
Homework 2

CSCI 567: Machine Learning

Seyed Mohammad Asghari Pari

Neural networks

Solution to Question 1.1:

Calculating $\frac{\partial l}{\partial \mathbf{u}}$: We first calculate $\frac{\partial l}{\partial \mathbf{u}_m}$ as follows,

$$\frac{\partial l}{\partial \mathbf{u}_m} = \frac{\partial l}{\partial \mathbf{h}_m} \frac{\partial \mathbf{h}_m}{\partial \mathbf{u}_m} = \frac{\partial l}{\partial \mathbf{h}_m} H(\mathbf{u}_m) \quad (1)$$

where the first equality is true because \mathbf{u}_m affects only \mathbf{h}_m and the last equality is true by the definition of Heaviside step function. From (1), we can write

$$\frac{\partial l}{\partial \mathbf{u}} = \frac{\partial l}{\partial \mathbf{h}} \cdot *H(\mathbf{u}). \quad (2)$$

Now, we calculate $\frac{\partial l}{\partial \mathbf{h}}$ and for that we first derive $\frac{\partial l}{\partial \mathbf{h}_m}$ as follows,

$$\frac{\partial l}{\partial \mathbf{h}_m} = \sum_{k=1}^K \frac{\partial l}{\partial \mathbf{a}_k} \frac{\partial \mathbf{a}_k}{\partial \mathbf{h}_m} = \sum_{k=1}^K \frac{\partial l}{\partial \mathbf{a}_k} W_{km}^{(2)} = [W^{(2)}]_{m\bullet}^\top \frac{\partial l}{\partial \mathbf{a}}, \quad (3)$$

where $[W^{(2)}]_{m\bullet}$ is the m -th column of matrix $W^{(2)}$. From (3), we can write

$$\frac{\partial l}{\partial \mathbf{h}} = (W^{(2)})^\top \frac{\partial l}{\partial \mathbf{a}}. \quad (4)$$

Finally, from (2) and (4), we get,

$$\frac{\partial l}{\partial \mathbf{u}} = \left((W^{(2)})^\top \frac{\partial l}{\partial \mathbf{a}} \right) \cdot *H(\mathbf{u}). \quad (5)$$

Calculating $\frac{\partial l}{\partial \mathbf{a}}$: We first calculate $\frac{\partial l}{\partial \mathbf{a}_k}$ as follows,

$$\frac{\partial l}{\partial \mathbf{a}_k} = \sum_{j=1}^K \frac{\partial l}{\partial \mathbf{z}_j} \frac{\partial \mathbf{z}_j}{\partial \mathbf{a}_k} = \sum_{j=1}^K -\frac{\mathbf{y}_j}{\mathbf{z}_j} \frac{\partial \mathbf{z}_j}{\partial \mathbf{a}_k} = \begin{bmatrix} \frac{\partial \mathbf{z}_1}{\partial \mathbf{a}_k} & \frac{\partial \mathbf{z}_2}{\partial \mathbf{a}_k} & \dots & \frac{\partial \mathbf{z}_K}{\partial \mathbf{a}_k} \end{bmatrix} \begin{bmatrix} -\frac{\mathbf{y}_1}{\mathbf{z}_1} \\ -\frac{\mathbf{y}_2}{\mathbf{z}_2} \\ \vdots \\ -\frac{\mathbf{y}_K}{\mathbf{z}_K} \end{bmatrix} = \frac{\partial \mathbf{z}}{\partial \mathbf{a}_k} \left(\mathbf{y} \cdot * \frac{1}{\mathbf{z}} \right) \quad (6)$$

where we have used $\frac{1}{\mathbf{z}}$ to denote $\begin{bmatrix} \frac{1}{\mathbf{z}_1} & \dots & \frac{1}{\mathbf{z}_K} \end{bmatrix}^\top$. Note that from (6) we can write,

$$\frac{\partial l}{\partial \mathbf{a}} = \frac{\partial \mathbf{z}}{\partial \mathbf{a}} \left(\mathbf{y} \cdot * \frac{1}{\mathbf{z}} \right). \quad (7)$$

Next, we calculate $\frac{\partial \mathbf{z}_j}{\partial \mathbf{a}_k}$ as follows: if $j = k$, then we have

$$\frac{\partial \mathbf{z}_k}{\partial \mathbf{a}_k} = \frac{e^{\mathbf{a}_k} \sum_{k'} e^{\mathbf{a}_{k'}} - e^{\mathbf{a}_k} e^{\mathbf{a}_k}}{\left(\sum_{k'} e^{\mathbf{a}_{k'}} \right)^2} = \mathbf{z}_k - \mathbf{z}_k \mathbf{z}_k, \quad (8)$$

and if $j \neq k$, then we have

$$\frac{\partial z_j}{\partial a_k} = \frac{-e^{a_k} e^{a_j}}{(\sum_{k'} e^{a_{k'}})^2} = -z_k z_j. \quad (9)$$

