

# Deep Learning: Assignment Two

Aditi Nair (asn264) and Akash Shah (ass502)

April 4, 2017

## 1 Batch Normalization

1. Let  $x_1, \dots, x_n$  be scalar features. Then we define the mean  $\mu_n$  as:

$$\mu_n = \frac{1}{n} \sum_{i=1}^n x_i$$

and the variance  $\sigma_n^2$  as:

$$\sigma_n^2 = \frac{1}{n} \sum_{i=1}^n (x_i - \mu_n)^2$$

Now to normalize the feature  $x_i$  (for  $1 \leq i \leq n$ ) to have zero mean and unit standard deviation, we compute:

$$\hat{x}_i = \frac{x_i - \mu_n}{\sigma_n}$$

Then over all  $\hat{x}_i$ , the expected value is 0:

$$\begin{aligned} E[x_1, \dots, x_n] &= \frac{1}{n} \sum_{i=1}^n \hat{x}_i = \frac{1}{n} \sum_{i=1}^n \frac{x_i - \mu_n}{\sigma_n} = \frac{1}{n\sigma_n} \sum_{i=1}^n (x_i - \mu_n) \\ &= \frac{1}{n\sigma_n} \left[ \left( \sum_{i=1}^n x_i \right) - n \cdot \mu_n \right] = \frac{1}{n\sigma_n} \left[ \sum_{i=1}^n x_i - \sum_{i=1}^n x_i \right] = 0 \end{aligned}$$

Then over all  $\hat{x}_i$  the variance is 1 since:

$$\begin{aligned} Var[x_1, \dots, x_n] &= \frac{1}{n} \sum_{i=1}^n (\hat{x}_i - \hat{\mu})^2 = \frac{1}{n} \sum_{i=1}^n (\hat{x}_i - 0)^2 \\ &= \frac{1}{n} \sum_{i=1}^n \hat{x}_i^2 = \frac{1}{n} \sum_{i=1}^n \left( \frac{x_i - \mu_n}{\sigma_n} \right)^2 \\ &= \frac{1}{n} \sum_{i=1}^n \frac{(x_i - \mu_n)^2}{\sigma_n^2} = \frac{1}{n \cdot \sigma_n^2} \sum_{i=1}^n (x_i - \mu_n)^2 \\ &= \frac{n}{n \cdot \sum_{i=1}^n (x_i - \mu_n)^2} \sum_{i=1}^n (x_i - \mu_n)^2 = 1 \end{aligned}$$

2. For scalar features  $x_1, \dots, x_n$  the output of the BN module can be written as:

$$y_i = BN_{\gamma, \beta}(x_i) = \gamma \hat{x}_i + \beta$$

with

$$\hat{x}_i = \frac{x_i - \mu_n}{\sqrt{\sigma_n^2 + \epsilon}}$$

$\mu_n$  and  $\sigma_n$  are defined as above. For numerical stability, the BN algorithm adds  $\epsilon$  to  $\sigma_n^2$  in the denominator of  $\hat{x}_i$  before taking the square root. At the  $k^{th}$  layer of a neural net, we write these variables are  $y_i^k$ ,  $\hat{x}_i^k$ ,  $\beta^k$  and  $\gamma^k$ .

**Now we write  $\frac{\delta E}{\delta \gamma^k}$  in terms of  $\frac{\delta E}{\delta y_i^k}$ :**

$$\frac{\delta E}{\delta \gamma^k} = \frac{\delta E}{\delta y_i^k} \cdot \frac{\delta y_i^k}{\delta \gamma^k}$$

$\gamma^k$  is a constant and  $y_i^k$  is a vector, so we compute  $\frac{\delta y_i^k}{\delta \gamma^k}$  as:

$$\frac{\delta y_i^k}{\delta \gamma^k} = \left\langle \frac{\delta(y_i^k)_1}{\delta \gamma^k}, \dots, \frac{\delta(y_i^k)_n}{\delta \gamma^k} \right\rangle^T$$

