

# Practical No. 2

Karol Pustelnik  
 kp446518@students.mimuw.edu.pl

March 20, 2022

## 1 The reasoning behind word2vec

a)

Show that the naive-softmax loss given in Equation (2) is equivalent to the cross-entropy loss between  $y$  and  $\hat{y}$ .

$$\begin{aligned}
 y_w &\rightarrow 0 \text{ for } w \neq o \\
 -\sum_{w \in \text{Voc}} y_w \cdot \log(\hat{y}_w) &= - (0 + 0 + \dots 1 \cdot \log(\hat{y}_o) + \dots 0) \\
 &= -\log(\hat{y}_o)
 \end{aligned}$$

b)

Compute the partial derivative of  $J_{\text{naive-softmax}}(v_c, o, U)$  with respect to  $v_c$ . Write your answer in terms of  $y$ ,  $\hat{y}$ , and  $U$ .

$$\begin{aligned}
 b) \quad \frac{-\sum_{w \in \text{Voc}} y_w \log(\hat{y}_w)}{\partial v_c} &= -\sum_{w \in \text{Voc}} y_w \cdot \frac{\partial \log(\hat{y}_w)}{\partial v_c} \\
 &= -\sum_{w \in \text{Voc}} y_w \cdot \frac{1}{\hat{y}_w} \cdot \underbrace{\hat{y}_w (1 - \hat{y}_w)}_{\frac{\partial \hat{y}_w}{\partial v_c}} \cdot \partial v_c^\top \\
 -\sum_{w \in \text{Voc}} (y_w - \hat{y}_w) \cdot \partial v_c^\top &= \sum_{w \in \text{Voc}} (\hat{y}_w - y_w) \cdot \partial v_c^\top \\
 &= U^\top (\hat{y} - y) \quad \text{in Python} \\
 &\quad \text{np.matmul(np.transpose(U), (\hat{y} - y))}
 \end{aligned}$$

c)

Compute the partial derivatives of  $J_{\text{naive-softmax}}(v_c, o, U)$  with respect to each of the outside word vectors,  $u_{ws}$ . Write your answer in terms of  $y$ ,  $\hat{y}$ , and  $v_c$ .

$$\begin{aligned}
 c) \quad \frac{\partial J}{\partial u_j} &= \frac{\partial -\sum_w y_w \log(\hat{y}_w)}{\partial u_j} \quad \text{pert when } w=j \\
 &= -\sum_{w \neq j} y_w \frac{1}{\hat{y}_w} \cdot \frac{\partial \hat{y}_w}{\partial u_j} - y_j \cdot \frac{1}{\hat{y}_j} \cdot \frac{\partial \hat{y}_j}{\partial u_j} \\
 1) \quad \frac{\partial \hat{y}_w}{\partial u_j} &= \frac{\exp(u_w^\top v_c)}{\sum_w \exp(u_w^\top v_c)}
 \end{aligned}$$

$$\begin{aligned}
 0 &= \frac{\exp(u_j^\top v_c) \cdot v_c \cdot \exp(u_j^\top v_c)}{\left(\sum_w \exp(u_w^\top v_c)\right)^2} \\
 2) \quad \frac{\partial \hat{y}_j}{\partial u_i} &= \frac{\exp(u_j^\top v_c)}{\sum_w \exp(u_w^\top v_c)} = \\
 &\frac{\exp(u_j^\top v_c) \cdot v_c \cdot \sum_{w \neq j} \exp(u_w^\top v_c) - \exp(u_j^\top v_c) \cdot v_c \cdot \exp(u_j^\top v_c)}{\left(\sum_w \exp(u_w^\top v_c)\right)^2}
 \end{aligned}$$

$$= \hat{y}_j \cdot v_c - (\hat{y}_i)^2 v_c = \hat{y}_j \cdot v_c (1 - \hat{y}_i)$$

Now I go back to derivative and I substitute

$$\begin{aligned} & - \sum_{w \neq j} y_w \frac{1}{\hat{y}_w} \cdot \frac{\partial \hat{y}_w}{\partial u_j} - y_j \cdot \frac{1}{\hat{y}_j} \cdot \frac{\partial \hat{y}_j}{\partial u_j} \\ &= - \sum_{w \neq j} y_w \frac{1}{\hat{y}_w} \cdot (\cancel{y_j \hat{y}_w v_c}) - y_j \cdot \frac{1}{\hat{y}_j} \cdot \cancel{y_j \cdot (1 - \hat{y}_j) \cdot v_c} \\ &= \sum_{w \neq j} y_w \hat{y}_j v_c + y_j \hat{y}_j v_c - y_j v_c \end{aligned}$$

$$\begin{aligned} &= \sum_w y_w \hat{y}_j v_c - y_j v_c \\ &= \hat{y}_j v_c - y_j v_c = (\hat{y}_j - y_j) v_c \end{aligned}$$

If we continue like that for every  $u_i$  we can create gradient:

$$(\hat{y} - y) \circ v_c$$

where  $\circ$  is outer product

in python:

$$\text{np.outer}((\hat{y} - y)), v_c)$$

d)

Please compute the derivative of  $\sigma(x)$  with respect to  $x$ .

$$\begin{aligned}
d) \quad g(x) &= \frac{e^x}{e^x + 1} \\
\frac{\partial g(x)}{\partial x} &= \frac{e^x \cdot (e^x + 1) - e^x \cdot e^x}{(e^x + 1)^2} = \frac{e^x}{(e^x + 1)^2} \\
&= g(x) \cdot \frac{1}{e^x + 1} = g(x) \cdot \left( \frac{e^x - e^x + 1}{e^x + 1} \right) \\
&= g(x) \cdot (1 - g(x))
\end{aligned}$$

e)

