
Convexified Convolution Neural Networks

Group 16

 -K.Sravanthi 

-M.Megha

-B.Kaushik

-R Harsha

Introduction:

Convolutional neural network (CNN) is a class of deep, feed-forward artificial neural networks that has successfully used in image classification, face recognition, speech recognition, text classification and game playing.

However, training a CNN is solving a non-convex optimization problem which is NP Hard.

Also, training a CNN via back-propagation has 2 drawbacks

1) Rate of convergence

2) Difficulty in understanding statistical properties of CNN

Due to non-convexity, the rate of convergence can be slowed down, which the author has tried to overcome by convexifying the network (which also gives us detailed information of the statistical properties).

These convexified networks are known as Convexified Convolutional Neural Networks (CCNNs). Also, the complexity of these networks is significantly lower than CNN (or fully connected neural network).

In this paper, the author has 2 main objectives:

- 1) Formulate CCNNs

- 2) Show that generalization error of CCNN model is comparable to the best of CNN models

Previous Works:

- Formulating neural networks as convex optimization method,using infinite parameters(Bengio et al)
- A method for learning multilayer latent variable models in which,for certain activation functions,it is a convex relaxation for learning fully connected neural networks(Aslan et al)
- Convolution Kernel Networks(Mairal et al)
- Scatnet method(Bruna and Mallat)
- Randomly initialized CNN can extract features as powerful as kernel methods(Daniel et al)

Notations used:

$\|A\|_*$ -nuclear norm-sum of singular values of A

$\|A\|_2$ -spectral norm -square root of max eigen values of $A^H A$ (A^H -conjugate transpose).

$$\langle u, v \rangle := \sum_{i=1}^{\infty} u_i v_i$$

The Basic CNN problem setup(2 Layer)

1) Split the input vector x into P patches $\{z_p(x)\}_{p=1}^P$.

2) Choose an activation function $\sigma: \mathbb{R} \rightarrow \mathbb{R}$, and collection of weight vectors

$$\{w_j\}_{j=1}^r$$

3) Define a filter function h_j $\{j \in [r]\}$ for each patch z . Same filters are applied to each patch (parameter sharing).

4) For each patch index $p \in [P]$, filter index $j \in [r]$, and output coordinate

$k \in [d_2]$, introduce a coefficient $\alpha_{k,j,p} \in \mathbb{R}$ which governs contribution of the filter h_j on patch $z_p(x)$ to output $f_k(x)$.

Final form of CNN is a vector containing all the individual functions where k th component is given by

$$f_k(\mathbf{x}) = \sum_{j=1}^r \sum_{p=1}^P \alpha_{k,j,p} h_j(z_p(\mathbf{x}))$$

5) For $f_k(\mathbf{x})$, we do empirical risk minimization by using a loss function (that is convex and L -Lipschitz). This optimization problem is non-convex (as $f_k(\mathbf{x})$ is non-linearly related to α and w).

6) We define some constraints on them

Therefore, we try to convexify this CNN.

We basically embed the non linear problem into a reproducing kernel hilbert space, which we reduce to a linear setting. Firstly, we show the convexifying for linear activation function

Reproducing Kernel Hilbert space(RKHS):

An RKHS is a hilbert space of functions in which point evaluation is a continuous linear functional.

An RKHS is associated with a kernel that reproduces every function in the space ,in the sense that for any x in the set, on which the functions are defined , "evaluation at x " can be performed by taking an inner product with a function determined by the kernel.

Such a reproducing kernel exists if and only if every evaluation functional is continuous.