For  $1 \leq j \leq n$ , we compute  $\frac{\delta(y_i^k)_j}{\delta \gamma^k}$  as:

$$\frac{\delta(y_i^k)_j}{\delta \gamma^k} = \frac{\delta(\gamma^k(\hat{x}_i)_j + \beta^k)}{\delta \gamma^k} = \hat{x}_{ij}$$

It follows that:

$$\frac{\delta y_i^k}{\delta \gamma^k} = \langle \hat{x}_{i1}, \dots, \hat{x}_{in} \rangle^T = \hat{x}_i$$

Finally we can write  $\frac{\delta E}{\delta \gamma^k}$  as:

$$\frac{\delta E}{\delta \gamma^k} = \frac{\delta E}{\delta y_i^k} \cdot \hat{x}_i$$

**Next we write  $\frac{\delta E}{\delta \beta^k}$  in terms of  $\frac{\delta E}{\delta y_i^k}$ :**

$$\frac{\delta E}{\delta \beta^k} = \frac{\delta E}{\delta y_i^k} \frac{\delta y_i^k}{\delta \beta^k}$$

$\beta^k$  is a constant and  $y_i^k$  is a vector, so we compute  $\frac{\delta y_i^k}{\delta \beta^k}$  as:

$$\frac{\delta y_i^k}{\delta \beta^k} = \left\langle \frac{\delta(y_i^k)_1}{\delta \beta^k}, \dots, \frac{\delta(y_i^k)_n}{\delta \beta^k} \right\rangle^T$$

For  $1 \leq j \leq n$ , we compute  $\frac{\delta(y_i^k)_j}{\delta \beta^k}$  as:

$$\frac{\delta(y_i^k)_j}{\delta \beta^k} = \frac{\delta(\gamma^k(\hat{x}_i)_j + \beta^k)}{\delta \beta^k} = 1$$

It follows that:

$$\frac{\delta y_i^k}{\delta \beta^k} = \langle 1, \dots, 1 \rangle^T$$

Finally we can write  $\frac{\delta E}{\delta \beta^k}$  as:

$$\frac{\delta E}{\delta \beta^k} = \frac{\delta E}{\delta y_i^k} \cdot \vec{1}$$

## 2 Convolution

1. Assuming a stride of 1 and no zero padding, we will get a  $3 \times 3$  matrix of 9 values.
2. SHOW WORK?

Assuming a stride of 1, no zero padding and an additive filter, the result of this convolution is:

$$\begin{bmatrix} 158 & 183 & 172 \\ 229 & 237 & 238 \\ 195 & 232 & 244 \end{bmatrix}$$

3. SKIPPED

## 3 Variants of Pooling

1. Three types of pooling are max-pooling, average-pooling and L-P pooling. Max-pooling is implemented in different dimensions by the `MaxPool1d`, `MaxPool2d` and `MaxPool3d` classes in PyTorch. Average-pooling is implemented in different dimensions by the `AvgPool1d`, `AvgPool2d`, and `AvgPool3d` classes in PyTorch. L-P pooling is implemented in by the `LPPool2d` module in PyTorch.
2. SKIPPED: Does this require the output of a full max-pooling operation on an image/matrix, or just one snapshot?
3. Max-pooling essentially applies a maximum function over (sometimes non-overlapping) regions of input matrices. For inputs with several layers of depth, max-pooling is generally applied separately for each layer's matrix. The pooling operation is used to reduce the size of the input matrices in each layer. Applying max-pooling will select only the maximum value in each region of the input matrix, discarding smaller values in each region. This can prevent overfitting by discarding information from the smaller values. In comparison, operations like average-pooling weigh all values in the region equally - and are therefore more conservative in discarding additional information. An average pooling operation in a region with many small values and one large value will have a relatively small output compared to the max pooling operation over the same region, which ignores the overall tendency of values in the region. Therefore the aggressive approach of the max pooling operation can be used as a regularizer.

