

PSet Assignment #1

1 Softmax

- a) Trivial
- b) Code

2 Neural Network Basics

- a) $\forall x \in \mathbb{R}, \sigma(x) = \frac{1}{1+e^{-x}}$
 $\forall x \in \mathbb{R} :$

$$\begin{aligned}\nabla \sigma(x) &= \frac{e^{-x}}{(1+e^{-x})^2} \\ &= \left(\frac{1}{\sigma(x)} - 1\right)\sigma(x)^2 \\ &= (1 - \sigma(x))\sigma(x) \\ &= \sigma(-x)\sigma(x)\end{aligned}$$

- b) $\mathbf{y} \in \mathbb{R}^n, \exists k \in \llbracket 1, n \rrbracket \quad \mathbf{y} = \mathbf{e}_k$. Therefore

$$\forall \theta \in \mathbb{R}^n \quad CE(\mathbf{y}, \hat{\mathbf{y}}) = -\log(\hat{y}_k) = -\theta_k + \log\left(\sum_{i=1}^n e^{\theta_i}\right)$$

Then $\forall \theta \in \mathbb{R}^n :$

$$\nabla_{\theta} CE(\mathbf{y}, \hat{\mathbf{y}}) = \hat{\mathbf{y}} - \mathbf{y}$$

- c)

$$\begin{aligned}\mathbf{h} &= \sigma(\mathbf{x}\mathbf{W}_1 + \mathbf{b}_1) \\ \hat{\mathbf{y}} &= \text{softmax}(\mathbf{h}\mathbf{W}_2 + \mathbf{b}_2)\end{aligned}$$

$\forall \mathbf{x} \in \mathbb{R}^n :$

$$\begin{aligned}
\nabla_{\mathbf{x}} CE(\mathbf{y}, \hat{\mathbf{y}}) &= \sum_{i=1}^n \frac{\partial CE(\mathbf{y}, \hat{\mathbf{y}})}{\partial x_i} \mathbf{e}_i \\
&= \nabla_{\theta} CE(\mathbf{y}, \hat{\mathbf{y}}) \cdot \left(\frac{\partial \theta}{\partial \mathbf{x}} \right)^{\top} \\
\forall i \in \llbracket 1, D_x \rrbracket, \forall j \in \llbracket 1, D_y \rrbracket \quad \left(\frac{\partial \theta_j}{\partial x_i} \right) &= \left(\frac{\partial (\sum_{l=1}^h h_l (W_2)_{lj} + (b_2)_j)}{\partial x_i} \right) \\
&= \left(\frac{\partial (\sum_{l=1}^h h_l (W_2)_{lj} + (b_2)_j)}{\partial x_i} \right) \\
&= \left(\frac{\partial (\sum_{l=1}^h \sigma(\sum_{m=1}^{D_x} x_m (W_1)_{ml} + (b_1)_l) (W_2)_{lj})}{\partial x_i} \right) \\
&= \left(\sum_{l=1}^h (W_1)_{il} \sigma' \left(\sum_{m=1}^{D_x} x_m (W_1)_{ml} + (b_1)_l \right) (W_2)_{lj} \right) \\
&= (\mathbf{W}_1)_{i \cdot} \cdot (\mathbf{h}' * (\mathbf{W}_2)_{\cdot j})
\end{aligned}$$

Then $\nabla_{\mathbf{x}} CE(\mathbf{y}, \hat{\mathbf{y}}) = (\hat{\mathbf{y}} - \mathbf{y}) \mathbf{W}_2^{\top} \mathbf{D} \mathbf{W}_1^{\top}$ where $\mathbf{D} = \text{diag}(\sigma'(\mathbf{x} \mathbf{W}_1 + \mathbf{b}_1))$

d) There are $H(1 + D_x) + D_y(1 + H)$ parameters.

e) Code

f) Code

g) Code

3 word2vec

a) $\hat{\mathbf{y}}_o = p(\mathbf{o} | \mathbf{c}) = \frac{\exp(\mathbf{u}_o^{\top} \mathbf{v}_c)}{\sum_{w=1}^W \exp(\mathbf{u}_w^{\top} \mathbf{v}_c)}.$

We can rewrite $\hat{\mathbf{y}}_o$ as follows :

$$\hat{\mathbf{y}}_o = \text{softmax}(\mathbf{U}^{\top} \mathbf{v}_c)_o \quad \text{where } \mathbf{U} = [\mathbf{u}_1, \dots, \mathbf{u}_W]$$

