

1 Introduction

2 MLP for binary classification

2.1 Implementing the gradient descent

2.1.1 Forward pass

First Layer

$$\begin{cases} a_{L,q}^{(1)} = (\mathbf{w}_{L,q}^{(1)})^T \cdot \mathbf{x}_L + b_{L,q}^{(1)}, & q = 1, \dots, h_1 \\ a_{R,q}^{(1)} = (\mathbf{w}_{R,q}^{(1)})^T \cdot \mathbf{x}_R + b_{R,q}^{(1)}, & q = 1, \dots, h_1 \end{cases}$$

Or in a more compact way:

$$\begin{cases} \mathbf{a}_L^{(1)} = \mathbf{W}_L^{(1)} \cdot \mathbf{x}_L + \mathbf{b}_L^{(1)} \\ \mathbf{a}_R^{(1)} = \mathbf{W}_R^{(1)} \cdot \mathbf{x}_R + \mathbf{b}_R^{(1)} \end{cases}$$
$$\begin{cases} \mathbf{z}_L^{(1)} = \mathbf{g}_1(\mathbf{a}_L^{(1)}) \\ \mathbf{z}_R^{(1)} = \mathbf{g}_1(\mathbf{a}_R^{(1)}) \end{cases}$$

Second Layer

$$a_{L,q}^{(2)} = (\mathbf{w}_{L,q}^{(2)})^T \cdot \mathbf{z}_L^{(1)} + b_{L,q}^{(2)}, \quad q = 1, \dots, h_2$$

$$a_{LR,q}^{(2)} = (\mathbf{w}_{LR,q}^{(2)})^T \cdot \begin{bmatrix} \mathbf{z}_L^{(1)} \\ \mathbf{z}_R^{(1)} \end{bmatrix} + b_{L,q}^{(2)}, \quad q = 1, \dots, h_2$$

$$a_{R,q}^{(2)} = (\mathbf{w}_{R,q}^{(2)})^T \cdot \mathbf{z}_R^{(1)} + b_{R,q}^{(2)}, \quad q = 1, \dots, h_2$$

More compactly,

$$\mathbf{a}_L^{(2)} = \mathbf{W}_L^{(2)} \cdot \mathbf{z}_L^{(1)} + \mathbf{b}_L^{(2)}$$

$$\mathbf{a}_{LR}^{(2)} = \mathbf{W}_{LR}^{(2)} \cdot \begin{bmatrix} \mathbf{z}_L^{(1)} \\ \mathbf{z}_R^{(1)} \end{bmatrix} + \mathbf{b}_{LR}^{(2)}$$

$$\mathbf{a}_R^{(2)} = \mathbf{W}_R^{(2)} \cdot \mathbf{z}_R^{(1)} + \mathbf{b}_R^{(2)}$$

And,

$$\mathbf{z}^{(2)} = \mathbf{g}_2(\mathbf{a}_L^{(2)}, \mathbf{a}_{LR}^{(2)}, \mathbf{a}_R^{(2)})$$

Third Layer

$$a_q^{(3)} = (\mathbf{w}_q^{(3)})^T \cdot \mathbf{z}^{(2)} + b_q^{(3)}, \quad q = 1, \dots, K$$

$$\mathbf{a}^{(3)} = \mathbf{W}^{(3)} \cdot \mathbf{z}^{(2)} + \mathbf{b}^{(3)}$$

2.2 backward pass

TODO

2.3 Results

(EXTREMELY good results on validation set, training set too easy)

2.4 Extrapolation

(creation of another binary test)

3 Multi-way classification

3.1 Linear multi-way classification

(squared error, tichinov regularizer, logistic error)

3.2 How to deal multis way with MLP?

3.3 Computation of gradient

3.3.1 Forward pass

The first and second layer or the forward bass are exactly the same as in the binary case. Only change the third laye r :

$$a_q^{(3)} = (\mathbf{w}_q^{(3)})^T \cdot \mathbf{z}^{(2)} + b_q^{(3)}, \quad q = 1, \dots, K$$

$$\mathbf{a}^{(3)} = \mathbf{W}^{(3)} \cdot \mathbf{z}^{(2)} + \mathbf{b}^{(3)}$$

