

IFT6135 – Homework #2

Convolutional neural networks and regularization - Theory

Jean-Philippe Reid
701300

April 15, 2018

1 Convolutions

Compute the full, valid, and same convolution (with kernel flipping) for the following 1D matrices :

$$(1, 2, 3, 4) \circledast (1, 0, 2) \tag{1.1}$$

It is important to note that the kernel needs to be rotated in order to compute the convolution.

Valid convolution: padding = 0

$$\begin{aligned} [1, 2, 3, 4] \circledast [1, 0, 2] &= [2 \times 1 + 0 \times 2 + 1 \times 3, 2 \times 2 + 0 \times 3 + 4 \times 1] \\ &= [5, 8] \end{aligned} \tag{1.2}$$

Same convolution: padding = 1

$$[0, 1, 2, 3, 4, 0] \circledast [1, 0, 2] = [2, 5, 8, 6] \tag{1.3}$$

Full convolution: padding = 2

$$[0, 0, 1, 2, 3, 4, 0, 0] \circledast [1, 0, 2] = [1, 2, 5, 8, 6, 8] \tag{1.4}$$

2 Convolutional Neural Networks

Let's assume the input of size $3 \times 256 \times 256$. The first layer convolves 64 8×8 kernels with the input, using a stride of 2 and no padding. The second layer subsamples the output of the first layer with a 5×5 non-overlapping max pooling. The third layer convolves 128 4×4 kernels with a stride of 1 and a zero-padding of size 1 on each border.

(a)

$$\begin{aligned} o &= \frac{i + 2p - k}{s} + 1 \\ &= \frac{256 - 8}{2} + 1 = 125 \end{aligned} \tag{2.1}$$

The size of the output of the first layer is $64 \times 125 \times 125$.

The second layer corresponds to a max pooling of 5×5 with non-overlapping. We find that,

$$\begin{aligned} o &= \frac{i + 2p - k}{s} + 1 \\ &= \frac{125 - 5}{5} + 1 = 25 \end{aligned} \tag{2.2}$$

The size of the second layer output is $64 \times 25 \times 25$.

The third and last layer convolves 128 4×4 kernels with a stride of 1 and padding of 1. We find that,

$$\begin{aligned} o &= \frac{i + 2p - k}{s} + 1 \\ &= \frac{25 + 2 - 4}{1} + 1 = 24 \end{aligned} \tag{2.3}$$

The size of the second layer output is $128 \times 24 \times 24 = 73728$ parameters.

(b) There are $128 \times 64 \times 4 \times 4$ weights in the last layer, i.e. 131072 weights.

3 Kernel configuration for CNNs

- (a) *For the first convolutional layer, the input shape is $3 \times 64 \times 64$ and the output shape is $64 \times 32 \times 32$. Provide a configuration of the convolution operation that would match the input/output shapes (i.e. provide the kernel, stride, padding and dilatation between kernel elements).*

$$o = \frac{i + 2p - k}{s} + 1 \quad (3.1)$$

where o , i , p , and s are the output, the input, the padding, the kernel and the stride parameters.

The dilatation has for main effect to increase the reach of the kernel filters, i.e.

$$\begin{aligned} \tilde{k} &= k + (k - 1)(d - 1) \\ &= k + kd - k - d + 1 \\ &= d(k - 1) + 1 \end{aligned} \quad (3.2)$$

For $d = 1$, $\tilde{k} = k$, obviously. But for $d = 0$, $\tilde{k} = 1$.

- (i) *Assuming no dilatation and kernel size of 8×8 .*

For the problem in hand, we have that $d = 1$ and $k = 8$.

$$\begin{aligned} o &= \frac{i + 2p - k}{s} + 1 \\ 32 &= \frac{64 + 2p - 8}{s} + 1 \\ 31 &= \frac{56 + 2p}{s} \\ 31s &= 56 + 2p \end{aligned} \quad (3.3)$$

If we set $s = 2$, we find that $p = 3$.

Answer: $k = 8, d = 1, s = 2$ and $p = 3$.

- (ii) *Assuming dilatation of 6 ($d = 7$) between kernel elements and stride of 2 ($s = 2$).*

We have that $d = 7$ and $s = 2$

$$\begin{aligned}
o &= \frac{i + 2p - \tilde{k}}{s} + 1 \\
&= \frac{i + 2p - d(k - 1) - 1}{s} + 1 \\
32 &= \frac{64 + 2p - 7(k - 1) - 1}{2} + 1 \\
32 &= \frac{64 + 2p - 7k + 7 - 1}{2} + 1 \\
31 &= \frac{70 + 2p - 7k}{2} \\
62 &= 70 + 2p - 7k \\
-8 &= 2p - 7k
\end{aligned} \tag{3.4}$$

If we set $k = 2$, we find that $p = 3$.

