

Deep learning

Unsupervised Feature Learning and Deep Learning ¹

LING KANGJIE²

September 26, 2017

¹<http://ufldl.stanford.edu/tutorial/>

²lingkangjie@Foxmail.com

Dedicated to somebody.

Contents

1	Introduction	3
1.1	Linear Regression	3
1.1.1	Problem Formulation and Function Minimization	3
1.1.2	MATLAB Implementation	4
1.2	Logistic Regression	6
1.2.1	Matlab Implementation	7
1.3	Debugging: Gradient Checking	8
1.3.1	Matlab Implementation	8
1.4	Softmax Regression	9
1.4.1	Cost Function and Its properties	9
1.4.2	Matlab Implementation	10
1.5	Neural Network and Backpropagation Algorithm	11
1.5.1	BackPropogation Implementation Using Matlab	13
1.6	Convolutional Neural Networks and pooing	16
1.6.1	Implement convolution and Pooling	17
1.7	Build up Your Own Convolutional Neural Network	18
2	Unsupervised Learning	27
2.1	Autoencoders	27
2.1.1	Sparse Autoencoder Implementation	28
2.2	PCA, PCA Whitening and ZCA Whitening on Images	33
2.3	Reconstruction ICA	39

List of Figures

List of Tables

1.1	The variable space	19
1.2	The variable space	23
2.1	The variable space	28
2.2	The variable space	30

Preface

This tutorial will teach you the main ideas of Unsupervised Feature Learning and Deep Learning. By working through it, you will also get to implement several feature learning/deep learning algorithms, get to see them work for yourself, and learn how to apply/adapt these ideas to new problems. This tutorial assumes a basic knowledge of machine learning (specifically, familiarity with the ideas of supervised learning, logistic regression, gradient descent). If you are not familiar with these ideas, we suggest you go to this Machine Learning course¹ and complete sections II, III, IV (up to Logistic Regression) first.

Supervised Learning and Optimization

- Linear Regression
- Logistic Regression
- Debugging: Gradient Checking
- Softmax Regression
- Debugging methods

Convolutional Neural Network

- Feature Extraction Using Convolution
- Pooling And Downsample
- Networks Architecture
- Classical CNNs' Models

Unsupervised learning

- Autoencoders
- PCA whitening
- Sparse Coding

¹<https://goo.gl/hSZRFq>

- Restricted Boltzmann Machine
- Deep Belief Networks

Reinforcement learning

In this section, it would cover the following problems:

- Introduction And Basic Equations
- Monte Carlo sampling
- Deep Q networks

Acknowledgements

Thanks my father, my mother and other person supported me in the past.

Lingkangjie

[`https://space.bilibili.com/8661791/#!/`](https://space.bilibili.com/8661791/#!/)

1

Introduction

“What I cannot create, I do not understand.”

– Richard Philip Feynman, 1918-1988

1.1 Linear Regression

1.1.1 Problem Formulation and Function Minimization

As a refresher, we will start by learning how to implement linear regression. We will introduce two main concepts: the procedure of machine learning and function optimization. Our goal in linear regression is to predict a target value \mathbf{y} starting from a vector of input values $\mathbf{x} \in \mathbb{R}^n$. For example, we might want to make predictions about the price of a house so that \mathbf{y} represents the price of the house in dollars and the elements \mathbf{x}_j of \mathbf{x} represent “features” that describe the house (such as its size and the number of bedrooms). Suppose that we are given many examples of houses where the features for the i 'th house are denoted $\mathbf{x}^{(i)}$ and the price is $\mathbf{y}^{(i)}$. For short, we will denote our goal is to find a function $\mathbf{y} = \mathbf{h}(\mathbf{x})$ so that we have $\mathbf{y}^{(i)} \approx \mathbf{h}(\mathbf{x}^{(i)})$ for each training example. Here, $\mathbf{h}(\mathbf{x}^{(i)})$ is our hypothesis.

For linear regression, we have $\mathbf{h}_\theta(\mathbf{x}) = \sum_j \theta_j \mathbf{x}_j = \boldsymbol{\theta}^T \mathbf{x}$. Here, $\mathbf{h}_\theta(\mathbf{x})$ represents a large family of functions parametrized by the choice of $\boldsymbol{\theta}$ parameter vector. (We call this space of functions a **hypothesis class**). In particular, we want to search for a choice of $\boldsymbol{\theta}$ that minimizes the cost function:

$$\mathbf{J}(\boldsymbol{\theta}) = \frac{1}{2} \sum_i (\mathbf{h}(\mathbf{x}^{(i)}) - \mathbf{y}^{(i)})^2 = \frac{1}{2} \sum_i (\boldsymbol{\theta}^T \mathbf{x}^{(i)} - \mathbf{y}^{(i)})^2$$

We use gradient descent method to optimize the cost function. The gradient of $\mathbf{J}(\boldsymbol{\theta})$ is $\nabla_{\boldsymbol{\theta}} \mathbf{J}(\boldsymbol{\theta})$, which is a vector:

$$\nabla_{\boldsymbol{\theta}} \mathbf{J}(\boldsymbol{\theta}) = \begin{bmatrix} \frac{\partial \mathbf{J}(\boldsymbol{\theta})}{\partial \theta_1} \\ \frac{\partial \mathbf{J}(\boldsymbol{\theta})}{\partial \theta_2} \\ \frac{\partial \mathbf{J}(\boldsymbol{\theta})}{\partial \theta_3} \\ \vdots \\ \frac{\partial \mathbf{J}(\boldsymbol{\theta})}{\partial \theta_n} \end{bmatrix}$$

Differentiating the cost function $\mathbf{J}(\boldsymbol{\theta})$ as given above with respect to a particular parameter θ_j gives us:

$$\frac{\partial \mathbf{J}(\boldsymbol{\theta})}{\partial \theta_j} = \sum_i \mathbf{x}_j^{(i)} (\mathbf{h}(\mathbf{x}^{(i)}) - \mathbf{y}^{(i)})$$

1.1.2 MATLAB Implementation

The file *ex1a_linreg.m* is to split and shuffle data, call *linear_regression.m* function and optimize it through minFunc toolbox, and plot the results, just as follow:

```

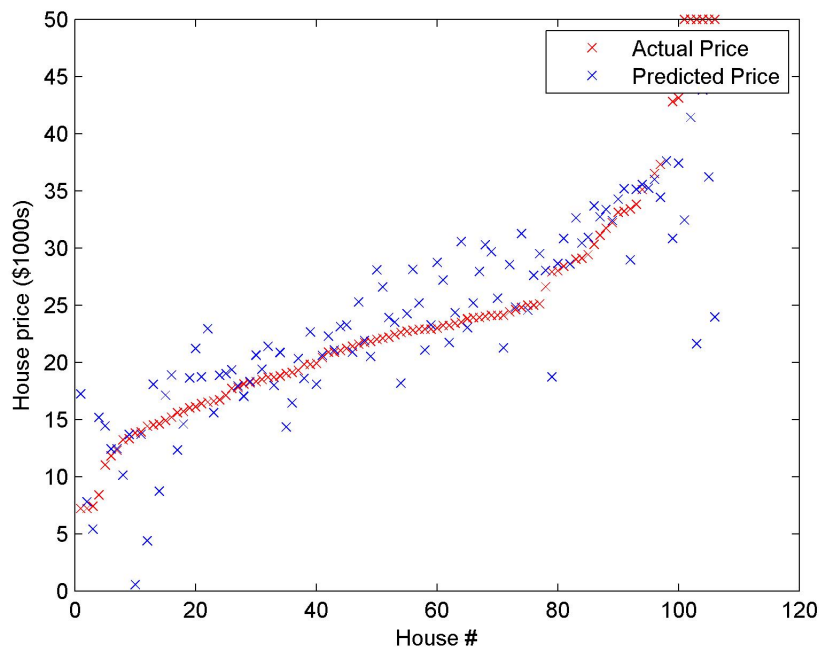
1  % the "housing.data" coming from http://archive.ics.uci.edu/ml.
2
3  clc;
4  clear;
5  close all;
6  addpath ../common
7  addpath ../common/minFunc_2012/minFunc
8  addpath ../common/minFunc_2012/minFunc/compiled
9
10 % Load housing data from file.
11 data = load('housing.data');
12 data=data'; % put examples in columns
13
14 % Include a row of 1s as an additional intercept feature.
15 data = [ ones(1,size(data,2)); data ];
16
17 % Shuffle examples.
18 data = data(:, randperm(size(data,2)));
19
20 % Split into train and test sets
21 % The last row of 'data' is the median home price.
22 train.X = data(1:end-1,1:400);
23 train.y = data(end,1:400);
24
25 test.X = data(1:end-1,401:end);
26 test.y = data(end,401:end);
27
28 [features_n,samples_n] = size(train.X);
29
30 % Initialize the coefficient vector theta to random values.
31 theta = rand(features_n,1);
32
33 % Run the minFunc optimizer with linear_regression.m as the objective.
34 tic;
35 options = struct('MaxIter', 200);
36 theta = minFunc(@linear_regression, theta, options, train.X, train.y);
37 fprintf('Optimization took %f seconds.\n', toc);
38
39 % Plot predicted prices and actual prices from training set.

```

```

40 actual_prices = train.y;
41 predicted_prices = theta'*train.X;
42
43 % Print out root-mean-squared (RMS) training error.
44 train_rms=sqrt(mean((predicted_prices - actual_prices).^2));
45 fprintf('RMS training error: %f\n', train_rms);
46
47 % Print out test RMS error
48 actual_prices = test.y;
49 predicted_prices = theta'*test.X;
50 test_rms=sqrt(mean((predicted_prices - actual_prices).^2));
51 fprintf('RMS testing error: %f\n', test_rms);
52
53 % Plot predictions on test data.
54 plot_prices=true;
55 if (plot_prices)
56     [actual_prices,I] = sort(actual_prices);
57     predicted_prices=predicted_prices(I);
58     plot(actual_prices, 'rx');
59     hold on;
60     plot(predicted_prices,'bx');
61     legend('Actual Price', 'Predicted Price');
62     xlabel('House #');
63     ylabel('House price (1000 dollars)');
64 end

```

The *linear_regression.m*

file is:

```

1 function [f,g] = linear_regression(theta, X,y)
2 %
3 % Arguments:

```

```

4      %   theta - A vector containing the parameter values to optimize.
5      %   X - The examples stored in a matrix.
6      %       X(i,j) is the i'th coordinate of the j'th example.
7      %   y - The target value for each example. y(j) is the target for exam
8      %
9
10     m=size(X,2);
11     n=size(X,1);
12
13     f=0;
14     g=zeros(size(theta));
15     % f: cost function sum , g:gradient sum
16     for row = 1:m
17         curr_x = X(:,row);
18         curr_y = y(:,row);
19         f = f + (1/2).*(theta'*curr_x-curr_y)^2;
20         g = g + curr_x.*(theta'*curr_x-curr_y);
21     end
22     % another version:vectorization
23     %   h = theta' * X - y;
24     %   f = (1/2) * h * h';
25     %   g = X*h';

```

1.2 Logistic Regression

In this section, We will use logistic regression to solve a classification problem. In linear regression, we tried to predict the value of $\mathbf{y}^{(i)}$ for the i 'th example $\mathbf{x}^{(i)}$ using a linear function $\mathbf{y} = \mathbf{h}_{\theta}(\mathbf{x}) = \boldsymbol{\theta}^T \mathbf{x}$. This is a continuous-valued problem. For 2-class problems, we use logistic regression which outputs denote the probability for each class as the form:

$$P(\mathbf{y} = 1|\mathbf{x}) = \mathbf{h}_{\theta}(\mathbf{x}) = \frac{1}{1 + \exp(-\boldsymbol{\theta}^T \mathbf{x})} = \sigma(\boldsymbol{\theta}^T \mathbf{x})$$

$$P(\mathbf{y} = 0|\mathbf{x}) = 1 - \mathbf{h}_{\theta}(\mathbf{x})$$

We use cross entropy to design the cost function:

$$\mathbf{J}(\boldsymbol{\theta}) = - \sum_i (\mathbf{y}^{(i)} \log(\mathbf{h}_{\theta}(\mathbf{x}^{(i)})) + (1 - \mathbf{y}^{(i)}) \log(1 - \mathbf{h}_{\theta}(\mathbf{x}^{(i)})))$$

