






Backprojection for Training Feedforward Neural Networks in the Input and Feature Spaces

Benyamin Ghojogh^(✉) , Fakhri Karray^(ID) , and Mark Crowley^(ID) 

Department of Electrical and Computer Engineering, University of Waterloo,
Waterloo, ON, Canada
{bghojogh,karray,mcrowley}@uwaterloo.ca

Abstract. After the tremendous development of neural networks trained by backpropagation, it is a good time to develop other algorithms for training neural networks to gain more insights into networks. In this paper, we propose a new algorithm for training feedforward neural networks which is fairly faster than backpropagation. This method is based on projection and reconstruction where, at every layer, the projected data and reconstructed labels are forced to be similar and the weights are tuned accordingly layer by layer. The proposed algorithm can be used for both input and feature spaces, named as backprojection and kernel backprojection, respectively. This algorithm gives an insight to networks with a projection-based perspective. The experiments on synthetic datasets show the effectiveness of the proposed method.

Keywords: Neural network · Backprojection · Kernel backprojection · Projection · Training

1 Introduction

In one of his recent seminars, Geoffrey Hinton mentioned that after all of the developments of neural networks [1] and deep learning [2], perhaps it is time to move on from backpropagation [3] to newer algorithms for training neural networks. Especially, now that we know why shallow [4] and deep [5] networks work very well and why local optima are fairly good in networks [6], other training algorithms can help improve the insights into neural nets. Different training methods have been proposed for neural networks, some of which are backpropagation [3], genetic algorithms [7, 8], and belief propagation as in restricted Boltzmann machines [9].

A neural network can be viewed from a manifold learning perspective [10]. Most of the spectral manifold learning methods can be reduced to kernel principal component analysis [11] which is a projection-based method [12]. Moreover, at its initialization, every layer of a network can be seen as a random projection [13]. Hence, a promising direction could be a projection view of training neural

networks. In this paper, we propose a new training algorithm for feedforward neural networks based on projection and *backprojection* (or so-called reconstruction). In the backprojection algorithm, we update the weights layer by layer. For updating a layer m , we project the data from the input, until the layer m . We also backproject the labels of data from the last layer to the layer m . The projected data and backprojected labels at layer m should be equal because in a perfectly trained network, projection of data by the entire layers should result in the corresponding labels. Thus, minimizing a loss function over the projected data and backprojected labels would correctly tune the layer's weights. This algorithm is proposed for both the input and feature spaces where in the latter, the kernel of data is fed to the network.

2 Backprojection Algorithm

2.1 Projection and Backprojection in Network

In a neural network, every layer without its activation function acts as a linear projection. Without the nonlinear activation functions, a network/autoencoder is reduced to a linear projection/principal component analysis [12]. If \mathbf{U} denotes the projection matrix (i.e., the weight matrix of a layer), $\mathbf{U}^\top \mathbf{x}$ projects \mathbf{x} onto the column space of \mathbf{U} . The reverse operation of projection is called reconstruction or backprojection and is formulated as $\mathbf{U}\mathbf{U}^\top \mathbf{x}$ which shows the projected data in the input space dimensionality (note that it is $\mathbf{U}\mathbf{f}^{-1}(\mathbf{f}(\mathbf{U}^\top \mathbf{x}))$ if we have a nonlinear function $\mathbf{f}(\cdot)$ after the linear projection). At the initialization, a layer acts as a random projection [13] which is a promising feature extractor according to the Johnson-Lindenstrauss lemma [14]. Fine tuning the weights using labels makes the features more useful for discrimination of classes.

