

The objective described in the negative sampling part is:-

$$\log \sigma(V_{w_0}^T V_{w_I}) + \sum E_{w_i \sim P_n(w)} [\log \sigma(-V_{w_i}^T V_{w_I})]$$

Here, ~~the~~  $w_0$  is a context word,  $w_I$  is the target word, and  $w_i$ 's are ~~the~~ random words chosen on the basis of the distribution  $P_n(w)$ , that gives:-

$$P_n(w_i) = \frac{f(w_i)}{\sum f(w_i)}, \quad f(w_i) \text{ is freq. of } w_i.$$

BASICALLY,

we take a real context word from the target word's neighborhood, and do gradient ascent for  $\log(\sigma(V_{w_0}^T V_{w_I})) = y$  (say).

$$\frac{\partial y}{\partial V_{w_0}} = \frac{1}{\sigma(V_{w_0}^T V_{w_I})} \cdot \sigma(V_{w_0}^T V_{w_I}) (1 - \sigma(V_{w_0}^T V_{w_I})) \cdot V_{w_I}$$

$$\frac{\partial y}{\partial V_{w_0}} = (1 - \sigma(V_{w_0}^T V_{w_I})) \cdot V_{w_I}$$

So, the parameter update for  $V_{w_0}$  =

$$V_{w_0} \leftarrow V_{w_0} + \alpha (1 - \sigma(V_{w_0}^T V_{w_I})) \cdot V_{w_I}$$

Same for the context-vector derivative and update:-

$$\frac{\partial y}{\partial V_{w_I}} = (1 - \sigma(V_{w_0}^T V_{w_I})) \cdot V_{w_0}$$

$$V_{w_I} \leftarrow V_{w_I} + \alpha (1 - \sigma(V_{w_0}^T V_{w_I})) V_{w_0}$$

$\alpha$  is learning rate.



Now, for the negative samples:-

The gradient steps will be in the opposite direction, because the expressions differ only by a  $-1$ .

For each 'fake', 'false' context word:-

~~Δw<sub>i</sub>~~

$$\Delta w_i = w_i - \alpha (1 - \sigma(v_{w_t}^T w_i)) w_i$$

$$w_t = w_t - \alpha (1 - \sigma(v_{w_t}^T w_i)) w_i$$