The gradient of $\mathbf{J}(\boldsymbol{\theta})$ and written in its vector form is:

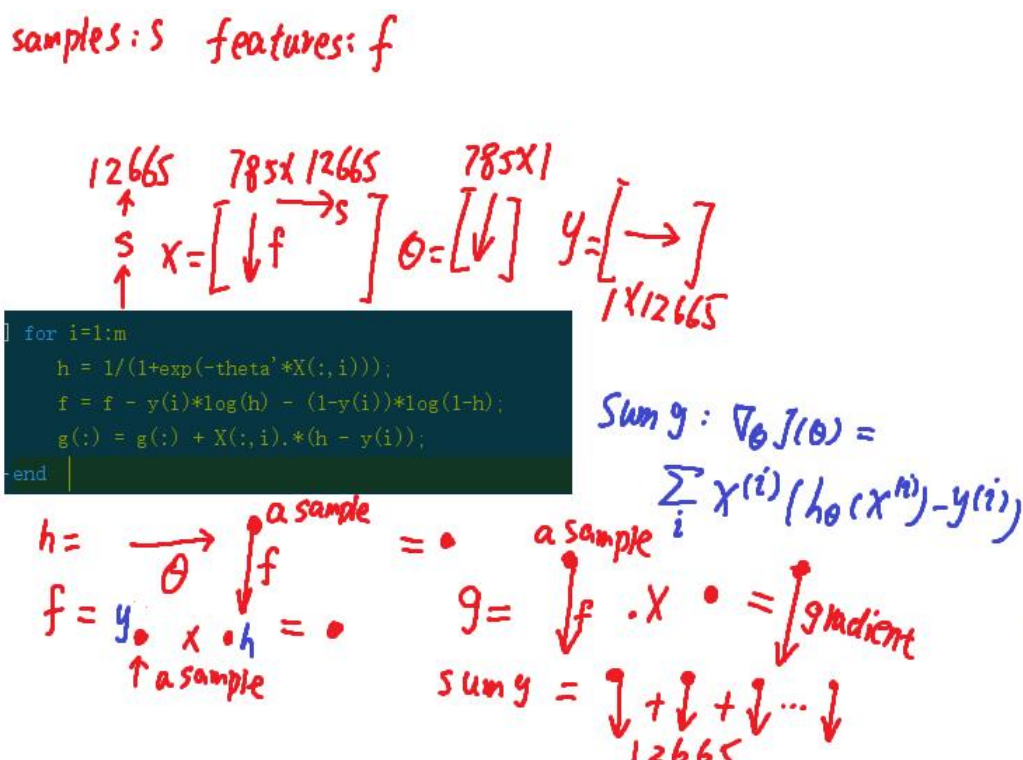
$$\nabla_{\boldsymbol{\theta}} \mathbf{J}(\boldsymbol{\theta}) = \sum_i \mathbf{x}^{(i)} (\mathbf{h}(\mathbf{x}^{(i)}) - \mathbf{y}^{(i)})$$

This is essentially the same as the gradient for linear regression except that the hypothesis function.

1.2.1 Matlab Implementation

There will involve 5 files:

1. *ex1b_logreg.m*: the main file, to optimize the cost function, train and test the model.
 2. *loadMNISTImages.m* *ex1_load_minst.m*: *ex1_load_minst.m* call *loadMNISTImages.m* function to load MNIST data, return train and test data set. Data form description: train data, a structure, containing 2 constants, X with 784*12665 and y with 1*12665, in which 12665 is train sample numbers. Test data set have the same form with 2115 train samples.
 3. *logistic_regression.m*: computing the cost function and its gradient. initializing the θ is a 785*1 random vector, after transpose, multiplied with X to get results passed to *sigmoid* function. Finally we get the computed feedforward function in which each element denote a probability correspondingly for a sample belongs to a certain class. Compared to the true results, we want to compute the cost function in order to get the gradient.
- NOTE: the shape of gradient matrix θ (in codes is *g*) is just the same as *x.f* is just a variable.



Vectorized version is:

```

1      h = 1./(1+exp(-theta' * X));
2      f = -y * log(h)' - (1-y) * (log(1-h)')
3      g = X * (h - y)'

```

θ (1xfeature) \times X (nxsamples) = [sample ... sample]

```

h = 1./(1+exp(-theta' * X));
f = -y * log(h)' - (1-y) * (log(1-h)')
g = X * (h - y)'

```

*f is a variable, $-y * \log(h)' * X * (h - y)'$*
 *$g = X * (h - y)'$*

4. *binary_classifier_accuracy.m*: to calculate the model classification accuracy. Multiplied \mathbf{X} input by the optimized $\boldsymbol{\theta}$ to get the output. If output is bigger than 0.5, that is to say classification correct.

1.3 Debugging: Gradient Checking

Sometimes, we need to check the gradient of objective function manually to insure that our optimization algorithm works correctly. Basically, according to derivative definition, for a extremely small ϵ , we have:

$$g(\boldsymbol{\theta}) \approx \frac{J(\boldsymbol{\theta} + \epsilon) - J(\boldsymbol{\theta} - \epsilon)}{2 \times \epsilon}$$

Because $\boldsymbol{\theta}$ is a vector, if we want to check the i -th element of \mathbf{g} vector, just:

$$g_i(\boldsymbol{\theta}) \approx \frac{J(\boldsymbol{\theta}^{i+}) - J(\boldsymbol{\theta}^{i-})}{2 \times \epsilon}$$

Where

$$\boldsymbol{\theta}^{i\pm} = \boldsymbol{\theta} \pm \epsilon \times \vec{e}_i$$

in which, \vec{e}_i is the i -th basis vector of the same dimension as $\boldsymbol{\theta}$, with a 1 in the i -th position and 0s everywhere else.

1.3.1 Matlab Implementation

We create a file called *grad_check.m*. In anywhere, We just pass the optimized function, initialized θ , number checks, and train data set to the *grad_check.m* function. The number checks mean that how many times we want to check the gradient of random selected dimension of θ . Using form likes:

```

1 error = grad_check(@logistic_regression_vec, ...
2   theta, 10, train.X, train.y)

1 function average_error = grad_check(fun, theta0, num_checks, varargin)
2   delta=1e-3;

```

```

3  sum_error=0;
4  fprintf(' Iter      i      err');
5  fprintf('      g      g_estimated      f\n')
6  for i=1:num_checks
7      T = theta0;
8      j = randsample(numel(T),1);
9      T0=T; T0(j) = T0(j)-delta;
10     T1=T; T1(j) = T1(j)+delta;
11
12     [f,g] = fun(T, varargin{:});
13     f0 = fun(T0, varargin{:});
14     f1 = fun(T1, varargin{:});
15     g_estimated = (f1-f0) / (2*delta);
16     error = abs(g(j) - g_estimated);
17     fprintf('% 5d % 6d % 15g % 15f % 15f % 15f\n', ...
18             i,j,error,g(j),g_estimated,f);
19     sum_error = sum_error + error;
20 end
21 average_error = sum_error/num_checks;

```

1.4 Softmax Regression

We use logistic regression classifier to distinguish between two kinds of hand-written digits. For multiple classes, we use softmax regression classifier. For $\mathbf{y}^{(i)} \in \mathfrak{R}^K$ (K classes), normalized the hypothesis distribution:

$$h_{\theta}(\mathbf{x}) = \frac{1}{\sum_{j=1}^K \exp(\boldsymbol{\theta}^{(j)T} \mathbf{x})} \begin{bmatrix} \exp(\boldsymbol{\theta}^{(1)T} \mathbf{x}) \\ \exp(\boldsymbol{\theta}^{(2)T} \mathbf{x}) \\ \vdots \\ \exp(\boldsymbol{\theta}^{(K)T} \mathbf{x}) \end{bmatrix}$$

Note that the $\boldsymbol{\theta}$ is a n -by- K matrix, n denoted the input feature dimensions of \mathbf{x} . Where $\boldsymbol{\theta} = [\boldsymbol{\theta}^{(1)} \boldsymbol{\theta}^{(2)} \dots \boldsymbol{\theta}^{(K)}]$

1.4.1 Cost Function and Its properties

From my point of view, softmax regression is a two layers neural networks, and the $\boldsymbol{\theta}$ can be seen as the weigh of the networks. In the equation below, $\mathbf{1}\{\bullet\}$ is the indicator function $\mathbf{1}\{\text{a true statement}\} = 1$, and $\mathbf{1}\{\text{a false statement}\} = 0$. That is to say when \bullet is true, the $\mathbf{1}\{\bullet\}$ outputs 1. Our cost function will be:

$$J(\boldsymbol{\theta}) = - \left[\sum_{i=1}^m \sum_{k=1}^K \mathbf{1}\{\mathbf{y}^{(i)}=k\} \log \frac{\exp(\boldsymbol{\theta}^{(k)T} \mathbf{x}^{(i)})}{\sum_{j=1}^K \exp(\boldsymbol{\theta}^{(j)T} \mathbf{x}^{(i)})} \right]$$

When $K = 2$, the cost function of softmax regression becomes:

$$J(\boldsymbol{\theta}) = - \sum_i (\mathbf{y}^{(i)} \log(h_{\boldsymbol{\theta}}(\mathbf{x}^{(i)})) + (1 - \mathbf{y}^{(i)}) \log(1 - h_{\boldsymbol{\theta}}(\mathbf{x}^{(i)})))$$

Which is the logistic regression cost function. The gradient of softmax regression cost function is a Hessian matrix:

$$\nabla_{\theta^{(k)}} \mathbf{J}(\boldsymbol{\theta}) = - \sum_{i=1}^m \left[\mathbf{x}^{(i)} \left(1\{\mathbf{y}^{(i)} = k\} - \frac{\exp(\boldsymbol{\theta}^{(k)T} \mathbf{x}^{(i)})}{\sum_{j=1}^K \exp(\boldsymbol{\theta}^{(j)T} \mathbf{x}^{(i)})} \right) \right]$$

in which, m is sample size, $\nabla_{\theta^{(k)}} \mathbf{J}(\boldsymbol{\theta})$ is itself a vector. Softmax regression has an unusual property that is it has a **redundant** set of parameters. In other words, suppose we take each of our parameter vectors $\boldsymbol{\theta}^{(j)}$, and subtract some fixed vector $\boldsymbol{\psi}$ from it, so that every $\boldsymbol{\theta}^{(j)}$ is now replaced with $\boldsymbol{\theta}^{(j)} - \boldsymbol{\psi}$ (for every $j = 1, 2, \dots, K$), which means $\boldsymbol{\theta}^{(j)}$ has redundant information. Let's see how it happens. For our hypothesis function $\mathbf{h}_{\boldsymbol{\theta}}(\mathbf{x})$, each element in it has the form:

$$\frac{\exp(\boldsymbol{\theta}^{(k)T} \mathbf{x}^{(i)})}{\sum_{j=1}^K \exp(\boldsymbol{\theta}^{(j)T} \mathbf{x}^{(i)})}$$

If we subtract a $\boldsymbol{\psi}$ from $\boldsymbol{\theta}^{(k)}$ in numerator and denominator at the same time, it, saying $\exp((\boldsymbol{\theta}^{(k)} - \boldsymbol{\psi})^T \mathbf{x})$. Because \exp function properties, the calculated hypothesis result is just the same as not been subtracted. In a extreme situation, suppose $\boldsymbol{\theta}^{(K)} = \boldsymbol{\psi}$, one can reduce $K \times n$ parameters to $(K - 1) \times n$ parameters needed to be optimized without harming the representational power of our hypothesis. When $K = 2$, the softmax regression becomes logistic regression.

$$\mathbf{h}_{\boldsymbol{\theta}}(\mathbf{x}) = \frac{1}{\exp(\boldsymbol{\theta}^{(1)T} \mathbf{x}) + \exp(\boldsymbol{\theta}^{(2)T} \mathbf{x})} \begin{bmatrix} \exp(\boldsymbol{\theta}^{(1)T} \mathbf{x}) \\ \exp(\boldsymbol{\theta}^{(2)T} \mathbf{x}) \end{bmatrix}$$

Tasking advantage of the properties that softmax regression hypothesis is overparameterized and setting $\boldsymbol{\psi} = \boldsymbol{\theta}^{(2)}$, we can subtract $\boldsymbol{\theta}^{(2)}$ from each of the two parameters, giving us

$$\begin{aligned} \mathbf{h}_{\boldsymbol{\theta}}(\mathbf{x}) &= \frac{1}{\exp((\boldsymbol{\theta}^{(1)} - \boldsymbol{\theta}^{(2)})^T \mathbf{x}) + 1} \begin{bmatrix} \exp((\boldsymbol{\theta}^{(1)} - \boldsymbol{\theta}^{(2)})^T \mathbf{x}) \\ 1 \end{bmatrix} \\ &= \begin{bmatrix} 1 - \frac{1}{\exp((\boldsymbol{\theta}^{(1)} - \boldsymbol{\theta}^{(2)})^T \mathbf{x}) + 1} \\ \frac{1}{\exp((\boldsymbol{\theta}^{(1)} - \boldsymbol{\theta}^{(2)})^T \mathbf{x}) + 1} \end{bmatrix} \\ &= \begin{bmatrix} 1 - \frac{1}{\exp(-\boldsymbol{\theta}'^T \mathbf{x}) + 1} \\ \frac{1}{\exp(-\boldsymbol{\theta}'^T \mathbf{x}) + 1} \end{bmatrix} \quad (\text{if we denote } \boldsymbol{\theta}^{(2)} - \boldsymbol{\theta}^{(1)} = -\boldsymbol{\theta}') \end{aligned}$$

From the above derivation, we find that softmax regression predicts the probability of one of the classes as $\frac{1}{\exp(-\boldsymbol{\theta}'^T \mathbf{x}) + 1}$, and that of the other class as $1 - \frac{1}{\exp(-\boldsymbol{\theta}'^T \mathbf{x}) + 1}$, same as logistic regression.

