
Paper Review on Kernel Analysis of Deep Networks

Yang Fu
Department of Statistics
University of Washington
Seattle, WA 98195
fuyang@uw.edu

1 Introduction

This is a paper review on *Kernel Analysis of Deep Networks*. In the original paper, depth and structure of deep neural networks are analysed via kernel principal component analysis. It verifies two hypotheses.

One is that as the network gets deeper, simpler and more accurate representations are obtained. Here simpleness is measured by the number of principal components used and accuracy is measured by the error rate of the softmax classifier built on those principal components.

The other is that the structure of the network controls how fast the representation of the task is formed layer after layer. It means that the error rate for networks with different structure will decrease in different paths as the number of principal components grows although the networks will have the same error rate as we control them to.

In this review, section 2 will prove that the kernel used in the paper is positive semi-definite. Section 3 will introduce the kernel-based method, explain how it is solved with a numerical algorithm and check the convergence of the algorithm. The last section repeats some of the experiments on the artificial data, MLP and CNN.

2 Kernel

In the paper, the Gaussian RBF kernel is used. In this section, it will be shown that RBF kernel is a positive semi-definite kernel.

2.1 Definition

A Gaussian RBF kernel with hyper-parameter $\sigma > 0$ on two sample $x, x' \in \mathbb{R}^d$ is defined as

$$k(x, x') = \exp\left(-\frac{\|x - x'\|^2}{\sigma^2}\right).$$

2.2 Positive semi-definite

Lemma 1. If $k_i, i = 1, 2, \dots$ are positive semi-definite kernels, then $\lim_{i \rightarrow \infty} k_i$ is also a positive semi-definite kernel if the limit exists.

Proof By the definition of positive semi-definite kernel, for all i ,

$$g_i = \sum_{k, \ell} \alpha_k \alpha_\ell k_i(x_k, x_\ell) \geq 0,$$
$$\therefore \lim_{i \rightarrow \infty} g_i \geq 0.$$

25

Lemma 2. If h is a positive semi-definite kernel, and $f : \mathbb{R}^d \rightarrow \mathbb{R}$ is a function. Then

$$k(x, x') = f(x)h(x, x')f(x')$$

is also a positive semi-definite kernel.

Proof $\forall \alpha_1, \dots, \alpha_n \in \mathbb{R}$,

$$\sum_{i,j} \alpha_i \alpha_j f(x_i) h(x_i, x_j) f(x_j) = \sum_{i,j} \beta_i \beta_j h(x_i, x_j) \geq 0,$$

where $\beta_i = \alpha_i f(x_i) \in \mathbb{R}$. ■

Theorem 3. Gaussian RBF kernel $k(\cdot, \cdot)$ is a positive semi-definite kernel.

Proof Without loss of generality, assume $\sigma = \sqrt{2}$. $\forall x_1, \dots, x_n \in \mathbb{R}^d$,

$$\begin{aligned} \because \|x_i - x_j\|^2 &= x_i^T x_i + x_j^T x_j - 2x_i^T x_j, \\ \therefore k(x_i, x_j) &= \exp\left(-\frac{x_i^T x_i}{2}\right) \exp(x_i^T x_j) \exp\left(-\frac{x_j^T x_j}{2}\right). \end{aligned}$$

By Taylor expansion,

$$\exp(x_i^T x_j) = \sum_{k=0}^{\infty} \frac{(x_i^T x_j)^k}{k!}.$$

Let $h_\ell(x, y) = \sum_{k=0}^{\ell} \frac{(x^T y)^k}{k!}$. Since for each ℓ , $h_\ell(x, y)$ is a finite sum of polynomial kernels, $h_\ell(x, y)$ is also positive semi-definite. By Lemma 1, $\exp(x_i^T x_j)$ is a positive semi-definite kernel and thus, by Lemma 2, RBF kernel is also positive semi-definite. ■

3 Kernel-based method

In the paper, Kernel PCA is used to analyse MLP and CNN. In this section, Kernel PCA will be introduced as well as the algorithm used to solve it.

