

통계학 2017/504/2
이재승

STAT 433 - Midterm Part I

1. During backpropagation, when the gradient passes backward through a sigmoid activation function, the gradient will always decrease in magnitude.

A. True
B. False

A: True.

2. Suppose that you find that your model's training error looks so good (potential overfitting). What can you do to address this issue? (Check all that apply)

A. Data augmentation
B. Dropout
C. Batch Normalization
D. RMSprop Optimizer

A, B, C

3. Which of the following is true?

A. Batch Normalization is an alternative method of dropout.
B. Batch Normalization makes training faster.
C. Batch Normalization is a non-linear transformation to give nonlinearity to the network.
D. Batch Normalization is standardizing the data before training neural network.

D

4. You want to make the weights sparse and smaller. How can you do that? Why?

Pruning을 하면 된다. Pruning은 Model의 Weights 중 중요도가 낮은 Weights의 연결을 제거하여 모델의 parameter를 줄이는 방법이다. 많은 Weights 값이 0이 되거나 작은 영향이 있다면, 상대적으로 영향력이 적은 Weights를 삭제하여도 결과에 주는 영향이 적을 것이다. 따라서 최종 결과에 영향력이 적은 Weights를 제거하면 작은 parameter로 가전적인 유사한 성능을 보여주는 모델을 만들 수 있는 것이다. 즉, pruning은 네트워크의 상층 크게 제한되지 않는 선에서 weights를 최대한 sparse하게 만드는 방법으로 정의할 수 있다. 이는 Neural Network를 실행하는데 필요한 계산 리소스를 줄이기 위해 수행한다.

Pruning에는 Magnitude pruning, sensitivity pruning 등 여러 방법이 있다. Pruning은 parameter가 줄어들기 때문에 복원 속도가 빨라지고, Regularization이 일반화 성능을 높이는 장점이 있다. Parameter를 줄여 모델의 속도를 빠르게 하기 위해 이용한다.

중요도가 낮은 Weights를 제거하여

5. $\tanh(x) = \frac{e^x - e^{-x}}{e^x + e^{-x}}$ has a similar performance as sigmoid function except that it is zero-centered. Write down $\tanh(x)$ in terms of $\sigma(x)$ where $\sigma(x) = 1/(1 + e^{-x})$. Show your work to get the full credit.

$$\sigma(x) = \frac{1}{1 + e^{-x}} \quad 1 - \sigma(x) = 1 - \frac{1}{1 + e^{-x}} = \frac{1}{1 + e^x} = \sigma(-x) \quad \therefore 1 - \sigma(x) = \sigma(-x)$$

$$\tanh(x) = \frac{e^x - e^{-x}}{e^x + e^{-x}} = \frac{e^x + e^{-x} - 2e^{-x}}{e^x + e^{-x}} = 1 + \frac{-2e^{-x}}{e^x + e^{-x}} = 1 - \frac{2}{e^{2x} + 1}$$

$$= 1 - 2\sigma(-2x)$$

$$= 1 - 2(1 - \sigma(2x))$$

$$= 1 - 2 + 2\sigma(2x)$$

$$= 2\sigma(2x) - 1$$

$$\therefore \tanh(x) = 2\sigma(2x) - 1$$

6. You have a single layer neural network for a binary classification with a sigmoid activation function as below. (X : $n \times m$ matrix, predicted \hat{y} & true label y : $1 \times m$)

$$z = WX + b$$

$$h = \sigma(z)$$

$$\hat{y} = h$$

$$L = - \sum_i y_i \log \hat{y}_i + (1 - y_i) \log(1 - \hat{y}_i)$$

What is $\frac{\partial L}{\partial W}$? Write your answer as a matrix-matrix multiplication.

$$\frac{\partial L}{\partial W} = \frac{\partial L}{\partial \hat{y}} \cdot \frac{\partial \hat{y}}{\partial h} \cdot \frac{\partial h}{\partial z} \cdot \frac{\partial z}{\partial W}$$

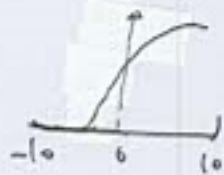
$$= \begin{bmatrix} \frac{\partial L}{\partial \hat{y}_1} & \frac{\partial L}{\partial \hat{y}_2} & \dots & \frac{\partial L}{\partial \hat{y}_m} \\ \frac{\partial L}{\partial \hat{y}_1} & \frac{\partial L}{\partial \hat{y}_2} & \dots & \frac{\partial L}{\partial \hat{y}_m} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial L}{\partial \hat{y}_1} & \frac{\partial L}{\partial \hat{y}_2} & \dots & \frac{\partial L}{\partial \hat{y}_m} \end{bmatrix} \cdot \begin{bmatrix} \frac{\partial \hat{y}_1}{\partial h_1} & \frac{\partial \hat{y}_2}{\partial h_2} & \dots & \frac{\partial \hat{y}_m}{\partial h_m} \\ \frac{\partial \hat{y}_1}{\partial h_1} & \frac{\partial \hat{y}_2}{\partial h_2} & \dots & \frac{\partial \hat{y}_m}{\partial h_m} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial \hat{y}_1}{\partial h_1} & \frac{\partial \hat{y}_2}{\partial h_2} & \dots & \frac{\partial \hat{y}_m}{\partial h_m} \end{bmatrix} \cdot \begin{bmatrix} \frac{\partial h_1}{\partial z_1} & \frac{\partial h_2}{\partial z_2} & \dots & \frac{\partial h_m}{\partial z_m} \\ \frac{\partial h_1}{\partial z_1} & \frac{\partial h_2}{\partial z_2} & \dots & \frac{\partial h_m}{\partial z_m} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial h_1}{\partial z_1} & \frac{\partial h_2}{\partial z_2} & \dots & \frac{\partial h_m}{\partial z_m} \end{bmatrix} \cdot \begin{bmatrix} \frac{\partial z_1}{\partial W_{11}} & \frac{\partial z_1}{\partial W_{12}} & \dots & \frac{\partial z_1}{\partial W_{1m}} \\ \frac{\partial z_2}{\partial W_{21}} & \frac{\partial z_2}{\partial W_{22}} & \dots & \frac{\partial z_2}{\partial W_{2m}} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial z_n}{\partial W_{n1}} & \frac{\partial z_n}{\partial W_{n2}} & \dots & \frac{\partial z_n}{\partial W_{nm}} \end{bmatrix}$$

