

Compression of Fully-Connected Layer in Neural Network by Kronecker Product

Shuchang Zhou

*State Key Laboratory of Computer Architecture
Institute of Computing Technology
Chinese Academy of Sciences, Beijing, China
and Megvii Inc.
shuchang.zhou@gmail.com*

Jia-Nan Wu

*Department of Computer Science & Engineering
Shanghai Jiao Tong University
800 Dong Chuan Road, Shanghai, China 200240
tankeco@sjtu.edu.cn*

July 23, 2015

1. Abstract

In this paper we propose and study a technique to reduce the number of parameters and computation time in fully-connected layers of neural networks using Kronecker product, at a mild cost of the prediction quality. The technique proceeds by replacing Fully-Connected layers with so-called Kronecker Fully-Connected layers, where the weight matrices of the FC layers are approximated by linear combinations of multiple Kronecker products of smaller matrices. In particular, given a model trained on SVHN dataset, we are able to construct a new KFC model with 73% reduction in total number of parameters, while the error only rises mildly. In contrast, using low-rank method can only achieve 35% reduction in total number of parameters given similar quality degradation allowance. If we only compare the KFC layer with its counterpart fully-connected layer, the reduction in the number of parameters exceeds 99%. The amount of computation is also reduced as we replace matrix product of the large matrices in FC layers with matrix products of a few smaller matrices in KFC layers. Further experiments on MNIST, SVHN and some Chinese Character recognition models also demonstrate effectiveness of our technique.

2. Introduction

Model approximation aims at reducing the number of parameters and amount of computation of neural network models, while keeping the quality of prediction results mostly the same.¹ Model approximation is important for real world application of neural network to satisfy the time and storage constraints of the applications.

1. In some circumstances, as less number of model parameters reduce the effect of overfitting, model approximation sometimes leads to more accurate predictions.

In general, given a neural network $\tilde{f}(\cdot; \tilde{\theta})$, we want to construct another neural network $f(\cdot; \theta)$ within some pre-specified resource constraint, and minimize the differences between the outputs of two functions on the possible inputs. An example setup is to directly minimize the differences between the output of the two functions:

$$\inf_{\theta} \sum_i d(f(x_i; \theta), \tilde{f}(x_i; \tilde{\theta})), \quad (1)$$

where d is some distance function and x_i runs over all input data.

The formulation 1 does not give any constraints between the structure of f and \tilde{f} , meaning that any model can be used to approximate another model. In practice, a structural similar model is often used to approximate another model. In this case, model approximation may be approached in a modular fashion w.r.t. to each layer.

2.1 Low Rank Model Approximation

Low rank approximation in linear regression dates back to Anderson (1951). In Sainath et al. (2013), Liao et al. (2013), Xue et al. (2013), Zhang et al. (2014b), Denton et al. (2014), low rank approximation of fully-connected layer is used; and Jaderberg et al. (2014), Rigamonti et al. (2013), Lebedev et al. (2014) considered low rank approximation of convolution layer. Zhang et al. (2014a) considered approximation of multiple layers with nonlinear activations.

We first outline the low rank approximation method below. The fully-connected layer widely used in neural network construction may be formulated as:

$$\mathbf{L}_a = h(\mathbf{L}_{a-1}\mathbf{M}_a + \mathbf{b}_a), \quad (2)$$

where \mathbf{L}_i is the output of the i -th layer of the neural network, \mathbf{M}_a is often referred to as “weight term” and \mathbf{b}_a as “bias term” of the a -th layer.

As the coefficients of the weight term in the fully-connected layers are organized into matrices, it is possible to perform low-rank approximation of these matrices to achieve an approximation of the layer, and consequently the whole model. Given Singular Value Decomposition of a matrix $\mathbf{M} = \mathbf{U}\mathbf{D}\mathbf{V}^*$, where \mathbf{U}, \mathbf{V} are unitary matrices and \mathbf{D} is a diagonal matrix with the diagonal made up of singular values of \mathbf{M} , a rank- k approximation of $\mathbf{M} \in \mathbb{F}^{m \times n}$ is:

$$\mathbf{M} \approx \mathbf{M}_k = \tilde{\mathbf{U}}\tilde{\mathbf{D}}\tilde{\mathbf{V}}^*, \text{ where } \tilde{\mathbf{U}} \in \mathbb{F}^{m \times k}, \tilde{\mathbf{D}} \in \mathbb{F}^{k \times k}, \tilde{\mathbf{V}} \in \mathbb{F}^{n \times k}, \quad (3)$$

where $\tilde{\mathbf{U}}$ and $\tilde{\mathbf{V}}$ are the first k -columns of the \mathbf{U} and \mathbf{V} respectively, and $\tilde{\mathbf{D}}$ is a diagonal matrix made up of the largest k entries of \mathbf{D} .

