleman, Kyle W. Willett and Joni Dambre. *Rotation-invariant convolutional neural networks for galaxy morphology prediction*. *Mon. Not. R. Astron. Soc.* 000, 1–20 (2014).
- [3] Aäron van den Oord, Ira Korshunova, Jeroen Burms, Jonas Degraeve, Lionel Pigou, Pieter Buteneers. *Classifying plankton with deep neural networks*. First prize of The National Data Science Bowl competition, March 2015.
- [4] Jonathan Krause, Timnit Gebru, Jia Deng, Li-Jia Li, Li Fei-Fei. ICPR, 2014, supported by an ONR MURI grant and the Yahoo! FREP program. *Learning Features and Parts for Fine-Grained Recognition*.
- [5] Hokuto Kagaya, Kiyoharu Aizawa, Makoto Ogawa. MM'14, November 3–7, 2014. *Food Detection and Recognition Using Convolutional Neural Network*
- [6] Evgeny A. Smirnov*, Denis M. Timoshenko, Serge N. Andrianov. *Comparison of Regularization Methods for ImageNet Classification with Deep Convolutional Neural Networks*. 2013 2nd AASRI Conference on Computational Intelligence and Bioinformatics
- [7] Yoshua Bengio, Ian Goodfellow, and Aaron Courville. *Draft book in preparation*. Deep learning, Jan 2016
- [8] Jiawei Han & Micheline Kamber. *Data Mining Concepts and Techniques*, 3rd edition, 2012
- [9] Jonathan Sachs. 1996-1999. *Digital Light & Color*
- [10] Rafael C. Gonzalez, Paul Wintz. *Digital Image Processing*, 2nd edition, 1987
- [11] Wikipedia. *RGB color model*
- [12] Stanford University. *Convolutional neural networks for visual recognition*
- [13] Mathworks. *Convolutional neural networks documentation*
- [14] Pippin online tutorial
- [15] Grzegorzgwardys. *Convolutional Neural Networks*
- [16] Eindhoven University of Technology – PARsE. *Convolutional Neural Networks*
- [17] Matthew D. Zeiler, Rob Fergus. Nov 2013. *Visualizing and Understanding Convolutional Networks*.
- [18] Dominik Scherer, Andreas Müller, and Sven Behnke. 20th International Conference on Artificial Neural Networks (ICANN), Thessaloniki, Greece, September 2010. *Evaluation of Pooling Operations in Convolutional Architectures for Object Recognition*.