2.2 Definitions

Let us have a training set $\mathcal{X} := \{\mathbf{x}_i \in \mathbb{R}^d\}_{i=1}^n$ and their one-hot encoded labels $\mathcal{Y} := \{\mathbf{y}_i \in \mathbb{R}^p\}_{i=1}^n$ where n , d , and p are the sample size, dimensionality of data, and dimensionality of labels, respectively. We denote the dimensionality or the number of neurons in layer m by d_m . By convention, we have $d_0 := d$ and $d_{n_\ell} = p$ where n_ℓ is the number of layers and p is the dimensionality of the output layer. Let the data after the activation function of the m -th layer be denoted by $\mathbf{x}^{(m)} \in \mathbb{R}^{d_m}$. Let the projected data in the m -th layer be $\mathbb{R}^{d_m} \ni \mathbf{z}^{(m)} := \mathbf{U}_m^\top \mathbf{x}^{(m-1)}$ where $\mathbf{U}_m \in \mathbb{R}^{d_{m-1} \times d_m}$ is the weight matrix of the m -th layer. Note that $\mathbf{x}^{(m)} = \mathbf{f}_m(\mathbf{z}^{(m)})$ where $\mathbf{f}_m(\cdot)$ is the activation function in the m -th layer. By convention, $\mathbf{x}^{(0)} := \mathbf{x}$. The data are projected and passed through the activation functions layer by layer; hence, $\mathbf{x}^{(m)}$ is calculated as:

$$\mathbb{R}^{d_m} \ni \mathbf{x}^{(m)} := \mathbf{f}_m(\mathbf{U}_m^\top \mathbf{f}_{m-1}(\mathbf{U}_{m-1}^\top \cdots \mathbf{f}_1(\mathbf{U}_1^\top \mathbf{x}))) = \mathbf{f}_m(\mathbf{U}_m^\top \mathbf{x}^{(m-1)}). \quad (1)$$

In a mini-batch gradient descent set-up, let $\{\mathbf{x}_i\}_{i=1}^b$ be a batch of size b . For a batch, we denote the outputs of activation functions at the m -th layer by $\mathbb{R}^{d_m \times b} \ni \mathbf{X}^{(m)} := [\mathbf{x}_1^{(m)}, \dots, \mathbf{x}_b^{(m)}]$.

Now, consider the one-hot encoded labels of batch, denoted by $\mathbf{y} \in \mathbb{R}^p$. We take the inverse activation function of the labels and then reconstruct or *backproject* them to the previous layer to obtain $\mathbf{y}^{(n_\ell-1)}$. We do similarly until the layer m . Let $\mathbf{y}^{(m)} \in \mathbb{R}^{d_m}$ denote the backprojected data at the m -th layer, calculated as:

$$\mathbf{y}^{(m)} := \mathbf{U}_{m+1} \mathbf{f}_{m+1}^{-1}(\mathbf{U}_{m+2} \mathbf{f}_{m+2}^{-1}(\cdots \mathbf{U}_{n_\ell} \mathbf{f}_{n_\ell}^{-1}(\mathbf{y}))) = \mathbf{U}_{m+1} \mathbf{f}_{m+1}^{-1}(\mathbf{y}^{(m+1)}). \quad (2)$$

By convention, $\mathbf{y}^{(n_\ell)} := \mathbf{y}$. The backprojected batch at the m -th layer is $\mathbb{R}^{d_m \times b} \ni \mathbf{Y}^{(m)} := [\mathbf{y}_1^{(m)}, \dots, \mathbf{y}_b^{(m)}]$. We use $\mathbf{X} \in \mathbb{R}^{d \times b}$ and $\mathbf{Y} \in \mathbb{R}^{p \times b}$ to denote the column-wise batch matrix and its one-hot encoded labels.

2.3 Optimization

In the backprojection algorithm, we optimize the layers' weights one by one. Consider the m -th layer whose loss we denote by \mathcal{L}_m :

$$\underset{\mathbf{U}_m}{\text{minimize}} \quad \mathcal{L}_m := \sum_{i=1}^b \ell(\mathbf{x}_i^{(m)} - \mathbf{y}_i^{(m)}) = \sum_{i=1}^b \ell(\mathbf{f}_m(\mathbf{U}_m^\top \mathbf{x}_i^{(m-1)}) - \mathbf{y}_i^{(m)}), \quad (3)$$