In this case approximation by SVD is optimal in the sense that the following holds Horn and Johnson (1991):

$$\mathbf{M}_r = \inf_{\mathbf{X}} \|\mathbf{X} - \mathbf{M}\|_F \text{ s.t. } \text{rank}(\mathbf{X}) \leq r, \quad (4)$$

and

$$\mathbf{M}_r = \inf_{\mathbf{X}} \|\mathbf{X} - \mathbf{M}\|_2 \text{ s.t. } \text{rank}(\mathbf{X}) \leq r. \quad (5)$$

The approximate fully connected layer induced by SVD is:

$$\mathbf{Z} = \mathbf{L}_{a-1}\tilde{\mathbf{U}} \quad (6)$$

$$\mathbf{L}_a = h(\mathbf{Z}\tilde{\mathbf{D}}\tilde{\mathbf{V}}^* + \mathbf{b}_a), \quad (7)$$

In the modular representation of neural network, this means that the original fully connected layer is now replaced by two consequent fully-connected layers.

However, the above post-processing approach only ensures getting an optimal approximation of \mathbf{M} under the rank constraint, while there is still no guarantee that such an approximation is optimal w.r.t. the input data. I.e., the optimum of the following may well be different from the rank- r approximation \mathbf{M}_r w.r.t. some given input \mathbf{X} :

$$\inf_{\mathbf{Z}} \|\mathbf{Z}\mathbf{X} - \mathbf{M}\mathbf{X}\|_F \text{ s.t. } \text{rank}(\mathbf{Z}) \leq r. \quad (8)$$

Hence it is often necessary for the resulting low-rank model $\tilde{\mathcal{M}}$ to be trained for a few more epochs on the input, which is also known as the ‘‘fine-tuning’’ process.

Alternatively, we note that the rank constraint can be enforced by the following structural requirement for $\mathbf{X} \in \mathbb{F}^{m \times n}$:

$$\text{rank}(\mathbf{X}) \leq r \Leftrightarrow \exists \mathbf{A} \in \mathbb{F}^{m \times r}, \mathbf{B} \in \mathbb{F}^{r \times n} \text{ s.t. } \mathbf{X} = \mathbf{A}\mathbf{B}. \quad (9)$$

In light of this, if we want to impose a rank constraint on a fully-connected layer $L(\mathbf{M}, g)$ in a neural network where $\mathbf{M} \in \mathbb{F}^{m \times n}$, we can replace that layer with two consecutive layers $L_1(\mathbf{B}, g_1)$ and $L_2(\mathbf{A}, g_2)$, where $g_1(x) = x$, $g_2 = g$, and $\mathbf{M} = \mathbf{A}\mathbf{B}$ where $\mathbf{A} \in \mathbb{F}^{m \times r}$, $\mathbf{B} \in \mathbb{F}^{r \times n}$, and then train the structurally constrained neural network on the training data.

As a third method, a regularization term inducing low rank matrices may be imposed on the weight matrices. In this case, the training of a k -layer model is modified to be:

$$\inf_{\theta} \sum_i f(x_i; \theta) + \sum_{j=1}^k r_j(\theta), \quad (10)$$

where r is the regularization term. For the weight term of the FC layers, conceptually we may use the matrix rank function as the regularization term. However, as the rank function is only well-defined for infinite-precision numbers, nuclear norm may be used as its convex proxy Jaderberg et al. (2014), Recht et al. (2010).

3. Model Approximation by Kronecker Product

Next we propose to use Kronecker product of matrices of particular shapes for model approximation in Section 3.1. We also outline the relationship between the Kronecker product approximation and low-rank approximation in Section 3.2.

