

# A PROOF OF CONVERGENCE OF MULTI-CLASS LOGISTIC REGRESSION NETWORK

MAREK RYCHLIK  
UNIVERSITY OF ARIZONA  
DEPARTMENT OF MATHEMATICS, 617 N SANTA RITA RD, P.O. BOX 210089  
TUCSON, AZ 85721-0089, USA

**ABSTRACT.** This paper revisits the special type of a neural network known under two names. In the statistics and machine learning community it is known as a multi-class logistic neural network. In the neural network community, it is simply the soft-max layer. The importance is underscored by its role in deep learning: as the last layer, whose output is actually the classification of the input patterns, such as images. Our exposition focuses on mathematically rigorous derivation of the key equation expressing the gradient. The fringe benefit of our approach is a fully vectorized expression, which is a basis of an efficient implementation. The second result of this paper is the positivity of the second derivative of the cross-entropy loss function as function of the weights. This result proves that optimization methods based on convexity may be used to train this network. As a corollary, we demonstrate that no  $L^2$ -regularizer is needed to guarantee convergence of gradient descent.

## 1. NOTATION AND DEFINITIONS

The multi-class logistic regression network is a neural network which takes an input vector  $\mathbf{x} \in \mathbb{R}^D$  and produces an activation vector  $\mathbf{a} \in \mathbb{R}^C$  by a linear transformation

$$\mathbf{a} = \mathbf{W} \mathbf{x},$$

where  $\mathbf{W} = [w_{jk}]$  is an  $C \times D$  matrix of weights. The vector space of such matrices will be denoted by  $L(\mathbb{R}^D, \mathbb{R}^C)$  and identified with the space of linear transformations

$$\mathbf{W} : \mathbb{R}^D \rightarrow \mathbb{R}^C.$$

The activations are subsequently transformed by the *soft-max function*  $\sigma : \mathbb{R}^C \rightarrow \mathbb{R}^C$  given by the formula

$$\sigma_i(\mathbf{u}) = \frac{e^{u_i}}{\sum_{j=1}^C e^{u_j}}.$$

Thus,

$$\mathbf{y} = \sigma(\mathbf{W} \mathbf{x})$$

The number  $C$  is the number of classes into which the input vectors will be classified into. For a *training vector*  $\mathbf{x}$  its classification is known and it is given by a target vector  $\mathbf{t} \in \{0, 1\}^C \subset \mathbb{R}^C$ , where  $t_k = 1$  iff vector  $\mathbf{x}$  belongs to class  $k$ . Thus

$$\sum_{k=1}^C t_k = 1.$$

In order to train the network, we need the *training set* consisting of  $N$  sample input vectors

$$\mathbf{x}^{(n)} \in \mathbb{R}^D, \quad n = 1, 2, \dots, N$$

and  $N$  corresponding target vectors

$$\mathbf{t}^{(n)} \in \mathbb{R}^C.$$

The loss function to be minimized is the *cross-entropy* loss function, given by the formula:

$$(1) \quad L(\mathbf{W}) = - \sum_{n=1}^N \sum_{i=1}^D t_i^{(n)} \log y_i^{(n)}.$$

The problem of training is the problem of finding the optimal weight matrix  $\widehat{\mathbf{W}}$  which minimizes  $L(\mathbf{W})$ :

$$\widehat{\mathbf{W}} = \arg \min_{\mathbf{W} \in L(\mathbb{R}^D, \mathbb{R}^C)} L(\mathbf{W}).$$

The natural question arises, whether the minimum exists and whether it is unique. It is known that the simplistic answer to this question is “no” because shifting the weights by a constant depending on  $n$  only does not change the value of  $L$ . More precisely, if

$$\tilde{\mathbf{W}} = \mathbf{W} + \mathbb{1} \mathbf{c}^\top$$

then  $L(\tilde{\mathbf{W}}) = L(\mathbf{W})$ , and, more strongly

$$\tilde{\mathbf{y}}^{(n)} = \mathbf{y}^{(n)}$$

where

$$\begin{aligned} \mathbf{y}^{(n)} &= \sigma \left( \mathbf{W} \mathbf{x}^{(n)} \right). \\ \tilde{\mathbf{y}}^{(n)} &= \sigma \left( \tilde{\mathbf{W}} \mathbf{x}^{(n)} \right). \end{aligned}$$

This is the consequence of the following identity: for every  $c \in \mathbb{R}$

$$\sigma(\mathbf{u} + c \mathbb{1}) = \sigma(\mathbf{u}).$$

This brief paper answers this question in the most satisfactory fashion, giving sufficient conditions for  $L$  to be a strictly convex function on the subspace of those weight matrices  $\mathbf{W}$  for which the sum of every column is 0, i.e.

$$\mathbb{1}^\top \mathbf{W} = 0.$$

This condition guarantees both the existence and uniqueness of the solution, and convergence of optimization algorithms which depend on strict convexity.

