

# Project Report for the AI (CS410 2017 Spring)

**Abstract**—this report is a solution of the classification problem under the dataset *Gene\_Chip\_Data*. As a basic problem in machine learning field, there are plenty of mature methods to solve such problem. This study proposed the LR (Logistic Regression) method and the SVM (Support Vector Machine) method to train a classifier for multiclass and binary classification problem. The model also adopted the preprocessing of PCA (principal component analysis) to reduce the dimension of original data. Besides, the study tried to implement a deep learning model called DBN (Deep Belief Network) to do the same work as LR and SVM. By solving the same problem with different methods, we can have simple discussion and comparison between those methods.

**Keywords**—Binary Classification, Multiclass Classification, Logistic Regression, Support Vector Machine, Deep Belief Network

## I. METHODS

This section reviews the usage of the mentioned methods, namely classification and preprocessing techniques.

### A. Principal Component Analysis

PCA is mathematically defined as an orthogonal linear transformation that transforms the data to a new coordinate system such that the greatest variance by some projection of the data comes to lie on the first coordinate (called the first principal component), the second greatest variance on the second coordinate, and so on [1].

Consider a data matrix  $\mathbf{X}$ , where each of the  $n$  rows represents a different repetition of the experiment, and each of the  $p$  columns gives a particular kind of feature.

Mathematically, the transformation is defined by a set of  $p$ -dimensional vectors of weights or loadings  $\mathbf{W}_{(k)} = (\mathbf{w}_1, \dots, \mathbf{w}_p)_{(k)}$  that map each row vector  $\mathbf{x}_{(i)}$  of  $\mathbf{X}$  to a new vector of principal component scores  $\mathbf{t}_{(i)} = (\mathbf{t}_1, \dots, \mathbf{t}_m)_{(i)}$ , given by

$$\mathbf{t}_{k(i)} = \mathbf{X}_{(i)} \cdot \mathbf{W}_{(k)} \quad \text{for } i = 1, \dots, n \quad k = 1, \dots, m$$

In such a way that the individual variables of  $\mathbf{t}$  considered over the data set successively inherit the maximum possible variance from  $\mathbf{x}$ , with each loading vector  $\mathbf{w}$  constrained to be a unit vector.

In order to maximize variance, the first loading vector  $\mathbf{w}_{(1)}$  thus has to satisfy

$$\mathbf{w}_{(1)} = \arg \max_{\|\mathbf{w}\|=1} \left\{ \sum_i (\mathbf{t}_1)_{(i)}^2 \right\} = \arg \max_{\|\mathbf{w}\|=1} \left\{ \sum_i (\mathbf{x}_{(i)} \cdot \mathbf{w})^2 \right\}$$

Equivalently, writing this in matrix form gives

$$\mathbf{w}_{(1)} = \arg \max_{\|\mathbf{w}\|=1} \{ \|\mathbf{X}\mathbf{w}\|^2 \} = \max_{\|\mathbf{w}\|=1} \{ \mathbf{w}^T \mathbf{X}^T \mathbf{X} \mathbf{w} \}$$

Since  $\mathbf{w}_{(1)}$  has been defined to be a unit vector, it equivalently also satisfies

$$\mathbf{w}_{(1)} = \arg \max \left\{ \frac{\mathbf{w}^T \mathbf{X}^T \mathbf{X} \mathbf{w}}{\mathbf{w}^T \mathbf{w}} \right\}$$

A standard result for a symmetric matrix such as  $\mathbf{X}^T \mathbf{X}$  is that the quotient's maximum possible value is the largest eigenvalue of the matrix, which occurs when  $\mathbf{w}$  is the corresponding eigenvector.