Below we measure the reduction in amount of computation by number of floating point operations. In particular, we will assume the computation complexity of two matrices of dimensions $M \times K$ and $K \times N$ to be $O(MKN)$, as many neural network implementations Bastien et al. (2012), Bergstra et al. (2010), Jia et al. (2014), Collobert et al. (2011) have not used algorithms of lower computation complexity for the typical inputs of the neural networks. Our analysis is mostly immune to the ‘‘hidden constant’’ problem in computation complexity analysis as the underlying computations of the transformed model may also be carried out by matrix products.

3.1 Weight Matrix Approximation by Kronecker Product

We next discuss how to use Kronecker product to approximate weight matrices of FC layers, leading to construction of a new kind of layer which we call Kronecker Fully-Connected layer.

The idea originates from the observation that for a matrix $\mathbf{M} \in \mathbb{F}^{m \times n}$ where the dimensions are not prime ², we have approximations like:

$$\mathbf{M} = \mathbf{M}_1 \otimes \mathbf{M}_2, \quad (11)$$

where $m = m_1 m_2$, $n = n_1 n_2$, $\mathbf{M}_1 \in \mathbb{F}^{m_1 \times n_1}$, $\mathbf{M}_2 \in \mathbb{F}^{m_2 \times n_2}$.

Any factors of m and n may be selected as m_1 and n_1 in the above formulation. However, in a Convolutional Neural Network, the input to a FC layer may be a tensor of order 4, which has some natural shape constraints that we will try to leverage in 3.1.1. Otherwise, when the input is a matrix, we do not have natural choices of m_1 and n_1 . We will explore heuristics to pick m_1 and n_1 in 3.1.7.

3.1.1 KRONECKER PRODUCT APPROXIMATION FOR FULLY-CONNECTED LAYER WITH 4D TENSOR INPUT

In a convolutional layer processing images, the input data \mathbf{L}_{a-1} may be a tensor of order 4 as \mathcal{T}_{nchw} where $n = 1, 2, \dots, N$ runs over N different instances of data, $c = 1, 2, \dots, C$ runs over C channels of the given images, $h = 1, 2, \dots, H$ runs over H rows of the images, and $w = 1, 2, \dots, W$ runs over W columns of the images. \mathcal{T} is often reshaped into a matrix before being fed into a fully connected layer as \mathbf{D}_{nj} , where $n = 1, 2, \dots, N$ runs over the N different instances of data and $j = 1, 2, \dots, CHW$ runs over the combined dimension of channel, height, and width of images. The weights of the fully-connected layer would then be a matrix \mathbf{M}_{jk} where $j = 1, 2, \dots, CHW$ and $k = 1, 2, \dots, K$ runs over output number of channels. I.e., the layer may be written as:

$$\mathbf{D} = \text{Reshape}(\mathbf{L}_{a-1}) \quad (12)$$

$$\mathbf{L}_a = h(\mathbf{D}\mathbf{M} + \mathbf{b}). \quad (13)$$

Though the reshaping transformation from \mathcal{T} to \mathbf{D} does not incur any loss in pixel values of data, we note that the dimension information of the tensor of order 4 is lost in the matrix representation. As a consequence, \mathbf{M} has $CHWK$ number of parameters.

Due to the shape of \mathbf{M} , we may propose a few kinds of structural constraint on \mathbf{M} by requiring \mathbf{M} to be Kronecker product of matrices of particular shapes.

3.1.2 FORMULATION I

In this formulation, we require $\mathbf{M} = \mathbf{M}_1 \otimes \mathbf{M}_2 \otimes \mathbf{M}_3$, where $\mathbf{M}_1 \in \mathbb{F}^{C \times K_1}$, $\mathbf{M}_2 \in \mathbb{F}^{H \times K_2}$, $\mathbf{M}_3 \in \mathbb{F}^{W \times K_3}$, and $K = K_1 K_2 K_3$. The number of parameters is reduced to $CK_1 + HK_2 + WK_3$. The underlying assumption for this model is that the transformation is invariant across rows and columns of the images.

3.1.3 FORMULATION II

In this formulation, we require $\mathbf{M} = \mathbf{M}_1 \otimes \mathbf{M}_2$, where $\mathbf{M}_1 \in \mathbb{F}^{C \times K_1}$, $\mathbf{M}_2 \in \mathbb{F}^{HW \times K_2}$, and $K = K_1 K_2$. The number of parameters is reduced to $CK_1 + HWK_2$. The underlying assumption for this model is that the channel transformation should be decoupled from the spatial transformation.