Answer: $k = 2, d = 7, s = 2$ and $p = 3$.

- (b) *For the second pooling layer, the input shape is $64 \times 32 \times 32$ and the output shape is $64 \times 8 \times 8$. Provide a configuration of the pooling operation that would match the input/output shapes, assuming no overlapping of pooling windows and no padding (i.e. what kernel size and stride to use).*

$$\begin{aligned}
o &= \frac{i + 2p - k}{s} + 1 \\
8 &= \frac{32 - k}{s} + 1 \\
7 &= \frac{32 - k}{s}
\end{aligned} \tag{3.5}$$

Answer: $k = 4, s = 4, p = 0$

- (c) *Without any padding and using the input from question 2 ($64 \times 32 \times 32$), what would have been the output size should the kernel be of size 8×8 , but with the stride 4×4 ?*

$$\begin{aligned}
o &= \frac{i + 2p - k}{s} + 1 \\
&= \frac{32 - 8}{4} + 1 \\
&= \frac{24}{4} + 1 = 7
\end{aligned} \tag{3.6}$$

Answer: output size = 7×7 .

- (d) *In the last convolutional layer, the input shape is $64 \times 8 \times 8$ and the output shape is $128 \times 4 \times 4$. Provide a configuration of the convolution operation that would match the input/output shapes (i.e. provide the kernel, stride, padding and dilatation between kernel elements).*

- (i) Let's assume that $p = 0$ and $d = 1$.

$$\begin{aligned}
 o &= \frac{i + 2p - k}{s} + 1 \\
 4 &= \frac{8 - k}{s} + 1 \\
 3 &= \frac{8 - k}{s}
 \end{aligned}
 \tag{3.7}$$

If we set $k = 2$, we find that $s = 2$.

Answer: $k = 2, d = 1, s = 2$ and $p = 0$.

- (ii) Let's assume that $p = 2$ and $d = 2$.

$$\begin{aligned}
 o &= \frac{i + 2p - \tilde{k}}{s} + 1 \\
 4 &= \frac{8 + 2p - d(k - 1) - 1}{s} + 1 \\
 4 &= \frac{8 + 4 - 2(k - 1) - 1}{s} + 1 \\
 4 &= \frac{8 + 4 - 2k + 2 - 1}{s} + 1 \\
 3 &= \frac{13 - 2k}{s}
 \end{aligned}
 \tag{3.8}$$

If we set $s = 1$, we find that $k = 5$.

Answer: $k = 5, d = 2, s = 1$ and $p = 2$.

(iii) Let's assume that $p = 1$ and $d = 1$.

$$\begin{aligned}
 o &= \frac{i + 2p - \tilde{k}}{s} + 1 \\
 4 &= \frac{8 + 2p - d(k - 1) - 1}{s} + 1 \\
 3 &= \frac{8 + 2 - k}{s} + 1 \\
 2 &= \frac{10 - k}{s} \\
 2 &= \frac{10 - k}{s}
 \end{aligned}
 \tag{3.9}$$

If we set $k = 4$, we find that $s = 3$.

Answer: $k = 4, d = 1, s = 3$ and $p = 1$.

4 Dropout as weight decay

- (a) When the dropout is activated, each element of the input data is conserved with probability p . Therefore, the data can be expressed as $R \odot X$ where R is a matrix of the same dimension as that X , and where $R_{ij} \sim \text{Ber}(p)$. The reader may interpret R_{ij} as a delta kronecker.

- (b) In this context, the loss is defined as

$$L(\Theta) = \mathbb{E}_{R \sim \text{Ber}(p)} [\|\mathbf{y} - (R \odot \mathbf{x}) \mathbf{w}\|^2]. \quad (4.1)$$