1.4.2 Matlab Implementation

It will involve three files:

1. *ex1c_softmax.m*: the main file.
2. *multi_classifier_accuracy.m*: to estimate the model.
3. *softmax_regression_vec.m*: to compute the cost function and its gradients. The unvectorized core codes are:

```

1  for i=1:m
2      label = y(i);
```

```

3     p_k_1 = exp(theta'*X(:,i));
4     p_norm = p_k_1 ./ sum(p_k_1);
5     f = f - log(p_norm(label));
6     for j=1:num_classes
7         g(:,j) = g(:,j) + X(:,i).*(p_norm(j));
8     end
9     g(:,label) = g(:,label) - X(:,i);
10  end
11  f = f./m;
12  g(:,num_classes) = [];
13  g=g./m;
14  % make gradient a vector for minFunc
15  g=g(:);

```

$X: 785 \times 60000$ $y: 1 \times 60000$
 $m = 60000$ $\theta: 785 \times 10$

$p_{k-1} = \downarrow_{10} 785 \quad \downarrow_{785} = \downarrow_{10} \Rightarrow p_{k-1} = \downarrow \rightarrow x \begin{bmatrix} \downarrow \downarrow \downarrow \end{bmatrix} = \begin{bmatrix} \downarrow \downarrow \downarrow \dots \end{bmatrix}$
 $p_{norm} = \begin{bmatrix} \downarrow \downarrow \downarrow \dots \end{bmatrix} \Rightarrow p_{norm} = \begin{bmatrix} \downarrow \downarrow \downarrow \dots \end{bmatrix}$

$f = - \sum_{i=1}^m \log(p_{norm}(label)) = - \sum_{i=1}^m \log(\frac{\exp(\theta^{(k)T} x^{(i)})}{\sum_{j=1}^K \exp(\theta^{(j)T} x^{(i)})}) = - \log(10) - \log(0) - \log(0) - \dots = - \log(10) = \text{a variable.}$

$m=1 \quad x^{(i)} \frac{\exp(\theta^{(k)T} x^{(i)})}{\sum_{j=1}^K \exp(\theta^{(j)T} x^{(i)})} = x^{(i)} \cdot 1_{\{y^{(i)}=K\}}$

$\begin{bmatrix} \downarrow \downarrow \downarrow \end{bmatrix} \cdot x \begin{bmatrix} \uparrow \uparrow \uparrow \end{bmatrix} \sim 784 \begin{bmatrix} \uparrow \uparrow \uparrow \end{bmatrix} \rightarrow -x(:,i)$
 $X(:,i) \quad p_{norm}(g)$

$g \uparrow label$

Note:

we would handle all 10 digit classes instead 2 classes as in the previous section. According to softmax regression properties, the last class corresponding to its gradient could be set 0, meaning the 10th class is left out of θ .

1.5 Neural Network and Backpropagation Algorithm

In this section, we will learn a model called Neural Network(NN), and see how does such model can help us to compute graph-based model. We will from another point of view to see softmax regression. First, I would introduce basic mathematical notations, and then briefly describe the network architecture.

1. $W_{i,j}^{(l)}$: the parameter(or weigh) associated with the connection between unit j in layer l , and unit i in layer $l + 1$. The network connection direction is from j to i .
2. $b_i^{(l)}$: the bias associated with unit i in layer $l + 1$.
3. a_i^l : activation of unit i in layer l . Especially, $a_i^{(1)} = x_i$.

4. \mathbf{W}, \mathbf{b} : combination of all the $\mathbf{W}_{i,j}^{(l)}$ and $\mathbf{b}_i^{(l)}$, are total parameters of network we basically want to learn.
5. m : training example size.
6. λ : weigh decay parameter.
7. $\delta_i^{(l)}$: error term, the core of our backpropagation algorithm, which measures how much that node(i -th node in layer l) responsible for any errors in our output.

The cost function(one-half squared-error cost function, with **L2** regularization) for the network to be:

$$\begin{aligned} J(\mathbf{W}, \mathbf{b}) &= \left[\frac{1}{m} \sum_{i=1}^m J(\mathbf{W}, \mathbf{b}; \mathbf{x}^{(i)}, \mathbf{y}^{(i)}) \right] + \frac{\lambda}{2} \sum_{l=1}^{n_{l-1}} \sum_{i=1}^{s_l} \sum_{j=1}^{s_{l+1}} \left(\mathbf{W}_{ji}^{(l)} \right)^2 \\ &= \left[\frac{1}{m} \sum_{i=1}^m \left(\frac{1}{2} \| \mathbf{h}_{\mathbf{W}, \mathbf{b}}(\mathbf{x}^{(i)}) - \mathbf{y}^{(i)} \|^2 \right) \right] + \frac{\lambda}{2} \sum_{l=1}^{n_{l-1}} \sum_{i=1}^{s_l} \sum_{j=1}^{s_{l+1}} \left(\mathbf{W}_{ji}^{(l)} \right)^2 \end{aligned}$$

As $\delta_i^{(l)}$ is a very important term in our backpropagation algorithm, we now formally describe how to compute it:

1. Perform a feedforward pass, computing The activations for layers **L₂**, **L₃**, and so on up to output layer **L_{n_l}**.
2. For each output unit i in layer n_l (the output layer), set

$$\delta_i^{(n_l)} = \frac{\partial}{\partial \mathbf{z}_i^{(n_l)}} \frac{1}{2} \| \mathbf{y} - \mathbf{h}_{\mathbf{W}, \mathbf{b}}(\mathbf{x}) \|^2 = -(\mathbf{y}_i - \mathbf{a}_i^{(n_l)}) \cdot \mathbf{f}'(\mathbf{z}_i^{(n_l)})$$

in which, $\mathbf{z}^{l+1} = \mathbf{W}^{(l)} \mathbf{a}^{(l)} + \mathbf{b}^{(l)}$, $\mathbf{a}^{(l+1)} = \mathbf{f}(\mathbf{z}^{(l+1)})$.

3. For $l = n_l - 1, n_l - 2, n_l - 3, \dots, 2$ For each node i in layer l , set

$$\delta_i^{(l)} = \left(\sum_{j=1}^{s_{l+1}} \mathbf{W}_{ji}^{(l)} \delta_j^{(l+1)} \right) \mathbf{f}'(\mathbf{z}_i^{(l)})$$

4. Compute the desired partial derivatives, which are given as:

$$\begin{aligned} \frac{\partial}{\partial \mathbf{W}_{ij}^{(l)}} J(\mathbf{W}, \mathbf{b}; \mathbf{x}, \mathbf{y}) &= \mathbf{a}_j^{(l)} \delta_i^{(l+1)} \\ \frac{\partial}{\partial \mathbf{b}_i^{(l)}} J(\mathbf{W}, \mathbf{b}; \mathbf{x}, \mathbf{y}) &= \delta_i^{(l+1)} \end{aligned}$$

Finally, we can re-write the algorithm using matrix-vectorial notation. We will use \bullet to denote the element-wise product operator. We also do The same for $\mathbf{f}'(\bullet)$ (so that $\mathbf{f}'([\mathbf{z}_1, \mathbf{z}_2]) = [\mathbf{f}'(\mathbf{z}_1), \mathbf{f}'(\mathbf{z}_2)]$). For step2, after vetorizing, we get (just eliminate i notation)

$$\boldsymbol{\delta}^{(n_l)} = -(\mathbf{y} - \mathbf{a}^{(n_l)}) \bullet \mathbf{f}'(\mathbf{z}^{(n_l)})$$

For step3, the equation can be re-write as:

$$\boldsymbol{\delta}^{(l)} = \left((\mathbf{W}^{(l)})^T \boldsymbol{\delta}^{(l+1)} \right) \bullet \mathbf{f}'(\mathbf{z}^{(l)})$$

For step4,

$$\nabla_{\mathbf{W}^{(l)}} J(\mathbf{W}, \mathbf{b}; \mathbf{x}, \mathbf{y}) = \delta^{(l+1)} (\mathbf{a}^{(l)})^T$$

$$\nabla_{\mathbf{b}^{(l)}} J(\mathbf{W}, \mathbf{b}; \mathbf{x}, \mathbf{y}) = \delta^{(l+1)}$$

If we use sigmoid activation, $f'(z_i^{(l)}) = a_i^{(l)}(1 - a_i^{(l)})$.

1.5.1 BackPropogation Implementation Using Matlab

Based on the previous section, we will train a neural network classifier to classify the 10 digits in the MNIST dataset. First, let's define a neural network We will work for it later. We have a network with a layer of 784 input units, a layer of 10 output units and a single hidden layer of 256 hidden units.

- *run_train.m*

The main script, to load MNIST data set, setup random initial weighs, make cost function optimized and compute accuracy on the test and train set.

- *initialize_weights.m*

Initialling the weights between each layer, the results are stored in **stack** variable. As the network has 3 layers, neurons of each layer are 784-256-1. The **stack** variable has such a form: **stack**<1×2 cell>, **stack**(1, 1) has variables of **W**<256×784 double> and **b**<256×1 double>, **stack**(1, 2) has variables of **W**<10×256 double> and **b**<10×1 double>.

- *params2stack.m stack2params.m* Converts a flattened parameter vector into a nice **stack** structure or flattens a stack to a vector for us to work with. For example, if we use **stack** = **params2stack(params, ei)**, the **params** is our object we want to convert, and the **ei** is an auxiliary variable(maybe a cell structure in matlab) which tells us the configuration of the network we need to satisfy after converted. The total size of **params** here is 203430, according to **ai** structure, we have $203430 = 256 \times 784 + 256 + 10 \times 256 + 10$. Because **minFunc** only can deal with parameters vector, we need such mechanism to convert back and forth.