2. In case any of m and n is prime, it is possible to add some extra dummy feature or output class to make the dimensions dividable.

3.1.4 FORMULATION III

In this formulation, we require $\mathbf{M} = \mathbf{M}_1 \otimes \mathbf{M}_2$, where $\mathbf{M}_1 \in \mathbb{F}^{CH \times K_1}$, $\mathbf{M}_2 \in \mathbb{F}^{W \times K_2}$, and $K = K_1 K_2$. The number of parameters is reduced to $CHK_1 + WK_2$. The underlying assumption for this model is that the transformation w.r.t. columns may be decoupled.

3.1.5 FORMULATION IV

In this formulation, we require $\mathbf{M} = \mathbf{M}_1 \otimes \mathbf{M}_2$, where $\mathbf{M}_1 \in \mathbb{F}^{CW \times K_1}$, $\mathbf{M}_2 \in \mathbb{F}^{H \times K_2}$, and $K = K_1 K_2$. The number of parameters is reduced to $CWK_1 + HK_2$. The underlying assumption for this model is that the transformation w.r.t. rows may be decoupled.

3.1.6 COMBINED FORMULATIONS

Note that the above four formulations may be linearly combined to produce more possible kinds of formulations. It would be a design choice with respect to trade off between the number of parameters, amount of computation and the particular formulation to select.

3.1.7 KRONECKER PRODUCT APPROXIMATION FOR MATRIX INPUT

For fully-connected layer whose input are matrices, there does not exist natural dimensions to adopt for the shape of smaller weight matrices in KFC. Through experiments, we find it possible to arbitrarily pick a decomposition of input matrix dimensions to enforce the Kronecker product structural constraint. We will refer to this formulation as KFCM.

Concretely, when input to a fully-connected layer is $\mathbf{X} \in \mathbb{F}^{N \times C}$ and the weight matrix of the layer is $\mathbf{W} \in \mathbb{F}^{C \times K}$, we can construct approximation of \mathbf{W} as:

$$\tilde{\mathbf{W}} = \mathbf{W}_1 \otimes \mathbf{W}_2 \approx \mathbf{W}, \quad (14)$$

where $C = C_1 C_2$, $K = K_1 K_2$, $\mathbf{W}_1 \in \mathbb{F}^{C_1 \times K_1}$ and $\mathbf{W}_2 \in \mathbb{F}^{C_2 \times K_2}$.

The computation complexity will be reduced from $O(NCK)$ to $O(NCK(\frac{1}{K_2} + \frac{1}{C_1})) = O(NC_2 C_1 K_1 + NC_2 K_1 K_2)$, while the number of parameters will be reduced from CK to $C_1 K_1 + C_2 K_2$.

Through experiments, we have found it sensible to pick $C_1 \approx \sqrt{C}$ and $K_1 \approx \sqrt{K}$.

As the choice of C_1 and K_1 above is arbitrary, we may use linear combination of Kronecker products if matrices of different shapes for approximation.

$$\tilde{\mathbf{W}} = \sum_{j=1}^J \mathbf{W}_{1j} \otimes \mathbf{W}_{2j} \approx \mathbf{W}, \quad (15)$$

where $\mathbf{W}_{1j} \in \mathbb{F}^{C_{1j} \times K_{1j}}$ and $\mathbf{W}_{2j} \in \mathbb{F}^{C_{2j} \times K_{2j}}$.

3.2 Relationship between Kronecker Product Constraint and Low Rank Constraint

It turns out that factorization by Kronecker product is closely related to the low rank approximation method. In fact, approximating a matrix \mathbf{M} with Kronecker product $\mathbf{M}_1 \otimes \mathbf{M}_2$ of two matrices may be casted into a Nearest Kronecker product Problem:

$$\inf_{\mathbf{M}_1, \mathbf{M}_2} \|\mathbf{M} - \mathbf{M}_1 \otimes \mathbf{M}_2\|_F. \quad (16)$$

An equivalence relation in the above problem is given in Van Loan and Pitsianis (1993), Van Loan (2000) as:

$$\arg \inf_{\mathbf{M}_1, \mathbf{M}_2} \|\mathbf{M} - \mathbf{M}_1 \otimes \mathbf{M}_2\|_F = \arg \inf_{\mathbf{M}_1, \mathbf{M}_2} \|\mathcal{R}(\mathbf{M}) - \text{vec } \mathbf{M}_1 (\text{vec } \mathbf{M}_2)^\top\|_F, \quad (17)$$

where $\mathcal{R}(\mathbf{M})$ is a matrix formed by a fixed reordering of entries \mathbf{M} .