## 2. RELATED READING

Many related techniques are described in [2, 3]. The softmax function use for pattern recognition is also discussed in [1], on pp. 215 and pp. 238–245.

## 3. THE CRITICAL POINTS OF $L$

In this section we summarize known conditions for  $\mathbf{W}$  to be a critical points of  $L$ . For simplicity of notation we assume  $N = 1$ , which allows us to avoid the superscript  $(n)$ . The summation over all samples yields the result for the entire training set. The manner in which we derive the formulas may differ from a typical derivation in machine learning texts, in that it relies upon the Chain Rule for *Fréchet derivatives*, rather than calculation of partials with respect of individual weights  $w_{jk}$ .<sup>1</sup>

---

<sup>1</sup>It should be remembered that Fréchet derivatives are closely related to *Jacobi matrices*, but they are **not the same**, as it will be quite apparent in our calculations. Fréchet derivative is a linear transformation, and upon the choice of a basis, or, in the case of  $\mathbb{R}^n$ , when the standard basis is used, this linear transformation is identified with the Jacobi matrix.

With this approach, we automatically derive *vectorized* formulas for the gradient of  $L$ , which subsequently leads to very efficient implementation of the training algorithm, as all critical operations are simply matrix products.

We also allow the target vector  $\mathbf{t} \in \mathbb{R}^C$ , without necessarily requiring that  $t_k \in \{0, 1\}$ , but we still require

$$\sum_{k=1}^C t_k = 1.$$

From the point of view of applications, this generalization is useful when the classification of the inputs is ambiguous. For instance, we can use several humans to classify the inputs, in which case the humans may classify the inputs differently. Then we could assign to  $t_k$  the fraction of humans who assign the input to class  $k$ .

We use the following representation of the loss function  $L$ :

$$L(\mathbf{W}) = Q_{\mathbf{t}}(-\log(\sigma(P_{\mathbf{x}}(\mathbf{W}))))$$

where

$$\begin{aligned} P_{\mathbf{x}}(\mathbf{W}) &= \mathbf{W} \mathbf{x}, \\ Q_{\mathbf{t}}(\mathbf{z}) &= \mathbf{t}^\top \mathbf{z}. \end{aligned}$$

is a vector-valued *linear operator on matrices*, consisting in multiplying a matrix by  $\mathbf{x}$  on the right. Thus

$$L = Q_{\mathbf{t}} \circ ((-\log) \circ \sigma) \circ P_{\mathbf{x}},$$

which is a composition involving linear operators  $Q_{\mathbf{t}}$  and  $P_{\mathbf{x}}$  and two non-linear transformations,  $\log$  (coordinatewise) and  $\sigma$ . Let  $\rho = (-\log) \circ \sigma$ . It will be beneficial to think of  $L$  as

$$L = Q_{\mathbf{t}} \circ \rho \circ P_{\mathbf{x}}$$

which is a composition with only one non-linear term.

**Remark 1.** *This form is particularly useful for calculating higher derivatives of  $L$ , specifically, the Hessian of  $L$ .*

The Chain Rule yields the Fréchet derivative of the composition:

$$DL = DQ_{\mathbf{t}} D\rho DP_{\mathbf{x}}$$

where the intermediate Fréchet derivatives are evaluated at respective intermediate values of the composition. The more detailed version of the formula is:

$$DL = (DQ_{\mathbf{t}} \circ \rho \circ P_{\mathbf{x}}) (D\rho \circ P_{\mathbf{x}}) DP_{\mathbf{x}}$$

It should be noted that for a function  $\mathbf{F} : \mathbb{R}^n \rightarrow \mathbb{R}^m$ ,

$$D\mathbf{F} : \mathbb{R}^n \rightarrow L(\mathbb{R}^n, \mathbb{R}^m),$$

i.e. it is an operator-valued function on  $\mathbb{R}^n$ .

