

A Network Architectures

Layer	CNN-3	CNN-9	CNN-18	CNN-45	CNN-60	CNN-69
Conv1.x	$[3 \times 3, 64] \times 1$	$[3 \times 3, 64] \times 3$	$[3 \times 3, 64] \times 6$	$[3 \times 3, 64] \times 15$	$[3 \times 3, 64] \times 20$	$[3 \times 3, 64] \times 23$
Pool1	2x2 Max Pooling, Stride 2					
Conv2.x	$[3 \times 3, 96] \times 1$	$[3 \times 3, 96] \times 3$	$[3 \times 3, 96] \times 6$	$[3 \times 3, 96] \times 15$	$[3 \times 3, 96] \times 20$	$[3 \times 3, 96] \times 23$
Pool2	2x2 Max Pooling, Stride 2					
Conv3.x	$[3 \times 3, 128] \times 1$	$[3 \times 3, 128] \times 3$	$[3 \times 3, 128] \times 6$	$[3 \times 3, 128] \times 15$	$[3 \times 3, 128] \times 20$	$[3 \times 3, 128] \times 23$
Pool3	2x2 Max Pooling, Stride 2					
Fully Connected	256	256	256	256	256	256

Table 5: Our plain CNN architectures with different convolutional layers. Conv1.x, Conv2.x and Conv3.x denote convolution units that may contain multiple convolution layers. E.g., $[3 \times 3, 64] \times 3$ denotes 3 cascaded convolution layers with 64 filters of size 3×3 .

Layer	ResNet-32 for Section 4.2	ResNet-32 for Section 4.3	ResNet-18 for Section 4.6
Conv0.x	N/A	N/A	$[7 \times 7, 256]$, Stride 2 3×3 , Max Pooling, Stride 2
Conv1.x	$[3 \times 3, 64] \times 1$ $[3 \times 3, 64] \times 5$	$[3 \times 3, 96] \times 1$ $[3 \times 3, 96] \times 5$	$[3 \times 3, 256] \times 2$
Conv2.x	$[3 \times 3, 96] \times 5$	$[3 \times 3, 192] \times 5$	$[3 \times 3, 512] \times 2$
Conv3.x	$[3 \times 3, 128] \times 5$	$[3 \times 3, 384] \times 5$	$[3 \times 3, 768] \times 2$
Conv4.x	N/A	N/A	$[3 \times 3, 1024] \times 2$
	Average Pooling		

Table 6: Our ResNet architectures with different convolutional layers. Conv0.x, Conv1.x, Conv2.x, Conv3.x and Conv4.x denote convolution units that may contain multiple convolutional layers, and residual units are shown in double-column brackets. Conv1.x, Conv2.x and Conv3.x usually operate on different size feature maps. These networks are essentially the same as [6], but some may have different number of filters in each layer. The downsampling is performed by convolutions with a stride of 2. E.g., $[3 \times 3, 64] \times 4$ denotes 4 cascaded convolution layers with 64 filters of size 3×3 , and S2 denotes stride 2.

B Experimental Details for Imagenet-2012

For the input data of the Imagenet-2012 experiment, we only use the minimum data augmentation. Specifically, we first resize the images to 256×256 resolution and then randomly crop patches of size 224×224 from the resized images. Besides that, we also randomly flip the image horizontally. For SphereResNet-18-v1, we use the cosine SphereConv and the cosine W-Softmax loss. For SphereResNet-18-v2, we use the cosine SphereConv and the softmax loss. Generally, we find that the standard softmax loss and all kinds of W-Softmax loss usually have similar empirical performance. Note that, we could obtain better performance by using the other SphereConvs (sigmoid SphereConv with $k = 0.3$ is a good choice), but it requires more GPU memory. Due to the width of our architecture and the limitation of GPU memory, the mini-batch size is set to 40 for all methods in the Imagenet-2012 experiment.

C More Discussions for Sphere-normalized Softmax Loss

The sphere-normalized softmax (S-Softmax) loss is essentially applying the SphereConv to the fully connected layer in the softmax loss². However, simply applying the SphereConv can not make such loss work, because this loss function is difficult to converge in practice. To address this, we rescale the logit output of the S-Softmax loss with a scaling factor s . Therefore, the output range is changed from $[-1, 1]$ to $[-s, s]$. Typically, setting s from 10 to 70 works pretty well in practice. We could also use the cross-validation strategy to set the hyperparameter s .