From (8) and (9), we can write,

$$\frac{\partial \mathbf{z}}{\partial \mathbf{a}} = \text{diag}(\mathbf{z}) - \mathbf{z}\mathbf{z}^\top, \quad (10)$$

where $\text{diag}(\mathbf{z})$ is a matrix whose diagonal is vector \mathbf{z} and all other entries are zero. Finally, from (7) and (10), we get

$$\boxed{\frac{\partial l}{\partial \mathbf{a}} = (\text{diag}(\mathbf{z}) - \mathbf{z}\mathbf{z}^\top) \left(\mathbf{y} \cdot \frac{\mathbf{1}}{\mathbf{z}} \right)}. \quad (11)$$

Calculating $\frac{\partial l}{\partial W^{(1)}}$: We first calculate $\frac{\partial l}{\partial W_{ij}^{(1)}}$ as follows,

$$\frac{\partial l}{\partial W_{ij}^{(1)}} = \frac{\partial l}{\partial \mathbf{u}_i} \frac{\partial \mathbf{u}_i}{\partial W_{ij}^{(1)}} = \frac{\partial l}{\partial \mathbf{u}_i} \mathbf{x}_j, \quad (12)$$

where the first equality is true because $W_{ij}^{(1)}$ affects only \mathbf{u}_i . From (12), we can write,

$$\boxed{\frac{\partial l}{\partial W^{(1)}} = \frac{\partial l}{\partial \mathbf{u}} \mathbf{x}^\top}. \quad (13)$$

Calculating $\frac{\partial l}{\partial \mathbf{b}^{(1)}}$: We first calculate $\frac{\partial l}{\partial \mathbf{b}_i^{(1)}}$ as follows,

$$\frac{\partial l}{\partial \mathbf{b}_i^{(1)}} = \frac{\partial l}{\partial \mathbf{u}_i} \frac{\partial \mathbf{u}_i}{\partial \mathbf{b}_i^{(1)}} = \frac{\partial l}{\partial \mathbf{u}_i}, \quad (14)$$

where the first equality is true because $\mathbf{b}_i^{(1)}$ affects only \mathbf{u}_i . From (14), we can write,

$$\boxed{\frac{\partial l}{\partial \mathbf{b}^{(1)}} = \frac{\partial l}{\partial \mathbf{u}}}. \quad (15)$$

Calculating $\frac{\partial l}{\partial W^{(2)}}$: We first calculate $\frac{\partial l}{\partial W_{ij}^{(2)}}$ as follows,

$$\frac{\partial l}{\partial W_{ij}^{(2)}} = \frac{\partial l}{\partial \mathbf{a}_i} \frac{\partial \mathbf{a}_i}{\partial W_{ij}^{(2)}} = \frac{\partial l}{\partial \mathbf{a}_i} \mathbf{h}_j, \quad (16)$$

where the first equality is true because $W_{ij}^{(2)}$ affects only \mathbf{a}_i . From (16), we can write,

$$\boxed{\frac{\partial l}{\partial W^{(2)}} = \frac{\partial l}{\partial \mathbf{a}} \mathbf{h}^\top}. \quad (17)$$

Solution to Question 1.2:

We know that the update equation in stochastic gradient descent for a variable X when we are going to minimize a function l is

$$X^{t+1} = X^t - \eta \frac{\partial l}{\partial X} \Big|_{X=X^t}, \quad (18)$$

where X^t is the value for variable X at the t -th iteration. Note that, (18) can be written as

$$X^{t+1} = X^0 - \sum_{s=0}^t \eta \frac{\partial l}{\partial X} \Big|_{X=X^s}. \quad (19)$$

Now, let X be any of variables $W^{(1)}$, $W^{(2)}$, $\mathbf{b}^{(1)}$. Since we have set the initial values for $W^{(1)}$, $W^{(2)}$, $\mathbf{b}^{(1)}$ to be zero matrices/vectors and we further know that $\frac{\partial l}{\partial W^{(1)}}$, $\frac{\partial l}{\partial W^{(2)}}$, and $\frac{\partial l}{\partial \mathbf{b}^{(1)}}$ are all zero matrices/vectors, we can easily conclude from (19) that $W^{(1)}$, $W^{(2)}$, $\mathbf{b}^{(1)}$ all stay zero matrices/vectors no matter how many iterations we perform. Therefore, no learning will happen on $W^{(1)}$, $W^{(2)}$, $\mathbf{b}^{(1)}$.

For the case of stochastic gradient descent with momentum, the update equation for velocity \mathbf{v} is

$$\mathbf{v}^{t+1} = \alpha \mathbf{v}^t - \eta \frac{\partial l}{\partial X} \Big|_{X=X^t}, \quad (20)$$

$$X^{t+1} = X^t + \mathbf{v}^{t+1}, \quad (21)$$

which can be written as

$$\mathbf{v}^{t+1} = \alpha^{t+1} \mathbf{v}^0 - \sum_{s=0}^t \eta \alpha^{t-s} \frac{\partial l}{\partial X} \Big|_{X=X^s}, \quad (22)$$

$$X^{t+1} = X^0 + \sum_{s=0}^{t+1} \mathbf{v}^s. \quad (23)$$

Again since \mathbf{v}^0 is zero vector, if let X be any of variables $W^{(1)}$, $W^{(2)}$, $\mathbf{b}^{(1)}$ no learning will happen on $W^{(1)}$, $W^{(2)}$, $\mathbf{b}^{(1)}$. The reason again is that we have set the initial values for $W^{(1)}$, $W^{(2)}$, $\mathbf{b}^{(1)}$ to be zero matrices/vectors and we further know that $\frac{\partial l}{\partial W^{(1)}}$, $\frac{\partial l}{\partial W^{(2)}}$, and $\frac{\partial l}{\partial \mathbf{b}^{(1)}}$ are all zero matrices/vectors. Therefore, \mathbf{v} stays zero and then from (23), X stays zero.