The  $k$ -th component can be found by subtracting the first  $k-1$  principal components from  $\mathbf{X}$ :

$$\hat{\mathbf{X}}_k = \mathbf{X} - \sum_{s=1}^{k-1} \mathbf{X} \mathbf{w}_{(s)} \mathbf{w}_{(s)}^T$$

And then finding the loading vector which extracts the maximum variance from this new data matrix

$$\mathbf{w}_{(k)} = \arg \max_{\|\mathbf{w}\|=1} \{ \|\hat{\mathbf{X}}_k \mathbf{w}\|^2 \} = \arg \max \left\{ \frac{\mathbf{w}^T \hat{\mathbf{X}}_k^T \hat{\mathbf{X}}_k \mathbf{w}}{\mathbf{w}^T \mathbf{w}} \right\}$$

It turns out that this gives the remaining eigenvectors of  $\mathbf{X}^T \mathbf{X}$ , with the maximum values for the quantity in brackets given by their corresponding eigenvalues. Thus the loading vectors are eigenvectors of  $\mathbf{X}^T \mathbf{X}$ .  $\mathbf{X}^T \mathbf{X}$  itself can be recognized as proportional to the empirical sample covariance matrix of the dataset  $\mathbf{X}$ .

The following is a detailed description of PCA using the covariance method as opposed to the correlation method [2]. And our study also use this method. But note that it is better to use the singular value decomposition (using standard software). The goal is to transform a given data set  $\mathbf{X}$  of dimension  $p$  to an alternative data set  $\mathbf{Y}$  of smaller dimension  $L$ :

$$\mathbf{Y} = \text{PCA}(\mathbf{X})$$

1) Calculate the empirical mean and the deviations from the mean

Find the empirical mean along each column  $j = 1, \dots, p$ . Place the calculate mean values into an empirical mean vector  $\mathbf{u}$  of dimensions  $p \times 1$

$$\mathbf{u}[j] = \frac{1}{n} \sum_{i=1}^n \mathbf{X}[i, j]$$

Mean subtraction is an integral part of the solution towards finding a principal component basis that minimizes the mean square error of approximating the data. Hence we should center the data by subtract the empirical mean vector  $\mathbf{u}$  from each row of the data matrix  $\mathbf{X}$  and store the result in the  $n \times p$  matrix  $\mathbf{B}$

$$\mathbf{B} = \mathbf{X} - \mathbf{h} \mathbf{u}^T$$

Where  $\mathbf{h}$  is an  $n \times 1$  column vector of all 1s

2) Find the covariance matrix

Find the  $p \times p$  empirical covariance matrix  $\mathbf{C}$  from the product of matrix  $\mathbf{B}$  with its transpose:

$$\mathbf{C} = \frac{1}{n} \mathbf{B} \cdot \mathbf{B}^T$$

3) Find the eigenvectors and eigenvalues of the covariance matrix

Compute the matrix  $\mathbf{V}$  of eigenvectors which diagonalizes the covariance matrix  $\mathbf{C}$ :

$$\mathbf{V}^{-1} \mathbf{C} \mathbf{V} = \mathbf{D}$$

Matrix  $\mathbf{D}$  will take the form of an  $p \times p$  diagonal matrix, where

$$\mathbf{D}[k, l] = \begin{cases} \lambda_k, & k = l \\ 0, & k \neq l \end{cases}$$

$\lambda_k$  is the  $k$ -th eigenvalue of the covariance matrix  $\mathbf{C}$ , and matrix  $\mathbf{V}$ , also of dimension  $p \times p$ , contains  $p$  column vectors, each of length  $p$ , which represent the  $p$  eigenvectors of the covariance matrix  $\mathbf{C}$

4) Rearrange the eigenvectors and eigenvalues

Sort the columns of the eigenvector matrix  $\mathbf{V}$  and eigenvalue matrix  $\mathbf{D}$  in order of decreasing eigenvalue. Make sure to maintain the correct pairings between the columns in each matrix.

