

IFT6390-fundamentals of machine learning

Assignment 3

Jonathan Guymont, Marzieh Mehdizadeh

$$x \tag{1}$$

Question 1

(a)

$$\begin{aligned} 0.5(\tanh(0.5x) + 1) &= 0.5\left(\frac{e^{0.5x} - e^{-0.5x}}{e^{0.5x} + e^{-0.5x}} + 1\right) \\ &= 0.5\left(\frac{1 - e^{-x}}{1 + e^{-x}} + \frac{1 + e^{-x}}{1 + e^{-x}}\right) \\ &= 0.5 \frac{2}{1 + e^{-x}} \\ &= \frac{1}{1 + e^{-x}} \end{aligned} \tag{2}$$

(b)

$$\begin{aligned} \log \text{sigmoid}(x) &= \log(1 + e^{-x})^{-1} \\ &= -\log(1 + e^{-x}) \\ &= -\text{softmax}(-x) \end{aligned} \tag{3}$$

(c)

$$\begin{aligned} \frac{d}{dx} \text{sigmoid}(x) &= \frac{d}{dx} (1 + e^{-x})^{-1} \\ &= -(1 + e^{-x})^{-2} \frac{d}{dx} 1 + e^{-x} \\ &= (1 + e^{-x})^{-2} e^{-x} \\ &= (1 + e^{-x})^{-1} \frac{e^{-x}}{1 + e^{-x}} \\ &= (1 + e^{-x})^{-1} \left(\frac{1 + e^{-x} - 1}{1 + e^{-x}}\right) \\ &= (1 + e^{-x})^{-1} \left(1 - \frac{1}{1 + e^{-x}}\right) \end{aligned} \tag{4}$$

(4)

$$\begin{aligned}
\tanh'(x) &= \left(\frac{e^x - e^{-x}}{e^x + e^{-x}} \right)' \\
&= \frac{(e^x + e^{-x})(e^x + e^{-x}) - (e^x - e^{-x})(e^x - e^{-x})}{(e^x + e^{-x})^2} \\
&= \frac{(e^x + e^{-x})^2 - (e^x - e^{-x})^2}{(e^x + e^{-x})^2} \\
&= 1 - \frac{(e^x - e^{-x})^2}{(e^x + e^{-x})^2} \\
&= 1 - \tanh^2(x)
\end{aligned} \tag{5}$$

(5) $\text{sign}(x) = -1 + 2 \cdot \mathbb{1}_{x>0}$

(6) $\text{abs}'(x) = \text{sign}(x)$

(7) $\text{rect}'(x) = \mathbb{1}_{x>0}$

(8) Let $f(\mathbf{x}) = \sum_{x_i \in \mathbf{x}} x_i^2$, then $f'(\mathbf{x}) = (f'_{x_1}, \dots, f'_{x_{|\mathbf{x}|}}) = (2x_1, \dots, 2x_{|\mathbf{x}|})$

(9) Let $f(\mathbf{x}) = \sum_{x_i \in \mathbf{x}} |x_i|$, then $f'(\mathbf{x}) = (f'_{x_1}, \dots, f'_{x_{|\mathbf{x}|}}) = (\text{sign}(x_1), \dots, \text{sign}(x_{|\mathbf{x}|}))$

Gradient Computation for Parameters optimizations in a neural net for multiclass classification

(1) The dimension of $\mathbf{b}^{(1)}$ is $d_h \times 1$. The formula of the preactivation vector is

$$\mathbf{h}^a = \mathbf{W}^{(1)} \mathbf{x} + \mathbf{b}^{(1)}$$

and the formula for obtaining the value of the element j is

$$\mathbf{h}_j^a = \mathbf{W}_{j,\cdot}^{(1)} \mathbf{x} + \mathbf{b}_j^{(1)} = \mathbf{b}_j^{(1)} + \sum_{i=1}^d \mathbf{W}_{j,i}^{(1)} x_i.$$

The output vector of the activation is given by

$$\mathbf{h}^s = \text{relu}(\mathbf{h}^a)$$

where $\text{relu}(\cdot)$ is applied element wise, i.e. $\mathbf{h}_j^s = \max(0, \mathbf{h}_j^a)$, $j = 1, \dots, d_h$.

(2) The dimension of $\mathbf{W}^{(2)}$ is $m \times d_h$ and the dimension of $\mathbf{b}^{(2)}$ is $m \times 1$.

$$\mathbf{o}^a = \mathbf{W}^{(2)} \mathbf{h}^s + \mathbf{b}^{(2)}$$

$$\mathbf{o}_k^a = \mathbf{W}_{k,\cdot}^{(2)} \mathbf{h}^s + \mathbf{b}_k^{(2)} = \mathbf{b}_k^{(2)} + \sum_{i=1}^{d_h} \mathbf{W}_{k,i}^{(2)} \mathbf{h}_i^s$$

for $k = 1, \dots, m$.