3.1 Kernel PCA

As a kernel counterpart of regular PCA, Kernel PCA aims at solving the optimization problem

$$\arg \max_{f \in \mathcal{H}} \frac{\langle f, \hat{S} f \rangle_{\mathcal{H}}}{\|f\|_{\mathcal{H}}^2},$$

where \mathcal{H} is RKHS, \hat{S} is the covariance operator defined as

$$\begin{aligned} \hat{S} &= \frac{1}{n} \sum_{i=1}^n [k(x_i, \cdot) - \hat{\mu}] \otimes [k(x_i, \cdot) - \hat{\mu}], \\ \langle f, \hat{S} g \rangle_{\mathcal{H}} &= \frac{1}{n} \sum_{i=1}^n \langle f, k(x_i, \cdot) - \hat{\mu} \rangle_{\mathcal{H}} \langle g, k(x_i, \cdot) - \hat{\mu} \rangle_{\mathcal{H}}, \end{aligned}$$

where $\hat{\mu} = \frac{1}{n} \sum k(x_i, \cdot)$ is the mean element.

After centering the kernel matrix, the optimization problem becomes

$$\arg \max_{f \in \mathcal{H}} \frac{\sum_i f^2(x_i)}{\|f\|_{\mathcal{H}}^2}.$$

40 Since the objective function takes $x_1, \dots, x_n, \|f\|_{\mathcal{H}}$ as inputs and is strictly decreasing in the
 41 $(i+1)$ -th variable $\|f\|_{\mathcal{H}}$, the Representer Theorem can be invoked and the solution has the form
 42 $f(x) = \sum_{i=1}^n \alpha_i k(x_i, x)$, where $\alpha_i \in \mathbb{R}$.

Thus, the optimization problem becomes

$$\arg \max_{\alpha \in \mathbb{R}^n} \frac{\alpha^T \tilde{K}^2 \alpha}{\alpha^T \tilde{K} \alpha},$$

43 where \tilde{K} is the centered kernel matrix.

Now we show how the problem is equivalent to finding the eigenvector correspond to the largest eigenvalue. By the Spectral Theorem,

$$\tilde{K} = \sum_{p=1} \lambda_p e_p e_p^T,$$

44 where e_p are unit-length eigenvectors and λ_p are eigenvalues in decreasing order.

45 Let $\alpha = \sum_{j=1}^n \beta_j e_j$, since e_i are orthogonal and unit-length,

$$\begin{aligned} \frac{\alpha^T \tilde{K}^2 \alpha}{\alpha^T \tilde{K} \alpha} &= \frac{\sum_j \beta_j^2 \lambda_j^2}{\sum_j \beta_j^2 \lambda_j} \equiv Obj, \\ \frac{\partial Obj}{\partial \beta_k} &= \frac{2\beta_k \lambda_k^2 (\sum_j \beta_j^2 \lambda_j) - 2\beta_k \lambda_k (\sum_j \beta_j^2 \lambda_j^2)}{(\sum_j \beta_j^2 \lambda_j)^2} = 0, \\ \Leftrightarrow \lambda_k &= \frac{\sum_j \beta_j^2 \lambda_j^2}{\sum_j \beta_j^2 \lambda_j}, k = 1, \dots, n. \end{aligned}$$

46 That is to say, the local maximums of the objective function are the eigenvalues. Thus, the universal
 47 maximum is the largest eigenvalue λ_1 , and

$$\begin{aligned} \frac{\alpha^T \tilde{K}^2 \alpha}{\alpha^T \tilde{K} \alpha} &= \lambda_1, \\ \Leftrightarrow \alpha^T \tilde{K} (\tilde{K} \alpha - \lambda_1 \alpha) &= 0. \end{aligned}$$

48 Thus, α is the eigenvector correspond to λ_1 . Normalize α so that $\alpha^T \tilde{K} \alpha = 1$ and we have $\alpha =$
 49 $\frac{1}{\sqrt{\lambda_1}} e_1$.