Case 1: Linear Activation Function

- 1) Let the activation function be $\sigma(t)=t$.
- 2) The filter h_j , when applied to the patch vector $z_p(x)$, outputs a euclidean inner product of the form $h_j(z_p(x)) = \langle z_p(x), w_j \rangle$ (h_j is already defined before).
- 3) Define a patch matrix ($P * d1\ dimensional$) $Z(x)$, which contains the P input patches.
- 4) Therefore, the final form can be rewritten as

$$f_k(x) = \sum_{j=1}^r \sum_{p=1}^P \alpha_{k,j,p} \langle z_p(x), w_j \rangle = \sum_{j=1}^r \alpha_{k,j}^T Z(x) w_j$$

$$\text{tr}(Z(x) (\sum_{j=1}^r w_j \alpha_{k,j}^T)) = \text{tr}(Z(x) A_k)$$

5) We see that every k th component of the function vector depends linearly on matrix parameter A_k . This has a rank of at most ' r ' (due to parameter sharing of CNN)

6) Vector A consists of all the matrix parameters, and we define a function ($f^A: \mathbb{R}^{d_1} \rightarrow \mathbb{R}^{d_2}$), which contains all the output components.

7) This model class is based on some constraints.

- If B_1 and B_2 are positive radii, then, $F_{\text{CNN}}(B_1, B_2)$ to be a set of functions f^A which satisfies constraints

$$\textbf{(C1)} \max_{j \in [r]} \|w_j\|_2 \leq B_1$$

(C2) $\text{rank}(A)=r$

Any matrix A satisfying the constraint C_1 and C_2 must have nuclear norm bounded as $\|A\|_* \leq B_1 B_2 r \sqrt{d_2}$.

Therefore we define the function class F_{CCNN} as :

$$F_{\text{CCNN}} := \{ f^A : \|A\|_* \leq B_1 B_2 r \sqrt{d_2} \}$$

which guarantees $F_{\text{CNN}} \subseteq F_{\text{CCNN}}$

7) We then try to minimize empirical risk minimization over F_{CCNN}

$$\hat{f}_{\text{ccnn}} := \arg \min_{f^A \in F_{\text{ccnn}}} \sum_{i=1}^n \mathcal{L}(f^A(x_i); y_i).$$

Output expressed as matrix product

$$f_k(x) = \text{tr} \left(\begin{array}{c} d_1 \\ \text{---} z_1(x) \text{---} \\ \vdots \\ \text{---} z_P(x) \text{---} \\ P \end{array} Z(x) \times \begin{array}{c} r \text{ (filters)} \\ \begin{array}{|c|} \hline w_1 \quad w_r \\ \hline \end{array} \times \begin{array}{c} P \text{ (patches)} \\ \begin{array}{|c|} \hline \alpha_{k,1} \\ \hline \alpha_{k,r} \\ \hline \end{array} \\ r \end{array} \right) A_k$$

The diagram illustrates the matrix product structure of the output $f_k(x)$. It is expressed as the trace of a product of three matrices. The first matrix is a vertical rectangle labeled $Z(x)$ with dimensions d_1 (width) and P (height). The second matrix is a vertical rectangle labeled w (filters) with dimensions d_1 (width) and r (height). The third matrix is a horizontal rectangle labeled α (patches) with dimensions r (width) and P (height). The product of the second and third matrices is enclosed in a dashed box labeled A_k . The trace operation tr is applied to the product of the three matrices.

Case 2: Non-linear activations: RKHS filters

For non-linear activation function, we relax the class of the CNN filters to an RKHS.

We define a positive semi-definite kernel function K . With this, and a sufficiently smooth activation function we can show that the filter 'h' is contained in RKHS induced by K .

Let $S := \{z_p(x_i) : p \in [P], i \in [n]\}$ be a set of patches in the training data set.

According to the representer theorem, function value for any patch $z_p(x_i) \in S$ can be represented as:

--Let $\mathbf{K} \in \mathbb{R}^{np \times np}$ be a symmetric kernel matrix. To optimize the number of calculations, we factorize $\mathbf{K} = \mathbf{Q}\mathbf{Q}^T$, where $\mathbf{Q} \in \mathbb{R}^{np \times m}$.