이것은 $\sigma(z) \cdot (1 - \sigma(z))$ 의 미분이다.

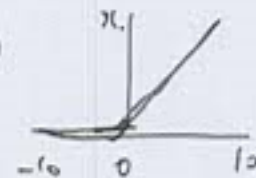
이것은 $\sigma(WX + b) \cdot (1 - \sigma(WX + b))$ 의 미분이다.

7. (continued from the above question) suppose that you apply ReLU activation before sigmoid activation. i.e., $\hat{y} = \sigma(\text{ReLU}(z))$. Then you classify the object by checking if $\hat{y} \geq 0.5$ or $\hat{y} < 0.5$. What will happen? Why?

$$\text{Sigmoid } \sigma(x) = \frac{1}{1 + e^{-x}}$$



$$\text{ReLU} : \max(0, x)$$



위와 같은 함수 그래프 형태를 생각했을 때, ReLU를 하면 Sigmoid를 적용하면 모든 prediction이 positive하게 될 것이다. 즉, $\hat{y} = \sigma(\text{ReLU}(z)) \geq 0.5$ 가 된다는 것을 의미한다. ReLU는 0 혹은 양수 값을 리턴하기 때문이, 이를 Sigmoid 함수에 넣으면 다 positive하게 되기 때문이다. 그래서 모든 $\hat{y} \geq 0.5$ 로 출력될 것이다.

8. Suppose that your classmate finds an activation function that is similar to ReLU such that

$$f(x) = \begin{cases} 1, & x \geq 0 \\ 0, & x < 0 \end{cases}$$

Will you use this? Why?

이것을 사용하지 않는다. 비록 이 함수가 non-linear하지만, discontinuous non-linear step function이긴 하지만, 미분할 수 없는 구간이 존재한다. 모든 $x < 0$ 구간에서 gradient가 0이다. 이는 곧 Backpropagation의 문제를 일으킬 것이다. 기울기가 거의 없는 구간을 야기하게 된다. 기울기가 거의 없는 구간은 최적화의 관점에서 굉장히 큰 문제를 일으킨다. 그래서 사용하지 않는다.

9. Provide two reasons why we are using convolutional layers instead of fully connected layers for image classification.

1) 이미지의 공간정보 소실을 막기 위해서다. Fully Connected Layer만으로 구성한 네트워크는 1차원 배열로만 이루어져 있다. 네트워크의 input이 이미지일때 Fully Connected Layer가 이미지를 학습하려면 (1,이러 (2,이러 가진 벡터 연관관계에 대한 정보가 있어야 하는데, 1차원 배열로 되면 (1,이 (1,1) ... (1,N) (2,이러 같이 두 벡터간의 연관관계에 대한 정보가 소실되어 사람이 파악할 수 있는 특징 파악이 어렵다. 그래서 1차원 배열로 늘여쓰기 편이 input을 특징 스페이스 안의 특징값들의 관계를 가진 노드로 하기 위해 Convolution Layer를 이용한다.

2) Convolutional Layer가 Fully Connected Layer보다 더 좁은 범위의 Feature를 추출기 때문이다. Fully Connected Layer의 Neuron은 이전 층의 모든 Neuron과 연결되기 위해 이 모든 Neuron 중 하나를 연결할 수 있다. 하지만 Convolutional Layer의 Neuron은 Convolutional kernel의 width를 따라 이전 계층의 좁은 feature와 연결된다. 이는 결과적으로 Convolutional Layer의 Neuron은 한 Neuron의 activation이 이전 layer의 대부분 Neuron activation이 intensive하다는 점에서 더 좁은 범위의 feature를 추출할 수 있다.

10. Consider to build a CNN for an image classification problem in which the layers are defined by the left column below. Fill the table below. Assume that width & height of the kernels (for Conv, Pool) are the same. Stride 1 Pad 1 for convolving layers. Stride 2 Pad 0 for Pooling layers. FC: a fully-connected layer.

Layer	Output Size		Layer		Number of parameters
	C	H/W	filters	kernel	
Input	3	32	-	-	0
Conv	16	32	16	3	448
ReLU	16	32	-	-	0
Pool	16	16	-	2	0
BatchNorm	16	16	-	-	32
Conv	16	16	16	3	2320
ReLU	16	16	-	-	0
Pool	16	8	-	2	0
Flatten	16 * 8 * 8	-	-	-	0
FC	10	-	-	-	10250