where $\ell(\cdot)$ is a loss function such as the squared ℓ_2 norm (or Mean Squared Error (MSE)), cross-entropy, etc. The loss \mathcal{L}_m tries to make the projected data $\mathbf{x}_i^{(m)}$ as similar as possible to the backprojected data $\mathbf{y}_i^{(m)}$ by tuning the weights \mathbf{U}_m . This is because the output of the network is supposed to be equal to the labels, i.e., $\mathbf{x}^{(n_\ell)} \approx \mathbf{y}$. In order to tune the weights for Eq. (3), we use a step of gradient descent. Using chain rule, the gradient is:

$$\mathbb{R}^{d_{m-1} \times d_m} \ni \frac{\partial \mathcal{L}_m}{\partial \mathbf{U}_m} = \sum_{i=1}^b \mathbf{vec}_{d_{m-1} \times d_m}^{-1} \left[\left(\frac{\partial \mathbf{z}_i^{(m)}}{\partial \mathbf{U}_m} \right)^\top \left(\frac{\partial \mathbf{f}_m(\mathbf{z}_i^{(m)})}{\partial \mathbf{z}_i^{(m)}} \right)^\top \frac{\partial \ell(\mathbf{f}_m(\mathbf{z}_i^{(m)}))}{\partial \mathbf{f}_m(\mathbf{z}_i^{(m)})} \right], \quad (4)$$

where we use the Magnus-Neudecker convention in which matrices are vectorized and $\mathbf{vec}_{d_{m-1} \times d_m}^{-1}$ is de-vectorization to $d_{m-1} \times d_m$ matrix. If the loss function is MSE or cross-entropy for example, the derivatives of the loss function w.r.t. the activation function, respectively, are:

$$\mathbb{R}^{d_m} \ni \frac{\partial \ell(\mathbf{f}_m(\mathbf{z}_i^{(m)}))}{\partial \mathbf{f}_m(\mathbf{z}_i^{(m)})} = 2(\mathbf{f}_m(\mathbf{z}_i^{(m)}) - \mathbf{y}_i^{(m)}), \text{ and} \quad (5)$$

$$\mathbb{R}^{d_m} \ni \frac{\partial \ell(\mathbf{f}_m(\mathbf{z}_i^{(m)}))}{\partial \mathbf{f}_m(\mathbf{z}_i^{(m)})} = - \left[\frac{\mathbf{y}_{i,j}^{(m)}}{\mathbf{f}_m(\mathbf{z}_{i,j}^{(m)})}, \forall j \in \{1, \dots, d_m\} \right]^\top, \quad (6)$$

where $\mathbf{y}_{i,j}^{(m)}$ and $\mathbf{z}_{i,j}^{(m)}$ are the j -th dimension of $\mathbf{y}_i^{(m)}$ and $\mathbf{z}_i^{(m)} = \mathbf{U}_m^\top \mathbf{x}_i^{(m-1)}$, respectively.

```

1 Procedure: UpdateLayerWeights( $\mathcal{U}, \mathbf{X}, \mathbf{Y}, m$ )
2 Input: weights:  $\mathcal{U} := \{\mathbf{U}_r\}_{r=1}^{n_\ell}$ , batch data:  $\mathbf{X} \in \mathbb{R}^{d \times b}$ , batch labels:
 $\mathbf{Y} \in \mathbb{R}^{p \times b}$ , layer:  $m \in [1, n_\ell]$ 
3  $\mathbf{X}^{(0)} := \mathbf{X}$ 
4 for layer  $r$  from 1 to  $(m - 1)$  do
5    $\mathbf{Z}^{(r)} := \mathbf{U}_r^\top \mathbf{X}^{(r-1)}$ 
6    $\mathbf{X}^{(r)} := \mathbf{f}_r(\mathbf{Z}^{(r)})$ 
7  $\mathbf{Y}^{(n_\ell)} := \mathbf{Y}$ 
8 for layer  $r$  from  $(n_\ell - 1)$  to  $m$  do
9    $\mathbf{Y}^{(r+1)} := \Pi(\mathbf{Y}^{(r+1)})$ 
10   $\mathbf{Y}^{(r)} := \mathbf{U}_{r+1} \mathbf{f}_{r+1}^{-1}(\mathbf{Y}^{(r+1)})$ 
11  $\mathbf{U}_m := \mathbf{U}_m - \eta (\partial \mathcal{L}_m / \partial \mathbf{U}_m)$ 
12 Return  $\mathbf{U}_m$ 

