

ASSIGNMENT 2: THEORY OF CNNs AND REGULARIZATION [IFT6135]

JOSEPH D. VIVIANO

1. CONVOLUTIONS

2. CONVOLUTIONAL NEURAL NETWORKS

(a) Since we have one zero padding in layer three, so the size of layer 3 is $128 \times 6 \times 6$. The image has 3 color, therefore the last layer is a fully connected layer of the size

$$3 \cdot 128 \cdot 6 \cdot 6 = 13824.$$

.

(b) For the last convolution the size of kernel is 4×4 and we have 128 kernels of 3 colors, so the number of parameters is as follows

$$\#W = 4 \cdot 4 \cdot 128 \cdot 3 = 6144.$$

.

3. KERNEL CONFIGURATIONS FOR CNNs

(a): i : We denote the dimension of input as $W_1 \times H_1$ and the dimension of output is $W_2 \times H_2$, the kernel size is $K \times K$, number of zero padding is P and stride size is S , then we easily can see that

$$(3.1) \quad W_2 = \frac{W_1 - K + 2P}{S} + 1,$$

by using this formula we get

$$32 = \frac{64 - 8 + 2P}{S} + 1$$

. So if we set $P = 3$ and $S = 2$, the this convolution operation works.

ii : If we have dilatation of size D , then

$$(3.2) \quad W_2 = \frac{W_1 - K + 2P + (W_1 - 1)D}{S} + 1,$$

So here we have

$$32 = \frac{64 - K + 2P + 63 \cdot 6}{2} + 1,$$

Date: March 2018.

. So if we set $K = 400$ and $P = 10$, then our convolution operation works.

(b): For the pooling layer if the the kernel size is $K = 4 \times 4$ with no overlapping which means the stride size is $S = 4$, Then the pooling operation works.

(c): Here we have $K = 8$, $W_1 = 32$ and $S = 4$. By replacing the values in the formula(3.1), we have

$$W_2 = \frac{32 - K}{4} + 1 = 7$$

. So the output is of the size 7×7 .

(d): *i* here we have $W_2 = 4$, $W_1 = 8$ and $P = 0$, so by using the formula in (3.1), we get

$$4 = \frac{8 - K + 0}{S} + 1,$$

. By setting $K = 2$ and $S = 2$, the operation will work.

ii : we have $W_2 = 4$, $W_1 = 8$, $P = 2$ and $D = 1$, so by applying (3.2), we obtain

$$4 = \frac{8 - K + 4 + 7}{S} + 1$$

by putting $K = 13$ and $S = 2$, the operation works.

iii : By substituting the values in (3.1), we have

$$4 = \frac{8 - K + 2}{S} + 1,$$

so we choose $K = 4$ and $S = 2$.

4. DROPOUT AS WEIGHT DECAY

5. DROPOUT AS A GEOMETRIC ENSEMBLE

Consider the case of a single linear layer model with a softmax output. Prove that weight scaling by 0.5 corresponds exactly to the inference of a conditional probability distribution proportional to a geometric mean over all dropout masks.

First, observe the single linear layer with softmax output with n input variables represented by the vector v with dropout mask d :

$$(5.1) \quad P(y = y|v; d) = \mathbf{softmax} \left(W^T(d \odot v) + b \right)_y$$

and the ensemble conditional probability distribution which represents the geometric mean over all dropout masks:

$$(5.2) \quad p_{ens}(y = y|v; d) \propto \left(\prod_{i=1}^N \hat{y}_v^{(i)} \right)^{\frac{1}{N}}.$$

Aren't they nice? Recall the alternative formulation of the softmax:

$$(5.3) \quad \mathbf{softmax}_i = \frac{e^{x_i}}{\sum_{k=1}^K e^{x_k}}$$

Which we now rewrite, subbing in our vector representation of the softmax and replacing e^x with $\exp(x)$:

$$(5.4) \quad \mathbf{softmax}_y = \frac{\exp(W_y^T(d \odot v) + b)}{\sum_{k=1}^K \exp(W_{y'}^T(d \odot v) + b)}$$

Now we show that the ensemble predictor is defined by re-normalizing the geometric mean over all the individual ensemble members' predictions:

$$(5.5) \quad P_{ens}(y = y|v) = \frac{\tilde{P}_{ens}(y = y|v)}{\sum y' \tilde{P}_{ens}(y = y'|v)}$$

Where each \tilde{P}_{ens} is the geometric mean over all dropout masks for a single y :

$$(5.6) \quad \tilde{P}_{ens}(y = y|v) = 2^n \sqrt{\prod_{d \in \{0,1\}^n} P(y = y|v; d)}.$$

Now we simply sub in our definition of *softmax* for P :

$$(5.7) \quad \tilde{P}_{ens}(y = y|v) = 2^n \sqrt{\prod_{d \in \{0,1\}^n} \frac{\exp(W_y^T(d \odot v) + b)}{\sum_{k=1}^K \exp(W_{y'}^T(d \odot v) + b)}}.$$

Since the denominator is a constant under this normalization scheme we ignore it and simplify:

$$(5.8) \quad \tilde{P}_{ens}(y = y|v) \propto 2^n \sqrt{\prod_{d \in \{0,1\}^n} \exp(W_y^T(d \odot v) + b)}$$

We convert the product to the sum by taking \exp of the entire equation:

$$(5.9) \quad \tilde{P}_{ens}(y = y|v) \propto \exp\left(\frac{1}{2^n} \sum_{d \in \{0,1\}^n} W_y^T(d \odot v) + b\right)$$

And finally the sum and exponent n cancel:

$$(5.10) \quad \tilde{P}_{ens}(y = y|v) \propto \exp\left(\frac{1}{2} W_y^T(d \odot v) + b\right)$$

Finally, we sub this back into our earlier formulation of the softmax to show that the weights W are scaled by $\frac{1}{2}$:

$$(5.11) \quad \mathbf{softmax}_y = \frac{\exp\left(\frac{1}{2} W_y^T(d \odot v) + b\right)}{\sum_{k=1}^K \exp\left(\frac{1}{2} W_{y'}^T(d \odot v) + b\right)}$$

Therefore, weight scaling by 0.5 is exactly equivalent to a conditional probability distribution proportional to a geometric mean over all dropout masks.

6. NORMALIZATION

UNIVERSITÉ DE MONTRÉAL
E-mail address: joseph@viviano.ca