The skip-gram model parameters are the center word vector and context word vector for each word in the vocabulary. In training, we learn the model parameters by maximizing the likelihood function (i.e., maximum likelihood estimation). This is equivalent to minimizing the following loss function:

$$-\sum_{t=1}^{T} \sum_{-m \le j \le m, \ j \ne 0} \log P(w^{(t+j)} \mid w^{(t)}). \tag{15.1.6}$$

When using stochastic gradient descent to minimize the loss, in each iteration we can randomly sample a shorter subsequence to calculate the (stochastic) gradient for this subsequence to update the model parameters. To calculate this (stochastic) gradient, we need to obtain the gradients of the log conditional probability with respect to the center word vector and the context word vector. In general, according to (15.1.4) the log conditional probability involving any pair of the center word w_c and the context word w_o is

$$\log P(w_o \mid w_c) = \mathbf{u}_o^{\mathsf{T}} \mathbf{v}_c - \log \left(\sum_{i \in \mathcal{V}} \exp(\mathbf{u}_i^{\mathsf{T}} \mathbf{v}_c) \right). \tag{15.1.7}$$

Through differentiation, we can obtain its gradient with respect to the center word vector \mathbf{v}_c as

$$\frac{\partial \log P(w_o \mid w_c)}{\partial \mathbf{v}_c} = \mathbf{u}_o - \frac{\sum_{j \in \mathcal{V}} \exp(\mathbf{u}_j^{\mathsf{T}} \mathbf{v}_c) \mathbf{u}_j}{\sum_{i \in \mathcal{V}} \exp(\mathbf{u}_i^{\mathsf{T}} \mathbf{v}_c)}$$

$$= \mathbf{u}_o - \sum_{j \in \mathcal{V}} \left(\frac{\exp(\mathbf{u}_j^{\mathsf{T}} \mathbf{v}_c)}{\sum_{i \in \mathcal{V}} \exp(\mathbf{u}_i^{\mathsf{T}} \mathbf{v}_c)} \right) \mathbf{u}_j \qquad (15.1.8)$$

$$= \mathbf{u}_o - \sum_{j \in \mathcal{V}} P(w_j \mid w_c) \mathbf{u}_j.$$

Training

Training continuous bag of words models is almost the same as training skip-gram models. The maximum likelihood estimation of the continuous bag of words model is equivalent to minimizing the following loss function:

$$-\sum_{t=1}^{T} \log P(w^{(t)} \mid w^{(t-m)}, \dots, w^{(t-1)}, w^{(t+1)}, \dots, w^{(t+m)}).$$
 (15.1.13)

Notice that

$$\log P(w_c \mid \mathcal{W}_o) = \mathbf{u}_c^{\mathsf{T}} \bar{\mathbf{v}}_o - \log \left(\sum_{i \in \mathcal{V}} \exp \left(\mathbf{u}_i^{\mathsf{T}} \bar{\mathbf{v}}_o \right) \right). \tag{15.1.14}$$

Through differentiation, we can obtain its gradient with respect to any context word vector $\mathbf{v}_{o_i}(i=1,\ldots,2m)$ as

$$\frac{\partial \log P(w_c \mid \mathcal{W}_o)}{\partial \mathbf{v}_{o_i}} = \frac{1}{2m} \left(\mathbf{u}_c - \sum_{j \in \mathcal{V}} \frac{\exp(\mathbf{u}_j^\top \bar{\mathbf{v}}_o) \mathbf{u}_j}{\sum_{i \in \mathcal{V}} \exp(\mathbf{u}_i^\top \bar{\mathbf{v}}_o)} \right) = \frac{1}{2m} \left(\mathbf{u}_c - \sum_{j \in \mathcal{V}} P(w_j \mid \mathcal{W}_o) \mathbf{u}_j \right).$$