- *supervised_dnn_cost.m*

Our core function, to slave cost and gradient computation. It contains three main components, forward prop, cost computation and gradients computed using back-propagation.

```

1 function [ cost, grad, pred_prob] = supervised_dnn_cost( theta, ei,..
2     data, labels, pred_only)
3 %SPNETCOSTSLAVE Slave cost function for simple phone net
4 %   Does all the work of cost / gradient computation
5 %   Returns cost broken into cross-entropy, weight norm, and prox reg
6 %       components (ceCost, wCost, pCost)
7
8 %% default values
9 po = false;
10 if exist('pred_only','var')
```

```

11     po = pred_only;
12 end;
13
14 %% reshape into network
15 stack = params2stack(theta, ei);
16 numHidden = numel(ei.layer_sizes) - 1;
17 hAct = cell(numHidden+1, 1); %hAct is <2*1 cell>
18 gradStack = cell(numHidden+1, 1); %gradStack is <2*1 cell>
19 %% forward prop
20 %%% YOUR CODE HERE %%%
21 m = size(data,2); %m=60000
22 m_inv = 1/m;
23 for l = 1:(numHidden+1)
24     if l==1
25         hAct{l} = data; % just store input data<784*60000>
26     else
27         % when run to here, l=2.
28         % stack{l-1}.W <256*784>, hAct{l-1}=hAct{1} <784*60000>
29         % after multiplying, <256*60000> meaning demension from
30         % 784 reduce to 256. stack{l-1}.b=stack{1} <256*1>
31         % repmat(stack{l-1}.b,[1,m]) is <256*60000>, meaning row direction
32         % repeats 1 times and column direction repeats 60000 times.
33         % because for a sample we have 256 dimensions in hidden layer, so
34         % we have 256 neurons, meaning has 256 bias for a sample.
35         % here hAct{l} is "L", l=2, not number 1.
36         hAct{l} = stack{l-1}.W * hAct{l-1} + ...
37             repmat(stack{l-1}.b,[1,m]);
38         if strcmp(ei.activation_fun,'logistic')
39             hAct{l} = sigmoid(hAct{l});
40         elseif strcmp(ei.activation_fun,'tanh')
41             hAct{l} = tanh(hAct{l});
42         elseif strcmp(ei.activation_fun,'rlu')
43             hAct{l}(hAct{l}<0) = 0;
44         end
45     end
46 end
47 % layer l=2, stack{l}.W <10,256>, hAct{l}< 256*60000>,
48 % stack{l}.b <10,1>
49 % repmat(stack{l}.b,[1,m])<10,60000>. From another perspective,
50 % the output layer mapping 256 demensions to 10 demensions.
51 % pred_prob <10*60000>
52 pred_prob = stack{l}.W * hAct{l} + repmat(stack{l}.b,[1,m]);
53 % perform element-by-element binary operation for two array.
54 % max(pred_prob,[],1) <1*60000>, here, "max" function
55 % produces the maximum values along the first
56 % dimension of "pred_prob".
57 % normalizes the output to range[0,1] to represent probability.
58 % bsxfun works like it, suppose

```



```

59 % A = [1 4 6 B=[2 2 2],after bsxfun(@minus,A,B) = [-1 2 4
60 %      5 7 3]                                     3 5 2]
61 % after sigmoid function, each neuron output
62 % a variable range[0,1],however
63 % for Softmax regression, we treat these outputs
64 % as weighed inputs,and put these outputs into
65 % exp power function and divided by its sum
66 % respectively for each sample

```

The cost The cost function is nearly identical to the softmax regression cost function. Note that instead of making predictions from the input data x the softmax function takes as input the final hidden layer of the network

$$h_{W,b}(x) \quad (1.1)$$

. The loss function is thus:

$$J(\theta) = - \left[\sum_{i=1}^m \sum_{k=1}^K 1\{y^{(i)}=k\} \log \frac{\exp(\theta^{(k)T} h_{W,b}(x^{(i)}))}{\sum_{j=1}^K \exp(\theta^{(j)T} h_{W,b}(x^{(i)}))} \right]$$

```

67 pred_prob = bsxfun(@minus,pred_prob,max(pred_prob,[],1));
68 pred_prob = exp(pred_prob);
69 pred_prob = bsxfun(@rdivide, pred_prob, sum(pred_prob));
70
71 %% return here if only predictions desired.
72 if po
73     cost = -1; ceCost = -1; wCost = -1; numCorrect = -1;
74     grad = [];
75     return;
76 end;
77
78 %% compute cost
79 %%% YOUR CODE HERE %%%
80 % groundTruth <10*60000>,meaning if groundTruth(j,i) not equal to 0,
81 % that j is the hand number label for input sample i
82 groundTruth = full(sparse(labels, 1:m, 1));
83 % groundTruth <10*60000>,pred_prob <10*60000>,after sum, <1*60000>
84 cost = -mean(sum(groundTruth .* log(pred_prob)));
85 % to add L2 regularization term.
86 ceCost = -(groundTruth - pred_prob);
87 wCost = 0;
88 % for all weighs in the network, just sum all weighs^2 up.Here wCost
89 % called weigh cost
90 for l=1:(numHidden+1)
91     wCost = wCost + sum(sum(stack{1}.W.^ 2));
92 end
93 wCost = 0.5 * ei.lambda * wCost;
94 cost = cost + wCost;
95
96 %% compute gradients using backpropagation

```

```

97  %%% YOUR CODE HERE %%%
98  delta = cell(numHidden+1,1);%delta<2*1 cell>
99  for l=numHidden+1:-1:1 %from back to front, numHidden = 1
100      % layer l is 2,2==1+1,output layer backpro
101      if(l==numHidden+1)
102          delta{l} = ceCost;
103          %the residual between groundTruth and computed results
104          % delta{2}<10*60000>,hAct{2}<256*60000>,actually is hidden
105          % layer (layer 2) output (activation output).
106          % I AM CONFUSED AT HERE.
107          gradStack{l}.W = delta{l} * hAct{l}' .* m_inv;
108          gradStack{l}.b = mean(delta{l},2);
109      else
110          % stack{2}.W'<60000*10>,delta{2}<10*60000>
111          delta{l} = stack{l+1}.W' * delta{l+1};
112          if strcmp(ei.activation_fun,'logistic')
113              delta{l} = delta{l} .* hAct{l+1} .* (1-hAct{l+1});
114          elseif strcmp(ei.activation_fun,'tanh')
115              delta{l} = delta{l} .* (1-hAct{l+1}.^2);
116          elseif strcmp(ei.activation_fun,'rlu')
117              delta{l}(hAct{l+1}<0) = 0;
118          end
119          % I AM doubt here,why not consider the W2 matrix?
120          % l = 1
121          gradStack{l}.W = delta{l} * hAct{l}' .* m_inv + ...
122              ei.lambda .* stack{l}.W;
123          gradStack{l}.b = mean(delta{l},2);
124      end
125  end
126
127  %% compute weight penalty cost and gradient for non-bias terms
128  %%% YOUR CODE HERE %%%
129
130  %% reshape gradients into vector
131  [grad] = stack2params(gradStack);
132  end
133
134  function sigm = sigmoid(x)
135      sigm = 1./(1+exp(-x));
136  end

```

1.6 Convolutional Neural Networks and pooling

If you have already been familiar with convolutional neural network, here I will emphasize some important details of it only. Basically, suppose input volume is a *width*height*depth* image, which width size is always the same with height, meaning that is a square shape.

Filter is $width_filter * height_filter * depth_filter$, which always is the case that $width_filter = height_filter, depth = depth_filter$. If we have K filters, then we will get K activation maps. For a activation map, suppose the padding size is for one border, it has output size $((width - width_filter + 2 * paddingsize) / stride + 1) \times ((height - height_filter + 2 * paddingsize) / stride + 1)$. In common settings, K is always the powers of 2, e.g. 32, 64, 128. For an element in activation, we get it by element-wise product of volume $width * height * depth$ and $width_filter * height_filter * depth_filter$. For a next time, we move the filter volume a stride, do the element-wise product, and get another element of the activation map. If we have a 1×1 filter with depth of D , the corresponding activation map size will not change compared to input volume but its depth to D .

Pooling, also known as downshampling, captures a feature-invariance property. Say we have a activation map $m \times m$ needed to be pool. If we implement $n \times n$ pooling, Then, our output has $(m \div n) \times (m \div n)$ regions. For a region, we take $n \times n$ mean (or maximum) feature pooling over $m \times m$ activation map to obtain the pooled convolved features.

1.6.1 Implement convolution and Pooling

There are some implementation tips. If we use `conv2(image, W)` MATLAB function, as MATLAB will first `flip W` before convolving `W` with `image`, we need to rotate it twice `rot90` before put it into `conv2` function, e.g.:

$$\begin{bmatrix} 1 & 2 & 3 \\ 4 & 5 & 6 \\ 7 & 8 & 9 \end{bmatrix} \xrightleftharpoons[\text{flip}]{\text{rotate-twice}} \begin{bmatrix} 9 & 8 & 7 \\ 6 & 5 & 4 \\ 3 & 2 & 1 \end{bmatrix}$$

- `cnnExercise.m` The main file to check our implementation whether it is right or not. To load MNISI data set, initial weigh and bias, call the `cnnConvolve.m` and `cnnPool.m` function.

- `cnnConvolve.m`
 $filterDim = 8$, meaning 8×8 filter. $numFilters = 100$, meaning has 100 filters. So W has shape of $< 8 \times 8 \times 100 >$. The b has shape of $< 100 \times 100 >$, but we only need $(numFilters, 1)$. The $images$ shape $< 28 \times 28 \times 8 >$, meaning we have 8 image samples of size $< 28 \times 28 >$, each image has one depth because it is a gray image not RGB. The $convolvedFeatures$ variable is a 4-D tensor, contains the activation maps, our convoluted results. For a activation map, it has dimension $28 - 8 + 1$, as our stride equals to one, and not padding. The $convolvedFeatures$ variable should have $numFilters \times numImages$ activation maps.

We apply 100 filters to a input sample (there total are 8 samples to run) shape $< 28 \times 28 >$ in turn and add bias(100 biases) corresponding to 100 filters, in other words, convolve every filter with every image to get the total $100 * 8$ activation maps. After a filters and an image convoluted, we copy the intermediate variable `convolvedImage` into `convolvedFeatures`

```
1 convolvedImage = conv2(im,filter,'valid');
2 convolvedImage = sigmoid(convolvedImage + b(filterNum,1));
3 convolvedFeatures(:, :, filterNum, imageNum) = convolvedImage;
```

- [*cnnPool.m*](#)

When we get the convolutional output, a 4-D tensor *convolvedFeatures* variable, we can perform the pooling. The pooling output is also a 4-D tensor with its shape $< 21/3 \times 21/3 \times 100 \times 8 >$ (*convolvedDim / poolDim, convolvedDim / poolDim, numFilters, numImage*). There are some tricks to perform pooling: as we want a computing operations compatibility environment, we use *conv2* operations again. We first convolute activation map with $< poolDim \times poolDim >$ matrix patch of 1s, and then downsample it with stride of *poolDim*, and each element in the result $< 7 \times 7 >$ patch will be divided by *poolDim * poolDim* to get final mean pooling.

```

1 for imageNum = 1:numImages
6   for filterNum = 1:numFilters
7     pooled = conv2(convolvedFeatures(:, :, filterNum, imageNum), ...
7     ones(poolDim, poolDim), 'valid');
7     pooledFeatures(:, :, filterNum, imageNum) =
7     pooled(1:poolDim:end, 1:poolDim:end) ./ (poolDim*poolDim);
8   end
10 end

```

1.7 Build up Your Own Convolutional Neural Network

Implement the CNN cost and gradient computation in this step. Your network will have two layers. The first layer is a convolutional layer followed by mean pooling and the second layer is a densely connected layer into softmax regression. The cost of the network will be the standard cross entropy between the predicted probability distribution over 10 digit classes for each image and the ground truth distribution.

- [*cnnTrain.m*](#)

Our main train file is to call other functions, train and test the CNN.

- [*cnnInitParams.m*](#)

Before we perform forward propagation, we need to initial parameters for a single layer convolutional neural. In this function, the parameters space are: *filterDim* = 9, *numFilters* = 20, *imageDim* = 28, *poolDim* = 2, *numClasses* = 10. So after initialing, filter volume has shape of $< 9 \times 9 \times 20 >$ (the filter volume is *Wc*) between input layer and hidden layer). After pooling, the activation map dimension *outDim* are $(imageDim - filterDim + 1) / poolDim = 10$, so hidden neurons are $outDim^2 \times numFilters = 2000$. The output layer size are 10, so *W* (matlab variable denoted by *Wd*) between hidden layer and output layer has size $< 10 \times 2000 >$. For a filter, we touch it a bias, so we totally get $< 20 \times 1 >$ biases *bc*. For output layer, each neuron contact a bias, so there are $< 10 \times 1 >$ biases *bd*. At last, flatten and concatenated all these initialled parameters together, [*cnnInitParams.m*](#) returns a θ containing $< 21650 \times 1 >$ ($9 \times 9 \times 20 + 10 \times 2000 + 20 + 10$) parameters.