We interpret $Q_{(i,p)}$ as a feature vector in place of original patch vector $z_p(x_i)$, and replace in the filter equation from earlier as:

$$h(z_p(x_i)) = \langle Q_{(i,p)}, w \rangle \quad \text{where} \quad w := \sum_{(i',p')} c_{i',p'} Q_{(i',p')}.$$

It suffices to know about the vector “w” to learn the filter “h”.

--We compute a patch matrix for the training data, and then test data. If input given is x ,

-The p -th row of this matrix is the feature vector for patch p , which is equal to $Q^T v(z_p(x)) \in \mathbb{R}^m$, where for any patch z , the vector $v(z)$ is defined as a np -dimensional vector whose (i, p) -th coordinate is equal to $K(z, z_p(x_i))$.

-Then compute the predictor.

Algorithm 1: Learning two-layer CCNNs

Input: Data $\{(x_i, y_i)\}_{i=1}^n$, kernel function \mathcal{K} , regularization parameter $R > 0$, number of filters r .

1. Construct a kernel matrix $K \in \mathbb{R}^{nP \times nP}$ such that the entry at column (i, p) and row (i', p') is equal to $\mathcal{K}(z_p(x_i), z_{p'}(x_{i'}))$. Compute a factorization $K = QQ^\top$ or an approximation $K \approx QQ^\top$, where $Q \in \mathbb{R}^{nP \times m}$.
2. For each x_i , construct patch matrix $Z(x_i) \in \mathbb{R}^{P \times m}$ whose p -th row is the (i, p) -th row of Q , where $Z(\cdot)$ is defined in Section 3.2.
3. Solve the following optimization problem to obtain a matrix $\hat{A} = (\hat{A}_1, \dots, \hat{A}_{d_2})$:

$$\hat{A} \in \underset{\|A\|_* \leq R}{\operatorname{argmin}} \tilde{\mathcal{L}}(A) \quad \text{where} \quad \tilde{\mathcal{L}}(A) := \sum_{i=1}^n \mathcal{L}\left(\left(\operatorname{tr}(Z(x_i)A_1), \dots, \operatorname{tr}(Z(x_i)A_{d_2})\right); y_i\right). \quad (12)$$

4. Compute a rank- r approximation $\tilde{A} \approx \hat{U}\hat{V}^\top$ where $\hat{U} \in \mathbb{R}^{m \times r}$ and $\hat{V} \in \mathbb{R}^{Pd_2 \times r}$.

Output: Return the predictor $\hat{f}_{\text{ccnn}}(x) := (\operatorname{tr}(Z(x)\hat{A}_1), \dots, \operatorname{tr}(Z(x)\hat{A}_{d_2}))$ and the convolutional layer output $H(x) := \hat{U}^\top(Z(x))^\top$.

Activation function and kernel functions used:

Kernel functions used - Gaussian kernel and inverse polynomial

Activation functions used -we have to choose such that the generalization error of F_{CCNN} is comparable to that of best CNN models.

We consider following types of activation functions:

- (a) arbitrary polynomial functions
- (b) sinusoid activation function $\sigma(t) := \sin(t)$
- (c) erf function $\sigma_{\text{erf}}(t) := 2/\sqrt{\pi} \int_0^t e^{-z^2} dz$ (approximation to sigmoid)
- (d) a smoothed hinge loss $\sigma_{\text{sh}}(t) := \int_{-\infty}^t 0.5 (\sigma_{\text{erf}}(z) + 1) dz$

Inverse polynomial kernel captures all the 4 categories (a,b,c,d,).

Gaussian kernel captures only the first 2 categories(a,b).

***sigmoid and ReLU functions are not used because they fail to converge quickly**

Generally,by choosing a hyperparameter R in the given algorithm,expected error is at most $\frac{R}{\sqrt{n}}$ where $R = C_\sigma(B_1)B_2r$

$$\mathbb{E}_{X,Y}[\mathcal{L}(\hat{f}_{ccnn}(X); Y)] \leq \inf_{f \in \mathcal{F}_{cnn}} \mathbb{E}_{X,Y}[\mathcal{L}(f(X); Y)] + \frac{c LC_\sigma(B_1)B_2r \sqrt{\log(nP) \mathbb{E}_X[\|K(X)\|_2]}}{\sqrt{n}},$$

where $c > 0$ is a universal constant.