5) Compute the cumulative energy content for each eigenvector

The eigenvalues represent the distribution of the source data's energy among each of the eigenvectors, where the eigenvectors form a basis for the data. The cumulative energy content  $g$  for the  $j$ -th eigenvector is the sum of the energy content across all of the eigenvalues from 1 through  $j$ :

$$g[j] = \sum_{k=1}^j \mathbf{D}[k, k]$$

6) Select a subset of the eigenvectors as basis vectors

Save the first  $L$  columns of  $\mathbf{V}$  as the  $p \times L$  matrix  $\mathbf{W}$ . Use the vector  $g$  as a guide in choosing an appropriate value for  $L$ . The goal is to choose a value of  $L$  as small as possible while achieving a reasonably high value of  $g$  on a percentage basis. For example, you may want to choose  $L$  so that the cumulative energy  $g$  is above a certain threshold, like 90 percent. In this case, choose the smallest value of  $L$  such that

$$\frac{g[L]}{g[p]} \geq 0.9$$

7) Project the data onto the new basis

$$\mathbf{Y} = \mathbf{X} \cdot \mathbf{W}$$

## B. Logistic Regression

In statistics, logistic regression is a regression model where the dependent variable is categorical—that is, where it can take only two values, 0 and 1. The binary logistic model is used to estimate the probability of a binary response based on one or more predictor variables (features). It allows one to say that the presence of a risk factor increases the probability of a given

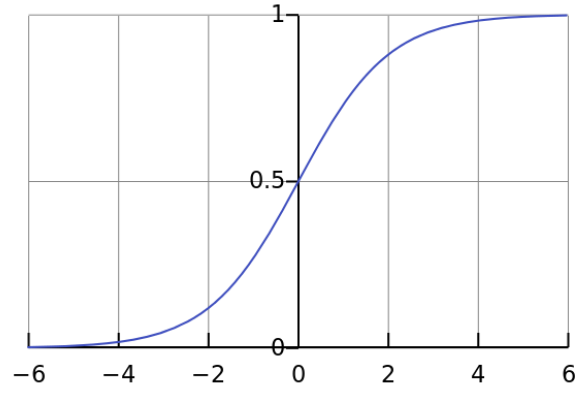


Figure 1. The standard logistic function  $\sigma(t)$

outcome by a specific percentage. In order to fit multiclass need, logistic regression can be binomial, ordinal or multinomial. Binomial or binary logistic regression deals with situations in which the observed outcome for a dependent variable can have only two possible types, “0” and “1”, as this leads to the most straightforward interpretation [3]. If a particular observed outcome for the dependent variable is the noteworthy possible outcome (referred to as a “success” or a “case”) it is usually coded as “1” and the contrary outcome (referred to as a “failure” or a “noncase”) as “0”. Logistic regression is used to predict the odds of being a case based on the values of the independent variables (predictors). The odds are defined as the probability that a particular outcome is a case divided by the probability that it is a noncase.

An explanation of logistic regression can begin with an explanation of the standard logistic function. The logistic function is useful because it can take any real input  $t$  ( $t \in \mathbf{R}$ ), whereas the output always takes values between 0 and 1 and hence is interpretable as a probability. The logistic function  $\sigma(t)$  is defined as follows:

$$\sigma(t) = \frac{1}{1 + e^{-t}}$$

A graph of the logistic function of the  $t$ -interval  $(-6, 6)$  is shown in Figure 1. Let us assume that  $t$  is a linear function of a single explanatory variable  $x$ . We can then express  $t$  as follows:

$$t = \mathbf{w}_0 + \mathbf{w}_1 x_1 + \dots + \mathbf{w}_n x_n = \mathbf{W}x$$

And the logistic function can now be written as:

$$F(x) = \frac{1}{1 + e^{-\mathbf{W}x}}$$

## C. Support Vector Machine

Support vector machines are used for clustering data into two categories according to maximum boundary geometry. SVMs are supervised learning models associated with learning algorithms and used to analyze data and recognize patterns. In the SVM training algorithm, new examples are assigned to one category or another as either nonlinear or linear binary classifiers obtained from a set of training examples. Generally, SVMs provide a suitable means of clustering data for small data sets especially, because the classification is based on support vectors, and data dimension-size ratio has no effect of model