- [*cnnCost.m*](#)

Our core function, to perform forward propagation, calculate cost, and compute back

propagation gradients, and make predictions if necessary. Our minibatch size is 256, so a *images* variable contains 256 pictures. The file *cnnParamsToStack.m* is to transform θ into four variables Wc , Wd , bc , bd . In other words, it unrolls θ . Our goal is to update these four parameters gradients.

Table 1.1: The variable space

Variable Name	Value	variable Name	Value
<i>filterDim</i>	9	<i>images</i>	$< 28 \times 28 \times 256 >$
<i>labels</i>	$< 256 \times 1 >$	<i>numClasses</i>	10
<i>numFilters</i>	20	<i>poolDim</i>	2
θ	$< 21650 \times 1 >$	λ	0.0001
<i>Wc</i>	$< 9 \times 9 \times 20 >$	<i>Wd</i>	$< 10 \times 2000 >$
<i>bc</i>	$< 20 \times 1 >$	<i>bd</i>	$< 10 \times 1 >$

STEP 1a: Forward Propagation

In forward propagation, we will compute CNN layer, and Softmax layer.

Listing 1.1: Forward Computing for Convolutional Layer

```

1 % dimension of convolved output
2 convDim = imageDim-filterDim+1;
3 % dimension of subsampled output
4 outputDim = (convDim)/poolDim;
5 % is a 4-D tensor with shape <20*20*20*256>,
6 % for storing convolved results.
7 activations = zeros(convDim,convDim,numFilters,numImages);
8 % is a 4-D tensor with shape <10*10*20*256>,
9 % for storing pooling results.
10 activationsPooled = zeros(outputDim,outputDim, ...
11 numFilters,numImages);
12 activations = cnnConvolve(filterDim, numFilters, ...
13 images, Wc, bc);
14 activationsPooled = cnnPool(poolDim, activations);
15 % Reshape activations into 2-d matrix for Softmax layer,
16 % hiddenSize=2000, activationsPooled <2000*256>.
17 activationsPooled = reshape(activationsPooled,[],numImages);

```

Listing 1.2: Forward Computing for Softmax Layer

```

1 % for storing probability that each image belongs to each class.
   probs <10*256>
2 probs = zeros(numClasses,numImages);
3 % Wd<10*2000>,activationsPooled<2000*256>, hidden neurons 2000,
   number images 256
4 probs = Wd * activationsPooled + repmat(bd,[1,numImages]);
5 % normalizing ,probs <10*256>
6 probs = bsxfun(@minus, probs, max(probs, [], 1));
7 probs = exp(probs);
8 probs = bsxfun(@rdivide, probs, sum(probs));

```

STEP 1b: Calculate Cost

$$J(\theta) = - \left[\sum_{i=1}^m \sum_{k=1}^K 1\{y^{(i)}=k\} \log \frac{\exp(\theta^{(k)T} \mathbf{h}_{W,b}(\mathbf{x}^{(i)}))}{\sum_{j=1}^K \exp(\theta^{(j)T} \mathbf{h}_{W,b}(\mathbf{x}^{(i)}))} \right]$$

Listing 1.3: Calculate Cost

```

1 cost = 0; % save results into cost
2 % groundTruth <10*256>
3 groundTruth = full(sparse(labels, 1:numImages, 1));
4 % numImages_inv = 1/number images = 1/256, cost is a variable, just a
   single number, neither vector nor matrix.
5 % probs <10*256>, groundTruth <10*256>, element wise. For a sample,
   compute 10 by 10 element wise to get one result, all the samples
   results added up. for details, see equation in file
   supervised_dnn_cost.m
6 cost = -numImages_inv*(groundTruth(:)'*log(probs(:))) ...
7 + (lambda/2.)*(sum(Wd(:).^2)+sum(Wc(:).^2));

```

STEP 1c: Backpropagation

Before computing gradients for the 4 variables mentioned above, we need to calculate backpropagate errors θ . we will need to propagate this error through the subsampling and convolutional layer. To do this upsampling quickly, we use kron function to upsample the error and propagate through the pooling layer. If X is 2-by-3, Y is another matrix, then $\text{kron}(X,Y)$ is $\begin{bmatrix} X(1,1)*Y & X(1,2)*Y & X(1,3)*Y \\ X(2,1)*Y & X(2,2)*Y & X(2,3)*Y \end{bmatrix}$ e.g:

$$\begin{bmatrix} 1 & 1 & 2 & 2 \\ 1 & 1 & 2 & 2 \\ 3 & 3 & 4 & 4 \\ 3 & 3 & 4 & 4 \end{bmatrix} \leftarrow \text{kron} \left(\begin{pmatrix} 1 & 2 \\ 3 & 4 \end{pmatrix}, \begin{pmatrix} 1 & 1 \\ 1 & 1 \end{pmatrix} \right)$$

For a normal BP, we have:

$$\delta^{(l)} = \left((\mathbf{W}^{(l)})^T \delta^{(l+1)} \right) \bullet \mathbf{f}'(\mathbf{z}^{(l)})$$

If the l -th layer is a convolutional and subsampling layer then the error is propagated through as

$$\delta^{(l)} = \text{upsample} \left((\mathbf{W}^{(l)})^T \delta^{(l+1)} \right) \bullet \mathbf{f}'(\mathbf{z}^{(l)})$$

In our experiment, the l is second layer, hidden layer, $l = 2$. $\mathbf{W}^{(l)}$ is between layer 2 and layer 3 (the output layer). $\delta^{(l+1)}$, $\delta^{(3)}$, the output layer propagated error. $\mathbf{f}'(\mathbf{z}_i^{(l)})$ is the derivative of the activation function, and $\mathbf{f}'(\mathbf{z}_i^{(l)}) = \mathbf{a}_i^{(l)}(1 - \mathbf{a}_i^{(l)})$. k indexes the filter number.

Listing 1.4: Calculate Delta Error for CNN Layer

```

1 %delta <10*256>
2 delta = -(groundTruth - probs);
3 % 4-D tensor, Wd is between layer 2 and layer 3.

```

```

4 delta_pool = reshape(Wd'*delta, outputDim, outputDim, ...
5 numFilters, numImages);
6 % delta_conv <20*20*20*256>
7 delta_conv = zeros(convDim,convDim,numFilters,numImages);
8 % upsampling the delta_pool to delta_conv
9 for i=1:numImages
10     for j=1:numFilters
11         % strictly speaking, here delta_conv variable is denoted
12         % the result of upsampled pooling. After kron, get<20*20>
13         delta_conv(:,:,j,i) = (1./poolDim^2) .* kron(...
14         squeeze(delta_pool(:,:,j,i)), ones(poolDim));
15     end
16 end
17 % To propagate pooling error through the convolutional layer,
18 % we simply need to multiply the incoming error by
19 % the derivative of the activation function as in the
20 % usual back propagation algorithm.
21 % activations, delta_conv: 4-D, element wise product
22 % activations variable comes from CNN forward path, and have
23 % not been put into Softmax layer
24 delta_conv = activations .* (1-activations) .* delta_conv;

```

STEP 1d: Gradient Calculation There are 4 parameters gradient we need to calculate. Compute the gradient for the densely connected weights and bias, $\mathbf{W_d}$ and $\mathbf{b_d}$, just as usual. In usual situation,

$$\nabla_{\mathbf{W}^{(l)}} J(\mathbf{W}, \mathbf{b}; \mathbf{x}, \mathbf{y}) = \delta^{(l+1)} (\mathbf{a}^{(l)})^T$$

$$\nabla_{\mathbf{b}^{(l)}} J(\mathbf{W}, \mathbf{b}; \mathbf{x}, \mathbf{y}) = \delta^{(l+1)}$$

To calculate the gradient w.r.t.(with respect to) to the filter maps, we rely on the border handling convolution operation again and flip the error matrix $\delta_k^{(l)}$ the same way we flip the filters in the convolutional layer. In layer 1 and layer 2 is to do convolutional and pooling operations, \mathbf{a}^1 is image input layer data. k denoted k -th filter

$$\nabla_{\mathbf{W}_k^{(l)}} J(\mathbf{W}, \mathbf{b}; \mathbf{x}, \mathbf{y}) = \sum_{i=1}^m (\mathbf{a}_i^{(l)}) * \text{rot90}(\delta_k^{(l+1)}, 2)$$

$$\nabla_{\mathbf{b}_k^{(l)}} J(\mathbf{W}, \mathbf{b}; \mathbf{x}, \mathbf{y}) = \sum_{a,b} (\delta_k^{(l+1)})_{a,b}$$

Attention: For a sample, we calculate its gradients as this way

$$\frac{\partial}{\partial \mathbf{W}_{ij}^{(l)}} J(\mathbf{W}, \mathbf{b}; \mathbf{x}, \mathbf{y}) = \mathbf{a}_j^{(l)} \delta_i^{(l+1)}$$

$$\frac{\partial}{\partial \mathbf{b}_i^{(l)}} J(\mathbf{W}, \mathbf{b}; \mathbf{x}, \mathbf{y}) = \delta_i^{(l+1)}$$

For a batch of sample, say we have m samples, the derivative of the overall cost function is:

$$\frac{\partial}{\partial \mathbf{W}_{ij}^{(l)}} J(\mathbf{W}, \mathbf{b}; \mathbf{x}, \mathbf{y}) = \left[\frac{1}{m} \sum_{i=1}^m \frac{\partial}{\partial \mathbf{W}_{ij}^{(l)}} J(\mathbf{W}, \mathbf{b}; \mathbf{x}^{(i)}, \mathbf{y}^{(i)}) \right] + \lambda \mathbf{W}_{ij}^{(l)}$$

$$\frac{\partial}{\partial \mathbf{b}_i^{(l)}} J(\mathbf{W}, \mathbf{b}; \mathbf{x}, \mathbf{y}) = \frac{1}{m} \sum_{i=1}^m \frac{\partial}{\partial \mathbf{b}_i^{(l)}} J(\mathbf{W}, \mathbf{b}; \mathbf{x}^{(i)}, \mathbf{y}^{(i)})$$

Listing 1.5: Gradient Calculation

```

1 % according the above two equations, for a batch of samples.
2 % delta <10*256> activationsPooled <2000*256>,
3 % Wd <10*2000> sum(delta, 2)<10,1>
4 Wd_grad = numImages_inv .* delta * ...
5 activationsPooled' + lambda .* Wd;
6 bd_grad = numImages_inv .* sum(delta, 2);
7 for i=1:numFilters
8     for j=1:numImages
9         % For a filter, we need to backpropagate error for
10        % that filter with each image and aggregate over images.
11        % Wc_grad contains 20 filters, meaning i=1:20.
12        % squeeze(images(:,:,j)) <28*28>
13        % squeeze(delta_conv(:,:,i,j)) <20*20>
14        Wc_grad(:,:,i) = Wc_grad(:,:,i) + conv2(...
15        squeeze(images(:,:,j)),...
16        rot90(squeeze(delta_conv(:,:,i,j)),2), ...
17        'valid');
18    end
19    % Wc_grad <9*9*20>, Wc <9*9*20>
20    Wc_grad(:,:,i) = numImages_inv .* ...
21    Wc_grad(:,:,i) + lambda .* Wc(:,:,i);
22
23    % temp 4-D tensor, sum(temp(:)) is a variable
24    % meaning, for a filter, we inspect all its
25    % images into which the filter has ever filtered.
26    % then we get this bais gradient represents all
27    % all images the filter has ever seen.
28    temp = delta_conv(:,:,i,:);
29    bc_grad(i) = numImages_inv .* sum(temp(:));
30 end
31 % Unroll gradient into grad vector for minFuncd
32 grad = [Wc_grad(:) ; Wd_grad(:) ; bc_grad(:) ; bd_grad(:)];

```

- *computeNumericalGradient.m*

Before we preform a large scale computation, it is a good habit to check gradient to insure our cost function calculated correctly. Once our code passes the gradient check we are ready to move onto training a real network on the full dataset. Make sure to switch the DEBUG boolean to false in order not to run the gradient check again.

Suppose we have already defined all the computing environment for gradient checking such as the number of images, labels, filter dimension et al (for more details to

see the *cnnTrain.m* file). Then, we just change a tiny value of a certain dimension of θ parameters (a vector contains all the network weights and biases) to see what happens in cost function J we have defined already. Compared the gradient J computed with our derivative definition computed, we can see how well our cost function works.

