

Word2Vec
skip-gram

$$\text{minimize } J(\theta) = -\frac{1}{T} \sum_{t=1}^T \sum_{\substack{-m \leq j \leq m \\ j \neq 0}} \log P(W_{t+j} | W_t)$$

all parameters
i.e. u_0, v_c

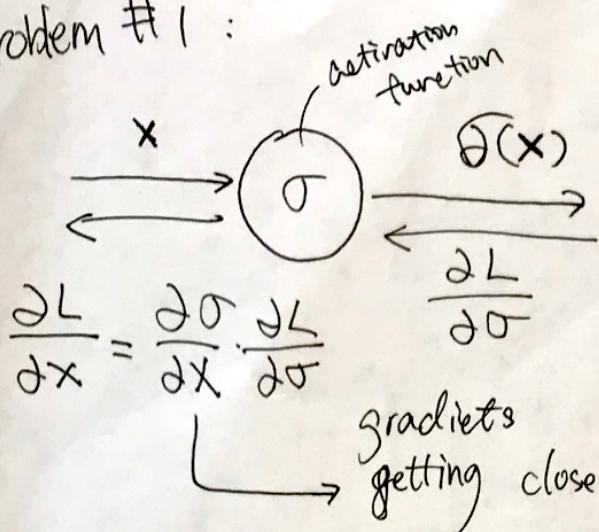
$$P(O|C) = \frac{\exp(u_0^T v_c)}{\sum_{W=1}^V \exp(u_W^T v_c)}$$

$$\begin{aligned} & \frac{\partial}{\partial v_c} \log \frac{\exp(u_0^T v_c)}{\sum_{W=1}^V \exp(u_W^T v_c)} \\ (\log \frac{x}{y}) &= \log x - \log y \\ &= \frac{\partial}{\partial v_c} \log \exp(u_0^T v_c) - \frac{\partial}{\partial v_c} \log \sum_{W=1}^V \exp(u_W^T v_c) \\ &\quad \text{chain rule } f(g(x)) \\ &\quad \frac{\partial \log x}{\partial x} = \frac{1}{x} \quad || \quad \frac{\partial f}{\partial x} = \frac{\partial f}{\partial g} \cdot \frac{\partial g}{\partial x} \\ &\quad \frac{1}{\sum_{W=1}^V \exp(u_W^T v_c)} \cdot \frac{\partial}{\partial v_c} \sum_{W=1}^V \exp(u_W^T v_c) \\ &\quad \frac{1}{\sum_{W=1}^V \exp(u_W^T v_c)} \cdot \sum_{X=1}^V \frac{\partial}{\partial v_c} \exp(u_X^T v_c) \\ &\quad \frac{1}{\sum_{W=1}^V \exp(u_W^T v_c)} \cdot \sum_{X=1}^V \exp(u_X^T v_c) \frac{\partial}{\partial v_c} u_X^T v_c \\ &\quad \text{chain rule again} \\ &\quad \frac{1}{\sum_{W=1}^V \exp(u_W^T v_c)} \cdot \sum_{X=1}^V \exp(u_X^T v_c) \cdot \frac{\partial}{\partial v_c} u_X^T v_c \\ &\quad \text{observed context word.} \end{aligned}$$

$$\begin{aligned} &= u_0 - \frac{\sum_{X=1}^V \exp(u_X^T v_c) \cdot u_X}{\sum_{W=1}^V \exp(u_W^T v_c)} P(X|C) \\ &= u_0 - \frac{\sum_{X=1}^V \frac{\exp(u_X^T v_c)}{\sum_{W=1}^V \exp(u_W^T v_c)}}{\sum_{X=1}^V \frac{\exp(u_X^T v_c)}{\sum_{W=1}^V \exp(u_W^T v_c)}} \cdot u_X \\ &= u_0 - \sum_{X=1}^V P(X|C) \cdot \bar{u}_X \end{aligned}$$

expectation of context word
by current model

Problem #1 :



e.g. tanh

$$\frac{\partial L}{\partial x} = \frac{\partial \sigma}{\partial x} \cdot \frac{\partial L}{\partial \sigma}$$

gradients
getting close to 0 \Rightarrow "saturated" neurons

Problem #2 : e.g.: ReLu

not zero-centered

\Rightarrow Input to a neuron is always positive ..

$$g(wx + b) = g(\sum_i w_i x_i + b)$$

\Rightarrow gradients on w will always be all positive or all negative

$$\begin{aligned}\frac{\partial L}{\partial w_i} &= \frac{\partial L}{\partial g} \cdot \cancel{\frac{\partial g}{\partial z}} \cdot \frac{\partial z}{\partial x} \\ &= \frac{\partial L}{\partial g} \cdot \frac{\partial g}{\partial z} \cdot x_i\end{aligned}$$