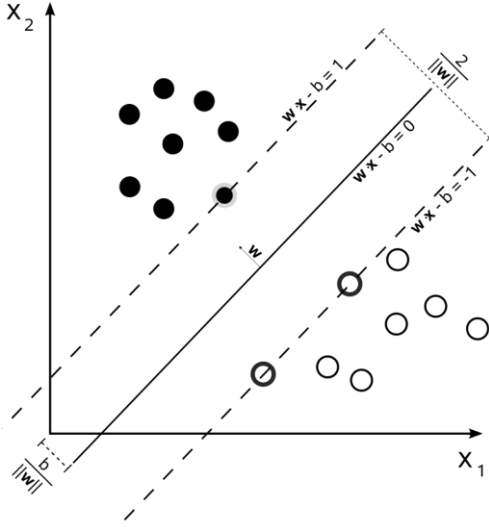


Figure 2. Support vector machine

complexity.

A graph of the SVM is shown in figure 2. Maximum-margin hyperplane and margins for an SVM trained with samples from two classes. Samples on the margin are called the support vectors.

Given a training dataset  $D(x_i, y_i)$ , where  $x_i$  denotes  $n$  observations of malware signatures,  $x_i \in \mathbb{R}^N, i = 1, \dots, N$ ; and  $y_i$  is the corresponding class label whose value is either 1 or -1 (i.e., malicious or benign), indicating the class to which the point  $x_i$  belongs,  $y_i \in \{1, -1\}$ , assigned to each observation  $x_i$ . Each behavioural signature  $x_i$  is of dimension  $d$  corresponding to the number of propositional variables.

$$D = \{(x_i, y_i) | x_i \in \mathbb{R}^N, y_i \in \{1, -1\}\}_{i=1}^N$$

A typical clustering problem is identifying the maximum margin of a hyperplane that divides the points exhibiting  $y_i = 1$  from those exhibiting  $y_i = -1$ . Any hyperplane can be written as the set of points,  $x$ , satisfying the following formula:

$$w \cdot x_i + b = 0 \quad \forall i$$

where  $\cdot$  denotes the inner product and  $W$  denotes the normal vector of the hyperplane. The parameter  $\frac{b}{||w||}$  determines the offset of a hyperplane from the origin along the normal vector  $W$ . Generally, a decision function  $D(x)$  is defined for clustering as  $D(x) = w \cdot x + b$ .

The basic idea is to find a set of support vectors which define the widest linear margin between two classes. A traditional linear SVM has a key drawback, i.e., assuming that the training data are linearly separable, it cannot hold true when applied to practical real cases. Bernhard et al.[4] suggested a novel approach to generating nonlinear classifiers by applying a kernel function to maximum-margin hyperplanes. The nonlinear classification algorithm is formally similar to the linear SVM, except that each inner product  $(\phi(x_i) \cdot \phi(x_j))$ , i.e) is replaced by a kernel function. This enables the algorithm to fit the maximum margin hyperplane in a transformed feature space.

SVM can also be used in multiclass classification problems by modifying the dimension of output value. If the dataset has  $n$  different labels, then we need to modify the output to hold  $n$

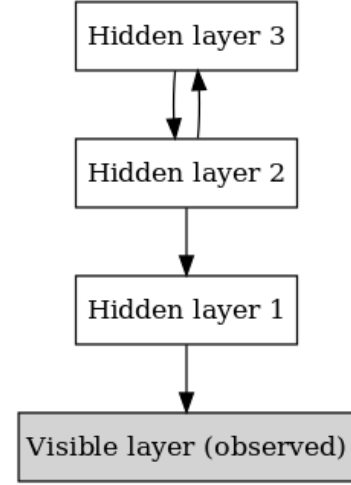


Figure 3. Schematic overview of a DBN

dimensions, making the corresponding dimension value be 1 and others be -1.