## 4 t-SNE

1. The crowding problem refers to the tendency of dimensionality reduction techniques to crowd points of varying similarity close together. Generally, SNE techniques model pairwise distances in high-dimensional space by a probability distribution  $P$ , and pairwise distances in the low-dimensional space by a probability distribution  $Q$ . In order to ensure that the low-dimensional representation is faithful to important properties of the original representation, SNE optimizes the KL-divergence between  $P$  and  $Q$  so that the pairwise distances in the low-dimensional space are as similar as possible to those in the original high-dimensional space.

Maarten and Hinton (2008) provide the following examples to illustrate the crowding problem. Given a set of points which lie on a two-dimensional manifold in high-dimensional space, it is fairly straightforward to effectively describe their pairwise distances with a two-dimensional map. On the other hand, suppose these points lie on a higher-dimensional manifold - then it becomes much more difficult to model pairwise distances in two dimensions. Maarten and Hinton provide an example on a 10-dimensional manifold, where it is possible to have ten points which are mutually equidistant - this is clearly impossible to map into a two-dimensional space.

More generally, the distribution of possible pairwise distances in the high-dimensional space is very

different than the distribution of possible pairwise distances in the two-dimensional space. Consider a set of points in high-dimensional space which are uniformly distributed around a central point. As the distance of the points from the central point grows, the volume of space these points could occupy also grows. However in two dimensions there is less area to accommodate points which are moderately far from the center than there is to accommodate points which are near to the center. Therefore, if we attempt to model small distances accurately in the two-dimensional space, we will be forced to place moderately large distances much further from the center point. Now, when trying to optimize the two-dimensional mapping, these too-far points will be pushed inward by the SNE objective function, effectively crowding all of the points together in the two-dimensional mapping.

t-SNE alleviates this by representing distances between points as probabilities using specific distributions. Distances in the high-dimensional spaces are converted to probabilities using a Gaussian distribution, whereas distances in the low-dimensional spaces are converted to probabilities using a distribution with a much heavier tail. This way, moderately far points are assigned larger distances in the lower-dimensional mapping compared to when the lower-dimensional distances are computed using a small-tailed distribution. Then the distances between pairs of moderately far points in higher dimensional space are closer to the distances between the same pairs of points in low dimensional space, and the SNE objective function will not push the low-dimensional representations closer together as much, ameliorating the crowding problem.

2. Let

$$C = KL(P||Q) = \sum_i \sum_j p_{ij} \log \frac{p_{ij}}{q_{ij}} = \sum_i \sum_j p_{ij} \log p_{ij} - p_{ij} \log q_{ij}$$

We want to compute  $\frac{\delta C}{\delta y_i}$ . Note that  $p_{ij}$  is constant with respect to  $y_i$  but  $q_{ij}$  is not. Specifically:

$$q_{ij} = \frac{(1 + \|y_i - y_j\|^2)^{-1}}{\sum_{k \neq l} (1 + \|y_k - y_l\|^2)^{-1}}$$

Following the derivation of van der Maaten and Hinton, we define two intermediate values:

$$d_{ij} = \|y_i - y_j\|$$

and

$$Z = \sum_{k \neq l} (1 + d_{kl}^2)^{-1}$$

Now we can rewrite  $q_{ij}$  as:

$$q_{ij} = \frac{(1 + d_{ij}^2)^{-1}}{\sum_{k \neq l} (1 + d_{kl}^2)^{-1}} = \frac{(1 + d_{ij}^2)^{-1}}{Z}$$

In the new notation, all terms are constant with respect to  $y_i$  except  $d_{ij}$  and  $d_{ji}$  for all  $j$ . Now, we can compute the partial derivative  $\frac{\delta C}{\delta y_i}$  in terms of the partial derivatives  $\frac{\delta C}{\delta d_{ij}}$  and  $\frac{\delta C}{\delta d_{ji}}$  for all  $d_{ij}$  and  $d_{ji}$ . In particular, we notice that the terms  $d_{ji}$  and  $d_{ij}$  appears once each for all possible value of  $j$ . That is:

$$\frac{\delta C}{\delta y_i} = \sum_j \frac{\delta C}{\delta d_{ij}} \frac{\delta d_{ij}}{\delta y_i} + \sum_j \frac{\delta C}{\delta d_{ji}} \frac{\delta d_{ji}}{\delta y_i}$$