```

Algorithm 1: Updating the weights of a layer in backprojection

For the activation functions in which the nodes are independent, such as linear, sigmoid, and hyperbolic tangent, the derivative of the activation function w.r.t. its input is a diagonal matrix:

$$\mathbb{R}^{d_m \times d_m} \ni \frac{\partial \mathbf{f}_m(\mathbf{z}_i^{(m)})}{\partial \mathbf{z}_i^{(m)}} = \mathbf{diag}\left(\frac{\partial \mathbf{f}_m(\mathbf{z}_{i,j}^{(m)})}{\partial \mathbf{z}_{i,j}^{(m)}}, \forall j \in \{1, \dots, d_m\}\right), \quad (7)$$

where $\mathbf{diag}(\cdot)$ makes a matrix with its input as diagonal.

The derivative of the projected data before the activation function (i.e., the input of the activation function) w.r.t. the weights of the layer is:

$$\mathbb{R}^{d_m \times (d_m d_{m-1})} \ni \frac{\partial \mathbf{z}_i^{(m)}}{\partial \mathbf{U}_m} = \frac{\partial \mathbf{U}_m^\top \mathbf{x}_i^{(m-1)}}{\partial \mathbf{U}_m} = \mathbf{I}_{d_m} \otimes \mathbf{x}_i^{(m-1)\top}, \quad (8)$$

where \otimes denotes the Kronecker product and \mathbf{I}_{d_m} is the $d_m \times d_m$ identity matrix.

The procedure for updating weights in the m -th layer is shown in Algorithm 1. Until the layer m , data is projected and passed through activation functions layer by layer. Also, the label is backprojected and passed through inverse activation functions until the layer m . A step of gradient descent is used to update the layer's weights where $\eta > 0$ is the learning rate. Note that the backprojected label at a layer may not be in the feasible domain of its inverse activation function. Hence, at every layer, we should project the backprojected label onto the feasible domain [15]. We denote projection onto the feasible set by $\Pi(\cdot)$.

2.4 Different Procedures

So far, we explained how to update the weights of a layer. Here, we detail updating the entire network layers. In terms of the order of updating layers, we can have three different procedures for a backprojection algorithm. One possible procedure is to update the first layer first and move to next layers one by one until

```

1 Procedure: Backprojection( $\mathcal{X}, \mathcal{Y}, b, e$ )
2 Input: training data:  $\mathcal{X}$ , training labels:  $\mathcal{Y}$ , batch size:  $b$ , number of epochs:
    $e$ 
3 Initialize  $\mathcal{U} = \{\mathbf{U}_r\}_{r=1}^{n_\ell}$ 
4 for epoch from 1 to  $e$  do
5   for batch from 1 to  $\lceil n/b \rceil$  do
6      $\mathbf{X}, \mathbf{Y} \leftarrow$  take batch from  $\mathcal{X}$  and  $\mathcal{Y}$ 
7     if procedure is forward then
8       for layer  $m$  from 1 to  $n_\ell$  do
9          $\mathbf{U}_m \leftarrow$  UpdateLayerWeights( $\{\mathbf{U}_r\}_{r=1}^{n_\ell}, \mathbf{X}, \mathbf{Y}, m$ )
10      else if procedure is backward then
11        for layer  $m$  from  $n_\ell$  to 1 do
12           $\mathbf{U}_m \leftarrow$  UpdateLayerWeights( $\{\mathbf{U}_r\}_{r=1}^{n_\ell}, \mathbf{X}, \mathbf{Y}, m$ )
13      else if procedure is forward-backward then
14        if batch index is odd then
15          for layer  $m$  from 1 to  $n_\ell$  do
16             $\mathbf{U}_m \leftarrow$  UpdateLayerWeights( $\{\mathbf{U}_r\}_{r=1}^{n_\ell}, \mathbf{X}, \mathbf{Y}, m$ )
17        else
18          for layer  $m$  from  $n_\ell$  to 1 do
19             $\mathbf{U}_m \leftarrow$  UpdateLayerWeights( $\{\mathbf{U}_r\}_{r=1}^{n_\ell}, \mathbf{X}, \mathbf{Y}, m$ )