#### D. Deep Belief Network

Deep belief network is a generative graphical model, or alternatively a type of deep neural network, composed of multiple layers of latent variables (hidden units), with connections between the layers but not between units within each layer. DBN can be viewed as a composition of simple, unsupervised networks of restricted Boltzmann machines (RBMs), where each sub-network's hidden layer serves as the visible layer for the next. This also leads to a fast, layer-by-layer unsupervised training procedure.

A scheme of a DBN is shown in figure 3, and the arrows represent directed connections in the graphical model that the net represents. The training algorithm for DBNs proceeds as follows [5]. Let  $X$  be a matrix of inputs, regarded as a set of feature vectors.

- 1) Train a RBM on  $X$  to obtain its weight matrix  $W$ . Use this as the weight matrix between the lower two layers of the network.
- 2) Transforms  $X$  by the RBM to produce new data  $X'$ , either by sampling or by computing the mean activation of the hidden units.
- 3) Repeat this procedure with  $X \leftarrow X'$  for the next pair of layers, until the top two layers of the network are reached.
- 4) Fine-tune all the parameters of this deep architecture with respect to a proxy for the DBN log-likelihood, or with respect to a supervised training criterion (after adding extra learning machinery to convert the learned representation into supervised predictions, e.g. the logistic regression classifier).

## II. RESULTS

In this section, the implementations and the results of the models above will be represented. The experiment environment is shown in the table 1:

Table 1: Experiment environment

Operating System	Microsoft Windows 10 Professional
RAM	8GB
Graphics Device	NVIDIA GeForce GT 755M
CPU	Inter i5-4200M 2.50GHz
Platform	Python 3.5 + tensorflow 1.0.1

In the *Gene\_Chip\_Data* dataset, file *microarray.oringinal.txt* is the data file, which saves 22283 rows and 5894 columns matrix data. Every row shows the feature and every column shows the observation. In another word, there are 5894 examples and every example has 22283 features. Another file *E-TBM-185.sdrf.txt* saves the labels for examples from *[material]* to *[disease status]*, etc. Our experiment uses the *[material]* labels as the binary classification standard and uses the *[disease status]* labels as the multiclass classification standard. There are 2 different label, *organism\_part* & *cell\_line* under the *[material]* and 194 kinds of disease status labels under *[disease status]*. And in order to test the performance of multiclass classifier and binary classifier, we conducted two kind of experiment for every methods. Multiclass test requires classifiers to tell all 194 kinds of status, and binary class test only requires classifiers to tell whether the observations comes from organism or cell.

As for the testing, we use the **k-fold cross-validation** method to test our models. The original data is divided into  $k$  groups and we make  $k-1$  groups be the training set and 1 group left be the test set. Then change the training set and test set in turn for  $k$  times so that every group has the chance to be tested. Finally we will calculate the average accuracy of  $k$  times cross-validations.

#### A. PCA & LR

Since the number of features in the dataset is more than 20000, we should use PCA first to reduce the dimension. We set the contribution ratio of PCA as 0.9 and we finally get 112 main features. Then we could use the logistic regression model to fit the reduced data. Besides, it's important to transfer the label to vectors. We used one-hot code to represent different kind of labels. For example, we could use  $[1, 0]$  to represent *organism\_part* and use  $[0, 1]$  to represent *cell\_line*. Figure 4 is the result of multiclass classification and binary classification. We trained the model for 100 iterations and every iteration has 100 batches.

#### B. PCA & SVM

In this experiment, we used soft margin SVM with  $C = 0.9$  rather than hard margin SVM. Same as PCA & LR, we first used PCA to reduce the dimension and then used SVM model to fit the data. However, instead of using original one-hot code to represent labels, we use the methods that similar to one-hot code, which uses -1 to replace 0 in one-hot code. For example, we use  $[1, -1]$  to represent *organism\_part* and use  $[-1, 1]$  to represent *cell\_line*. Figure 5 is the result of multiclass and binary classification. We trained the model for 100 iterations

and every iteration has 100 batches.