Note the right-hand side of formula 17 is a rank-1 approximation of matrix $\mathcal{R}(\mathbf{M})$, hence has a closed form solution. However, the above approximation is only optimal w.r.t. the parameters of the weight matrices, but not w.r.t. the prediction quality over input data.

Similarly, though there are iterative algorithms for rank-1 approximation of tensor Friedland et al. (2013), Lathauwer et al. (2000), the optimality of the approximation is lost once input data distribution is taken into consideration.

Hence in practice, we only use the Kronecker Product constraint to construct KFC layers and optimize the values of the weights through the training process on the input data.

3.3 Extension to Sum of Kronecker Product

Just as low-rank approximation may be extended beyond rank-1 to arbitrary number of ranks, one could extend the Kronecker Product approximation to Sum of Kronecker Product approximation. Concretely, one not the following decomposition of \mathbf{M} :

$$\mathbf{M} = \sum_{i=1}^{\text{rank}(\mathcal{R}(\mathbf{M}))} \mathbf{A}_i \otimes \mathbf{B}_i. \quad (18)$$

Hence it is possible to find k -approximations:

$$\mathbf{M} \approx \sum_{i=1}^k \mathbf{A}_i \otimes \mathbf{B}_i. \quad (19)$$

We can then generalize Formulation I-IV in 3.1 to the case of sum of Kronecker Product.

We may further combine the multiple shape formulation of 15 to get the general form of KFC layer:

$$\mathbf{M} \approx \sum_{j=1}^J \sum_{i=1}^k \mathbf{A}_{ij} \otimes \mathbf{B}_{ij}. \quad (20)$$

where $\mathbf{A}_{ij} \in \mathbb{F}^{C_{1j} \times K_{1j}}$ and $\mathbf{B}_{ij} \in \mathbb{F}^{C_{2j} \times K_{2j}}$.

4. Empirical Evaluation of Kronecker product method

We next empirically study the properties and efficacy of the Kronecker product method and compare it with some other common low rank model approximation methods.

To make a fair comparison, for each dataset, we train a convolutional neural network with a fully-connected layer as a baseline. Then we replace the fully-connected layer with different layers according to different methods and train the new network until quality metrics stabilizes. We then compare KFC method with low-rank method and the baseline model in terms of number of parameters and prediction quality. We do the experiments based on implementation of KFC layers in TheanoBergstra et al. (2010), Bastien et al. (2012) framework.

As the running time may depend on particular implementation details of the KFC and the Theano work, we do not report running time below. However, there is no noticeable slow down in our experiments and the complexity analysis suggests that there should be significant reduction in amount of computation.

4.1 MNIST

The MNIST dataset LeCun et al. (1998) consists of 28×28 grey scale images of handwritten digits. There are 60000 training images and 10000 test images. We select the last 10000 training images as validation set.

Our baseline model has 8 layers and the first 6 layers consist of four convolutional layers and two pooling layers. The 7th layer is the fully-connected layer and the 8th is the softmax output. The input of the fully-connected layer is of size $32 \times 3 \times 3$, where 32 is the number of channel and 3 is the side length of image patches (the mini-batch size is omitted). The output of the fully-connected layer is of size 256, so the weight matrix is of size 288×256 .

CNN training is done with Adam Kingma and Ba (2014) with weight decay of 0.0001. Dropout Hinton et al. (2012) of 0.5 is used on the fully-connected layer and KFC layer. $y = |\tanh x|$ is used as activation function. Initial learning rate is $1e - 4$ for Adam.

Test results are listed in Table 1. The number of layer parameters means the number of parameters of the fully-connected layer or its counterpart layer(s). The number of model parameters is the number of the parameters of the whole model. The test error is the min-validation model's test error.