```

Algorithm 2: Backprojection

we reach the last layer. Repeating this procedure for the batches results in the *forward procedure*. In an opposite direction, we can have the *backward procedure* where, for each batch, we update the layers from the last layer to the first layer one by one. If we have both directions of updating, i.e., forward update for a batch and backward update for the next batch, we call it the *forward-backward procedure*. Algorithm 2 shows how to update the layers in different procedures of the backprojection algorithm. Note that in this algorithm, an updated layer impacts the update of next/previous layer. One alternative approach is to make updating of layers dependent only on the weights tuned by previous mini-batch. In that approach, the training of layers can be parallelized within mini-batch.

3 Kernel Backprojection Algorithm

Suppose $\phi : \mathcal{X} \rightarrow \mathcal{H}$ is the pulling function to the feature space. Let t denote the dimensionality of the feature space, i.e., $\phi(\mathbf{x}) \in \mathbb{R}^t$. Let the matrix-form of \mathcal{X} and \mathcal{Y} be denoted by $\mathbb{R}^{d \times n} \ni \check{\mathbf{X}} := [\mathbf{x}_1, \dots, \mathbf{x}_n]$ and $\mathbb{R}^{p \times n} \ni \check{\mathbf{Y}} := [\mathbf{y}_1, \dots, \mathbf{y}_n]$. The kernel matrix [16] for the training data $\check{\mathbf{X}}$ is defined as $\mathbb{R}^{n \times n} \ni \check{\mathbf{K}} := \check{\Phi}(\check{\mathbf{X}})^\top \check{\Phi}(\check{\mathbf{X}})$ where $\mathbb{R}^{t \times n} \ni \check{\Phi}(\check{\mathbf{X}}) := [\phi(\mathbf{x}_1), \dots, \phi(\mathbf{x}_n)]$. We normalize the kernel matrix [17] as $\check{\mathbf{K}}(i, j) := \check{\mathbf{K}}(i, j) / [\check{\mathbf{K}}(i, i)\check{\mathbf{K}}(j, j)]^{1/2}$ where $\check{\mathbf{K}}(i, j)$ denotes the (i, j) -th element of the kernel matrix.

According to representation theory [18], the projection matrix $\mathbf{U}_1 \in \mathbb{R}^{d \times d_1}$ can be expressed as a linear combination of the projected training data. Hence, we have $\mathbb{R}^{t \times d_1} \ni \Phi(\mathbf{U}_1) = \Phi(\check{\mathbf{X}}) \Theta$ where every column of $\Theta := [\theta_1, \dots, \theta_{d_1}] \in \mathbb{R}^{n \times d_1}$ is the vector of coefficients for expressing a projection direction as a linear combination of projected training data. The projection of the pulled data is $\mathbb{R}^{d_1 \times n} \ni \Phi(\mathbf{U}_1)^\top \Phi(\check{\mathbf{X}}) = \Theta^\top \Phi(\check{\mathbf{X}})^\top \Phi(\check{\mathbf{X}}) = \Theta^\top \check{\mathbf{K}}$.

In the kernel backprojection algorithm, in the first network layer, we project the pulled data from the feature space with dimensionality t to another feature space with dimensionality d_1 . The projections of the next layers are the same as in backprojection. In other words, *kernel backprojection applies backprojection in the feature space rather than the input space*. In a mini-batch set-up, we use the columns of the normalized kernel corresponding to the batch samples, denoted by $\{\mathbf{k}_i \in \mathbb{R}^n\}_{i=1}^b$. Therefore, the projection of the i -th data point in the batch is $\mathbb{R}^{d_1} \ni \Theta^\top \mathbf{k}_i$. In kernel backprojection, the dimensionality of the input is n and the kernel vector \mathbf{k}_i is fed to the network as input. If we replace the \mathbf{x}_i by \mathbf{k}_i , Algorithms 1 and 2 are applicable for kernel backprojection.