- (c)

$$\begin{aligned} \mathbb{E}_{R \sim \text{Ber}(p)} [\|\mathbf{y} - (R \odot \mathbf{x}) \mathbf{w}\|^2] &= \dots \\ &= \mathbb{E} [\|\mathbf{y}\|^2 - 2\mathbf{y}^\top R \odot \mathbf{x} \mathbf{w} + \mathbf{w}^\top (R \odot \mathbf{x})^\top (R \odot \mathbf{x}) \mathbf{w}] \\ &= \|\mathbf{y}\|^2 - 2\mathbf{y}^\top p \mathbf{x} \mathbf{w} + \mathbf{w}^\top \mathbb{E} [(R \odot \mathbf{x})^\top (R \odot \mathbf{x})] \mathbf{w} \end{aligned} \quad (4.2)$$

where

$$\begin{aligned} \mathbb{E} [(R \odot \mathbf{x})^\top (R \odot \mathbf{x})] &= \mathbb{E} \left[\sum_k b_{ik} x_{ik} x_{jk} b_{jk} \right] \\ &= \sum_k \mathbb{E} [b_{ik} x_{ik} x_{jk} b_{jk}] \end{aligned} \quad (4.3)$$

Let's consider the case where $i = j$,

$$\begin{aligned} \mathbb{E} [b_{ik} x_{ik} x_{ik} b_{ik}] &= \mathbb{E} [b_{ik}^2 x_{ik}^2] \\ &= \mathbb{E} [b_{ik}^2] x_{ik}^2 \\ &= [\text{Var} [b_{ik}] + \mathbb{E} [b_{ik}]^2] x_{ik}^2 \\ &= [p(1-p) + p^2] x_{ik}^2. \end{aligned} \quad (4.4)$$

Then, the case where $i \neq j$,

$$\begin{aligned} \mathbb{E} [b_{ik} x_{ik} x_{jk} b_{jk}] &= \mathbb{E} [b_{ik}] \mathbb{E} [b_{jk}] x_{ik} x_{jk} \\ &= p^2 x_{ik} x_{jk} \end{aligned} \quad (4.5)$$

Equation 4.2 then becomes

$$\begin{aligned} \mathbb{E}_{R \sim \text{Ber}(p)} [\|\mathbf{y} - (R \odot \mathbf{x}) \mathbf{w}\|^2] &= \dots \\ &= \|\mathbf{y}\|^2 - 2\mathbf{y}^\top p \mathbf{x} \mathbf{w} + \mathbf{w}^\top \mathbb{E} [(R \odot \mathbf{x})^\top (R \odot \mathbf{x})] \mathbf{w} \\ &= \|\mathbf{y}\|^2 - 2\mathbf{y}^\top p \mathbf{x} \mathbf{w} + \mathbf{w}^\top [p(1-p) \text{diag}(\mathbf{x}^\top \mathbf{x}) + p^2 \mathbf{x}^\top \mathbf{x}] \mathbf{w} \end{aligned} \quad (4.6)$$

We need to find the minimize the above to obtain the ideal w . We find that

$$\begin{aligned}
\frac{\partial}{\partial w} L(\Theta) &= -y^\top p \mathbf{x} + w p (1 - p) \text{diag}(\mathbf{x}^\top \mathbf{x}) + p^2 w \mathbf{x}^\top \mathbf{x} \\
0 &= -y^\top \mathbf{x} + w [(1 - p) \text{diag}(\mathbf{x}^\top \mathbf{x}) + p \mathbf{x}^\top \mathbf{x}] \\
y^\top \mathbf{x} &= \bar{w} p \left[\frac{(1 - p)}{p} \Gamma^2 + \mathbf{x}^\top \mathbf{x} \right]
\end{aligned} \tag{4.7}$$

$$\boxed{\bar{w}^\star = \left[\frac{(1 - p)}{p} \Gamma^2 + \mathbf{x}^\top \mathbf{x} \right]^{-1} y \mathbf{x}^\top} \tag{4.8}$$

We defined in the above equations $\Gamma^2 = \text{diag}(\mathbf{x}^\top \mathbf{x})$ and $\bar{w} = pw$.

5 Dropout as geometric ensemble

Note: for sake of transparency, we refer the reader to section 7.12 of I. Goodfellow *et al.* for more details.

Let's consider the case of a single linear layer model with softmax output,

$$P(y = j | v) = \hat{y}_j = \text{softmax}(\mathbf{W}^\top \mathbf{v} + \mathbf{b})_j, \quad (5.1)$$

where \mathbf{v} is the input data.

The effect of dropout is to introduce a mask on the vector \mathbf{W} , transforming a significant fraction of the weights to zero. Since the mask is different for every evaluation of the model, we may represent mathematically this procedure by introducing a matrix mask \mathbf{m} of the same dimension as \mathbf{x} which set to zero some of the inputs,

$$\mathbf{x} \rightarrow \mathbf{m} \odot \mathbf{x}. \quad (5.2)$$

The reader may appreciate better the above expression by considering the non vectorize equation, given by $x_{ij} \rightarrow m_{ij}x_{i,j}$ where $m_{i,j}$ is equivalent to the kronecker operator $\delta_{i,j}$.

