

ASSIGNMENT 2: THEORY OF CNNs AND REGULARIZATION [IFT6135]

JOSEPH D. VIVIANO

1. CONVOLUTIONS

$$\begin{aligned}\mathbf{x} &= (\dots, x_{-1}, x_0, x_1, x_2, x_3, x_4, \dots) = (\dots, 0, 1, 2, 3, 4, 0, \dots) \\ \mathbf{w} &= (\dots, w_{-1}, w_0, w_1, w_2, w_3, w_4, \dots) = (\dots, 0, 1, 0, 2, 0, 0, \dots)\end{aligned}$$

$$(1.1) \quad [\mathbf{x} * \mathbf{w}](i) = \sum_{k=-\infty}^{\infty} x_k w_{i-k}$$

$$[\mathbf{x} * \mathbf{w}](0) = x_0 w_0 = 1 \cdot 1 = 1$$

$$[\mathbf{x} * \mathbf{w}](1) = x_0 w_1 + x_1 w_0 = 1 \cdot 0 + 2 \cdot 1 = 2$$

$$[\mathbf{x} * \mathbf{w}](2) = x_0 w_2 + x_1 w_1 + x_2 w_0 = 1 \cdot 2 + 2 \cdot 0 + 3 \cdot 1 = 5$$

$$[\mathbf{x} * \mathbf{w}](3) = x_0 w_3 + x_1 w_2 + x_2 w_1 + x_3 w_0 = 1 \cdot 0 + 2 \cdot 2 + 3 \cdot 0 + 4 \cdot 1 = 8$$

$$[\mathbf{x} * \mathbf{w}](4) = x_2 w_2 + x_3 w_1 = 3 \cdot 2 + 4 \cdot 0 = 6$$

$$[\mathbf{x} * \mathbf{w}](5) = x_3 w_2 = 4 \cdot 2 = 8$$

Let's $y_i = [\mathbf{x} * \mathbf{w}](i)$. Thus we have $\mathbf{y} = (1, 2, 5, 8, 6, 8)$ where the null elements are not shown.

2. CONVOLUTIONAL NEURAL NETWORKS

(a) NB: zero padding in layer three, so size of layer 3 is $128 \times 6 \times 6$. The image is RGB (i.e., 3 channels), therefore the last layer is a fully connected layer of the size

$$3 \times 128 \times 6 \times 6 = 13824$$

.

(b) The last convolution has a kernel size 4×4 and there are 128 filters with 3 channels, so

$$n_params = 4 \times 4 \times 128 \times 3 = 6144$$

.

3. KERNEL CONFIGURATIONS FOR CNNs

(a): i : input is $W_1 \times H_1$ and output is $W_2 \times H_2$. Kernel size is K , zero padding is P and stride is S . Therefore

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

plugging our numbers in, we get

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

. Either $P = 3$ and $S = 2$ would produce a proper convolution.

ii : Dilatation size is D ,

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

So plugging in,

$$32 = \frac{64 - K + 2P + 63.6}{2} + 1.$$

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

(b): If the kernel size of the pooling layer is $K = 4 \times 4$ with no overlap, and the stride size is $S = 4$, the pooling operation works.

(c): $K = 8$, $W_1 = 32$ and $S = 4$, we plug them in and presto

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

The output is 7×7 .

(d): i $W_2 = 4$, $W_1 = 8$ and $P = 0$, plugging in, we get

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

Therefore $K = 2$ and $S = 2$ are appropriate.

ii : $W_2 = 4$, $W_1 = 8$, $P = 2$ and $D = 1$. Plugging in, we get

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

So $K = 13$ and $S = 2$ are appropriate.

iii :

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

so $K = 4$ and $S = 2$ are appropriate.

4. DROPOUT AS WEIGHT DECAY

(a) Let $\tilde{X} = X \odot \delta$ where $\delta = (\delta_1, \dots, \delta_n)$,

$$(4.1) \quad \delta_i = \begin{cases} \mathbf{0}, & p \\ \mathbf{1}, & 1 - p, \end{cases}$$

and $\mathbf{0}$ and $\mathbf{1}$ are vecteur of dimension $1 \times d$ of elements 0 and 1 respectively.

(b) Let L_{MSE} be the cost function. The general formula for the cost is

$$(4.2) \quad L_{MSE}(w) = \frac{1}{n} \sum_{i=1}^n (\hat{y}^{(i)} - y^{(i)})^2$$

where $\hat{y}^{(i)}$ is the the prediction and if we add dropout we have

$$(4.3) \quad L_{MSE}(w) = \frac{1}{n} \sum_{i=1}^n (x^{(i)} \delta_i w - y^{(i)})^2$$

(c) In this part, we will assume that the probability of dropping a input unit is $1 - p$. The expected value of the prediction is

$$(4.4) \quad \mathbb{E}[\hat{y}^{(i)}] = \mathbb{E}[X \delta_i w] = Xwp.$$

This mean that the expected value of the mean square error when using dropout is obtain by replacing the parameters vector w by wp . If we use matrices to express the expected cost we have

$$(4.5) \quad \begin{aligned} L_{MSE}(w) &= (Xwp - y)^\top (Xwp - y) \\ &= p^2 w^\top X^\top Xw - 2pw^\top X^\top y + y^\top y \end{aligned}$$