Therefore, if $\mathbf{y} = \mathbf{e}_o$, then $\forall \mathbf{v}_c \in \mathbb{R}^n$

$$\begin{aligned}
\nabla_{\mathbf{v}_c} CE(\mathbf{y}, \hat{\mathbf{y}}) &= \mathbf{U} \nabla_{\mathbf{U}^{\top} \mathbf{v}_c} CE(\mathbf{y}, \hat{\mathbf{y}}) \\
&= \mathbf{U} (\hat{\mathbf{y}} - \mathbf{y})
\end{aligned}$$

b) • If $w \neq o$:

$$\begin{aligned}
\nabla_{\mathbf{u}_w} CE(\mathbf{y}, \hat{\mathbf{y}}) &= \left(\frac{\partial (\mathbf{U}^{\top} \mathbf{v}_c)}{\partial \mathbf{u}_w} \right)^{\top} (\hat{\mathbf{y}} - \mathbf{y}) \\
&= \left(\sum_{i=1}^W \frac{\partial ((\mathbf{u}_i^{\top} \mathbf{v}_c) \mathbf{e}_i)}{\partial \mathbf{u}_w} \right)^{\top} (\hat{\mathbf{y}} - \mathbf{y}) \\
&= \hat{\mathbf{y}}_w \cdot \mathbf{v}_c
\end{aligned}$$

- If $w = o$:

$$\nabla_{\mathbf{u}_o} CE(\mathbf{y}, \hat{\mathbf{y}}) = (\hat{\mathbf{y}}_o - 1) \cdot \mathbf{v}_c$$

c) Now $J_{\text{neg-sample}}(\mathbf{o}, \mathbf{v}_c, \mathbf{U}) = -\log(\sigma(\mathbf{u}_o^\top \mathbf{v}_c)) - \sum_{k=1}^K \log(\sigma(-\mathbf{u}_k^\top \mathbf{v}_c))$

- $\forall \mathbf{v}_c \in \mathbb{R}^n$

$$\begin{aligned} \nabla_{\mathbf{v}_c} CE(\mathbf{y}, \hat{\mathbf{y}}) &= -\frac{\sigma(-\mathbf{u}_o^\top \mathbf{v}_c)\sigma(\mathbf{u}_o^\top \mathbf{v}_c)}{\sigma(\mathbf{u}_o^\top \mathbf{v}_c)} \mathbf{u}_o + \sum_{k=1}^K \frac{\sigma(-\mathbf{u}_k^\top \mathbf{v}_c)\sigma(\mathbf{u}_k^\top \mathbf{v}_c)}{\sigma(-\mathbf{u}_k^\top \mathbf{v}_c)} \mathbf{u}_k \\ &= -\sigma(-\mathbf{u}_o^\top \mathbf{v}_c) \mathbf{u}_o + \sum_{k=1}^K \sigma(\mathbf{u}_k^\top \mathbf{v}_c) \mathbf{u}_k \end{aligned}$$

- If $w \neq o$:

$$\nabla_{\mathbf{u}_w} CE(\mathbf{y}, \hat{\mathbf{y}}) = \sigma(\mathbf{u}_w^\top \mathbf{v}_c) \mathbf{v}_c$$

- If $w = o$:

$$\nabla_{\mathbf{u}_o} CE(\mathbf{y}, \hat{\mathbf{y}}) = -\sigma(-\mathbf{u}_o^\top \mathbf{v}_c) \mathbf{v}_c$$

d) $J_{\text{skip-gram}}(\text{word}_{c-m\dots c+m}) = \sum_{-m \leq j \leq m, j \neq 0} F(\mathbf{w}_{c+j}, \mathbf{v}_c)$

- $\forall \mathbf{v}_c \in \mathbb{R}^n$

$$\nabla_{\mathbf{v}_c} J = \sum_{-m \leq j \leq m, j \neq 0} \frac{\partial F(\mathbf{w}_{c+j}, \mathbf{v}_c)}{\partial \mathbf{v}_c}$$

•

$$\nabla_{\mathbf{u}_w} J = \sum_{-m \leq j \leq m, j \neq 0} \frac{\partial F(\mathbf{w}_{c+j}, \mathbf{v}_c)}{\partial \mathbf{u}_w}$$

$$J_{\text{CBOW}}(\text{word}_{c-m\dots c+m}) = F(\mathbf{w}_c, \hat{\mathbf{v}}) \text{ with } \hat{\mathbf{v}} = \sum_{-m \leq j \leq m, j \neq 0} \mathbf{v}_{c+j}$$

•

$$\nabla_{\mathbf{v}_c} J = \vec{0}$$

•

$$\nabla_{\mathbf{u}_w} J = \frac{\partial F(\mathbf{w}_c, \hat{\mathbf{v}})}{\partial \mathbf{u}_w}$$

e) Code

f) Code

g) Code

h) Code

4 Sentiment Analysis

- a) Code
- b) Blah blah blah ... OverFitting ... Blah blah blah more Bias for less Variance
... Blah blah blah
- c) Code
- d) Code