²The softmax loss is defined as the combination of the last fully connected layer, the softmax function and the cross-entropy loss.

D Proofs of Lemmas

D.1 Proof of Lemma 1

The gradient is

$$\nabla \mathcal{G}(\mathbf{U}, \mathbf{V}) = \begin{bmatrix} \nabla_{\mathbf{U}} \mathcal{G}(\mathbf{U}, \mathbf{V}) \\ \nabla_{\mathbf{V}} \mathcal{G}(\mathbf{U}, \mathbf{V}) \end{bmatrix} = \begin{bmatrix} (\mathbf{U}\mathbf{V}^\top - \mathbf{F})\mathbf{V} \\ (\mathbf{V}\mathbf{U}^\top - \mathbf{F}^\top)\mathbf{U} \end{bmatrix}$$

The Hessian matrix is

$$\begin{aligned} \nabla^2 \mathcal{G}(\mathbf{U}, \mathbf{V}) &= \begin{bmatrix} \nabla_{\mathbf{U}}^2 \mathcal{G}(\mathbf{U}, \mathbf{V}) & \nabla_{\mathbf{U}, \mathbf{V}}^2 \mathcal{G}(\mathbf{U}, \mathbf{V}) \\ \nabla_{\mathbf{V}, \mathbf{U}}^2 \mathcal{G}(\mathbf{U}, \mathbf{V}) & \nabla_{\mathbf{V}}^2 \mathcal{G}(\mathbf{U}, \mathbf{V}) \end{bmatrix} \\ &= \begin{bmatrix} \mathbf{V}^\top \mathbf{V} \otimes \mathbf{I}_n & (\mathbf{U}\mathbf{V}^\top - \mathbf{F}) \otimes \mathbf{I}_k + \mathbf{U} \boxtimes \mathbf{V} \\ (\mathbf{V}\mathbf{U}^\top - \mathbf{F}^\top) \otimes \mathbf{I}_k + \mathbf{V} \boxtimes \mathbf{U} & \mathbf{U}^\top \mathbf{U} \otimes \mathbf{I}_m \end{bmatrix}, \end{aligned} \quad (12)$$

where \mathbf{I}_n is an $n \times n$ identity matrix for any integer n , given matrices $\mathbf{A} \in \mathbb{R}^{n \times r}$ and $\mathbf{B} \in \mathbb{R}^{m \times k}$ with $\mathbf{A}_{:,i}$ denoting the i -th column of \mathbf{A} , $\mathbf{A} \boxtimes \mathbf{B} \in \mathbb{R}^{nk \times mr}$ is defined as

$$\mathbf{A} \boxtimes \mathbf{B} = \begin{bmatrix} \mathbf{A}_{:,1} \mathbf{B}_{:,1}^\top & \mathbf{A}_{:,2} \mathbf{B}_{:,1}^\top & \cdots & \mathbf{A}_{:,r} \mathbf{B}_{:,1}^\top \\ \mathbf{A}_{:,1} \mathbf{B}_{:,2}^\top & \mathbf{A}_{:,2} \mathbf{B}_{:,2}^\top & \cdots & \mathbf{A}_{:,r} \mathbf{B}_{:,2}^\top \\ \vdots & \vdots & \ddots & \vdots \\ \mathbf{A}_{:,1} \mathbf{B}_{:,k}^\top & \mathbf{A}_{:,2} \mathbf{B}_{:,k}^\top & \cdots & \mathbf{A}_{:,r} \mathbf{B}_{:,k}^\top \end{bmatrix}.$$

At a global optimum, we have $\mathbf{U}\mathbf{V}^\top = \mathbf{F}$. Then it is easy to see that for any real c , if $\tilde{\mathbf{U}} = c\mathbf{U}$ and $\tilde{\mathbf{V}} = \mathbf{V}/c$, then we have

$$\nabla^2 \mathcal{G}(\tilde{\mathbf{U}}, \tilde{\mathbf{V}}) = \begin{bmatrix} \frac{1}{c^2} \mathbf{V}^\top \mathbf{V} \otimes \mathbf{I}_n & \mathbf{U} \boxtimes \mathbf{V} \\ \mathbf{V} \boxtimes \mathbf{U} & c^2 \mathbf{U}^\top \mathbf{U} \otimes \mathbf{I}_m \end{bmatrix}$$