Listing 1.6: gradient checking

```

1 function numgrad = computeNumericalGradient(J, theta)
2 % Initialize numgrad with zeros
3 numgrad = zeros(size(theta));
4 epsilon = 1e-4;
5 for i =1:length(numgrad)
6     oldT = theta(i);
7     theta(i)=oldT+epsilon;
8     pos = J(theta);
9     theta(i)=oldT-epsilon;
10    neg = J(theta);
11    numgrad(i) = (pos-neg)/(2*epsilon);
12    theta(i)=oldT;
13    if mod(i,100)==0
14        fprintf('Done with %d\n',i);
15    end;
16 end;
17 end

```

- *computeNumericalGradient.m*

Batch methods, such as limited memory BFGS, which use the full training set to compute the next update to parameters at each iteration tend to converge very well to local optima. They are also straight forward to get working provided a good off the shelf implementation (e.g. `minFunc`) because they have very few hyper-parameters to tune. However, often in practice computing the cost and gradient for the entire training set can be very slow and sometimes intractable on a single machine if the dataset is too big to fit in main memory. Another issue with batch optimization methods is that they don't give an easy way to incorporate new data in an 'online' setting. Stochastic Gradient Descent (SGD) addresses both of these issues by following the negative gradient of the objective after seeing only a single or a few training examples. The use of SGD in the neural network setting is motivated by the high cost of running back propagation over the full training set. SGD can overcome this cost and still lead to fast convergence. We will use SGD with momentum to accelerate computing and avoid local minimums to some extent.

Table 1.2: The variable space

Variable Name	Value	variable Name	Value
<i>funObj</i>	cost function handle	<i>data</i>	$< 28 \times 28 \times 60000 >$
<i>labels</i>	$< 60000 \times 1 >$	<i>peochs</i>	3
<i>alpha</i>	0.1000	<i>momentum</i>	0.9500
θ	$< 21650 \times 1 >$	<i>minibatch</i>	256

We obey the following equations to update our parameter θ

$$\mathbf{v} = \gamma \mathbf{v} + \alpha \nabla_{\theta} J(\theta; \mathbf{x}^{(i)}, \mathbf{y}^{(i)})$$

$$\theta = \theta - \mathbf{v}$$

In above equation \mathbf{v} is the current velocity vector which is of the same dimension as the parameter vector θ . Finally $\lambda \in (0, 1]$ determines for how many iterations the previous gradients are incorporated into the current update.

Listing 1.7: SGD with Momentum

```

1 % Setup for momentum
2 mom = 0.5;
3 momIncrease = 20;
4 velocity = zeros(size(theta));
5 it = 0;
6 for e = 1:epochs % 1:3
7 % randomly permute indices of data for quick minibatch sampling
8     rp = randperm(m);
9     for s=1:minibatch:(m-minibatch+1) % s = 1:256:(60000-256), total
        234 iterations
10         it = it + 1;
11         % increase momentum after momIncrease iterations
12         % for a epoch, suppose after 19 iterations, when in 20-th
13         % iterations, we update mom to 0.9500, and for rest of iterations,
14         % we use the 0.9500 momentum all the time
15         if it == momIncrease % momIncrease = 20
16             mom = options.momentum;
17         end;
18         % get next randomly selected minibatch
19         mb_data = data(:, :, rp(s:s+minibatch-1)); %mb_data <28*28*256>
20         mb_labels = labels(rp(s:s+minibatch-1)); % mb_labels <256*1>
21         % evaluate the objective function on the next minibatch
22         % cost: a variable, grad:<21650*1>
23         [cost grad] = funObj(theta, mb_data, mb_labels);
24         % Update velocity and theta
25         % velocity, theta: <21650*1>
26         velocity = mom.* velocity + alpha .* grad;
27         theta = theta - velocity;
28         fprintf('Epoch %d: Cost on iteration %d is %f\n', e, it, cost);
29     end;
30     % anneal learning rate by factor of two after each epoch
31     % we are total 3 epochs, for a epoch, we observe all data set once
32     alpha = alpha/2.0;
33 end;
34 opttheta = theta;

```

The last work is to test our model. As we access the optimized $\theta = \text{opttheta}$ from the file `computeNumericalGradient.m` output. We just put it to `cnnCost.m` file to compute forward output.

Listing 1.8: Test

```
1 [~,cost,preds]=cnnCost(opttheta,testImages,testLabels,numClasses,...
2 filterDim,numFilters,poolDim,true);
3 % testLabels <10000*1>,compare two data sets element by
4 % element to check our model accuracy.
5 acc = sum(preds==testLabels)/length(preds);
6 % Accuracy should be around 97.4% after 3 epochs
7 fprintf('Accuracy is %f\n',acc);
```

Listing 1.9: cnnCost.m

```
1 if pred
2     % probs <10*10000>,10000 test samples.preds<1*10000>
3     % suppose preds(1,1)=7,denoted the first test sample output 7.
4     [~,preds] = max(probs,[],1);
5     preds = preds';
6     grad = 0;
7     return;
8 end;
```

Congratulations, you've successfully implemented a Convolutional Neural Network!

2

Unsupervised Learning

2.1 Autoencoders

In last section of chapter 1, we pour great efforts to build up a CNN from scratch. In this section, we will introduce the first kind of unsupervised learning methods, autoencoders. To begin with, we will introduce some Sparse Autoencoder (SAE) notations. The SAE is consist of 3 layers: input layer, hidden layer and output layer. The number of neurons in input layer is the same as output layer. If the number of neurons in hidden layer are less than input layer, then meaning we constraint our network to learn some compression representations of inputs. It's very useful when the input features are correlated, because this algorithm will be able to discover some of those correlations. When there does not exit structure in the data, say inputs coming from i.i.d.(independent and identically distributed) Gaussian independent of the other features, it would very difficult for Autoencoder to learn such a compression task. If the number of neurons in hidden layer are more than input neurons, and if we impose a **sparsity** constraint on the hidden units, then the autoencoder will still discover interesting structure in the data, even if the number of hidden units is large.

- \mathbf{x} Input features for a training example, $\mathbf{x} \in \mathfrak{R}^n$
- \mathbf{y} Output or target values, a vector, $\mathbf{y} = \mathbf{x}$
- $(\mathbf{x}^{(i)}, \mathbf{y}^{(i)})$ The i -th training example
- $\mathbf{h}_{\mathbf{W}, \mathbf{b}}(\mathbf{x})$ Output of our hypothesis on input \mathbf{x} , using parameters \mathbf{W}, \mathbf{b} . It is a vector of the same dimension as the target value \mathbf{y}
- $\mathbf{W}_{ij}^{(l)}$ The parameter associated with connection between unit j in layer l , and unit i in layer $l + 1$
- $\mathbf{b}_i^{(l)}$ The bias term associated with unit i in layer $l + 1$. Can also be thought of as the parameter associated with the connection between the bias unit in layer l and unit i in layer $l + 1$
- $\boldsymbol{\theta}$ Our parameter vector, can be seen as the results of taking the parameters \mathbf{W}, \mathbf{b} and unrolling them into a long column vector
- $\mathbf{a}_i^{(l)}$ Activation(output) of unit i in layer l of the networks. $\mathbf{a}_i^{(1)} = \mathbf{x}_i$

- $f(\cdot)$ The activation function. Here we use $f(z) = \tanh(z)$
- $z_i^{(l)}$ Total weighted sum of inputs to unit i in layer l , $a_i^{(l)} = f(z_i^{(l)})$
- α Learning rate parameter
- s_l Number of units in layer l (not counting the bias unit).
- n_l Number layers in the network, just say how many layers in the network. Layer L_1 is usually the input layer, and layer L_{n_l} is the output layer
- λ Weight decay parameter
- \hat{x} the reconstruction outputs, same meaning as $h_{W,b}(x)$
- ρ Sparsity parameter, which specifies our desired level of sparsity
- $\hat{\rho}_i$ The average activation of hidden unit i (in the sparse autoencoder)
- β Weight of the sparsity penalty term (in the sparse autoencoder objective)

2.1.1 Sparse Autoencoder Implementation

STEP 1: Generate training set `sampleIMAGES.m`

Variable `IMAGES` contains 10 images with the shape of 512×512 , each image have been preprocessed using whitening method. Whitening is a preprocessing step which removes redundancy in the input, by causing adjacent pixels to become less correlated. We will randomly select an image, and sample randomly 8×8 image patch 1000 times for that image, so that we have 10000 patches for each patch has shape of 8×8 , and reshaped it into a 64×10000 matrix. BUT in codes, we have a trick to implement it. The function `sampleIMAGES.m` will return 10000 normalized image patches with shape of 64×10000 .

Table 2.1: The variable space

Variable Name	Value	variable Name	Value
<code>IMAGES</code>	$<512*512*10>$	<code>bumpatches</code>	10000
<code>patches</code>	$<64 \times 10000>$	<code>patchsize</code>	8
<code>idx</code>	from 0 to 10000		

Listing 2.1: backbone of `sampleIMAGES`

```

1 idx = 0;
2 img_rnd = randperm(10);
3 patch_x_rnd = randperm(512-7);
4 patch_y_rnd = randperm(512-7);
5 % consider the IMAGES shape particularity, we divide the sample process
  into
6 % two same parts but with different random characteristic.
7 for i=1:10
8     for j=1:500
9         idx = idx+1;
10        patches(:,idx) = reshape(IMAGES(patch_x_rnd(j):(patch_x_rnd(j)+7), ...

```

```

11     patch_y_rnd(j):(patch_y_rnd(j)+7),img_rnd(i)), 64, 1);
12     end
13 end
14 % we use the same code again to fill up the other half of 'patches'
    variable.
15 % note: idx here is from 5001 to 10000
16 img_rnd = randperm(10);
17 patch_x_rnd = randperm(512-7);
18 patch_y_rnd = randperm(512-7);
19
20 for i=1:10
21     for j=1:500
22         idx = idx+1;
23         patches(:,idx) = reshape(IMAGES(patch_x_rnd(j):(patch_x_rnd(j)+7),
            ....
24         patch_y_rnd(j):(patch_y_rnd(j)+7),img_rnd(i)), 64, 1);
25     end
26 end
27
28 patches = normalizeData(patches);
29 function patches = normalizeData(patches)
30 % Squash data to [0.1, 0.9] since we use sigmoid as the activation
31 % function in the output layer
32
33 % Remove DC (mean of images).
34 patches = bsxfun(@minus, patches, mean(patches));
35
36 % Truncate to +/-3 standard deviations and scale to -1 to 1
37 % pstd: a variable
38 pstd = 3 * std(patches(:));
39 patches = max(min(patches, pstd), -pstd) / pstd;
40
41 % Rescale from [-1,1] to [0.1,0.9]
42 patches = (patches + 1) * 0.4 + 0.1;
43
44 end

```

In our main file *train.m*, we can visualize the image patches result using the file *display_network.m*. A visualization version maybe like following picture.

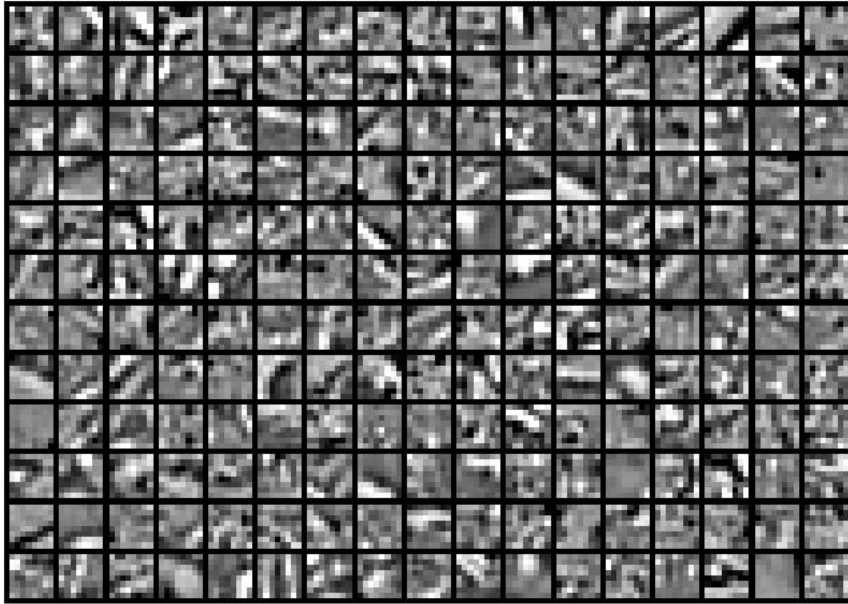
Listing 2.2: Visualizing Input image patches

```

1 patches = sampleIMAGES;
2 % size(patches,2) =10000 randi(10000,200,1) = <200*1>
3 % randi(imax,m,n) means generate a m-by-n array of integer values from the
    uniform
4 % distribution on the set 1:imax. display_network(A) function display a A
    matrix (n*m) which
5 % has m images with the shape of <squareroot of n, squareroot of n> (
    default)

```

```
6 display_network(patches(:,randi(size(patches,2),200,1)),8);
```



Then

we use file `initializeParameters.m` to initialize parameters. As the number of input is just the same as the number of output, we only need two parameters to call this function. We put $hiddenSize = 25$ and $visibleSize = 64$, get the initialized parameters θ with shape $< 3289 \times 1 >$ ($64 \times 25 + 64 + 25 \times 62 + 25$, each bias associates a neuron, input layer has bias and output layer has bias).

