

# word2vec gradients

Tambet Matiisen

October 6, 2015

## 1 Softmax loss and gradients

Let's denote

$$x_i = \mathbf{w}_i^T \hat{\mathbf{r}}$$

$x_i$  is a scalar and can be considered as (unnormalized) "similarity" of vectors  $\mathbf{w}_i$  and  $\hat{\mathbf{r}}$ . Vectors  $\mathbf{w}_i$  and  $\hat{\mathbf{r}}$  are both  $D \times 1$ -dimensional. Full matrix  $\mathbf{w}$  is  $V \times D$  dimensional.

$$\hat{y}_i = Pr(word_i | \hat{\mathbf{r}}, \mathbf{w}) = \frac{e^{x_i}}{\sum_{j=1}^V e^{x_j}}$$

$\hat{y}_i$  turns similarity into probability - probability that predicted word  $\hat{\mathbf{r}}$  is the word  $\mathbf{w}_i$ .  $V$  is the total number of words in vocabulary.

$$J = CE(\mathbf{y}, \hat{\mathbf{y}}) = - \sum_{i=1}^V y_i \log(\hat{y}_i)$$

$J$  is the cross-entropy loss, where  $\mathbf{y}$  is the vector of target values and  $\hat{\mathbf{y}}$  is the vector of predictions. In  $\mathbf{y}$  only  $y_k$  is 1 and all others are 0s. The loss could have been written also as  $J = -\log(\hat{y}_k)$ , but gradients are much easier to write with above notation.

### 1.1 Gradients with respect to $\hat{\mathbf{r}}$

$$\frac{\partial J}{\partial \hat{\mathbf{r}}} = \sum_{i=1}^V \frac{\partial J}{\partial x_i} \frac{\partial x_i}{\partial \hat{\mathbf{r}}}$$

The sum is needed, because all  $x_i$ -s depend on  $\hat{\mathbf{r}}$  and  $\hat{y}_i$ -s (which depend on  $x_i$ -s) are summed together in  $J$ .

We know cross-entropy loss derivative with respect to softmax input from task 2b:

$$\frac{\partial J}{\partial x_i} = \hat{y}_i - y_i$$

Also

$$\frac{\partial x_i}{\partial \hat{\mathbf{r}}} = \mathbf{w}_i$$

Altogether

$$\frac{\partial J}{\partial \hat{\mathbf{r}}} = \sum_{i=1}^V (\hat{y}_i - y_i) \mathbf{w}_i$$

Instead of  $\mathbf{w}_i$  we can write it in terms of full matrix  $\mathbf{w}$ .

$$\frac{\partial J}{\partial \hat{\mathbf{r}}} = \mathbf{w}^T (\hat{\mathbf{y}} - \mathbf{y})$$

## 1.2 Gradients with respect to $\mathbf{w}_i$

$$\frac{\partial J}{\partial \mathbf{w}_i} = \frac{\partial J}{\partial x_i} \frac{\partial x_i}{\partial \mathbf{w}_i}$$

There is no sum this time, because only one  $x_i$  depends on  $w_i$ . Because

$$\frac{\partial x_i}{\partial \mathbf{w}_i} = \hat{\mathbf{r}}$$

the result is

$$\frac{\partial J}{\partial \mathbf{w}_i} = (\hat{y}_i - y_i) \hat{\mathbf{r}}$$

In terms of full matrix  $\mathbf{w}$  that is

$$\frac{\partial J}{\partial \mathbf{w}} = (\hat{\mathbf{y}} - \mathbf{y}) \hat{\mathbf{r}}^T$$

## 2 Negative sampling loss and gradient

Let's start with notation again.  $x_i$  will stay the same.

$$x_i = \mathbf{w}_i^T \hat{\mathbf{r}}$$

For probability let's use  $p_i$  this time, because now we are using *sigmoid* function:

$$p_i = \sigma(x_i)$$

For negative sampling let's use different notation:

$$z_k = -\mathbf{w}_k^T \hat{\mathbf{r}}$$

$$q_k = \sigma(z_k)$$

Negative sampling loss becomes then:

$$J = J(\hat{\mathbf{r}}, \mathbf{w}_i, \mathbf{w}_{1..K}) = -\log p_i - \sum_{k=1}^K \log q_k$$

## 2.1 Gradient with respect to $\hat{\mathbf{r}}$

$$\frac{\partial J}{\partial \hat{\mathbf{r}}} = \frac{\partial J}{\partial p_i} \frac{\partial p_i}{\partial x_i} \frac{\partial x_i}{\partial \hat{\mathbf{r}}} + \sum_{k=1}^K \frac{\partial J}{\partial q_k} \frac{\partial q_k}{\partial z_k} \frac{\partial z_k}{\partial \hat{\mathbf{r}}}$$

Here  $i$  is a parameter to loss function  $J$  (the target word  $\mathbf{w}_i$ ), so we can leave it free. Then derivatives are:

$$\begin{aligned} \frac{\partial J}{\partial p_i} &= -\frac{1}{p_i} & \frac{\partial J}{\partial q_k} &= -\frac{1}{q_k} \\ \frac{\partial p_i}{\partial x_i} &= p_i(1 - p_i) & \frac{\partial q_k}{\partial z_k} &= q_k(1 - q_k) \\ \frac{\partial x_i}{\partial \hat{\mathbf{r}}} &= \mathbf{w}_i & \frac{\partial z_k}{\partial \hat{\mathbf{r}}} &= -\mathbf{w}_k \end{aligned}$$

Altogether

$$\frac{\partial J}{\partial \hat{\mathbf{r}}} = -(1 - p_i)\mathbf{w}_i + \sum_{k=1}^K (1 - q_k)\mathbf{w}_k$$

## 2.2 Gradient with respect to $\mathbf{w}_j$

$$\frac{\partial J}{\partial \mathbf{w}_j} = \frac{\partial J}{\partial p_i} \frac{\partial p_i}{\partial x_i} \frac{\partial x_i}{\partial \mathbf{w}_j} + \sum_{k=1}^K \frac{\partial J}{\partial q_k} \frac{\partial q_k}{\partial z_k} \frac{\partial z_k}{\partial \mathbf{w}_j}$$

Notice that we have to use new index  $j$ , because the same word may (or may not) appear in positive part as  $w_i$  and in negative part as  $w_k$ .

Only new derivatives we need to take here are

$$\frac{\partial x_i}{\partial \mathbf{w}_j} = \delta_{ij} \hat{\mathbf{r}} \quad \frac{\partial z_k}{\partial \mathbf{w}_j} = -\delta_{kj} \hat{\mathbf{r}}$$

Where  $\delta_{xy}$  is Kronecker delta, which is 1 if  $x = y$  and 0 otherwise. Altogether

$$\frac{\partial J}{\partial \mathbf{w}_j} = -(1 - p_i)\delta_{ij} \hat{\mathbf{r}} + \sum_{k=1}^K (1 - q_k)\delta_{kj} \hat{\mathbf{r}}$$

We can simplify it further, if we write it in terms of  $\mathbf{w}_i$  and  $\mathbf{w}_k$ :

$$\frac{\partial J}{\partial \mathbf{w}_i} = -(1 - p_i)\hat{\mathbf{r}} \quad \frac{\partial J}{\partial \mathbf{w}_k} = (1 - q_k)\hat{\mathbf{r}}$$

So we have to use the left term to update target vector weights and right term to update all negative word weights. Pay attention, that the target word may also occur in negative sample - we want the model learn, that the same word shouldn't occur next to itself. In this case both updates must be applied to target word vector.