3.3.2 Backward pass

Third Layer The goal is the minimize $E_2(\mathbf{w}) = \frac{1}{2} \sum_{i=1}^N \|\mathbf{a}^{(3)}(\mathbf{x}_i) - \tilde{\mathbf{t}}_i\|^2$. We compute the gradient for one i that is we minimize $E_{2,i}$ that we will refer to as E_i

$$r_q^{(3)} = \frac{\partial E_i}{\partial a_q^{(3)}} = 2(a_q^{(3)} - \tilde{t}_q), \quad q = 1, \dots, K$$

$$\mathbf{r}^{(3)} = [r_q^{(3)}] = 2(\mathbf{a}^{(3)} - \tilde{\mathbf{t}})$$

$$\begin{cases} \nabla_{\mathbf{W}^{(3)}} E_i = \mathbf{r}^{(3)} \cdot \nabla_{\mathbf{W}^{(3)}} a^{(3)} = \mathbf{r}^{(3)} \cdot (\mathbf{z}^{(2)})^T \\ \nabla_{\mathbf{b}^{(3)}} E_i = \mathbf{r}^{(3)} \cdot \nabla_{\mathbf{b}^{(3)}} a^{(3)} = \mathbf{r}^{(3)} \end{cases}$$

Second Layer First, let's compute the derivatives (copy from mlp_implementation.docx)

$$\begin{aligned} r_{L,q}^{(2)} &= \frac{\partial E_i}{\partial a_{L,q}^{(2)}} = \sum_j^K \frac{\partial E_i}{\partial a_j^{(3)}} \cdot \frac{\partial a_j^{(3)}}{\partial a_{L,q}^{(2)}} \\ &= \sum_j^K r_j^{(3)} \frac{\partial[(\mathbf{w}_j^{(3)})^T \cdot \mathbf{g}_2(\mathbf{a}_L^{(2)}, \mathbf{a}_{LR}^{(2)}, \mathbf{a}_R^{(2)}) + b_j^{(3)}]}{\partial a_{L,q}^{(2)}} \\ &= \sum_j^K r_j^{(3)} \frac{\partial[\sum_{\mathbf{k}} \mathbf{w}_{j,k}^{(3)} g_2(\mathbf{a}_{L,k}^{(2)}, \mathbf{a}_{LR,k}^{(2)}, \mathbf{a}_{R,k}^{(2)})]}{\partial a_{L,q}^{(2)}} \\ &= \sum_j^K r_j^{(3)} \frac{\partial[\mathbf{w}_{j,q}^{(3)} g_2(\mathbf{a}_{L,q}^{(2)}, \mathbf{a}_{LR,q}^{(2)}, \mathbf{a}_{R,q}^{(2)})]}{\partial a_{L,q}^{(2)}} \\ &= \sum_j^K r_j^{(3)} w_{j,q}^{(3)} g_2'(\mathbf{a}_{L,q}^{(2)}, \mathbf{a}_{LR,q}^{(2)}, \mathbf{a}_{R,q}^{(2)}) \end{aligned}$$

Vectorizing gives:

$$\begin{cases} \mathbf{r}_L^{(2)} = \left(\text{diag} \left(g_2'(\mathbf{a}_L^{(2)}, \mathbf{a}_{LR}^{(2)}, \mathbf{a}_R^{(2)}) \right) \right) \cdot (\mathbf{W}^{(3)})^T \cdot \mathbf{r}^{(3)} \\ \mathbf{r}_{LR}^{(2)} = \left(\text{diag} \left(g_2'(\mathbf{a}_L^{(2)}, \mathbf{a}_{LR}^{(2)}, \mathbf{a}_R^{(2)}) \right) \right) \cdot (\mathbf{W}^{(3)})^T \cdot \mathbf{r}^{(3)} \\ \mathbf{r}_R^{(2)} = \left(\text{diag} \left(g_2'(\mathbf{a}_L^{(2)}, \mathbf{a}_{LR}^{(2)}, \mathbf{a}_R^{(2)}) \right) \right) \cdot (\mathbf{W}^{(3)})^T \cdot \mathbf{r}^{(3)} \end{cases}$$

So now the gradients are:

$$\begin{cases} \nabla_{\mathbf{W}_L^{(2)}} E_i = \mathbf{r}_L^{(2)} \left(\nabla_{\mathbf{W}_L^{(2)}} a_L^{(2)} \right)^T = \mathbf{r}_L^{(2)} \cdot (\mathbf{z}_L^{(1)})^T \\ \nabla_{\mathbf{W}_{LR}^{(2)}} E_i = \mathbf{r}_{LR}^{(2)} \left(\nabla_{\mathbf{W}_{LR}^{(2)}} a_{LR}^{(2)} \right)^T = \mathbf{r}_{LR}^{(2)} \cdot \begin{bmatrix} \mathbf{z}_L^{(1)} \\ \mathbf{z}_R^{(1)} \end{bmatrix}^T \\ \nabla_{\mathbf{W}_R^{(2)}} E_i = \mathbf{r}_R^{(2)} \left(\nabla_{\mathbf{W}_R^{(2)}} a_R^{(2)} \right)^T = \mathbf{r}_R^{(2)} \cdot (\mathbf{z}_R^{(1)})^T \end{cases}$$

