

ELEG 5491: Homework #1

Due on Tuesday, February 14, 2017, 3:30pm (in class)

Xiaogang Wang

Problem 1

[15 points]

Cross entropy is often used as the objective function when training neural networks in classification problems. Suppose the training set includes N training pairs $\mathcal{D} = \{(\mathbf{x}_i^{(\text{train})}, y_i^{(\text{train})})\}_{i=1}^N$, where $\mathbf{x}_i^{(\text{train})}$ is a training sample and $y_i^{(\text{train})} \in \{1, \dots, c\}$ is its class label. \mathbf{z}_i is the output of the network given input $\mathbf{x}_i^{(\text{train})}$ and the nonlinearity of the output layer is softmax. \mathbf{z}_i is a c dimensional vector, $z_{i,k} \in [0, 1]$ and $\sum_{k=1}^c z_{i,k} = 1$. Please write the objective function of cross entropy and show that it is equivalent to the negative log-likelihood on the training set, assuming the training samples are independent.

Answer: Denote the training label y_i by an one-hot target vector $\mathbf{t}_i \in \{0, 1\}^c$, where $\mathbf{t}_{i,y_i} = 1$ and $\mathbf{t}_i^T \mathbf{t}_i = 1$. The cross-entropy between the two distributions \mathbf{t}_i and \mathbf{z}_i is thus

$$L(\mathbf{t}_i, \mathbf{z}_i) = - \sum_{j=1}^c \mathbf{t}_{i,j} \log \mathbf{z}_{i,j} = - \log \mathbf{z}_{i,y_i} = - \log P(y = y_i | \mathbf{x}_i), \quad (1)$$

which is the negative log-likelihood of the i -th training sample. As the training samples are independent, it can be extended to the whole training set.

Problem 2

[25 points]

x_1 and x_2 are two input variables, and y is the target variable to be predicted. The network structure is shown in Figure 1(a). $h_{11} = f_{11}(x_1)$, $h_{12} = f_{12}(x_2)$, and $y = g(h_{11}, h_{12})$.

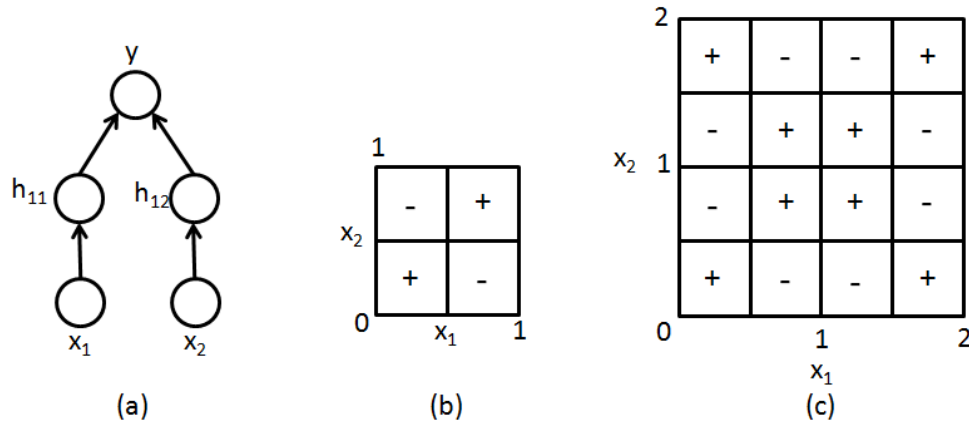


Figure 1: Problem 2

- Assuming $x_1 \in [0, 1]$ and $x_2 \in [0, 1]$, in order to obtain the decision regions in Figure 1(b), decide functions f_{11} , f_{12} , and g . [5 points]

Answer:

A possible solution could be

$$f_{11}(x_1) = \begin{cases} 1 & \text{if } x_1 \geq 0.5 \\ -1 & \text{if } x_1 < 0.5 \end{cases}, \quad f_{12}(x_2) = \begin{cases} 1 & \text{if } x_2 \geq 0.5 \\ -1 & \text{if } x_2 < 0.5 \end{cases}, \quad \text{and } g(h_{11}, h_{12}) = h_{11}h_{12}. \quad (2)$$

- Now we extend the range of x_1 and x_2 to $[0, 2]$. Please add one more layer to Figure 1(a) in order to obtain the decision regions in Figure 1(c). [5 points]

Answer:

Note that the decision regions are symmetric with respect to $x_1 = 1$ and $x_2 = 1$. Then a possible solution could be adding a layer z between h and x that makes

$$z_i = 1 - |x_i - 1| = \begin{cases} 2 - x_i & \text{if } x_i \geq 1 \\ x_i & \text{if } x_i < 1 \end{cases}, \quad \text{for } i = 1, 2. \quad (3)$$

- Although the decision boundaries in Figure 1(c) look complicated, there exist regularity and global structure. Please explain such regularity and global structure. Based on your observation, draw the decision boundaries when the range of x_1 and x_2 are extended to $[0, 4]$. [5 points]

Answer:

As pointed out in the previous answer, we can extend the symmetry property by “unfolding” Figure 1(c) again to $[0, 4]^2$, which results in the decision regions shown in Figure 2.

- Following the question above and assuming the range of x_1 and x_2 is extended to $[0, 2^n]$, draw the network structure and the transform function in each layer, in order to obtain the decision regions with the same regularity and global structure in Figure 1 (b) and (c). The complexity of computation units should be $O(n)$. [5 points]

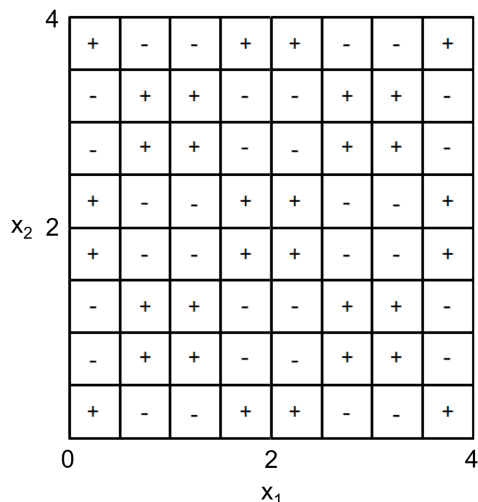


Figure 2: Answer to problem 2

Answer: We could insert n consecutive layers $z^{(1)}, \dots, z^{(n)}$ between x and h , such that

$$z_i^{(j)} = \begin{cases} 2^j - z_i^{(j-1)} & \text{if } z_i^{(j-1)} \geq 2^{j-1} \\ z_i^{(j-1)} & \text{if } z_i^{(j-1)} < 2^{j-1} \end{cases}, \quad \text{for } j = 1, \dots, n, \text{ and } i = 1, 2, \quad (4)$$

and let $z^{(0)} = x$ for convenience.

- Assuming the range of x_1 and x_2 is $[0, 2^n]$ and only one hidden layer is allowed, specify the network structure and transform functions. [5 points]

Answer: The pattern of decision regions also has a period of 2 along each axis, meaning that the decision regions for x are the same as $(x \bmod 2)$. Thus formally, one possible hidden layer function could be

$$h_{1i}(x_i) = f_{1i}(1 - |x_i - 2\lfloor x_i/2 \rfloor - 1|), \quad \text{for } i = 1, 2, \quad (5)$$

where f_{1i} is the same as defined in Eq. (2).

Problem 3

[20 points]

Figure 3 shows two convolutional neural networks. What is the receptive field of a neuron as the output the pooling layer in (a)? What is the receptive field of a neuron as the output the second convolutional layer in (b)? Justify your answers. We assume the the stride is 1 in the convolutional layer and the stride is equal to the size of the pooling region in the pooling layer.

Answer: See Figure 4 for illustrations in one dimension. It can be seen that the receptive field is (a) 8×8 (b) 13×13 .

Problem 4

[30 points]

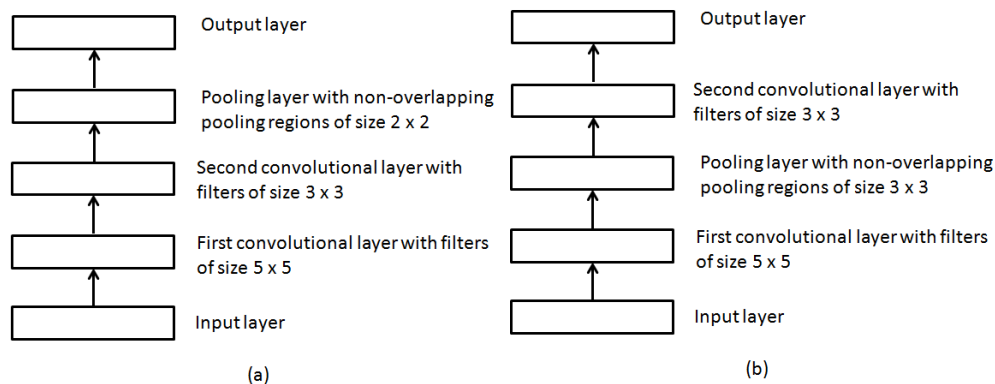


Figure 3:

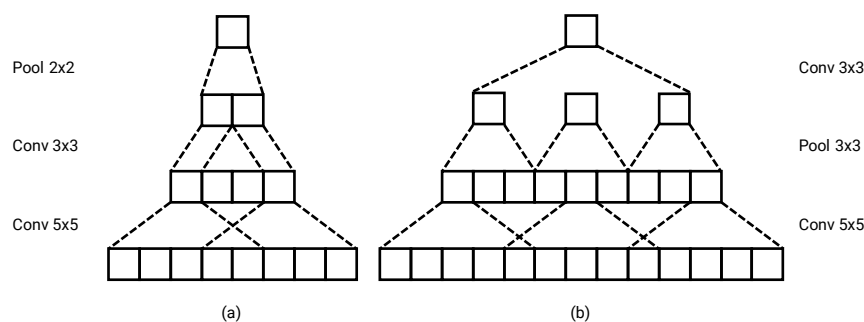


Figure 4:

Equivariance is an appealing property when design neural network operations. It means that transforming the input image (*e.g.*, translation) will also transform the output feature maps similarly after certain operations.

Formally, denote the image coordinate by $\mathbf{x} \in \mathbb{Z}^2$, and the pixel values at each coordinate by a function $f : \mathbb{Z}^2 \mapsto \mathbb{R}^K$, where K is the number of image channels. A convolution filter can also be formulated as a function $w : \mathbb{Z}^2 \mapsto \mathbb{R}^K$. Note that f and w are zero outside the image and filter kernel region, respectively. The convolution operation (correlation indeed for simplicity) is thus defined by

$$[f * w](\mathbf{x}) = \sum_{\mathbf{y} \in \mathbb{Z}^2} \sum_{k=1}^K f_k(\mathbf{y}) w_k(\mathbf{y} - \mathbf{x}). \quad (6)$$

1. **[15 pts]** Let $L_{\mathbf{t}}$ be the translation $\mathbf{x} \mapsto \mathbf{x} + \mathbf{t}$ on the image or feature map, *i.e.*, $[L_{\mathbf{t}}f](\mathbf{x}) = f(\mathbf{x} - \mathbf{t})$. Prove that convolution has equivariance to translation:

$$[[L_{\mathbf{t}}f] * w](\mathbf{x}) = [L_{\mathbf{t}}[f * w]](\mathbf{x}), \quad (7)$$

which means that first translating the input image then doing the convolution is equivalent to first convolving with the image and then translating the output feature map.

Proof.

$$\begin{aligned}
 [[L_{\mathbf{t}}f] * w](\mathbf{x}) &= \sum_{\mathbf{y} \in \mathbb{Z}^2} \sum_{k=1}^K [L_{\mathbf{t}}f]_k(\mathbf{y}) w_k(\mathbf{y} - \mathbf{x}) \\
 &= \sum_{\mathbf{y} \in \mathbb{Z}^2} \sum_{k=1}^K f_k(\mathbf{y} - \mathbf{t}) w_k(\mathbf{y} - \mathbf{x}) \\
 &= \sum_{\mathbf{y} \in \mathbb{Z}^2} \sum_{k=1}^K f_k(\mathbf{y}) w_k(\mathbf{y} + \mathbf{t} - \mathbf{x}) \\
 &= \sum_{\mathbf{y} \in \mathbb{Z}^2} \sum_{k=1}^K f_k(\mathbf{y}) w_k(\mathbf{y} - (\mathbf{x} - \mathbf{t})) \\
 &= [f * w](\mathbf{x} - \mathbf{t}) = [L_{\mathbf{t}}[f * w]](\mathbf{x}).
 \end{aligned}$$

□

2. [15 pts] Let $L_{\mathbf{R}}$ be the 90°-rotation on the image or feature map, where

$$\mathbf{R} = \begin{bmatrix} \cos(\pi/2) & -\sin(\pi/2) \\ \sin(\pi/2) & \cos(\pi/2) \end{bmatrix}, \quad (8)$$

then $[L_{\mathbf{R}}f](\mathbf{x}) = f(\mathbf{R}^{-1}\mathbf{x})$. However, convolution is not equivariant to rotations, *i.e.*, $[L_{\mathbf{R}}f] * w \neq L_{\mathbf{R}}[f * w]$, which is illustrated by Figure 5 ((a) is not equivalent to (b) rotated by 90°). In order to establish the equivalence, the filter also needs to be rotated (*i.e.* (b) is equivalent to (c) in Figure 5). Prove that:

$$[[L_{\mathbf{R}}f] * w](\mathbf{x}) = L_{\mathbf{R}}[f * [L_{\mathbf{R}^{-1}}w]](\mathbf{x}). \quad (9)$$