In the test phase, we normalize the kernel over the matrix $[\check{\mathbf{X}}, \mathbf{x}_t]$ where $\mathbf{x}_t \in \mathbb{R}^d$ is the test data point. Then, we take the portion of normalized kernel which correspond to the kernel over the training versus test data, denoted by $\mathbb{R}^n \ni \mathbf{k}_t := \Phi(\check{\mathbf{X}})^\top \Phi(\mathbf{x}_t)$. The projection at the first layer is then $\mathbb{R}^{d_1} \ni \Theta^\top \mathbf{k}_t$.

4 Experiments

Datasets: For experiments, we created two synthetic datasets with 300 data points each, one for binary-class and one for three-class classification (see Figs. 1 and 2). For more difficulty, we set different variances for the classes. The data were standardized as a preprocessing. For this conference short-paper, we limit ourselves to introduction of this new approach and small synthetic experiments. Validation on larger real-world datasets is ongoing for future publication.

Neural Network Settings: We implemented a neural network with three layers whose number of neurons are $\{15, 20, p\}$ where $p = 1$ and $p = 3$ for the binary and ternary classification, respectively. In different experiments, we used MSE loss for the middle layers and MSE or cross-entropy losses for the last layer. Moreover, we used Exponential Linear Unit (ELU) [19] or linear functions for activation functions of the middle layers while sigmoid or hyperbolic tangent (tanh) were used for the last layer. The derivative and inverse of these activation functions are as the following:

$$\text{ELU: } f(z) = \begin{cases} e^z - 1, & z \leq 0 \\ z, & z > 0 \end{cases}, \quad f'(z) = \begin{cases} e^z, & z \leq 0 \\ 1, & z > 0 \end{cases}, \quad f^{-1}(y) = \begin{cases} \ln(y + 1), & y \leq 0 \\ y, & y > 0 \end{cases},$$

$$\text{Linear: } f(z) = z, \quad f'(z) = 1, \quad f^{-1}(y) = y,$$

$$\text{Sigmoid: } f(z) = \frac{1}{1 + e^{-z}}, \quad f'(z) = f(1 - f), \quad f^{-1}(y) = \ln\left(\frac{y}{1 - y}\right),$$

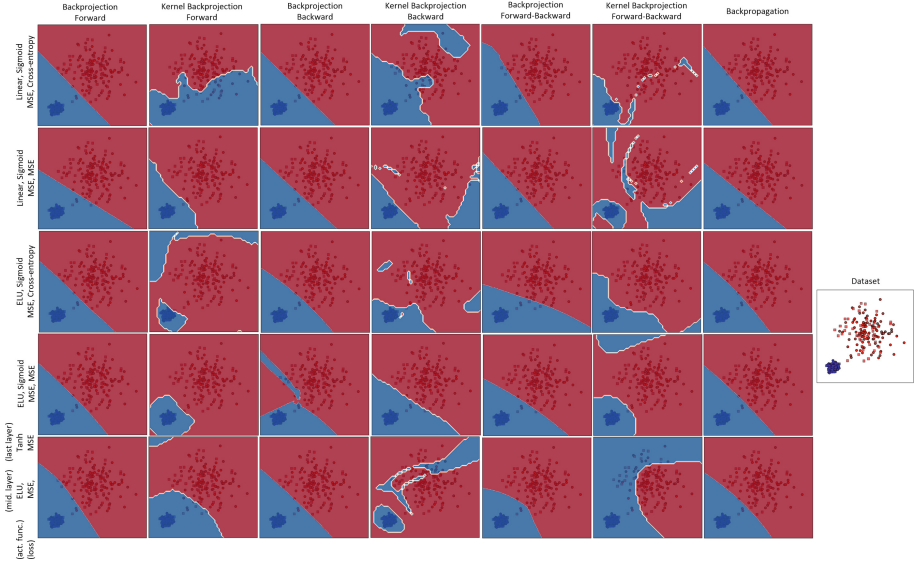


Fig. 1. Discrimination of two classes by different training algorithms with various activation functions and loss functions. The label for each row indicates the activation functions and the loss functions for the middle then the last layers. (Color figure online)