We have that at a global optimal point, $\nabla^2 \mathcal{G}(\mathbf{U}, \mathbf{V})$ is a positive semidefinite matrix with the smallest eigenvalue equal to 0. Specifically, due to the existence of the invariance, i.e., $\mathbf{U}\mathbf{V}^\top = \mathbf{U}\mathbf{R}(\mathbf{V}\mathbf{R})^\top$ for any orthogonal matrix $\mathbf{R} \in \mathbb{R}^{r \times r}$, there are $r(r-1)/2$ number of eigenvectors of $\nabla^2 \mathcal{G}(\mathbf{U}, \mathbf{V})$ at $\mathbf{U}\mathbf{V}^\top = \mathbf{F}$ corresponding to 0 eigenvalue [10]. Then for any $c > 1$, we have

$$\begin{aligned} \text{Tr}(\nabla^2 \mathcal{G}(\tilde{\mathbf{U}}, \tilde{\mathbf{V}})) &= \frac{1}{c^2} \text{Tr}(\mathbf{V}^\top \mathbf{V} \otimes \mathbf{I}_n) + c^2 \text{Tr}(\mathbf{U}^\top \mathbf{U} \otimes \mathbf{I}_m) \\ &\geq \frac{c^2}{2} (\text{Tr}(\mathbf{V}^\top \mathbf{V} \otimes \mathbf{I}_n) + \text{Tr}(\mathbf{U}^\top \mathbf{U} \otimes \mathbf{I}_m)) = \frac{c^2}{2} \text{Tr}(\nabla^2 \mathcal{G}(\mathbf{U}, \mathbf{V})). \end{aligned}$$

This indicates that the largest eigenvalue of $\nabla^2 \mathcal{G}(\tilde{\mathbf{U}}, \tilde{\mathbf{V}})$ is on the order of $\Theta(c^2)$ times the largest eigenvalue of $\nabla^2 \mathcal{G}(\mathbf{U}, \mathbf{V})$ following the perturbation bound analysis [15] and \mathbf{U} and \mathbf{V} are balanced. Using a similar idea, the smallest nonzero eigenvalue of $\nabla^2 \mathcal{G}(\tilde{\mathbf{U}}, \tilde{\mathbf{V}})$ is no greater than the smallest nonzero eigenvalue of $\nabla^2 \mathcal{G}(\mathbf{U}, \mathbf{V})$, which results in our claim on the restricted condition number.

D.2 Proof of Lemma 2

The gradient of $\mathcal{G}_S(\mathbf{U}, \mathbf{V})$ is

$$\nabla \mathcal{G}_S(\mathbf{U}, \mathbf{V}) = \begin{bmatrix} \nabla_{\mathbf{U}} \mathcal{G}_S(\mathbf{U}, \mathbf{V}) \\ \nabla_{\mathbf{V}} \mathcal{G}_S(\mathbf{U}, \mathbf{V}) \end{bmatrix} \quad \text{with}$$

$$\begin{aligned} \nabla_{\mathbf{U}} \mathcal{G}_S(\mathbf{U}, \mathbf{V}) &= \mathbf{D}_{\mathbf{U}} (\mathbf{D}_{\mathbf{U}} \mathbf{U} \mathbf{V}^\top \mathbf{D}_{\mathbf{V}} - \mathbf{F}) \mathbf{D}_{\mathbf{V}} \mathbf{V} - (\mathbf{D}_{\mathbf{U}}^3 (\mathbf{D}_{\mathbf{U}} \mathbf{U} \mathbf{V}^\top \mathbf{D}_{\mathbf{V}} - \mathbf{F}) \otimes_k (\mathbf{U} \mathbf{V}^\top \mathbf{D}_{\mathbf{V}})) \odot \mathbf{U}, \\ \nabla_{\mathbf{V}} \mathcal{G}_S(\mathbf{U}, \mathbf{V}) &= \mathbf{D}_{\mathbf{V}} (\mathbf{D}_{\mathbf{V}} \mathbf{V} \mathbf{U}^\top \mathbf{D}_{\mathbf{U}} - \mathbf{F}^\top) \mathbf{D}_{\mathbf{U}} \mathbf{U} - (\mathbf{D}_{\mathbf{V}}^3 (\mathbf{D}_{\mathbf{V}} \mathbf{V} \mathbf{U}^\top \mathbf{D}_{\mathbf{U}} - \mathbf{F}^\top) \otimes_k (\mathbf{V} \mathbf{U}^\top \mathbf{D}_{\mathbf{U}})) \odot \mathbf{V}, \end{aligned}$$