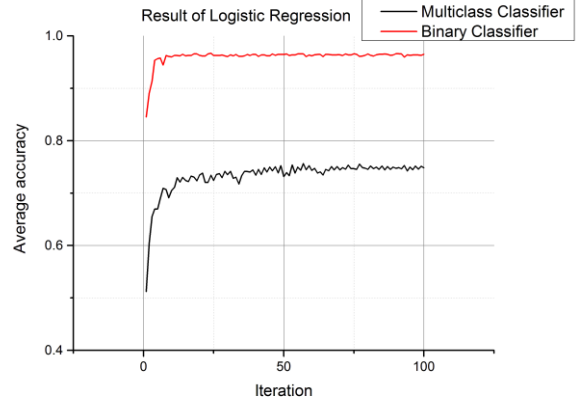


Figure 4. Result of LR

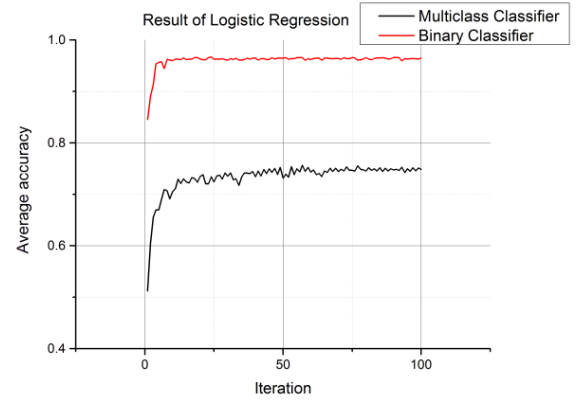


Figure 5. Result of SVM

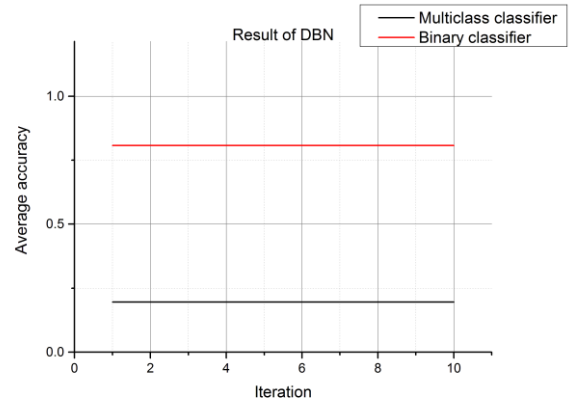


Figure 6. Result of DBN

#### C. DBN

When we focus on Figure 6, we can find that DBN model has so perfect convergence that the accuracy of every epoch is almost the same. We think the reason is pre-training of deep learning. However, because the number of observations under every label is much too insufficient that it can't satisfy the standard of deep learning, it performs poor in classification problems. The accuracy of multiclass classification problem is only about 0.2 and the accuracy of binary classification

problem is about 0.8, which is far away from LR and SVM. It may do well if there are enough data.

### III. DISCUSSIONS

In this section, we decide to discuss the results shown in last section and make some comparisons between different classifiers. At first we want to claim the reason why we used the *[material]* label as the standard of binary classification problem rather than continuously using *[disease status]*. In the dataset, there are 4752 observations with *organism\_part* label and 1142 observations with *cell\_line* label. The ratio is about 4:1. We would have decided to use healthy and unhealthy as the standard. However, there are only 85 healthy observations and 4917 of all observations are unhealthy, where the ratio is about 58:1. The incommensurable ratio will cause big error in statistics. So we changed the standard.