50 The Kernel PCA algorithm can be summerized as follow,

Algorithm 1 Kernel PCA

Input: Kernel Matrix K ;

Output: Principal component α ;

- 1: Center kernel matrix $K = (I - \frac{1}{n} \mathbf{1}\mathbf{1}^T) K (I - \frac{1}{n} \mathbf{1}\mathbf{1}^T)$;
 - 2: Compute eigenvectors e_j , eigenvalues λ_j of K ;
 - 3: The j -th principal component is the normalized eigenvector $\alpha_{(j)} = \frac{1}{\sqrt{\lambda_j}} e_j$.
-

51 3.2 Power Method

52 The algorithm used to solve eigen decomposition is called the Power Method.

Algorithm 2 Power Method

Input: Matrix A ;**Output:** Eigenvector v , eigenvalues λ ;1: Randomly initialize $v^{(0)}$;2: **repeat**3: $z^{(t)} = Av^{(t-1)}$;4: $v^{(t)} = \frac{z^{(t)}}{\|z^{(t)}\|_2}$;5: $\lambda^{(t)} = (v^{(t)})^T Av^{(t)}$;6: **until** Max iteration7: Project A orthogonal to previous v : $A = (I - vv^T)A(I - vv^T)$;8: Redo previous steps to get other eigenvectors and values.

The convergence property of the Power Method can be examined. Let A be a diagonalizable matrix and u_1, \dots, u_n be the eigenvectors. Since the eigenvectors form a basis of \mathbb{R}^n , the initialized vector v_0 in the power method can be written as

$$v^{(0)} = \sum_{i=1}^n a_i u_i,$$

53 where a_i are scalars. Multiply both sides by A^t , we have

$$\begin{aligned} A^t v^{(0)} &= \sum_{i=1}^n a_i A^t u_i \\ &= \sum_{i=1}^n a_i \lambda_i^t u_i \\ &= a_1 \lambda_1^t \left(u_1 + \sum_{i=2}^n \frac{a_i}{a_1} \left(\frac{\lambda_i}{\lambda_1} \right)^t u_i \right), \end{aligned}$$

54 where λ_1 is the dominant eigenvalue. Thus, $\left(\frac{\lambda_i}{\lambda_1} \right)^t \rightarrow 0$ and $A^t v^{(0)} \rightarrow a_1 \lambda_1^t u_1$.

55 In order to scale things up, Simultaneous Power Method is used, which is the same as (non-
56 simultaneous) Power Method in nature. Except that, instead of a vector, it multiplies A by multiple
57 vectors put in a matrix Q , where $Q^T Q = I$. And each step the vectors are normalized with QR
58 decomposition. The convergence property is the same as the non-simultaneous version.

59 4 Experiments

60 In the experiments on MNIST-10K, the deep networks are trained on the 10000 samples until a
61 training error rate of 2.5% is reached, same as the original paper. Then, kernel principal components
62 of different layers of the networks are estimated and fed into a softmax classifier with no penalty to
63 get the error rates. For each layer (0,1,2), each network (MLP, CNN) and each number of principal
64 components (1-10), a grid of hyper-parameters of the RBF kernel is tested on and the hyper-parameter
65 that gives the lowest error rate of the softmax classifier is chosen.

66 The parameters of the deep networks are exactly the same as the original paper except that ReLU is
67 used instead of sigmoid to get a shorter training time.

68 4.1 Artificial data

69 The experiment on the artificial data coincides with the original paper's conclusion. Fig 1 shows that
70 it does take more principal components for the more distorted data to be correctly classified.