Please repeat parts (b) and (c), computing the partial derivatives of  $J_{negsample}$  with respect to  $v_c$ , with respect to  $u_o$ , and with respect to a negative sample  $u_k$ . Please write your answers in terms of the vectors  $u_o$ ,  $v_c$ , and  $u_k$ , where  $k \in [1, K]$ .

$$\begin{aligned}
c) \quad \frac{\partial J}{\partial v_c} &\approx -\frac{1}{\overline{p}(u_o^\top v_c)} \cdot \cancel{g(u_o^\top v_c)} \cdot (1 - \cancel{g(u_o^\top v_c)}) \cdot u_o^\top \\
&\quad - \sum_{s=1}^k (g(u_{ws}^\top v_c) - 1) \cdot u_{ws}^\top = \\
&= (\cancel{g(u_o^\top v_c)} - 1) u_o^\top + \sum_{s=1}^k (\cancel{g(u_{ws}^\top v_c)} - 1) u_{ws}^\top \\
&\quad \text{with respect to } u_o: \quad \frac{\partial J}{\partial u_o} = (\cancel{g(u_o^\top v_c)} - 1) v_c \\
&\quad \text{with respect to } u_{ws}: \quad \frac{\partial J}{\partial u_{ws}} = (\cancel{g(u_{ws}^\top v_c)} - 1) v_c
\end{aligned}$$

1. Why this loss function is much more efficient to compute than the naive-softmax loss?

- Because the computations are based on much smaller sample.

f)

f)

1)  $\frac{\sum_{\substack{1 \leq j \leq m \\ j \neq i}} \partial J(v_c, w_{t+j}, u)}{\partial u}$

2)  $\frac{\sum_{\substack{1 \leq j \leq m \\ j \neq i}} \partial J(v_c, w_{t+j}, u)}{\partial v_c}$

3) 0

## 2 Building your version of word2vec

In the picture, I can see that there is almost no clustering of words. Everything is scattered. We also have words "swietny" and "zly" close to each other which is weird. I think that the selected seed was very unfortunate and led to weird results.

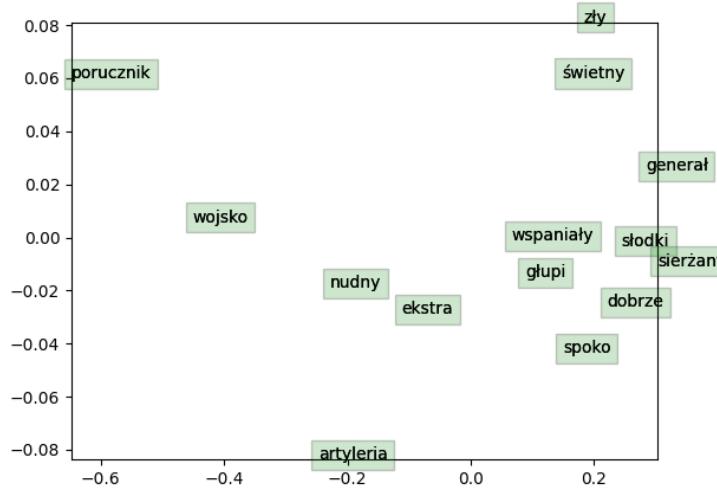


Figure 1: Visualization of 2D word embeddings

### 3 Optimizing Neural Networks

a)

1. Briefly explain (you can just give an intuition) how using  $m$  stops the updates from varying as much and why this low variance may be helpful to learning, overall.
  - Adding momentum to parameter updates is like increasing the speed of the ball that rolls down the hill. A faster ball is better because:
    - \* it oscillates less (because it has high momentum)
    - \* it reaches local (global) minima faster (because more momentum means more speed)
    - \* it is less likely to be stuck at bad local minima because it has enough speed to leave it

b)

1. Since Adam divides the update  $\sqrt{v}$ , which of the model parameters will get larger updates? Why might this help with by learning?
  - With  $\sqrt{v}$ , Adam normalizes the parameter updates, reducing high gradients parameters learning rate and increasing low gradients low gradients parameters learning rate. The normalization helps parameters learn at a similar pace and helps with faster convergence.

c) \*

1. Why L2 regularization can be helpful in low-data setting?
  - Because if our dataset is small, we have a high chance of overfitting. L2 regularization adds noise so its harder for the model to learn the data.
2. What is their proposal for improving the Adam update scheme?
  - Authors suggest "decoupling the weight decay from the optimization steps taken w.r.t the loss function". It means that AdamW adjusts the weight decay term to appear in the gradient update.

d)

1. Why should we apply dropout during training but not during evaluation?
  - Dropout makes neurons output wrong values on purpose. Because neurons are disabled randomly the network will have different outputs every (sequences of) activation.
2. How dropout can be seen from the perspective of Bayesian theory?\*

- Based on the paper, ”a neural network with arbitrary depth and non-linearities, with dropout applied before every weight layer, is mathematically equivalent to an approximation to the probabilistic deep Gaussian process. This theory gives tools to model uncertainty with dropout NNs – extracting information from existing models that has been thrown away so far. This mitigates the problem of representing uncertainty in deep learning without sacrificing either computational complexity or test accuracy.”