In Cut-96 method, we use 96 output neurons instead of 256 in fully-connected layer. In the LowRank-96 method, we replace the fully-connected layer with two fully-connected layer where the first FC layer output size is 96 and the second FC layer output size is 256. In the KFC-II method, we replace the fully-connected layer with KFC layer using formulation II with $K_1 = 64$ and $K_2 = 4$. In the KFC-Combined method, we replace the fully-connected layer with KFC layer and linear combined the formulation II, III and IV ($K_1 = 64, K_2 = 4$ in formulation II, $K_1 = 128, K_2 = 2$ in formulation III and IV).

Table 1: Comparison of using Low-Rank method and using KFC layers on MNIST dataset

Methods	# of Layer Params(%Reduction)	# of Model Params(%Reduction)	Test Error
Baseline	74.0K	99.5K	0.51%
Cut-96	27.8K(62.5%)	51.7K(48.1%)	0.58%
LowRank-96	52.6K(39.0%)	78.1K(21.6%)	0.54%
KFC-II	2.1K(97.2%)	27.7K(72.2%)	0.76%
KFC-Combined	27.0K(63.51%)	52.5K(47.2%)	0.57%

4.2 Street View House Numbers

The SVHN dataset Netzer et al. (2011) is a real-world digit recognition dataset consisting of photos of house numbers in Google Street View images. The dataset comes in two formats and we consider the second format: 32-by-32 colored images centered around a single character. There are 73257 digits for training, 26032 digits for testing, and 531131 less difficult samples

which can be used as extra training data. To build a validation set, we randomly select 400 images per class from training set and 200 images per class from extra training set as Sermanet et al. (2012), Goodfellow et al. (2013) did.

Here we use a similar but larger neural network as used in MNIST to be the baseline. The input of the fully-connected layer is of size $256 \times 5 \times 5$. The fully-connected layer has 256 output neurons. Other implementation details are not changed. Test results are listed in Table 3. In the Cut- N method, we use N output neurons instead of 256 in fully-connected layer. In the LowRank- N method, we replace the fully-connected layer with two fully-connected layer where the first FC layer output size is N and the second FC layer output size is 256. In the KFC-II method, we replace the fully-connected layer with KFC layer using formulation II with $K_1 = 64$ and $K_2 = 4$. In the KFC-Combined method, we replace the fully-connected layer with KFC layer and linear combined the formulation II, III and IV ($K_1 = 64, K_2 = 4$ in formulation II, $K_1 = 128, K_2 = 2$ in formulation III and IV). In the KFC-Rank N method, we use KFC formulation II with $K_1 = 64, K_2 = 2$ and extend it to rank N with as described above.

Table 2: Comparison of using Low-Rank method and using KFC layers on SVHN dataset

Methods	# of Layer Params(%Reduction)	# of Model Params(%Reduction)	Test Error
Baseline	1.64M	2.20M	2.57%
Cut-128	0.82M(50.0%)	1.38M(37.3%)	2.79%
Cut-64	0.41(25.0%)	0.97(55.9%)	3.19%
LowRank-128	0.85M(48.2%)	1.42M(35.7%)	3.02%
LowRank-64	0.43M(73.7%)	0.99M(55.1%)	3.67%
KFC-II	0.016M(99.0%)	0.58M(73.7%)	3.33%
KFC-Combined	0.34M(79.3%)	0.91M(58.6%)	2.60%
KFC-Rank10	0.17M(89.6%)	0.73M(66.8%)	3.19%

4.3 Chinese Character Recognition

We also evaluate application of KFC to a Chinese character recognition model. Our experiments are done on a private dataset for the moment and may extend to other established Chinese character recognition datasets like HCL2000(Zhang et al. (2009)) and CASIA-HWDB(Liu et al. (2013)).

For this task we also use a convolutional neural network. The distinguishing feature of the neural network is that following the convolution and pooling layers, it has two FC layers, one with 1536 hidden size, and the other with more than 6000 hidden size.

The two FC layers happen to be different type. The 1st FC layer accepts tensor as input and the 2nd FC layer accepts matrix as input. We apply KFC-I formulation to 1st FC and KFCM to 2nd FC.

It can be seen KFC can significantly reduce the number of parameters. However, in case of “KFC and KFCM (rank=1)”, this also leads to serious degradation of prediction quality. However, by increasing the rank from 1 to 10, we are able to recover most of the lost prediction quality. Nevertheless, the rank-10 model is still very small compared to the baseline model.