Learning multi-layer CCNNs:

For models with more than 1 hidden layer, the CCNN theory doesn't apply.

We need to compute the filters explicitly unlike in the 2-layer case.

We approximate the parametric matrix A as a rank r approximation $A \approx \hat{U} \hat{V}^T$ and set the j th filter h_j to the mapping

$$z \mapsto \langle \hat{U}_j, Q^\dagger v(z) \rangle$$

We apply all the r filters to all patches of the input. The resulting output is :

$$H(x) := \hat{U}^T Z(x)^T \quad \{r * P \text{ vector}\}$$

For multilayer CCNNs , two important tasks are average pooling and processing of multi channel inputs.

Average pooling reduces the output dimension of convolution layer. Because of the linearity of CCNN, pooling can be done before convolution.

Processing multi channel inputs: a multi channel patch vector is defined as a concatenation of patch vector for each channel. we form the feature matrix $Z(x)$ using these concatenated patch vectors and continue the process as in the previously given algorithm.

Algorithm 2: Learning multi-layer CCNNs

Input: Data $\{(x_i, y_i)\}_{i=1}^n$, kernel function \mathcal{K} , number of layers m , regularization parameters R_1, \dots, R_m , number of filters r_1, \dots, r_m .

Define $H_1(x) = x$.

For each layer $s = 2, \dots, m$:

- Train a two-layer network by Algorithm 1, taking $\{(H_{s-1}(x_i), y_i)\}_{i=1}^n$ as training examples and R_s, r_s as parameters. Let H_s be the output of the convolutional layer and \hat{f}_s be the predictor.

Output: Predictor \hat{f}_m and the top convolutional layer output H_m .

Cifar-10 dataset results

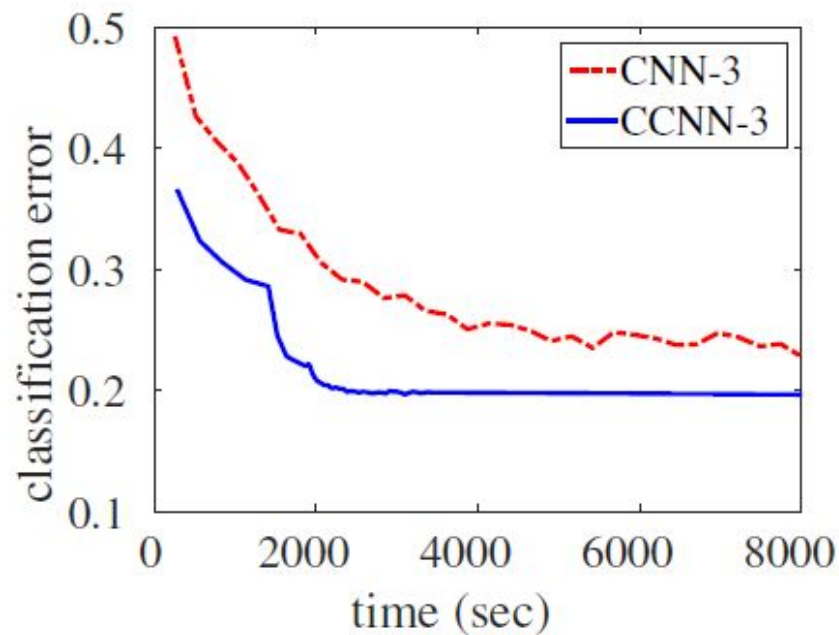


Fig : The convergence of CNN-3 and CCNN-3 on the CIFAR-10 dataset.

CNN-3 represents convolution neural networks with 3 hidden layers.

THANK YOU