

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

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

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

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

- c)

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

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

$$\nabla_{\mathbf{x}} CE(\mathbf{y}, \hat{\mathbf{y}}) = \nabla_{\boldsymbol{\theta}} CE(\mathbf{y}, \hat{\mathbf{y}}) \left(\frac{\partial \boldsymbol{\theta}}{\partial \mathbf{x}} \right)$$

where $\frac{\partial \boldsymbol{\theta}}{\partial \mathbf{x}}$ is the Jacobian Matrix :

$$\begin{aligned} \forall i \in \llbracket 1, D_y \rrbracket, \forall j \in \llbracket 1, D_x \rrbracket \quad & \left(\frac{\partial \boldsymbol{\theta}}{\partial \mathbf{x}} \right)_{ij} = \left(\frac{\partial \theta_i}{\partial x_j} \right) \\ & = \left(\frac{\partial (\sum_{l=1}^h h_l (W_2)_{li} + (b_2)_i)}{\partial x_j} \right) \\ & = \left(\frac{\partial (\sum_{l=1}^h \sigma(\sum_{m=1}^{D_x} x_m (W_1)_{ml} + (b_1)_l) (W_2)_{li})}{\partial x_j} \right) \\ & = \left(\sum_{l=1}^h (W_1)_{jl} \sigma' \left(\sum_{m=1}^{D_x} x_m (W_1)_{ml} + (b_1)_l \right) (W_2)_{li} \right) \\ & = (\mathbf{W}_1)_{j \cdot} \cdot (\mathbf{h}'^\top * (\mathbf{W}_2)_{\cdot i}) \\ & = (\mathbf{W}_1)_{j \cdot} \cdot (\mathbf{D} \mathbf{W}_2)_{\cdot i} \text{ where } \mathbf{D} = \text{diag}(\sigma'(\mathbf{x} \mathbf{W}_1 + \mathbf{b}_1)) \\ & = ((\mathbf{W}_1 \mathbf{D} \mathbf{W}_2)^\top)_{ij} \end{aligned}$$

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)

$$\begin{aligned} \nabla_{\mathbf{W}_2} CE(\mathbf{y}, \hat{\mathbf{y}}) &= \nabla_{\boldsymbol{\theta}} CE(\mathbf{y}, \hat{\mathbf{y}}) \left(\frac{\partial \boldsymbol{\theta}}{\partial \mathbf{W}_2} \right) \\ &= \mathbf{h}^\top (\hat{\mathbf{y}} - \mathbf{y}) \end{aligned}$$

$$\begin{aligned} \nabla_{\mathbf{b}_2} CE(\mathbf{y}, \hat{\mathbf{y}}) &= \nabla_{\boldsymbol{\theta}} CE(\mathbf{y}, \hat{\mathbf{y}}) \left(\frac{\partial \boldsymbol{\theta}}{\partial \mathbf{b}_2} \right) \\ &= (\hat{\mathbf{y}} - \mathbf{y}) \end{aligned}$$

$$\begin{aligned} \nabla_{\mathbf{W}_1} CE(\mathbf{y}, \hat{\mathbf{y}}) &= \nabla_{\boldsymbol{\theta}} CE(\mathbf{y}, \hat{\mathbf{y}}) \left(\frac{\partial \boldsymbol{\theta}}{\partial \mathbf{W}_1} \right) \\ &= \nabla_{\boldsymbol{\theta}} CE(\mathbf{y}, \hat{\mathbf{y}}) \left(\frac{\partial \boldsymbol{\theta}}{\partial \mathbf{h}} \right) \left(\frac{\partial \mathbf{h}}{\partial \mathbf{W}_1} \right) \\ &= \nabla_{\boldsymbol{\theta}} CE(\mathbf{y}, \hat{\mathbf{y}}) \mathbf{W}_2^\top \left(\frac{\partial \mathbf{h}}{\partial \mathbf{W}_1} \right) \\ &= \mathbf{x}^\top \nabla_{\boldsymbol{\theta}} CE(\mathbf{y}, \hat{\mathbf{y}}) \mathbf{W}_2^\top \mathbf{D} \end{aligned}$$

$$\nabla_{b_1} CE(\mathbf{y}, \hat{\mathbf{y}}) = \nabla_{\theta} CE(\mathbf{y}, \hat{\mathbf{y}}) \mathbf{W}_2^{\top} \mathbf{D}$$

3 word2vec

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

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

$$\hat{\mathbf{y}}_o = \text{softmax}(\mathbf{v}_c \mathbf{U}^{\top})_o \quad \text{where } \mathbf{U} = \begin{pmatrix} \mathbf{u}_1 \\ \vdots \\ \vdots \\ \vdots \\ \mathbf{u}_W \end{pmatrix}$$

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

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

b) $\forall w \in \llbracket 1, W \rrbracket$:

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

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

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

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

• If $w \neq o$:

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

- If $w = o$:

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

Conclusion : $\forall w \in \llbracket 1, W \rrbracket$:

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

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

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

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

•

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

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

•

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

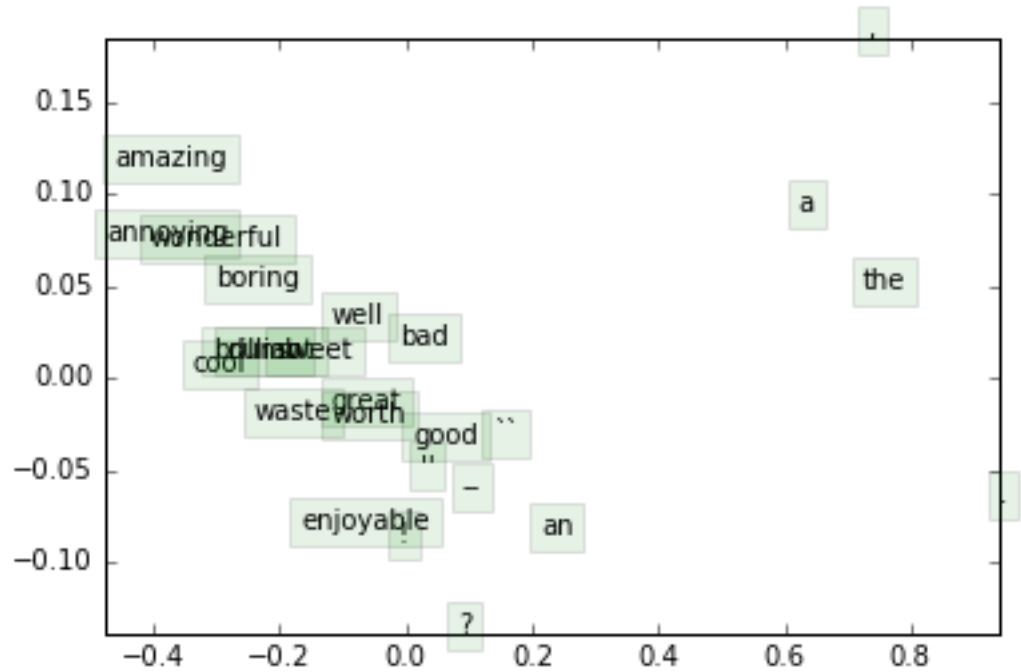
•

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

e) Code

f) Code

g) Code



h)

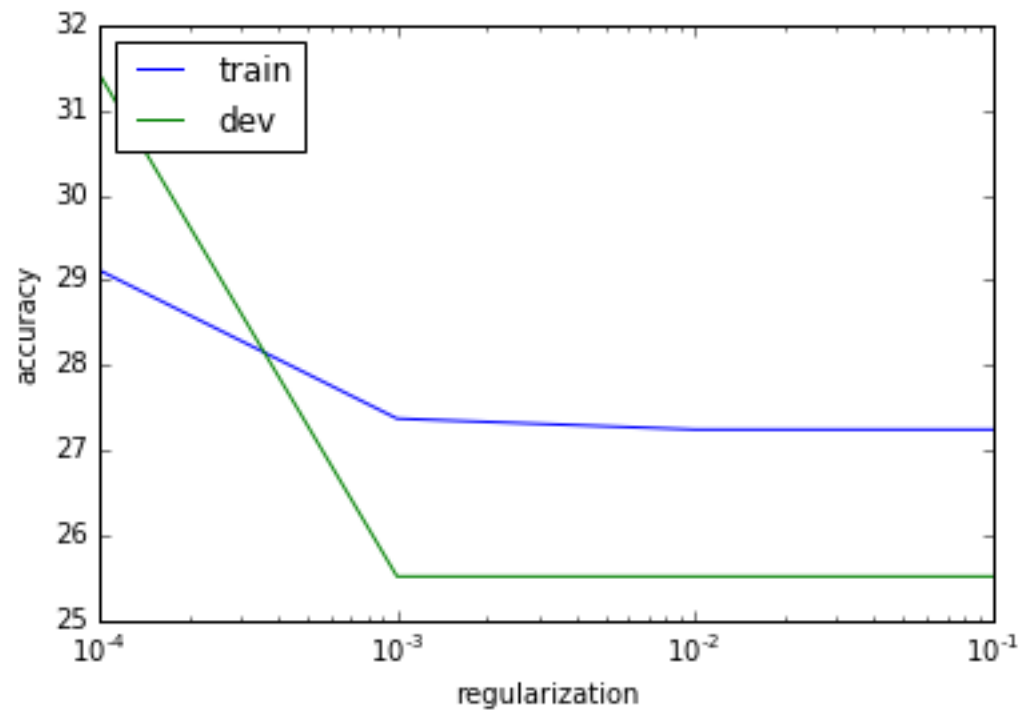
4 Sentiment Analysis

a) $J = \frac{1}{N} \sum_{i=1}^N CE(\hat{\mathbf{y}}_i, \mathbf{y}_i) + \frac{\lambda}{2} \|\mathbf{W}\|^2$ with $\hat{\mathbf{y}}_i = \text{softmax}(\mathbf{x}_i \mathbf{W})$
 $\forall \mathbf{W}$:

$$\begin{aligned} \nabla_{\mathbf{W}} J &= \frac{1}{N} \sum_{i=1}^N \nabla_{\mathbf{W}} CE(\hat{\mathbf{y}}_i, \mathbf{y}_i) + \frac{\lambda}{2} \nabla_{\mathbf{W}} \|\mathbf{W}\|^2 \\ &= \frac{1}{N} \sum_{i=1}^N \mathbf{x}_i^\top (\hat{\mathbf{y}}_i - \mathbf{y}_i) + \lambda \mathbf{W} \end{aligned}$$

b) Blah blah blah ... OverFitting ... Blah blah blah more Bias for less Variance
 ... Blah blah blah

c) I selected $\lambda \in \{0.1, 0.01, 0.01, 0.001, 0.0001\}$



d)