Table 3: Effect of using KFC layers on a Chinese recognition dataset

Methods	%Reduction of 1st FC Layer Params	%Reduction of 2nd FC Layer Params	%Reduction of Total Params	Test Error
Baseline	0%	0%	0%	10.6%
KFC-II	99.3%	0%	36.0%	11.6%
KFC-KFCM- rank1	98.7%	99.9%	94.5%	21.8%
KFC-KFCM- rank-10	93.3%	99.1%	91.8%	13.0%

5. Conclusion and Future Work

In this paper, we propose and study methods for approximating the weight matrices of fully-connected layers with sums of Kronecker product of smaller matrices, resulting in a new type of layer which we call Kronecker Fully-Connected layer. We consider both the cases when input to the fully-connected layer is a tensor of order 4 and when the input is a matrix. We have found that using the KFC layer can significantly reduce the number of parameters and amount of computation in experiments on MNIST, SVHN and Chinese character recognition.

As future work, we note that when weight parameters of a convolutional layer is a tensor of order 4 as $\mathcal{T} \in \mathbb{F}^{K \times C \times H \times W}$, it can be represented as a collection of $H \times W$ matrices \mathbf{T}_{hw} . We can then approximate each matrix by Kronecker products as $\mathbf{T}_{hw} = \mathbf{A}_{hw} \otimes \mathbf{B}_{hw}$ following KFCM formulation, and apply the other techniques outlined in this paper. It is also noted that the Kronecker product technique may also be applied to other neural network architectures like Recurrent Neural Network, for example approximating transition matrices with linear combination of Kronecker products.

References

- Theodore Wilbur Anderson. Estimating linear restrictions on regression coefficients for multivariate normal distributions. *The Annals of Mathematical Statistics*, pages 327–351, 1951.
- Frédéric Bastien, Pascal Lamblin, Razvan Pascanu, James Bergstra, Ian J. Goodfellow, Arnaud Bergeron, Nicolas Bouchard, and Yoshua Bengio. Theano: new features and speed improvements. Deep Learning and Unsupervised Feature Learning NIPS 2012 Workshop, 2012.
- James Bergstra, Olivier Breuleux, Frédéric Bastien, Pascal Lamblin, Razvan Pascanu, Guillaume Desjardins, Joseph Turian, David Warde-Farley, and Yoshua Bengio. Theano: a CPU and GPU math expression compiler. In *Proceedings of the Python for Scientific Computing Conference (SciPy)*, June 2010. Oral Presentation.
- Ronan Collobert, Koray Kavukcuoglu, and Clément Farabet. Torch7: A matlab-like environment for machine learning. In *BigLearn, NIPS Workshop*, number EPFL-CONF-192376, 2011.
- Emily L Denton, Wojciech Zaremba, Joan Bruna, Yann LeCun, and Rob Fergus. Exploiting linear structure within convolutional networks for efficient evaluation. In *Advances in Neural Information Processing Systems*, pages 1269–1277, 2014.