$$\text{Tanh: } f(z) = \frac{e^z - e^{-z}}{e^z + e^{-z}}, \quad f'(z) = 1 - f^2, \quad f^{-1}(y) = 0.5 \ln\left(\frac{1+y}{1-y}\right),$$

where in the inverse functions, we bound the output values for computational reasons in computer. Mostly, a learning rate of $\eta = 10^{-4}$ was used for backprojection and backpropagation and $\eta = 10^{-5}$ was used for kernel backprojection.

Comparison of Procedures: The performance of different forward, backward, and forward-backward procedures in backprojection and kernel backprojection are illustrated in Fig. 1. In these experiments, the Radial Basis Function (RBF) kernel was used in kernel backprojection. Although the performance of these procedures are not identical but all of them are promising discrimination of classes. This shows that all three proposed procedures work well for backprojection in the input and feature spaces. In other words, the algorithm is fairly robust to the order of updating layers.

Comparison to Backpropagation: The performances of backprojection, kernel backprojection, and backpropagation are compared in the binary and ternary classification, shown in Figs. 1 and 2, respectively. In Fig. 2, the linear kernel was used. In Fig. 1, most often, kernel backprojection considers a spherical class around the blue (or even red) class which is because of the choice of RBF kernel. Comparison to backpropagation in the two figures shows that backprojection's performance nearly matches that of backpropagation.

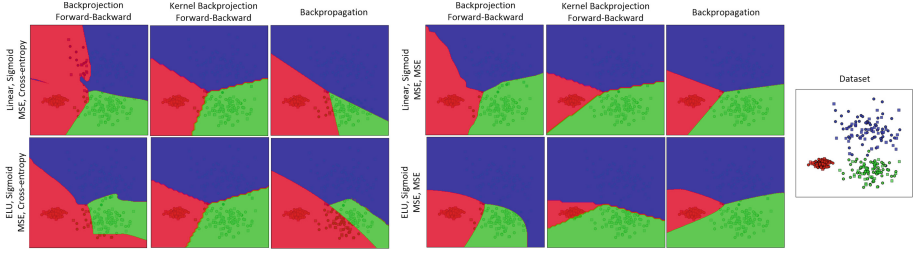


Fig. 2. Discrimination of three classes by different training algorithms with various activation functions and loss functions.

In the different experiments, the mean time of every epoch was often 0.08, 0.11, and 0.2s for backprojection, kernel backprojection, and backpropagation, respectively, where the number of epochs were fairly similar in the experiments. This shows that backprojection is *faster* than backpropagation. This is because backpropagation updates the weights one by one while backprojection updates layer by layer.

5 Conclusion and Future Direction

In this paper, we proposed a new training algorithm for feedforward neural network named backprojection. The proposed algorithm, which can be used for both the input and feature spaces, tries to force the projected data to be similar to the backprojected labels by tuning the weights layer by layer. This training algorithm, which is moderately faster than backpropagation in our initial experiments, can be used with either forward, backward, or forward-backward procedures. It is noteworthy that adding a penalty term for weight decay [20] to Eq. (3) can regularize the weights in backprojection [21]. Moreover, batch normalization can be used in backprojection by standardizing the batch at the layers [22]. This paper concentrated on feedforward neural networks. As a future direction, we can develop backprojection for other network structures such as convolutional networks [23] and carry more expensive validation experiments on real-world data.