Note that after each iteration of SGD, we perform the entry-wise normalization for both \mathbf{U} and \mathbf{V} , which means $\mathbf{D}_{\mathbf{U}} = \mathbf{I}_n$ and $\mathbf{D}_{\mathbf{V}} = \mathbf{I}_m$. Then the gradient of $\mathcal{G}_S(\mathbf{U}, \mathbf{V})$ is

$$\nabla \mathcal{G}_S(\mathbf{U}, \mathbf{V}) = \begin{bmatrix} \nabla_{\mathbf{U}} \mathcal{G}_S(\mathbf{U}, \mathbf{V}) \\ \nabla_{\mathbf{V}} \mathcal{G}_S(\mathbf{U}, \mathbf{V}) \end{bmatrix} = \begin{bmatrix} (\mathbf{U}\mathbf{V}^\top - \mathbf{F})\mathbf{V} - ((\mathbf{U}\mathbf{V}^\top - \mathbf{F}) \otimes_k (\mathbf{U}\mathbf{V}^\top)) \odot \mathbf{U} \\ (\mathbf{V}\mathbf{U}^\top - \mathbf{F}^\top)\mathbf{U} - ((\mathbf{V}\mathbf{U}^\top - \mathbf{F}^\top) \otimes_k (\mathbf{V}\mathbf{U}^\top)) \odot \mathbf{V} \end{bmatrix},$$

where given matrices $\mathbf{A}, \mathbf{B} \in \mathbb{R}^{n \times m}$ with $\mathbf{A}_{:,i}$ denoting the i -th column of \mathbf{A} , $\mathbf{A} \odot \mathbf{B} \in \mathbb{R}^{n \times m}$ is the Hadamard (pointwise) product, and the operation $\mathbf{A} \otimes_k \mathbf{B} \in \mathbb{R}^{n \times k}$ is defined as

$$\mathbf{A} \otimes_k \mathbf{B} = \begin{bmatrix} \mathbf{A}_{1,:} \mathbf{B}_{1,:}^\top \\ \mathbf{A}_{2,:} \mathbf{B}_{2,:}^\top \\ \vdots \\ \mathbf{A}_{n,:} \mathbf{B}_{n,:}^\top \end{bmatrix} \mathbf{1}_{1 \times k},$$

where $\mathbf{1}_{1 \times k}$ is a $1 \times k$ vector with all entries equal to 1.