(3)

$$\mathbf{o}_k^s = \frac{\exp(\mathbf{o}_k^a)}{\sum_{k=1}^m \exp(\mathbf{o}_k^a)} \quad (6)$$

They are all positive because $\exp: \mathbb{R} \mapsto \mathbb{R}^+$. Also a sum of positive number is positive. And the ratio of a positive number over a positive number is also positive.

$$\sum_{k=1}^m \mathbf{o}_k^s = \frac{1}{\sum_{k=1}^m \exp(\mathbf{o}_k^a)} \sum_{k=1}^m \exp(\mathbf{o}_k^a) = 1$$

(4) Let $Z = \sum_{k=1}^m \exp(\mathbf{o}_k^a)$. Then $\mathbf{o}^s = \frac{1}{Z}(\exp(\mathbf{o}_1^a), \dots, \exp(\mathbf{o}_m^a))^\top$ and

$$L(\mathbf{x}, y) = -\log \text{onehot}_m(y)(\exp(\mathbf{o}_1^a(\mathbf{x}))/Z, \dots, \exp(\mathbf{o}_m^a(\mathbf{x}))/Z)^\top = -\log \text{onehot}_m(y)\mathbf{o}^s(\mathbf{x})$$

where $\text{onehot}_m(y)$ is a $1 \times m$ onehot representation for y .

(5) \hat{R} is an estimation of the expected value of the loss function (minus the loglikelihood in our case)

$$\hat{R} = \frac{1}{n} \sum_{i=1}^n L(\mathbf{x}^{(i)}, y^{(i)}) = \frac{1}{n} \sum_{i=1}^n -\log \text{onehot}_m(y^{(i)})\mathbf{o}^s(\mathbf{x}^{(i)})$$

The set of trainable parameters is $\boldsymbol{\theta} = \{W^{(1)}, W^{(2)}, b^{(1)}, b^{(2)}\}$. The number of scalar parameters is $n_\theta = d_h \cdot d + d_h + m \cdot d_h + m = d_h(d+1) + m(d_h+1)$.

Optimization problem. First we need to initialize the parameters properly. To find the parameters that minimize the loss, we need to compute the derivative of the loss function w.r.t each parameters. Then we update each parameters by moving them in the opposite direction of their gradient (since we minimize). We repeat this step until a stopping criterion is met (e.g. maximum number of iteration is reached when using early stopping).

(6)

Algorithm 1 Pseudocode for Batch Gradient Descent

Require: Step size η

Require: Initial parameter $\boldsymbol{\omega}_0$

Require: Number of iterations T

for $i = 1$ to T **do**

 Compute gradient $\mathbf{g}_t = \frac{1}{m} \nabla_{\boldsymbol{\omega}} \sum_i L(\mathbf{x}^{(i)}, \mathbf{y}^{(i)})$

 Apply update: $\boldsymbol{\omega}_t = \boldsymbol{\omega}_{t-1} - \eta \mathbf{g}_t$

end for

(7)

$$\begin{aligned}
\frac{\partial}{\partial \boldsymbol{o}^a} L &= \frac{\partial}{\partial \boldsymbol{o}^a} -\log \text{onehot}_m(y) \boldsymbol{o}^s \\
&= -\frac{1}{\text{onehot}_m(y) \boldsymbol{o}^s} \frac{\partial}{\partial \boldsymbol{o}^a} \text{onehot}_m(y) \boldsymbol{o}^s \\
&= -\frac{1}{\text{onehot}_m(y) \boldsymbol{o}^s} \frac{\partial}{\partial \boldsymbol{o}^a} \text{onehot}_m(y) \text{softmax}(\boldsymbol{o}^a) \\
&= -\frac{1}{\text{onehot}_m(y) \boldsymbol{o}^s} \text{onehot}_m(y) \text{softmax}(\boldsymbol{o}^a) (\text{onehot}_m(y) - \text{softmax}(\boldsymbol{o}^a)) \\
&= -(\text{onehot}_m(y) - \text{softmax}(\boldsymbol{o}^a))
\end{aligned} \tag{7}$$

(8)

```
import numpy as np

def grad_oa(x, y):
    onehot = np.zeros(m)
    onehot[y] = 1
    return os(x) - onehot
```

Neural Network Implementation and Experiment