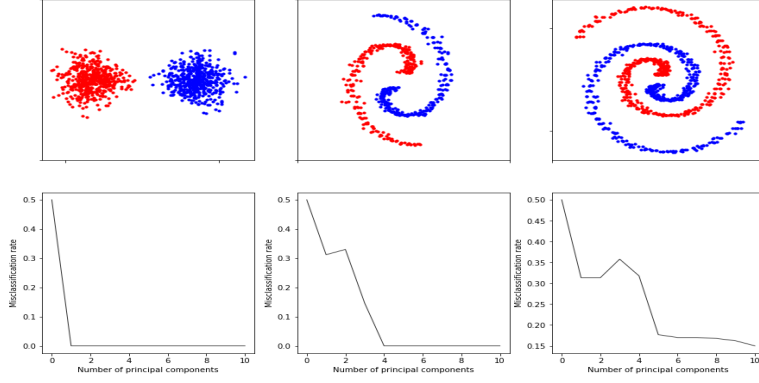


Figure 1: Results on artificial data

4.2 MNIST-10K

The experiment in the original paper to confirm the first hypothesis is repeated on MNIST-10K. Fig 2 shows that the result is very similar to the one in the original paper and the hypothesis is confirmed. Better representations are built layer after layer.

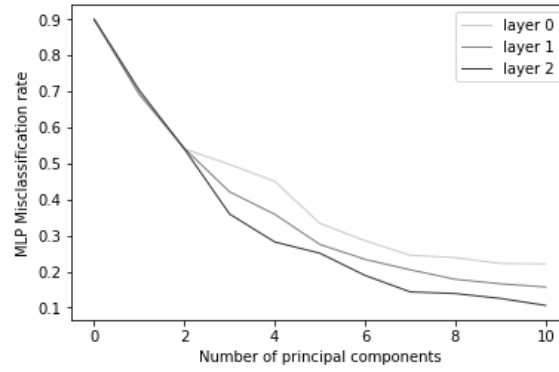


Figure 2: Layer-wise evolution of the error (MLP)

The experiment to confirm the second hypothesis is also repeated except for the pretrained multilayer perceptron (PMLP) and transfer learning part. The result is less clear compared to the original paper. In the original paper, CNN and MLP have very different paths for layer-wise evolution of the representation. Fig 3 shows that the path is not that different.

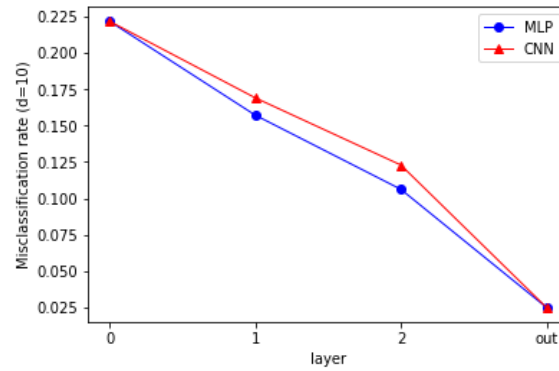


Figure 3: Layer-wise evolution of the error obtained for each training procedure for $d = 10$

79 **5 Conclusion**

80 The original paper proposes a new way to analyse deep networks via kernel PCA and certain
81 experiments is carried out to test two hypotheses.

82 The first hypothesis that as the network gets deeper, simpler and more accurate representations are
83 obtained is confirmed while the second hypothesis about the structure of the networks is less confirmed
84 to me in this paper review. Experiments on more different datasets and networks are needed in order
85 to make that conclusion.

86 **References**

87 [1] Montavon, G., Braun, M.L. & Müller, K.R. (2011) Kernel analysis of deep networks. *Journal of Machine*
88 *Learning Research* **12**(3):2563-2581.

89 [2] Convolutional Neural Networks (2018, January 27). Retrieved from [https://www.tensorflow.org/](https://www.tensorflow.org/tutorials/deep_cnn)
90 [tutorials/deep_cnn](https://www.tensorflow.org/tutorials/deep_cnn)

91 [3] Power Iteration (2017, June 26). Retrieved from http://mlwiki.org/index.php/Power_Iteration