Consequently, the Hessian matrix is

$$\begin{aligned} \nabla^2 \mathcal{G}_S(\mathbf{U}, \mathbf{V}) &= \begin{bmatrix} \nabla_{\mathbf{U}}^2 \mathcal{G}_S(\mathbf{U}, \mathbf{V}) & \nabla_{\mathbf{U}, \mathbf{V}}^2 \mathcal{G}_S(\mathbf{U}, \mathbf{V}) \\ \nabla_{\mathbf{V}, \mathbf{U}}^2 \mathcal{G}_S(\mathbf{U}, \mathbf{V}) & \nabla_{\mathbf{V}}^2 \mathcal{G}_S(\mathbf{U}, \mathbf{V}) \end{bmatrix} \text{ with} \\ \nabla_{\mathbf{U}}^2 \mathcal{G}_S(\mathbf{U}, \mathbf{V}) &= \mathbf{V}^\top \mathbf{V} \otimes \mathbf{I}_n - \text{diag}(\text{vec}((\mathbf{U}\mathbf{V}^\top - \mathbf{F}) \otimes_k (\mathbf{U}\mathbf{V}^\top))) \\ &\quad - \begin{bmatrix} \text{diag}(\mathbf{U}_{:,1} \odot ((2\mathbf{U}\mathbf{V}^\top - \mathbf{F})\mathbf{V}_{:,1})) & \cdots & \text{diag}(\mathbf{U}_{:,1} \odot ((2\mathbf{U}\mathbf{V}^\top - \mathbf{F})\mathbf{V}_{:,k})) \\ \vdots & \ddots & \vdots \\ \text{diag}(\mathbf{U}_{:,k} \odot ((2\mathbf{U}\mathbf{V}^\top - \mathbf{F})\mathbf{V}_{:,1})) & \cdots & \text{diag}(\mathbf{U}_{:,k} \odot ((2\mathbf{U}\mathbf{V}^\top - \mathbf{F})\mathbf{V}_{:,k})) \end{bmatrix} \\ \nabla_{\mathbf{U}, \mathbf{V}}^2 \mathcal{G}_S(\mathbf{U}, \mathbf{V}) &= \mathbf{I}_k \otimes (\mathbf{U}\mathbf{V}^\top - \mathbf{F}) + \mathbf{U} \boxtimes \mathbf{V} \\ &\quad - \begin{bmatrix} (2\mathbf{U}\mathbf{V}^\top - \mathbf{F}) \odot ((\mathbf{U}_{:,1} \odot \mathbf{U}_{:,1})\mathbf{1}_{1 \times m}) & \cdots & (2\mathbf{U}\mathbf{V}^\top - \mathbf{F}) \odot ((\mathbf{U}_{:,1} \odot \mathbf{U}_{:,k})\mathbf{1}_{1 \times m}) \\ \vdots & \ddots & \vdots \\ (2\mathbf{U}\mathbf{V}^\top - \mathbf{F}) \odot ((\mathbf{U}_{:,k} \odot \mathbf{U}_{:,1})\mathbf{1}_{1 \times m}) & \cdots & (2\mathbf{U}\mathbf{V}^\top - \mathbf{F}) \odot ((\mathbf{U}_{:,k} \odot \mathbf{U}_{:,k})\mathbf{1}_{1 \times m}) \end{bmatrix} \\ \nabla_{\mathbf{V}, \mathbf{U}}^2 \mathcal{G}_S(\mathbf{U}, \mathbf{V}) &= (\nabla_{\mathbf{U}, \mathbf{V}}^2 \mathcal{G}_S(\mathbf{U}, \mathbf{V}))^\top \\ \nabla_{\mathbf{V}}^2 \mathcal{G}_S(\mathbf{U}, \mathbf{V}) &= \mathbf{U}^\top \mathbf{U} \otimes \mathbf{I}_n - \text{diag}(\text{vec}((\mathbf{V}\mathbf{U}^\top - \mathbf{F}^\top) \otimes_k (\mathbf{V}\mathbf{U}^\top))) \\ &\quad - \begin{bmatrix} \text{diag}(\mathbf{V}_{:,1} \odot ((2\mathbf{V}\mathbf{U}^\top - \mathbf{F}^\top)\mathbf{U}_{:,1})) & \cdots & \text{diag}(\mathbf{V}_{:,1} \odot ((2\mathbf{V}\mathbf{U}^\top - \mathbf{F}^\top)\mathbf{U}_{:,k})) \\ \vdots & \ddots & \vdots \\ \text{diag}(\mathbf{V}_{:,k} \odot ((2\mathbf{V}\mathbf{U}^\top - \mathbf{F}^\top)\mathbf{U}_{:,1})) & \cdots & \text{diag}(\mathbf{V}_{:,k} \odot ((2\mathbf{V}\mathbf{U}^\top - \mathbf{F}^\top)\mathbf{U}_{:,k})) \end{bmatrix} \end{aligned}$$

Then we have $\lambda_i(\nabla^2 \mathcal{G}_S(\tilde{\mathbf{U}}, \tilde{\mathbf{V}})) = \lambda_i(\nabla^2 \mathcal{G}_S(\mathbf{U}, \mathbf{V}))$ for all $i \in [(n+m)k] = \{1, 2, \dots, (n+m)k\}$ by noticing that we normalize the data as $\frac{\mathbf{U}_{i,j}}{\|\mathbf{U}_{i,:}\|_2}$ for all $i \in [n]$ and $\frac{\mathbf{V}_{i,j}}{\|\mathbf{V}_{i,:}\|_2}$ for all $i \in [m]$. This finishes the proof.