

COMPARISON OF UNSUPERVISED PRE-TRAINING METHODS

Jifu Zhao

University of Illinois at Urbana-Champaign
Department of Nuclear, Plasma, and Radiological Engineering
Urbana, Illinois 61801, USA

1. INTRODUCTION

As one important part of artificial intelligence (AI), machine learning has been widely used in our daily life, such as natural language processing (NLP), zip code recognition, self-driving cars and so on. Among those applications, supervised learning is the most common and perhaps the most powerful form of machine learning [1]. For supervised learning, the choice of appropriate representations of original features is very important [2]. For a long time, most machine learning applications require us to specify the input \mathbf{x} and target \mathbf{y} . The choice of feature representation \mathbf{x} is usually determined by experts. However, sometimes the input \mathbf{x} may be just the raw feature without pre-processing.

Researchers have been working on discovering the hidden structure of the data for years. In this process, several methods have been developed and found to be useful in some domains, such as Principal Component Analysis (PCA), Independent Component Analysis (ICA), Nonnegative Matrix Factorization (NMF), Restricted Boltzmann Machine (RBM), Auto-encoders and its variants and so on. These methods work in an unsupervised way and the output from these methods are thought to be a transformation of the raw input features. For supervised learning, some supervised learning mechanism will be built based on the output from these pre-training methods rather than the raw features. As expected, a good representation can help the supervised learning algorithms work better.

In this work, we focus on the comparison of several unsupervised pre-training methods. More specifically, we will focus on PCA, kernel PCA and auto-encoders.

2. METHODS

2.1. PCA

Principal Component Analysis (PCA) is one of the most widely used pre-training methods for many applications. The goal of PCA is to find a lower dimensional linear space such that the variance of the orthogonal projection of the original data is maximized [3]. PCA is widely used for dimensionality reduction, feature extraction and data visualization. (Since

this part has been covered in class, we won't cover too much derivations about PCA. Instead, we will focus on its application.)

2.2. Kernel PCA

As a linear transformation, PCA works well in some applications where the input data can be well represented by linear transformation. However, in a lot of applications, this assumption doesn't hold. We need some non-linear transformation for these scenarios. One direct idea is to map the input \mathbf{x} into some high dimensional space $\phi(\mathbf{x})$ then perform PCA in this high dimensional space. However, sometimes the dimensionality of the new space is so high that the computation PCA becomes very expensive.

For this situation, Schölkopf [4, 5] proposed the method called kernel PCA where the computation can be performed in the original feature space rather than the new space using so called kernel trick. In this way, through kernel PCA, a non-linear transformation makes it possible for us to find some complex representations.

2.3. Auto-encoders

Auto-encoders use the input as the target, trying to find a reconstruction $r(\mathbf{x}) = g(h(\mathbf{x}))$ through so called encoder $h(\cdot)$ and decoder $g(\cdot)$ [2]. In most cases, an auto-encoder is a multiple layer deep neural network which has a bottleneck structure. An 5 layer auto-encoder system is shown in Figure 1.

Auto-encoder system can be trained using back-propagation methods. However, it works well only if we have a good initialization for the auto-encoder systems [6]. Hinton and Salakhutdinov [6] then proposed an effective way of initialization which utilize a stack of RBMs for pre-training. After unrolling, the whole auto-encoder system can be fine-tuned.

3. DATA

In this work, we plan to use the MNIST dataset [7]. The MNIST dataset contains a huge amount of handwritten digits. It has 60000 training examples (more precisely, 55000 training examples, 5000 validation examples) and 10000 test ex-

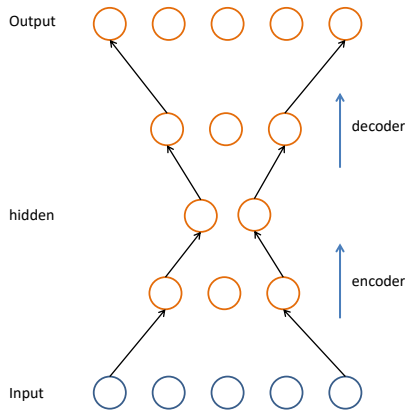


Fig. 1. Illustration of auto-encoder system.

amples. It has 10 classes which corresponding to digits from 0 to 9. Each example is a 28×28 pixel grayscale image. The pixel values range from 0 to 1. A subset of 36 examples is shown in Figure 2.