If you are already familiar with CNN and BP, then you know that our core function or our algorithm soul is to achieve cost function and compute its gradient. Don't worry, you are unnecessarily afraid for it. In my point of view, it is much more easier than CNN's backpropagation, but still has some impressed equations we need to understand.

- `sparseAutoencoderCost.m`

Table 2.2: The variable space

Variable Name	Value	variable Name	Value
θ	$< 3289 \times 1 >$	$visibleSize$	64
$hiddenSize$	25	λ	1.0000e-04
$sparsityParam$	0.0100	beta	3
$data$	$< 64 \times 10000 >$	$W1$	$< 25 \times 64 >$
$W2$	$< 64 \times 25 >$	$b1$	$< 25 \times 1 >$
$b2$	$< 64 \times 1 >$		

We use sigmoid function for activation function, and matrices $W^{(1)} \in \Re^{s_1 \times s_2}$, $W^{(2)} \in \Re^{s_2 \times s_3}$, vectors $b^{(1)} \in \Re^{s_2}$, $b^{(2)} \in \Re^{s_3}$. The objective $J_{sparse}(W, b)$ contains 3 terms, corresponding to squared error term, the weight decay term, and the sparsity penalty.

We write $\mathbf{a}_j^{(2)}(\mathbf{x})$ to denote the activation of this hidden unit (unit j in layer 2) when the network is given a specific input \mathbf{x} . Then the average activation of hidden unit j (averaged over the training set) given as follow:

$$\hat{\rho}_j = \frac{1}{m} \sum_{i=1}^m [\mathbf{a}_j^{(2)}(\mathbf{x}^{(i)})]$$

If we use KL divergence (KL-divergence is a standard function for measuring how different two different distributions are.) and sum up all $\hat{\rho}_j$ in hidden layer, (in codes, *mean_act1* is $\hat{\rho}_j$ for all j , ρ is *sarsityParam*=0.0100), then

$$\sum_{j=1}^{s_2} \rho \log \frac{\rho}{\hat{\rho}_j} + (1 - \rho) \log \frac{1 - \rho}{1 - \hat{\rho}_j}$$

and also can be written as

$$\sum_{j=1}^{s_2} \mathbf{KL}(\rho || \hat{\rho}_j)$$

For *W2grad*

$$\begin{aligned} \nabla_{\mathbf{W}^{(l)}} \mathbf{J}(\mathbf{W}, \mathbf{b}; \mathbf{x}, \mathbf{y}) &= \delta^{(l+1)} (\mathbf{a}^{(l)})^T \\ \mathbf{W}^{(l)} &= \mathbf{W}^{(l)} - \alpha \left[\left(\frac{1}{m} \Delta \mathbf{W}^{(l)} \right) + \lambda \mathbf{W}^{(l)} \right] \end{aligned}$$

normally, we have

$$\boldsymbol{\delta}^{(l)} = \left((\mathbf{W}^{(l)})^T \boldsymbol{\delta}^{(l+1)} \right) \bullet \mathbf{f}'(\mathbf{z}^{(l)})$$

now instead compute

$$\boldsymbol{\delta}_i^{(2)} = \left(\left(\sum_{j=1}^{s_3} \mathbf{W}_{ji}^{(2)} \boldsymbol{\delta}_j^{(3)} \right) + \beta \left(\frac{\rho}{\hat{\rho}_i} + \frac{1 - \rho}{1 - \hat{\rho}_i} \right) \right) \mathbf{f}'(\mathbf{z}_i^{(2)})$$

In implementation, we will transmit *delta1* to *delta2*. Here *delta1* is contacted layer 2, these codes may not will be denoted. Because for $\delta^{(l)}$ should from n_l (contacted layer n_l) to 2 (contacted layer 2). After all, we just know, we use the previous δ (suppose coming from layer l) and multiply the weigh \mathbf{W} between layer l and $l - 1$.

Listing 2.3: backbone of spareAutoencoderCost.m

```

1 % numImages=10000
2 numImages = size(data,2);
3 numImages_inv = 1./numImages;
4 % for each image, we bind it a vector of bias vector b.
5 % note: we use all samples to compute, not a batch
6 % activations1 <25*10000> mean_act1 <25*1>
7 activations1 = sigmoid(W1*data+repmat(b1,[1,numImages]));
8 output = sigmoid(W2*activations1+repmat(b2,[1,numImages]));
9 mean_act1 = mean(activations1,2);
10 % squared_error <64*10000>
11 squared_error = 0.5.*(output - data).^2;
12 % squared error term, cost is a variable
13 cost = numImages_inv .* sum(squared_error(:));
14 % the weight decay term, cost is a variable
15 cost = cost + 0.5 .* lambda .* (sum(W1(:).^2) + sum(W2(:).^2));
16 % the sparsity penalty term, cost is a variable, after kl_div,
17 % will return <25*1>, and we just need to sum it up
18 cost = cost + beta .* sum(kl_div(sparsityParam, mean_act1));
19 % delta2 <64*10000>
20 delta2 = -(data-output) .* output .* (1-output);
21 % gradient of bias in layer 2 connected layer 3: b2grad <64*1>
22 b2grad = mean(delta2,2);
23 W2grad = numImages_inv .* delta2 * activations1' + lambda .* W2;
24 % delta_sparsity <25*10000>, we use repmat to generate 10000 samples
    delta_sparsity,
25 % saying that for each sample, we really hope it has such a sparsity.
26 delta_sparsity = repmat(beta.*(-sparsityParam./mean_act1+(1-
    sparsityParam)./(1-mean_act1)), ...
27 [1,numImages]);
28 % delta2 <64*10000>, W2 <64*25>, delta1 <25*10000>
29 % delta1 has 25 rows, in hidden layer
30 delta1 = (W2' * delta2 + delta_sparsity) .* activations1 .* (1-
    activations1);
31 b1grad = mean(delta1,2);
32 % data in layer 1, is activation
33 W1grad = numImages_inv .* delta1 * data' + lambda .* W1;
34 grad = [W1grad(:) ; W2grad(:) ; b1grad(:) ; b2grad(:)];

```

At last, we use [computeNumericalGradient.m](#) to check our cost function whether its computing process is correct or not. Then we use [minFunc.m](#) to optimize our cost function to get the *opttheta* variable. Finally, we need to visualize *W1*. Because θ is consist of 4 parameters, *W1*, *W2*, *b1*, *b2* successively.

Listing 2.4: Visualization

```

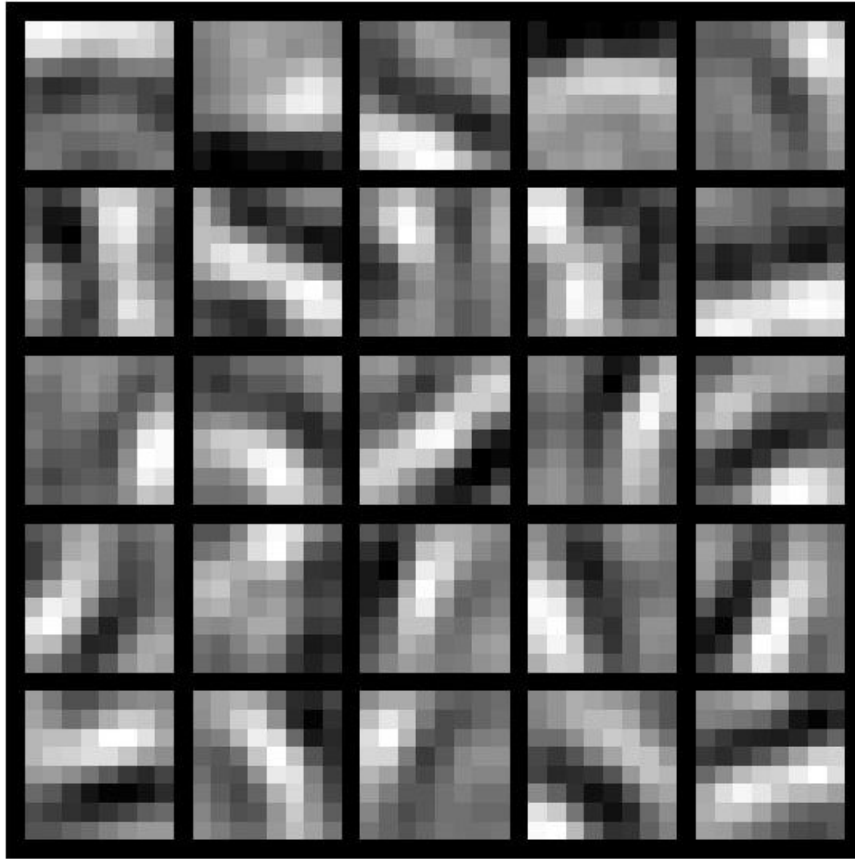
1 % W1 <25*64>
2 W1 = reshape(opttheta(1:hiddenSize*visibleSize), hiddenSize,
    visibleSize);

```

```

3 display_network(W1', 12);
4 print -djpeg weights.jpg    % save the visualization to a file

```



2.2 PCA, PCA Whitening and ZCA Whitening on Images

For images, the input will be somewhat redundant, because the values of adjacent pixels in an image are highly correlated. We need a way to do dimensionality reduction and make features uncorrelated. The goal of whitening is to make the input less redundant; more formally, our desiderata are that our learning algorithms sees a training input where (i) the features are less correlated with each other, and (ii) the features all have the same variance. In detail, in order for PCA to work well, informally we require that (i) The features have approximately zero mean, and (ii) The different features have similar variances to each other. With natural images, (ii) is already satisfied even without variance normalization, and so we won't perform any variance normalization. In this section, we only need a file [pca_gen.m](#)

1. For each image $\mathbf{x}^{(i)}$, we might normalize its the intensity as follows:

$$\mu^{(i)} := \frac{1}{n} \sum_{j=1}^n \mathbf{x}_j^{(i)}$$

$$\mathbf{x}_j^{(i)} := \mathbf{x}_j^{(i)} - \mu^{(i)}$$

for all j .

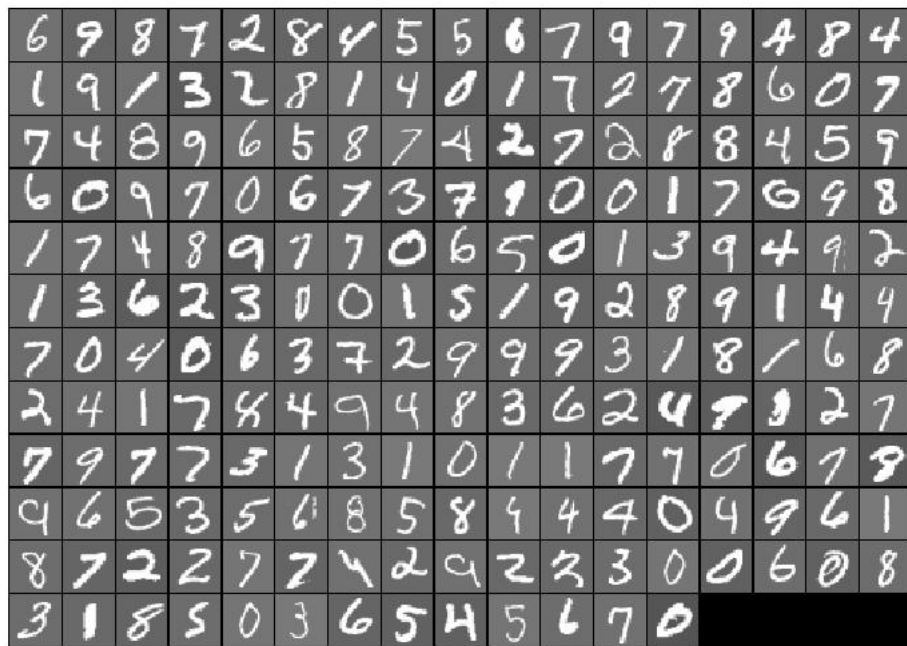
Here i is to denoted images index, and j is to denoted the i -th image j -th pixel intensity.