Thus, Eq. 5.1 becomes

$$P(y = j | \mathbf{v}, \mathbf{m}) = \text{softmax}(\mathbf{W}^\top (\mathbf{m} \odot \mathbf{x}) + \mathbf{b})_j. \quad (5.3)$$

The predictive distribution is defined as the normalized geometrical mean,

$$P_{\text{ensemble}}(y = j | \mathbf{v}) = \frac{\tilde{P}_{\text{ensemble}}(y = j | \mathbf{v})}{\sum_{j'} \tilde{P}_{\text{ensemble}}(y = j' | \mathbf{v})} \quad (5.4)$$

where $\tilde{P}_{\text{ensemble}}(y = j | v)$ is defined as the geometrical mean

$$\tilde{P}_{\text{ensemble}}(y = j | \mathbf{v}) = \left[\prod_{\mathbf{m}} P(y = j | \mathbf{v}, \mathbf{m}) \right]^{1/2^d}. \quad (5.5)$$

In the above expression, d is the number of units that are put to zero due to dropout.

Putting Eqs.5.3 and 5.5 together, we find that

$$\begin{aligned}
P_{\text{ensemble}}(y = j \mid \mathbf{v}) &\propto \left[\prod_{\mathbf{m}} \text{softmax}(\mathbf{W}^\top (\mathbf{m} \odot \mathbf{x}) + \mathbf{b})_j \right]^{1/2^d} \\
&= \left[\prod_{\mathbf{m}} \exp(\mathbf{W}_{j,:}^\top (\mathbf{m} \odot \mathbf{x}) + \mathbf{b}_j) \right]^{1/2^d} \\
&= \exp \left(\sum_{\mathbf{m}} \mathbf{W}_{j,:}^\top (\mathbf{m} \odot \mathbf{x}) + \mathbf{b}_j \right)^{1/2^d} \\
&= \exp \left(\frac{1}{2^d} \sum_{\mathbf{m}} \mathbf{W}_{j,:}^\top (\mathbf{m} \odot \mathbf{x}) + \mathbf{b}_j \right) \\
&= \exp \left(\sum_{\mathbf{m}} \frac{1}{2} \mathbf{W}_{j,:}^\top (\mathbf{m} \odot \mathbf{x}) + \mathbf{b}'_j \right) \\
&= \exp \left(\frac{1}{2} \mathbf{W}_{j,:}^\top \mathbf{x} + \mathbf{b}'_j \right) \tag{5.6}
\end{aligned}$$

where $\mathbf{W}_{j,:}$ represent the j^{th} line of W . This result implies that the scaling with a factor of 0.5 is equivalent to the predictive probability distribution of the ensemble.

6 Normalization

(a) Let us suppose a vector \mathbf{x} and \mathbf{y} related through the affine relation

$$\mathbf{y} = W^T \mathbf{x} + b. \quad (6.1)$$

The batch normalization is defined as

$$\text{BN}(\mathbf{y}) = \gamma \odot \frac{(\mathbf{y} - \mathbb{E}[\mathbf{y}])}{\text{Var}[\mathbf{y}]} + \beta \quad (6.2)$$

For the problem in hand, we have that

$$\begin{aligned} \mathbb{E}[\mathbf{y}] &= \mathbb{E}[(W\mathbf{x} + b)] \\ &= W^T \mathbb{E}[\mathbf{x}] + \mathbb{E}[b] \\ &= b \end{aligned} \quad (6.3)$$

$$\begin{aligned} \text{Var}[\mathbf{y}] &= \text{Var}[(W\mathbf{x} + b)] \\ &= W^T \text{Var}[\mathbf{x}] W + \text{Var}[b] \\ &= W^T W \\ &= \|W\| \end{aligned} \quad (6.4)$$

By combining Eqs. 6.2, 6.3 and 6.4, we find that

$$\begin{aligned} \text{BN}(\mathbf{y}) &= \gamma \odot \frac{(\mathbf{y} - \mathbb{E}[\mathbf{y}])}{\text{Var}[\mathbf{y}]} + \beta \\ &= \gamma \odot \frac{(\mathbf{y} - b)}{\|W\|} + \beta \\ &= \gamma \odot \frac{W^T}{\|W\|} \mathbf{x} + \beta \\ &= \gamma \odot \left(\frac{W}{\|W\|} \right)^T \mathbf{x} + \beta \end{aligned} \quad (6.5)$$

which corresponds to the weight normalization.