5	0	4	1	9	2
1	3	1	4	3	5
3	6	1	7	2	8
6	9	4	0	9	1
1	2	4	3	2	7
3	8	6	9	0	5

Fig. 2. Examples of MNIST dataset.

4. EXPERIMENTS

This work is an extension of the assignment 4, where we are required to compare the classification accuracy for raw image, PCA, RBM and stacked RBMs. In this work, we will use similar procedure, but focus on kernel PCA and auto-encoder

systems.

Our work can be roughly divided into two parts. The first part will be focused on the comparison of PCA and kernel PCA. More specifically, we will conduct experiments on PCA and kernel PCA with polynomial kernel and rbf kernel. Through choosing different number of transformed features, we will build a supervised multiclass classifier to classify the dataset into 10 classes. Through calculating the error rate of different methods, we hope to compare the performance of different methods.

The second part will focused on the auto-encoders. Except the study on number of hidden nodes and number of layers, we will also study the influence of pre-training using stacks of RBMs as described by Hinton and Salakhutdinov [6]. Then the same classification strategy will be implemented.

Since the goal of this work is to study different pre-training methods, the final classifier won't be very complicated, a multiclass Support Vector Machine (SVM) or a two-layer multi-class logistic regression neural network will be used.

5. RESOURCE FEASIBILITY

With a modern computer (i.e., Intel i7 processor, 16 RAM), most of the work described above can be implemented without any problem. However, for kernel PCA part, the kernel matrix K has the dimensionality of $n \times n$, where n is the number of given samples. In this work, for the MNIST dataset, there will be 60000 training examples. The dimensionality of K will be 60000×60000 . To store this matrix in memory, it will need more than 25 GB, which is not possible for our current computer. For this part, we will implement the algorithm on Amazon Web Service (AWS).

6. TIMELINE

We expect to finish the theoretical derivation part before the Thanksgiving. Then focus on actual programming part after the Thanksgiving (11/24/2016). For programming part, we will implement kernel PCA through Python3 + Numpy. For auto-encoder part, we will implement it through TensorFlow + TFLearn. We expect the conduct the actual experiment part from 11/28/2016 and finish all parts before 12/07/2016.

7. REFERENCES

- [1] Yann LeCun, Yoshua Bengio, and Geoffrey Hinton, “Deep learning,” *Nature*, vol. 521, no. 7553, pp. 436–444, 2015.
- [2] Yoshua Bengio et al., “Deep learning of representations for unsupervised and transfer learning.,” *ICML Unsupervised and Transfer Learning*, vol. 27, pp. 17–36, 2012.
- [3] Christopher Bishop, *Pattern Recognition and Machine Learning (Information Science and Statistics)*, 1st edn. 2006. corr. 2nd printing edn, Springer, New York, 2007.
- [4] Bernhard Schölkopf, Alexander Smola, and Klaus-Robert Müller, “Nonlinear component analysis as a kernel eigenvalue problem,” *Neural computation*, vol. 10, no. 5, pp. 1299–1319, 1998.
- [5] Jerome Friedman, Trevor Hastie, and Robert Tibshirani, *The elements of statistical learning*, vol. 1, Springer series in statistics Springer, Berlin, 2001.
- [6] Geoffrey E Hinton and Ruslan R Salakhutdinov, “Reducing the dimensionality of data with neural networks,” *Science*, vol. 313, no. 5786, pp. 504–507, 2006.
- [7] Yann LeCun, Léon Bottou, Yoshua Bengio, and Patrick Haffner, “Gradient-based learning applied to document recognition,” *Proceedings of the IEEE*, vol. 86, no. 11, pp. 2278–2324, 1998.