First, we compute  $\frac{\delta d_{ij}}{\delta y_i}$ . For  $x, y \in \mathbb{R}^n$ , notice that:

$$\|x - y\| = \sqrt{(x_1 - y_1)^2 + \dots + (x_n - y_n)^2}$$

Then it follows that:

$$\frac{\delta(\|x - y\|)}{\delta x} = \left\langle \frac{\delta(\|x - y\|)}{\delta x_1}, \dots, \frac{\delta(\|x - y\|)}{\delta x_n} \right\rangle^T$$

$$\begin{aligned}
&= \left\langle \frac{2(x_1 - y_1)}{2\|x - y\|}, \dots, \frac{2(x_n - y_n)}{2\|x - y\|} \right\rangle^T \\
&= \frac{1}{\|x - y\|} \left\langle x_1 - y_1, \dots, x_n - y_n \right\rangle^T = \frac{x - y}{\|x - y\|}
\end{aligned}$$

Since  $d_{ij} = \|y_i - y_j\|$ , it follows that:

$$\frac{\delta d_{ij}}{\delta y_i} = \frac{\delta(\|y_i - y_j\|)}{\delta y_i} = \frac{y_i - y_j}{\|y_i - y_j\|}$$

Moreover since  $d_{ji} = \|y_j - y_i\| = \|y_i - y_j\| = d_{ij}$ , we can say that:

$$\frac{\delta d_{ji}}{\delta y_i} = \frac{y_i - y_j}{\|y_i - y_j\|}$$

Finally:

$$\begin{aligned}
\frac{\delta C}{\delta y_i} &= \sum_j \frac{\delta C}{\delta d_{ij}} \frac{y_i - y_j}{\|y_i - y_j\|} + \sum_j \frac{\delta C}{\delta d_{ji}} \frac{y_i - y_j}{\|y_i - y_j\|} \\
&= \sum_j \left( \frac{\delta C}{\delta d_{ij}} + \frac{\delta C}{\delta d_{ji}} \right) \frac{y_i - y_j}{\|y_i - y_j\|} = 2 \sum_j \frac{\delta C}{\delta d_{ij}} \frac{y_i - y_j}{\|y_i - y_j\|} = 2 \sum_j \frac{\delta C}{\delta d_{ij}} \frac{y_i - y_j}{d_{ij}}
\end{aligned}$$

Next, we need to compute  $\frac{\delta C}{\delta d_{ij}}$ . Recall that:

$$C = \sum_i \sum_j p_{ij} \log p_{ij} - p_{ij} \log q_{ij}$$

and only  $q_{ij}$  is non-constant with respect to  $d_{ij}$ .

Therefore:

$$\begin{aligned}
\frac{\delta C}{\delta d_{ij}} &= \frac{\delta \left( \sum_i \sum_j p_{ij} \log p_{ij} - p_{ij} \log q_{ij} \right)}{\delta d_{ij}} \\
&= \frac{-\delta \left( \sum_i \sum_j p_{ij} \log q_{ij} \right)}{\delta d_{ij}}
\end{aligned}$$

When optimizing t-SNE, we define  $p_{ii} = q_{ii} = 0$  for all  $i$ . Then we can say:

$$\frac{\delta C}{\delta d_{ij}} = \frac{-\delta \left( \sum_{k \neq l} p_{kl} \log q_{kl} \right)}{\delta d_{ij}}$$

Following, the derivation of van der Maaten and Hinton, we observe that:

$$\log q_{kl} = \log \frac{q_{kl} Z}{Z} = \log q_{kl} Z - \log Z$$

It follows that:

$$\frac{\delta C}{\delta d_{ij}} = \frac{-\delta \left( \sum_{k \neq l} p_{kl} (\log q_{kl} Z - \log Z) \right)}{\delta d_{ij}}$$