Now we take the derivatives with respect to w

$$(4.6) \quad \frac{\partial}{\partial w} L_{MSE}(w) = 2p^2 X^\top Xw - 2pX^\top y$$

And making the dervatives equal to zero

$$(4.7) \quad \begin{aligned} 2p^2 X^\top Xw^* - 2pX^\top y &= 0 \\ X^\top Xw^* p &= X^\top y \\ w^* p &= (X^\top X)^{-1} X^\top y \end{aligned}$$

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 :

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

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 :

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

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

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

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 $\frac{1}{2}$:

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

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

(a) *Show batchnorm and weightnorm are the same when you only have one layer and input feature x .*

To normalize the minibatch of activations B , we do

$$(6.1) \quad B' = \frac{B - \mu}{\sigma}$$

Where μ is the mean of B , and σ is the standard deviation of B (with a small positive value added for numerical stability).

We can replace B with $w^\top x$, where w is our weight matrix to see:

$$(6.2) \quad B' = \frac{w^\top x}{\sqrt{\text{Var}[w^\top x]}} - \frac{\mathbb{E} w^\top x}{\sqrt{\text{Var}[w^\top x]}}$$

Now notice the following about the unit vector

$$(6.3) \quad \frac{u}{\|u\|} = \frac{w^\top}{\|w\|}$$

And

$$(6.4) \quad g = \frac{\|w\|_2}{\sqrt{\text{Var}[w^\top x]}}$$

Therefore if we assume x and w are independent and that x has 0 mean:

$$(6.5) \quad B' = \frac{\|w\|_2}{\sqrt{\text{Var}[w^\top x]}} \frac{w^\top}{\|w\|_2} x - \frac{\mathbb{E} w^\top x}{\sqrt{\text{Var}[w^\top x]}}$$

$$(6.6) \quad B' = g \frac{u}{\|u\|} x - \frac{\mathbb{E} w^\top x}{\sqrt{\text{Var}[w^\top x]}}$$

The expectation and standard deviation are constant under these conditions, so we get:

$$(6.7) \quad B' = g \frac{u}{\|u\|} x - c$$

But we can ignore c for this question.

(b) *Show the gradients of L with respect to u can be expressed as $sW^* \nabla_w L$.*

From above:

$$(6.8) \quad B' = g \frac{u}{\|u\|} x$$

With weightnorm, we explicitly reparameterize the model to perform gradient descent in the new parameters g and u directly. By decoupling the norm of the weight vector g and it's direction $\frac{u}{\|u\|}$, we can speed up convergence dramatically.

If we differentiate through the above with respect to some new parameters v , we get:

$$(6.9) \quad \nabla_g L = \frac{\nabla_u L \cdot u}{\|u\|}, \nabla_v L = \frac{g}{\|u\|} \nabla_u L - \frac{g \nabla_g L}{\|u\|^2} u$$

Where $\nabla_u L$ is the gradient with respect to the weights.

Let's sub in $\nabla_g L$ into $\nabla_v L$ to get:

$$(6.10) \quad \nabla_v L = \frac{g}{\|u\|} \nabla_u L - \frac{g \frac{\nabla_u L \cdot u}{\|u\|}}{\|u\|^2} u$$

Or,

$$(6.11) \quad \nabla_v L = \frac{g}{\|u\|} \nabla_u L - \frac{g \nabla_u L}{\|u\|^3} u^\top u$$

This leads us to the formulation:

$$(6.12) \quad \nabla_v L = \frac{g}{\|u\|} M_u \nabla_u L$$

where,

$$(6.13) \quad M_u = Id - \frac{u^\top u}{\|u\|^2}$$

where M_u is a projection matrix that projects onto the complement of the u vector, and Id is the identity matrix.

(c) *Explain a graph of different learning rates.*

Let λ be the learning rate. During learning we update u via at step k using $v_k \leftarrow v_k - \lambda \nabla_v L$.

As we said in the previous question, the matrix M_u project onto the complement of u . Therefore, $\nabla_v L$ is equal to a constant times M_u , i.e., $u \perp \lambda \nabla_v L$.

Since our update v is proportional to w , the update must be orthogonal to v and the norm increases by the Pythagorean theorem, which states that for any two orthogonal vectors v and v' the new weight vector must have the norm

$$(6.14) \quad \|v'\| = \sqrt{\|v\|^2 + c^2 \|v\|^2}$$

if

$$(6.15) \quad c = \|\lambda \nabla_v L\| / \|v\|$$

This tells us a few things. If the norm of the gradients is small, $\sqrt{1 + c^2}$ is close to 1 and the norm of v stops increasing. As the norm of the gradients grow, the norm of v will also grow. Also, the norm of the updated parameter is proportional to the absolute value of the learning rate. These observations explain the graph.