- Shmuel Friedland, Volker Mehrmann, Renato Pajarola, and SK Suter. On best rank one approximation of tensors. *Numerical Linear Algebra with Applications*, 20(6):942–955, 2013.
- Ian J Goodfellow, David Warde-Farley, Mehdi Mirza, Aaron Courville, and Yoshua Bengio. Maxout networks. *arXiv preprint arXiv:1302.4389*, 2013.
- Geoffrey E Hinton, Nitish Srivastava, Alex Krizhevsky, Ilya Sutskever, and Ruslan R Salakhutdinov. Improving neural networks by preventing co-adaptation of feature detectors. *arXiv preprint arXiv:1207.0580*, 2012.
- Roger A Horn and Charles R Johnson. *Topics in matrix analysis*. Cambridge university press, 1991.
- Max Jaderberg, Andrea Vedaldi, and Andrew Zisserman. Speeding up convolutional neural networks with low rank expansions. *arXiv preprint arXiv:1405.3866*, 2014.
- Yangqing Jia, Evan Shelhamer, Jeff Donahue, Sergey Karayev, Jonathan Long, Ross Girshick, Sergio Guadarrama, and Trevor Darrell. Caffe: Convolutional architecture for fast feature embedding. In *Proceedings of the ACM International Conference on Multimedia*, pages 675–678. ACM, 2014.
- Diederik Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*, 2014.
- Lieven De Lathauwer, Bart De Moor, and Joos Vandewalle. On the best rank-1 and rank-(r_1, r_2, \dots, r_n) approximation of higher-order tensors. *SIAM J. Matrix Anal. Appl.*, 21(4):1324–1342, March 2000. ISSN 0895-4798. doi: 10.1137/S0895479898346995. URL <http://dx.doi.org/10.1137/S0895479898346995>.
- Vadim Lebedev, Yaroslav Ganin, Maksim Rakhuba, Ivan V. Oseledets, and Victor S. Lempitsky. Speeding-up convolutional neural networks using fine-tuned cp-decomposition. *CoRR*, abs/1412.6553, 2014. URL <http://arxiv.org/abs/1412.6553>.
- Yann LeCun, Léon Bottou, Yoshua Bengio, and Patrick Haffner. Gradient-based learning applied to document recognition. *Proceedings of the IEEE*, 86(11):2278–2324, 1998.
- Hank Liao, Erik McDermott, and Andrew Senior. Large scale deep neural network acoustic modeling with semi-supervised training data for youtube video transcription. In *Automatic Speech Recognition and Understanding (ASRU), 2013 IEEE Workshop on*, pages 368–373. IEEE, 2013.
- Cheng-Lin Liu, Fei Yin, Da-Han Wang, and Qiu-Feng Wang. Online and offline handwritten chinese character recognition: benchmarking on new databases. *Pattern Recognition*, 46(1):155–162, 2013.
- Yuval Netzer, Tao Wang, Adam Coates, Alessandro Bissacco, Bo Wu, and Andrew Y Ng. Reading digits in natural images with unsupervised feature learning. In *NIPS workshop on deep learning and unsupervised feature learning*, volume 2011, page 5. Granada, Spain, 2011.
- Benjamin Recht, Maryam Fazel, and Pablo A. Parrilo. Guaranteed minimum-rank solutions of linear matrix equations via nuclear norm minimization. *SIAM Rev.*, 52(3):471–501, August 2010. ISSN 0036-1445. doi: 10.1137/070697835. URL <http://dx.doi.org/10.1137/070697835>.

- Roberto Rigamonti, Amos Sironi, Vincent Lepetit, and Pascal Fua. Learning separable filters. In *Computer Vision and Pattern Recognition (CVPR), 2013 IEEE Conference on*, pages 2754–2761. IEEE, 2013.
- Tara N Sainath, Brian Kingsbury, Vikas Sindhwani, Ebru Arisoy, and Bhuvana Ramabhadran. Low-rank matrix factorization for deep neural network training with high-dimensional output targets. In *Acoustics, Speech and Signal Processing (ICASSP), 2013 IEEE International Conference on*, pages 6655–6659. IEEE, 2013.
- Pierre Sermanet, Sandhya Chintala, and Yann LeCun. Convolutional neural networks applied to house numbers digit classification. In *Pattern Recognition (ICPR), 2012 21st International Conference on*, pages 3288–3291. IEEE, 2012.
- Charles F Van Loan. The ubiquitous kronecker product. *Journal of computational and applied mathematics*, 123(1):85–100, 2000.
- Charles F Van Loan and Nikos Pitsianis. *Approximation with Kronecker products*. Springer, 1993.
- Jian Xue, Jinyu Li, and Yifan Gong. Restructuring of deep neural network acoustic models with singular value decomposition. In *INTERSPEECH*, pages 2365–2369, 2013.
- Honggang Zhang, Jun Guo, Guang Chen, and Chunguang Li. Hcl2000-a large-scale handwritten chinese character database for handwritten character recognition. In *Document Analysis and Recognition, 2009. ICDAR'09. 10th International Conference on*, pages 286–290. IEEE, 2009.
- Xiangyu Zhang, Jianhua Zou, Xiang Ming, Kaiming He, and Jian Sun. Efficient and accurate approximations of nonlinear convolutional networks. *arXiv preprint arXiv:1411.4229*, 2014a.
- Yu Zhang, Ekapol Chuangsuwanich, and James Glass. Extracting deep neural network bottleneck features using low-rank matrix factorization. In *Proc. ICASSP*, 2014b.