Notice that  $q_{kl}$  and  $Z$  are non-constant with respect to  $d_{ij}$  and that:

$$q_{kl} Z = \frac{(1 + d_{kl}^2)^{-1}}{Z} Z = (1 + d_{kl}^2)^{-1}$$

Then we can compute  $\frac{\delta C}{\delta d_{ij}}$  as:

$$\frac{\delta C}{\delta d_{ij}} = - \sum_{k \neq l} p_{kl} \left( \frac{1}{q_{kl} Z} \frac{\delta(1 + d_{kl}^2)^{-1}}{\delta d_{ij}} - \frac{1}{Z} \frac{\delta Z}{\delta d_{ij}} \right)$$

Now we notice that  $\frac{\delta(1 + d_{kl}^2)^{-1}}{\delta d_{ij}}$  is 0 unless  $k = i$  and  $l = j$ . Then we can write:

$$\begin{aligned} \frac{\delta C}{\delta d_{ij}} &= \frac{-p_{ij}}{q_{ij} Z} \cdot -(1 + d_{ij}^2)^{-2} \cdot 2d_{ij} - \sum_{k \neq l} -p_{kl} \frac{1}{Z} \cdot -(1 + d_{ij}^2)^{-2} \cdot 2d_{ij} \\ &= \frac{2p_{ij}}{q_{ij} Z} \cdot (1 + d_{ij}^2)^{-2} \cdot d_{ij} - 2 \sum_{k \neq l} \frac{p_{kl}}{Z} \cdot (1 + d_{ij}^2)^{-2} d_{ij} \end{aligned}$$

Since  $q_{ij} Z = (1 + d_{ij}^2)^{-1}$  we can simplify the left summand further.

$$\begin{aligned} \frac{\delta C}{\delta d_{ij}} &= 2p_{ij} \cdot (1 + d_{ij}^2)^{-1} \cdot d_{ij} - 2 \sum_{k \neq l} \frac{p_{kl}}{Z} \cdot (1 + d_{ij}^2)^{-2} d_{ij} \\ &= 2p_{ij} \cdot (1 + d_{ij}^2)^{-1} \cdot d_{ij} - \frac{2}{Z} \sum_{k \neq l} p_{kl} \cdot (1 + d_{ij}^2)^{-2} d_{ij} \end{aligned}$$

Since  $\sum_{k \neq l} p_{kl} = 1$ :

$$\frac{\delta C}{\delta d_{ij}} = 2p_{ij}(1 + d_{ij}^2)^{-1}d_{ij} - \frac{2}{Z}(1 + d_{ij}^2)^{-2}d_{ij}$$

Since  $q_{ij} = \frac{(1 + d_{ij}^2)^{-1}}{Z}$ , we can simplify the right summand as:

$$\begin{aligned} \frac{\delta C}{\delta d_{ij}} &= 2p_{ij}(1 + d_{ij}^2)^{-1}d_{ij} - 2q_{ij}(1 + d_{ij}^2)^{-1}d_{ij} \\ &= 2(p_{ij} - q_{ij})(1 + d_{ij}^2)^{-1}d_{ij} \end{aligned}$$

Now we can substitute this into our expression for  $\frac{\delta C}{\delta y_i}$ :

$$\begin{aligned} \frac{\delta C}{\delta y_i} &= 2 \sum_j \frac{\delta C}{\delta d_{ij}} \frac{y_i - y_j}{d_{ij}} = 2 \sum_j \left( 2(p_{ij} - q_{ij})(1 + d_{ij}^2)^{-1}d_{ij} \right) \frac{y_i - y_j}{d_{ij}} \\ &= 4 \sum_j (p_{ij} - q_{ij})(1 + d_{ij}^2)^{-1}(y_i - y_j) \end{aligned}$$

## 5 Sentence Classification

### 5.1 ConvNet

### 5.2 RNN

### 5.3 Extra credit experiments of fastText

### 5.4 Extra credit question

## 6 Language Modeling