(b)

$$\begin{aligned}
\nabla_{\mathbf{u}} L &= \frac{\partial L}{\partial \mathbf{u}} \\
&= \frac{\partial L}{\partial \mathbf{w}} \frac{\partial \mathbf{w}}{\partial \mathbf{u}} \\
&= \frac{\partial L}{\partial \mathbf{w}} \frac{\partial}{\partial \mathbf{u}} \left(g \frac{\mathbf{u}}{\|\mathbf{u}\|} \right) \\
&= \frac{\partial L}{\partial \mathbf{w}} \left(\frac{g}{\|\mathbf{u}\|} - g \frac{\mathbf{u}}{\|\mathbf{u}\|^2} \frac{\mathbf{u}^\top}{\|\mathbf{u}\|} \right) \\
&= \frac{g}{\|\mathbf{u}\|} \left(1 - \frac{\mathbf{u}\mathbf{u}^\top}{\|\mathbf{u}\|^2} \right) \nabla_{\mathbf{w}} L
\end{aligned} \tag{6.6}$$

$1 - \frac{\mathbf{u}\mathbf{u}^\top}{\|\mathbf{u}\|^2}$ is defined as the orthogonal complement projection matrix.

$$\begin{aligned}
1 - \frac{\mathbf{u}\mathbf{u}^\top}{\|\mathbf{u}\|^2} &= 1 - \frac{\mathbf{u}\mathbf{u}^\top}{\mathbf{u}^\top \mathbf{u}} \\
&= 1 - \mathbf{u} (\mathbf{u}^\top \mathbf{u})^{-1} \mathbf{u}^\top \\
&= 1 - \mathbf{P} \\
&= W^*
\end{aligned} \tag{6.7}$$

where $\mathbf{P} = \mathbf{u} (\mathbf{u}^\top \mathbf{u})^{-1} \mathbf{u}^\top$ is defined as the projection matrix. Therefore, we find that

$$\begin{aligned}
\nabla_{\mathbf{u}} L &= \frac{g}{\|\mathbf{u}\|} \left(1 - \frac{\mathbf{u}\mathbf{u}^\top}{\|\mathbf{u}\|^2} \right) \nabla_{\mathbf{w}} L \\
&= \frac{g}{\|\mathbf{u}\|} W^* \nabla_{\mathbf{w}} L
\end{aligned} \tag{6.8}$$

Taken from Tim Salimans *et al.* (2016) *Weight normalization has two effects: it scales the weight gradient, but it also projects the gradient away from the current weight vector. Both effects help to bring the covariance matrix of the gradient closer to identity and benefit optimization.*

(c) The update of \mathbf{u} is given by

$$\mathbf{u}' \leftarrow \mathbf{u} + \alpha \nabla_{\mathbf{u}} L \tag{6.9}$$

where $\nabla_{\mathbf{u}} L \propto W^* \nabla_{\mathbf{w}} L$.

It is clear that $\nabla_{\mathbf{u}} L$ is orthogonal to \mathbf{w} since we project away from it due to the orthogonal complement projection matrix W^* . Since \mathbf{w} is proportional to \mathbf{u} , we find

that $\nabla_{\mathbf{u}}L \perp \mathbf{u}$. From the Pythagorean theorem, we find that

$$\begin{aligned}
\mathbf{u}' &= \mathbf{u} + \alpha \nabla_{\mathbf{u}}L \\
\|\mathbf{u}'\| &= (\|\mathbf{u}\|^2 + \alpha^2 \|\nabla_{\mathbf{u}}L\|^2)^{1/2} \\
&= (\|\mathbf{u}\|^2 + \alpha^2 k^2 \|\mathbf{u}\|^2)^{1/2} \\
&= \|\mathbf{u}\| (1 + \alpha^2 k^2)^{1/2} > \|\mathbf{u}\|
\end{aligned} \tag{6.10}$$

We see from the above equation that $\|\mathbf{u}'\|$ grows monotonically with the number of updates. Further, it also explain its evolution with the learning rate, i.e. the growth increases with α .

It is important to note that we assumed that $\nabla_{\mathbf{u}}L \propto k\mathbf{u}$ in the above equation. This is a consequence of Eq. 6.5. This is not a necessary condition, though.