Proof.

$$\begin{aligned}
 [[L_{\mathbf{R}}f] * w](\mathbf{x}) &= \sum_{\mathbf{y} \in \mathbb{Z}^2} \sum_{k=1}^K [L_{\mathbf{R}}f]_k(\mathbf{y}) w_k(\mathbf{y} - \mathbf{x}) \\
 &= \sum_{\mathbf{y} \in \mathbb{Z}^2} \sum_{k=1}^K f_k(\mathbf{R}^{-1}\mathbf{y}) w_k(\mathbf{y} - \mathbf{x}) \\
 &= \sum_{\mathbf{y} \in \mathbb{Z}^2} \sum_{k=1}^K f_k(\mathbf{y}) w_k(\mathbf{R}\mathbf{y} - \mathbf{x}) \\
 &= \sum_{\mathbf{y} \in \mathbb{Z}^2} \sum_{k=1}^K f_k(\mathbf{y}) w_k(\mathbf{R}(\mathbf{y} - \mathbf{R}^{-1}\mathbf{x})) \\
 &= \sum_{\mathbf{y} \in \mathbb{Z}^2} \sum_{k=1}^K f_k(\mathbf{y}) ([L_{\mathbf{R}^{-1}}w]_k(\mathbf{y} - \mathbf{R}^{-1}\mathbf{x})) \\
 &= [f * [L_{\mathbf{R}^{-1}}w]](\mathbf{R}^{-1}\mathbf{x}) = L_{\mathbf{R}}[f * [L_{\mathbf{R}^{-1}}w]](\mathbf{x}).
 \end{aligned}$$

□

3. [optional] To make convolution equivariant to rotations, we need to extend the definition of convolution and transformation. Recall a group (G, \otimes) in algebra is a set G , together with an binary operation \otimes , which satisfies four requirements:

Closure $a \otimes b \in G, \forall a, b \in G$.

Associativity $(a \otimes b) \otimes c = a \otimes (b \otimes c), \forall a, b, c \in G$.

Identity element There exists a unique $e \in G, e \otimes a = a \otimes e = a, \forall a \in G$.

Inverse element $\forall a \in G, \exists a^{-1} \in G, a \otimes a^{-1} = a^{-1} \otimes a = e$.

We can formulate 90°-rotation and translation by a group (G, \otimes) consisting of

$$\mathbf{g}(r, u, v) = \begin{bmatrix} \cos(r\pi/2) & -\sin(r\pi/2) & u \\ \sin(r\pi/2) & \cos(r\pi/2) & v \\ 0 & 0 & 1 \end{bmatrix}, \quad (10)$$

where $r \in \{0, 1, 2, 3\}$ and $(u, v) \in \mathbb{Z}^2$. $G = \{\mathbf{g}\}$ and \otimes is matrix multiplication. Translation is a special case of G when $r = 0$ (i.e. $\mathbf{g}(0, u, v)$) and rotation is a special case of G when $u = v = 0$ (i.e. $\mathbf{g}(r, 0, 0)$).

A key concept is to extend the definition of both the feature f and the filter w to G . Imagine the feature map is duplicated four times with rotation of 0°, 90°, 180°, and 270°. Then $f(\mathbf{g})$ is the feature values at particular rotated pixel coordinate, and the convolution operation becomes

$$[f * w](\mathbf{g}) = \sum_{\mathbf{h} \in G} \sum_{k=1}^K f_k(\mathbf{h}) w_k(\mathbf{g}^{-1} \mathbf{h}). \quad (11)$$

A rotation-translation $u \in G$ on the feature map is thus $[L_{\mathbf{u}}f](\mathbf{g}) = f(\mathbf{u}^{-1}\mathbf{g})$. Prove that under such extensions, the convolution is equivariant to rotation-translation:

$$[[L_{\mathbf{u}}f] * w](\mathbf{g}) = [L_{\mathbf{u}}[f * w]](\mathbf{g}). \quad (12)$$

Briefly explain how to implement this group convolution with traditional convolution and by rotating the feature map or filter.

Answer: The proof is similar to question 1 but here we extend the original \mathbb{Z}^2 group with 90° rotations.