And for b :

$$\begin{cases} \nabla_{\mathbf{b}_L^{(2)}} E_i = \mathbf{r}_L^{(2)} \left(\nabla_{\mathbf{b}_L^{(2)}} a_L^{(2)} \right)^T = \mathbf{r}_L^{(2)} \\ \nabla_{\mathbf{b}_{LR}^{(2)}} E_i = \mathbf{r}_{LR}^{(2)} \left(\nabla_{\mathbf{b}_{LR}^{(2)}} a_{LR}^{(2)} \right)^T = \mathbf{r}_{LR}^{(2)} \\ \nabla_{\mathbf{b}_R^{(2)}} E_i = \mathbf{r}_R^{(2)} \left(\nabla_{\mathbf{b}_R^{(2)}} a_R^{(2)} \right)^T = \mathbf{r}_R^{(2)} \end{cases}$$

First Layer

$$\begin{aligned}
r_{L,q}^{(1)} &= \frac{\partial E_i}{\partial a_{L,q}^{(1)}} = \sum_j^{h_2} \frac{\partial E_i}{\partial a_{L,j}^{(2)}} \cdot \frac{\partial a_{L,j}^{(2)}}{\partial a_{L,q}^{(1)}} + \sum_j^{h_2} \frac{\partial E_i}{\partial a_{LR,j}^{(2)}} \cdot \frac{\partial a_{LR,j}^{(2)}}{\partial a_{L,q}^{(1)}} \\
&= \sum_j^{h_2} r_{L,j}^{(2)} \frac{\partial a_{L,j}^{(2)}}{\partial a_{L,q}^{(1)}} + \sum_j^{h_2} r_{LR,j}^{(2)} \frac{\partial a_{LR,j}^{(2)}}{\partial a_{L,q}^{(1)}}
\end{aligned}$$

$$\begin{aligned}
r_{R,q}^{(1)} &= \frac{\partial E_i}{\partial a_{R,q}^{(1)}} = \sum_j^{h_2} \frac{\partial E_i}{\partial a_{R,j}^{(2)}} \cdot \frac{\partial a_{R,j}^{(2)}}{\partial a_{R,q}^{(1)}} + \sum_j^{h_2} \frac{\partial E_i}{\partial a_{LR,j}^{(2)}} \cdot \frac{\partial a_{LR,j}^{(2)}}{\partial a_{R,q}^{(1)}} \\
&= \sum_j^{h_2} r_{R,j}^{(2)} \frac{\partial a_{R,j}^{(2)}}{\partial a_{R,q}^{(1)}} + \sum_j^{h_2} r_{LR,j}^{(2)} \frac{\partial a_{LR,j}^{(2)}}{\partial a_{R,q}^{(1)}}
\end{aligned}$$

So, for a (p, j) pair, we have

$$\begin{aligned}
\frac{\partial a_{L,j}^{(2)}}{\partial a_{L,q}^{(1)}} &= \frac{\partial \left[\mathbf{w}_{L,j}^{(2)} \cdot \mathbf{g}_1(\mathbf{a}_L^{(1)}) + b_{L,j}^{(2)} \right]}{\partial a_{L,q}^{(1)}} \\
&= \frac{\partial \left[\sum_k w_{L,j,k}^{(2)} \cdot g_1(a_{L,k}^{(1)}) \right]}{\partial a_{L,q}^{(1)}} \\
&= \frac{\partial \left[w_{L,j,q}^{(2)} \cdot g_1(a_{L,q}^{(1)}) \right]}{\partial a_{L,q}^{(1)}} \\
&= w_{L,j,q}^{(2)} \cdot g_1'(a_{L,q}^{(1)})
\end{aligned}$$