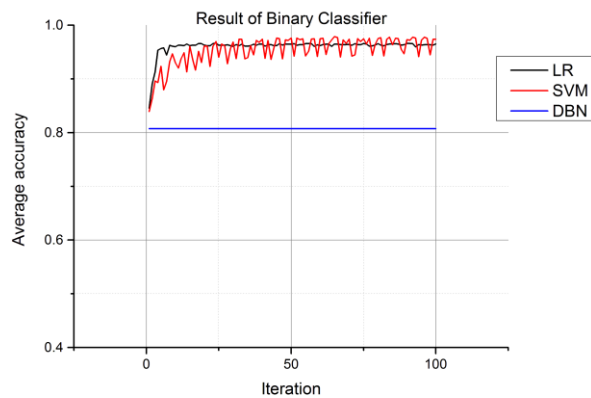


Figure 7. Result of Binary Classifier

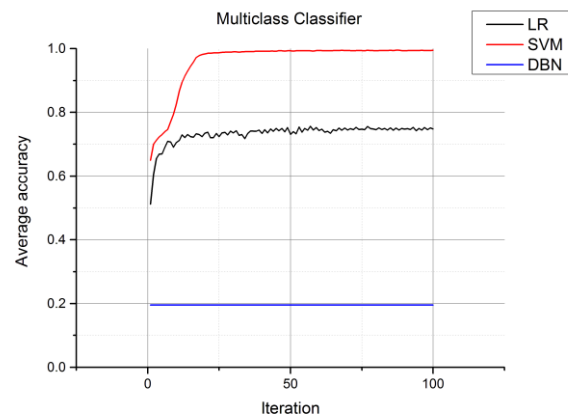


Figure 8. Result of Multiclass Classifier

#### A. LR

When we focus on Figure 4, we can find that LR model did well in binary classification problems and the accuracy reached 0.96. It also had an excellent performance in convergence. However, it didn't have a high accuracy in multiclass classification problems. And it had little fluctuation in the graph.

#### B. SVM

When we focus on Figure 5, it's obvious that VM models perfectly fitted multiclass classification problem that its convergence is excellent and can reach 0.99. However, in the aspect of binary classification problem, its bad convergence is also obvious. It fluctuates drastically around the accuracy of 0.95. But it's undeniable that SVM models have higher accuracy than LR.

#### C. DBN

When we focus on Figure 6, we can find that DBN model has so perfect convergence that the accuracy of every epoch is almost the same. We think the reason is pre-training of deep learning. However, because the number of observations under every label is much too insufficient that it can't satisfy the standard of deep learning, it performs poor in classification problems. The accuracy of multiclass classification problem is only about 0.2 and the accuracy of binary classification problem is about 0.8, which is far away from LR and SVM. It may do well if there are enough data.

#### D. Comparison

According to the comparison in figure 7 and figure 8, we can draw the conclusion that LR performs better when there are fewer target labels while SVM performs better when there are more target labels, and DBN has absolute advantages in convergence while it doesn't have high accuracy because of the restriction of data quantities.

### REFERENCES

- [1] Jolliffe I.T. Principal Component Analysis, Series: Springer Series in Statistics, 2nd ed., Springer, NY, 2002, XXIX, 487 p. 28 illus. ISBN 978-0-387-95442-4
- [2] "Engineering Statistics Handbook Section 6.5.5.2". Retrieved 19 January 2015.
- [3] Hosmer, David W.; Lemeshow, Stanley (2000). Applied Logistic Regression (2nd ed.). Wiley. ISBN 0-471-35632-8.
- [4] E.B. Bernhard, M.G. Isabelle, V. Vapnik, and N. Vladimir (1992). A training algorithm for optimal margin classifiers. In Haussler, David (editor); 5th Annual ACM Workshop on COLT, pages 144–152, Pittsburgh, PA, ACM Press.
- [5] Hinton, G. E.; Osindero, S.; Teh, Y. W. (2006). "A Fast Learning Algorithm for Deep Belief Nets" (PDF). Neural Computation. 18 (7): 1527–1554. doi:10.1162/neco.2006.18.7.1527. PMID 16764513.