Solution to Question 1.3:

If $\mathbf{a} = W^{(2)}\mathbf{u} + \mathbf{b}^{(2)}$, then since $\mathbf{u} = W^{(1)}\mathbf{x} + \mathbf{b}^{(1)}$ we have

$$\mathbf{a} = W^{(2)}(W^{(1)}\mathbf{x} + \mathbf{b}^{(1)}) + \mathbf{b}^{(2)} = W^{(2)}W^{(1)}\mathbf{x} + W^{(2)}\mathbf{b}^{(1)} + \mathbf{b}^{(2)}. \quad (24)$$

Hence, we have

$$\boxed{\mathbf{U} = W^{(2)}W^{(1)}, \quad \mathbf{v} = W^{(2)}\mathbf{b}^{(1)} + \mathbf{b}^{(2)}}. \quad (25)$$

Kernel methods

Solution to Question 2.1:

In this question, we have

$$J(\mathbf{w}) = \sum_{n=1}^N l(\mathbf{w}^\top \phi(\mathbf{x}_n), y_n) + \frac{\lambda}{2} \|\mathbf{w}\|_2^2. \quad (26)$$

If we take the derivative of $J(\mathbf{w})$ with respect to \mathbf{w} and set it to zero, we get

$$\frac{\partial J(\mathbf{w})}{\partial \mathbf{w}} = 0, \quad \Rightarrow \quad \sum_{n=1}^N \frac{\partial l(s_n, y_n)}{\partial s_n} \frac{\partial s_n}{\partial \mathbf{w}} \Big|_{s_n=\mathbf{w}^\top \phi(\mathbf{x}_n)} + \lambda \mathbf{w} = 0 \quad (27)$$

$$\Rightarrow \quad \sum_{n=1}^N \frac{\partial l(s_n, y_n)}{\partial s_n} \Big|_{s_n=\mathbf{w}^\top \phi(\mathbf{x}_n)} \phi(\mathbf{x}_n) + \lambda \mathbf{w} = 0 \quad (28)$$

$$\Rightarrow \quad \mathbf{w} = \sum_{n=1}^N \frac{-1}{\lambda} \frac{\partial l(s_n, y_n)}{\partial s_n} \Big|_{s_n=\mathbf{w}^\top \phi(\mathbf{x}_n)} \phi(\mathbf{x}_n) \quad (29)$$

$$\Rightarrow \quad \boxed{\mathbf{w} = \sum_{n=1}^N \alpha_n \phi(\mathbf{x}_n)}, \quad (30)$$

where we have denoted $\frac{-1}{\lambda} \frac{\partial l(s_n, y_n)}{\partial s_n} \Big|_{s_n = \mathbf{w}^\top \phi(\mathbf{x}_n)}$ by α_n .

Solution to Question 2.2:

Note that from (30), we can write

$$\mathbf{w} = \Phi^\top \boldsymbol{\alpha}, \quad (31)$$

where $\Phi = [\phi(\mathbf{x}_1)^\top \ \dots \ \phi(\mathbf{x}_N)^\top]^\top$. Now, by substituting (31) in (26), we get

$$J(\boldsymbol{\alpha}) = \sum_{n=1}^N l(\boldsymbol{\alpha}^\top \Phi \phi(\mathbf{x}_n), y_n) + \frac{\lambda}{2} \|\Phi^\top \boldsymbol{\alpha}\|_2^2 = \sum_{n=1}^N l(\boldsymbol{\alpha}^\top \Phi \phi(\mathbf{x}_n), y_n) + \frac{\lambda}{2} \boldsymbol{\alpha}^\top \Phi \Phi^\top \boldsymbol{\alpha} \quad (32)$$

which can be written as,

$$J(\boldsymbol{\alpha}) = \sum_{n=1}^N l(\boldsymbol{\alpha}^\top [K]_{n\bullet}, y_n) + \frac{\lambda}{2} \boldsymbol{\alpha}^\top K \boldsymbol{\alpha}, \quad (33)$$

where $K = \Phi \Phi^\top$ and $[K]_{n\bullet}$ is the n -th column of matrix K , that is,

$$[K]_{n\bullet} = \begin{bmatrix} \phi(x_1)^\top \phi(x_n) \\ \vdots \\ \phi(x_N)^\top \phi(x_n) \end{bmatrix}. \quad (34)$$