$$\begin{aligned}
\frac{\partial a_{LR,j}^{(2)}}{\partial a_{L,q}^{(1)}} &= \frac{\partial \left[\mathbf{w}_{LR,j}^{(2)} \cdot \begin{bmatrix} \mathbf{g}_1(\mathbf{a}_L^{(1)}) \\ \mathbf{g}_1(\mathbf{a}_R^{(1)}) \end{bmatrix} + b_{L,j}^{(2)} \right]}{\partial a_{L,q}^{(1)}} \\
&= \frac{\partial \left[\sum_{k=1}^{h_1} w_{LR,j}^{(2)}(k) \cdot g_1(a_{L,k}^{(1)}) + \sum_{k=1}^{h_1} w_{LR,j}^{(2)}(h_1 + k) \cdot g_1(a_{R,k}^{(1)}) \right]}{\partial a_{L,q}^{(1)}} \\
&= w_{LR,j}^{(2)}(q) \cdot g_1'(a_{L,q}^{(1)})
\end{aligned}$$

$$\begin{aligned}
r_{L,q}^{(1)} &= \sum_j^{h_2} r_{L,j}^{(2)} \frac{\partial a_{L,j}^{(2)}}{\partial a_{L,q}^{(1)}} + \sum_j^{h_2} r_{LR,j}^{(2)} \frac{\partial a_{LR,j}^{(2)}}{\partial a_{L,q}^{(1)}} \\
&= \sum_j^{h_2} r_{L,j}^{(2)} \cdot w_{L,j,q}^{(2)} \cdot g'_1(a_{L,q}^{(1)}) + \sum_j^{h_2} r_{LR,j}^{(2)} \cdot w_{LR,j}^{(2)}(q) \cdot g'_1(a_{L,q}^{(1)}) \\
&= g'_1(a_{L,q}^{(1)}) \cdot \left(\mathbf{W}_L^{(2)} \right)_q^T \cdot \mathbf{r}_L^{(2)} + g'_1(a_{L,q}^{(1)}) \cdot \left(\mathbf{W}_{LR}^{(2)} \right)_q^T \cdot \mathbf{r}_{LR}^{(2)}
\end{aligned}$$

In summary after vectorization we have:

$$\begin{cases} \mathbf{r}_L^{(1)} &= \text{diag} \left(g'_1(\mathbf{a}_L^{(1)}) \right) \cdot \left(\mathbf{W}_L^{(2)} \right)^T \cdot \mathbf{r}_L^{(2)} + \text{diag} \left(g'_1(\mathbf{a}_L^{(1)}) \right) \cdot \left(\mathbf{W}_{LR}^{(2)}(:, 1 : h_1) \right)^T \cdot \mathbf{r}_{LR}^{(2)} \\ \mathbf{r}_R^{(1)} &= \text{diag} \left(g'_1(\mathbf{a}_R^{(1)}) \right) \cdot \left(\mathbf{W}_R^{(2)} \right)^T \cdot \mathbf{r}_R^{(2)} + \text{diag} \left(g'_1(\mathbf{a}_R^{(1)}) \right) \cdot \left(\mathbf{W}_{LR}^{(2)}(:, h_1 + 1 : 2 \times h_1) \right)^T \cdot \mathbf{r}_{LR}^{(2)} \end{cases}$$

Now that we have computed the residual variables, we can easily find the gradients for the first layer. Indeed,

$$\begin{cases} \nabla_{\mathbf{W}_L^{(1)}} E_i = \mathbf{r}_L^{(1)} \cdot \left(\nabla_{\mathbf{W}_L^{(1)}} a_L^{(1)} \right)^T = \mathbf{r}_L^{(1)} \cdot (\mathbf{x}_L)^T \\ \nabla_{\mathbf{W}_R^{(1)}} E_i = \mathbf{r}_R^{(1)} \cdot \left(\nabla_{\mathbf{W}_R^{(1)}} a_R^{(1)} \right)^T = \mathbf{r}_R^{(1)} \cdot (\mathbf{x}_R)^T \end{cases}$$

And,

$$\begin{cases} \nabla_{\mathbf{b}_L^{(1)}} E_i = \mathbf{r}_L^{(1)} \cdot \left(\nabla_{\mathbf{b}_L^{(1)}} a_L^{(1)} \right)^T = \mathbf{r}_L^{(1)} \\ \nabla_{\mathbf{b}_R^{(1)}} E_i = \mathbf{r}_R^{(1)} \cdot \left(\nabla_{\mathbf{b}_R^{(1)}} a_R^{(1)} \right)^T = \mathbf{r}_R^{(1)} \end{cases}$$

3.4 Results and optimization

3.5 Discussion

4 Conclusion