References

1. Fausett, L.: Fundamentals of Neural Networks: Architectures, Algorithms, and Applications. Prentice-Hall Inc, Upper Saddle River (1994)
2. Goodfellow, I., Bengio, Y., Courville, A.: Deep Learning. MIT Press, Cambridge (2016)
3. Rumelhart, D.E., Hinton, G.E., Williams, R.J.: Learning representations by back-propagating errors. *Nature* **323**(6088), 533–536 (1986)
4. Soltanolkotabi, M., Javanmard, A., Lee, J.D.: Theoretical insights into the optimization landscape of over-parameterized shallow neural networks. *IEEE Trans. Inf. Theory* **65**(2), 742–769 (2018)

5. Allen-Zhu, Z., Li, Y., Liang, Y.: Learning and generalization in overparameterized neural networks, going beyond two layers. In: *Advances in Neural Information Processing Systems*, pp. 6155–6166 (2019)
6. Feizi, S., Javadi, H., Zhang, J., Tse, D.: Porcupine neural networks: (almost) all local optima are global. arXiv preprint [arXiv:1710.02196](https://arxiv.org/abs/1710.02196) (2017)
7. Montana, D.J., Davis, L.: Training feedforward neural networks using genetic algorithms. In: *International Joint Conference on Artificial Intelligence (IJCAI)*, vol. 89, pp. 762–767 (1989)
8. Leung, F.H.F., Lam, H.K., Ling, S.H., Tam, P.K.S.: Tuning of the structure and parameters of a neural network using an improved genetic algorithm. *IEEE Trans. Neural Networks* **14**(1), 79–88 (2003)
9. Hinton, G.E., Salakhutdinov, R.R.: Reducing the dimensionality of data with neural networks. *Science* **313**(5786), 504–507 (2006)
10. Hauser, M., Ray, A.: Principles of Riemannian geometry in neural networks. In: *Advances in Neural Information Processing Systems*, pp. 2807–2816 (2017)
11. Ham, J.H., Lee, D.D., Mika, S., Schölkopf, B.: A kernel view of the dimensionality reduction of manifolds. In: *International Conference on Machine Learning* (2004)
12. Ghojogh, B., Crowley, M.: Unsupervised and supervised principal component analysis: tutorial. arXiv preprint [arXiv:1906.03148](https://arxiv.org/abs/1906.03148) (2019)
13. Karimi, A.H.: Exploring new forms of random projections for prediction and dimensionality reduction in big-data regimes. Master’s thesis, University of Waterloo (2018)
14. Achlioptas, D.: Database-friendly random projections: Johnson-Lindenstrauss with binary coins. *J. Comput. Syst. Sci.* **66**(4), 671–687 (2003)
15. Parikh, N., Boyd, S.: Proximal algorithms. *Found. Trends® Optim.* **1**(3), 127–239 (2014)
16. Hofmann, T., Scholkopf, B., Smola, A.J.: Kernel methods in machine learning. *Ann. Stat.* **36**, 117–1220 (2008)
17. Ah-Pine, J.: Normalized kernels as similarity indices. In: Zaki, M.J., Yu, J.X., Ravindran, B., Pudi, V. (eds.) *PAKDD 2010. LNCS (LNAI)*, vol. 6119, pp. 362–373. Springer, Heidelberg (2010). https://doi.org/10.1007/978-3-642-13672-6_36
18. Alperin, J.L.: *Local Representation Theory: Modular Representations as an Introduction to the Local Representation Theory of Finite Groups*, vol. 11. Cambridge University Press, Cambridge (1993)
19. Clevert, D.A., Unterthiner, T., Hochreiter, S.: Fast and accurate deep network learning by exponential linear units (ELUs). In: *International Conference on Learning Representations (ICLR)* (2016)
20. Krogh, A., Hertz, J.A.: A simple weight decay can improve generalization. In: *Advances in Neural Information Processing Systems*, pp. 950–957 (1992)
21. Ghojogh, B., Crowley, M.: The theory behind overfitting, cross validation, regularization, bagging, and boosting: tutorial. arXiv preprint [arXiv:1905.12787](https://arxiv.org/abs/1905.12787) (2019)
22. Ioffe, S., Szegedy, C.: Batch normalization: Accelerating deep network training by reducing internal covariate shift. arXiv preprint [arXiv:1502.03167](https://arxiv.org/abs/1502.03167) (2015)
23. LeCun, Y., Bottou, L., Bengio, Y., Haffner, P.: Gradient-based learning applied to document recognition. *Proc. IEEE* **86**(11), 2278–2324 (1998)