Listing 2.5: Zero-mean the data

```

1 % x is out input images set,<784*60000>, avg <1*60000>, a random
2 % selection of 200 samples for visualization as follows
3 avg = mean(x,1);
4 x = x - repmat(avg,size(x,1),1);

```



- Next we need to compute the matrix Σ defined as follows:

$$\Sigma = \frac{1}{m} \sum_{i=1}^m (\mathbf{x}^{(i)})(\mathbf{x}^{(i)})^T$$

If \mathbf{x} has zero mean, then Σ is exactly the covariance matrix of \mathbf{x} . In matlab, we can use $\text{sigma} = x * x' / \text{size}(x,2)$, which x have m samples with shape $\langle n*m \rangle$. (what is the relationship between the equation and matlab code? I am confused about it.) When we use Matlab *svd* function, then the matrix U will contain the eigenvectors of Sigma (one eigenvector per column, sorted in order from top to bottom eigenvector), and the diagonal entries of the matrix S will contain the corresponding eigenvalues (also sorted in decreasing order). The matrix V will be equal to transpose of U , and can be safely ignored. The *svd* function actually computes the singular vectors and singular values of a matrix, which for the special case of a symmetric positive semi-definite matrix—which is all that we're concerned with here—is equal to its

eigenvectors and eigenvalues. \mathbf{x}_{rot} (The subscript **rot** comes from the observation that this corresponds to a rotation (and possibly reflection) of the original data.) is another version of original data \mathbf{x} . If we get \mathbf{x}_{rot} , then we can go from the rotated vectors \mathbf{x}_{rot} back to the original data \mathbf{x} as follows: $\mathbf{x} = \mathbf{U}\mathbf{x}_{rot}$

Listing 2.6: Zero-mean the data

```

1 xRot = zeros(size(x)); % We need to compute this
2 % sigma <784*784>
3 sigma = x * x' / size(x,2);
4 % u<784*784>, s<784*784>,v<784*784>
5 [u,s,v] = svd(sigma);
6 % xRot <784*10000>
7 xRot = u' * x; % rotated version of the data x
8 % If we want to a dimensioned version of original data, saying we
9 % want to keep only k dimensions, use the follows statement,
10 % where k is the number of eigenvectors to keep
11 % xTilde = U(:,1:k)' * x;
```

3. How can we insure that our xRot variable is right? Well, covariance matrix for the xRot should be a diagonal matrix with non-zero entries only along the main diagonal. So, we need to compute the covariance matrix of xRot and visualize it.

Listing 2.7: Check your implementation of PCA

```

1 covar = zeros(size(x, 1)); % You need to compute this
2 covar = xRot * xRot' / size(xRot,2);
3 % Visualise the covariance matrix. You should see a line across the
4 % diagonal against a blue background.
5 figure('name','Visualisation of covariance matrix');
6 % you will see a coloured diagonal line against a blue background,
7 % however, with so many samples, the diagonal line may not be
  apparent.
8 % you can check and see the concrete covar matrix in Matlab variables
  space
9 imagesc(covar);
```

4. Now we need to decide how many dimensions should be retain. We will pick k to be as small as possible, but so that at least 99% of the variance is retained.

Listing 2.8: Find k, the number of components to retain

```

1 k = 0; % Set k accordingly
2 % s contains our all eigenvalues
3 total = sum(s(:));
4 % mat <784*784>
5 mat = zeros(size(s,1));
6 % size(s,1) = 784
7 for i=1:size(s,1)
8     mat(i,i) = 1;
9     % s is a diagonal matrix, only diagonal has data.
```

```

10 % In face, we just sum k elements in s in turn by diagonal and
11 % compare the result to total sum to see the percentage
12 res = mat * s;
13 if (sum(res(:)) / total) >= 0.99
14     break;
15 end
16 end
17 k = i;

```

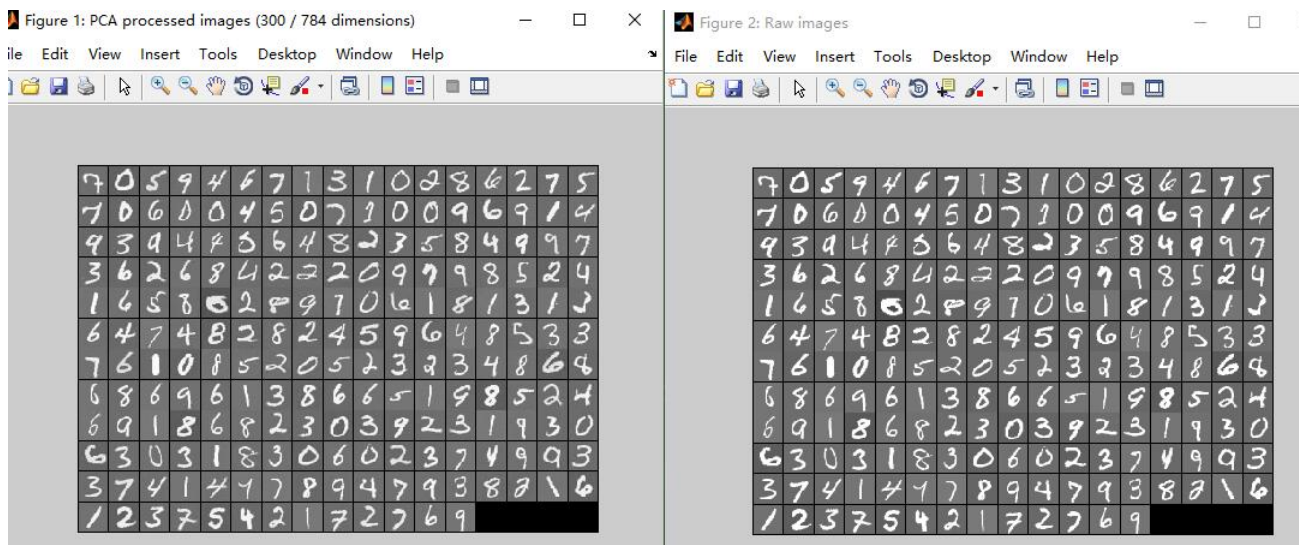
5. We will recover our data with k dimensions as follows:

Listing 2.9: Implement PCA with dimension reduction

```

1 xHat = zeros(size(x)); % You need to compute this
2 u_k = u(:,1:k); % <784*300>
3 xHat = u_k * u_k' * x; % x <784*60000>
4 figure('name',['PCA processed images',sprintf('%d / %d dimensions)',
5     k, size(x, 1)], '');
6 display_network(xHat(:,randsel));
7 figure('name','Raw images');
8 display_network(x(:,randsel))

```



6. PCA whitening and regularisation.

Listing 2.10: Implement PCA with whitening and regularisation

```

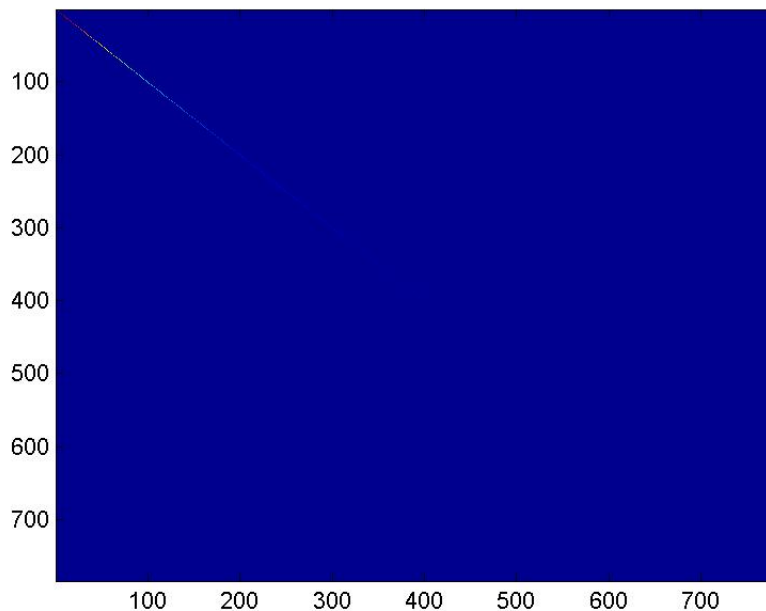
1 % Without regularisation (set epsilon to 0 or close to 0),
2 epsilon = 0.1;
3 xPCAWhite = zeros(size(x));
4 % As we have compute, xRot = u' * x; xPCAWhite <784*60000>
5 xPCAWhite = diag(1./sqrt(diag(s)+epsilon)) * xRot;
6
7 % use the dimensioned x to do pca whitening
8 % xPCAWhite_rd = diag(1./sqrt(diag(s)+epsilon)) * u' * xHat;
9

```

```

10 % Check our implementation of PCA whitening with and without
    regularisation.
11 % PCA whitening without regularisation results a covariance matrix
12 % that is equal to the identity matrix. PCA whitening with
    regularisation
13 % results in a covariance matrix with diagonal entries starting
    close to
14 % 1 and gradually becoming smaller
15 covar = xPCAWhite * xPCAWhite' / size(xPCAWhite,2);
16 figure('name','Visualisation of covariance matrix');
17 imagesc(covar);

```

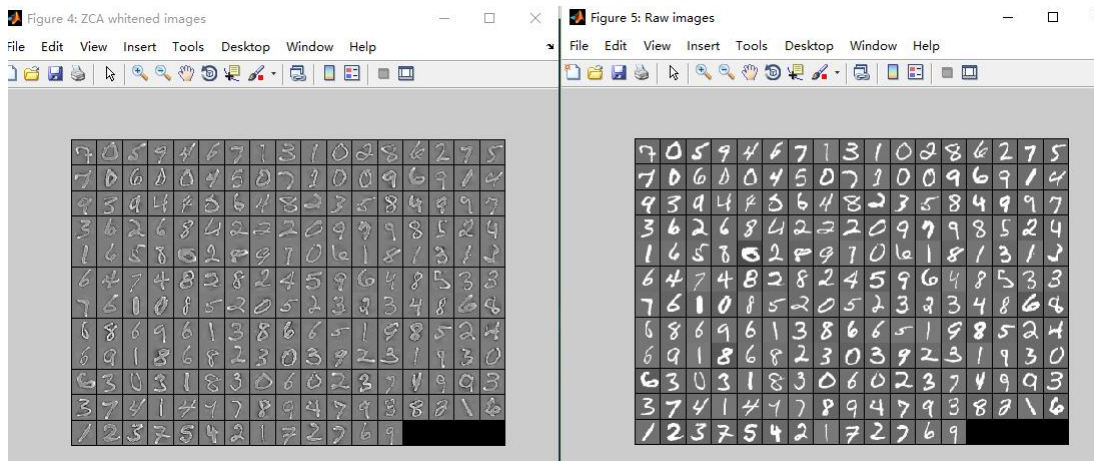


Listing 2.11: Implement ZCA whitening

```

7.
1 xZCAWhite = zeros(size(x));
2 % xZCAWhite <784*60000>
3 xZCAWhite = u * xPCAWhite;
4 % Visualise the data, and compare it to the raw data.
5 % You should observe that the whitened images have enhanced edges.
6 figure('name','ZCA whitened images');
7 display_network(xZCAWhite(:,randsel));
8 figure('name','Raw images');
9 display_network(x(:,randsel));

```

2.3 Reconstruction ICA

Independent Component Analysis (ICA) allows us to generate sparse representations of whitened data. But ICA, has difficulties arise when the number of features (rows of \mathbf{W} matrix), exceed the dimensionality of input, \mathbf{x} .

2.4 Self-raught Learning

I have not completed it, but I know its principles.