Proof.

$$\begin{aligned}
 [[L_{\mathbf{u}}f] * w](\mathbf{g}) &= \sum_{\mathbf{h} \in G} \sum_{k=1}^K [L_{\mathbf{u}}f]_k(\mathbf{h}) w_k(\mathbf{g}^{-1}\mathbf{h}) \\
 &= \sum_{\mathbf{h} \in G} \sum_{k=1}^K f_k(\mathbf{u}^{-1}\mathbf{h}) w_k(\mathbf{g}^{-1}\mathbf{h}) \\
 &= \sum_{\mathbf{h} \in G} \sum_{k=1}^K f_k(\mathbf{h}) w_k(\mathbf{g}^{-1}\mathbf{u}\mathbf{h}) \\
 &= \sum_{\mathbf{h} \in G} \sum_{k=1}^K f_k(\mathbf{h}) w_k((\mathbf{u}^{-1}\mathbf{g})^{-1}\mathbf{h}) \\
 &= [f * w](\mathbf{u}^{-1}\mathbf{g}) = [L_{\mathbf{u}}[f * w]](\mathbf{g}).
 \end{aligned}$$

□

To implement this group convolution, we just need to

- (a) duplicate and rotate the input image by $\{0^\circ, 90^\circ, 180^\circ, 270^\circ\}$, which leads to four times of input channels;
- (b) duplicate and rotate the original filters by $\{0^\circ, 90^\circ, 180^\circ, 270^\circ\}$ (the parameters are shared among duplicates). Note that each filter now receives four times of input channels, thus having four times of parameters.

The output feature maps directly serve as the input channels of the next group convolution layer and no need to be rotated again. If you are interested in group equivariance in CNNs, please refer to [1] Taco S. Cohen, *et. al.*, Group Equivariant Convolutional Networks, *ICML*, 2016.

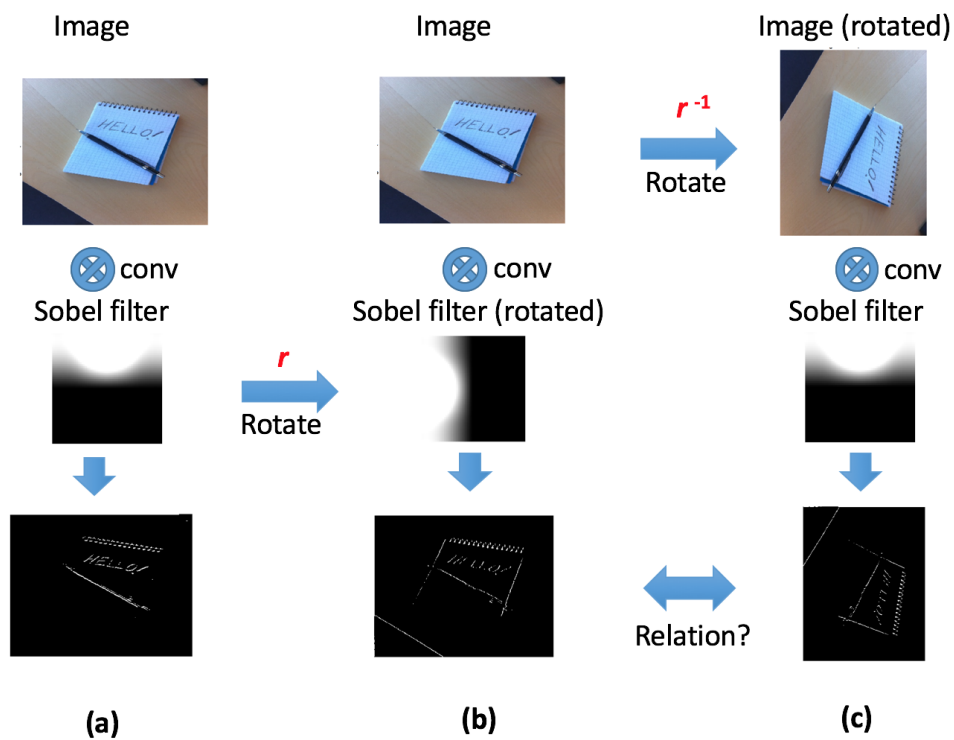


Figure 5: Equivariance relationship between convolution and rotation. (a) An image is convolved with a Sobel filter to detect horizontal edges. (b) The filter is rotated counterclockwise and then convolves the original image. (c) The image is first rotated clockwise, then it is convolved with the filter.