The Fréchet derivative of a linear transformation is the transformation itself. Thus

$$(2) \quad DQ_{\mathbf{t}} = Q_{\mathbf{t}}$$

$$(3) \quad DP_{\mathbf{x}} = P_{\mathbf{x}}$$

regardless of the argument.

We state one more useful formula. For any vectorized scalar function, evaluated elementwise, i.e. function  $\mathbf{F} : \mathbb{R}^n \rightarrow \mathbb{R}^n$  given by:

$$\mathbf{F}(\mathbf{x}) = (f(x_1), f(x_2), \dots, f(x_n))$$

is a diagonal matrix:

$$D\mathbf{F}(\mathbf{x}) = \begin{bmatrix} f'(x_1) & 0 & \dots & 0 \\ 0 & f'(x_2) & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \dots & f'(x_n) \end{bmatrix} = \text{diag}(f'(\mathbf{x})).$$

where  $\text{diag}$  is a MATLAB-like operator converting a vector to a diagonal matrix. We make an observation that  $\text{diag}$  itself is a linear operator from vectors to matrices.

In summary,

$$(4) \quad DL(\mathbf{W}) = Q_{\mathbf{t}} D\rho(\mathbf{a}) P_{\mathbf{x}}$$

Let us compute the derivative of  $\rho$ . Firstly,

$$\rho(\mathbf{u}) = -\mathbf{u} + \left( \log \left( \sum_{k=1}^C e^{u_i} \right) \right) \mathbb{1}$$

where  $\mathbb{1} = (1, 1, \dots, 1)$  is a (column) vector of 1's. This implies easily, using a combination of techniques already mentioned:

$$\begin{aligned} D\rho(\mathbf{u})\mathbf{h} &= -\mathbf{h} + \left( \frac{1}{\sum_{k=1}^C e^{u_i}} (e^{\mathbf{u}})^{\top} \mathbf{h} \right) \mathbb{1} \\ &= -\mathbf{h} + (\boldsymbol{\sigma}(\mathbf{u})^{\top} \mathbf{h}) \mathbb{1} = -\mathbf{h} + \mathbb{1} (\boldsymbol{\sigma}(\mathbf{u})^{\top} \mathbf{h}) = (-I + \mathbb{1} \boldsymbol{\sigma}(\mathbf{u})^{\top}) \mathbf{h} \end{aligned}$$

where  $I$  is the  $C \times C$  identity matrix. Thus, in short,

$$D\rho(\mathbf{u}) = -I + \mathbb{1} \boldsymbol{\sigma}(\mathbf{u})^{\top}.$$

Finally, carefully looking at the formula (4) and equating matrix product with composition of linear transformations (this convention should be familiar from linear algebra), we obtain

$$\begin{aligned} DL(\mathbf{W})\mathbf{V} &= \mathbf{t}^{\top} (-I + \mathbb{1} \boldsymbol{\sigma}(\mathbf{W}\mathbf{x})^{\top}) \mathbf{V} \mathbf{x} \\ &= -\mathbf{t}^{\top} \mathbf{V} \mathbf{x} + (\mathbf{t}^{\top} \mathbb{1}) \mathbf{y}^{\top} \mathbf{V} \mathbf{x} \\ &= -\mathbf{t}^{\top} \mathbf{V} \mathbf{x} + \mathbf{y}^{\top} \mathbf{V} \mathbf{x} \\ &= -(\mathbf{t} - \mathbf{y})^{\top} \mathbf{V} \mathbf{x}. \end{aligned}$$

We note that we used  $\mathbf{t}^{\top} \mathbb{1} = 1$  because  $\sum_{k=1}^C t_k = 1$  was assumed.

This formula generalizes easily to  $N \geq 1$ , when  $L$  is given by (1):

$$DL(\mathbf{W})\mathbf{V} = - \sum_{n=1}^N \left( \mathbf{t}^{(n)} - \mathbf{y}^{(n)} \right)^{\top} \mathbf{V} \mathbf{x}^{(n)}.$$

In order to use methods such as gradient descent, we need to find the gradient  $\nabla L(\mathbf{W})$  from the formula for  $DL(\mathbf{W})$ . We should note that the gradient belongs to the same vector space as the argument  $\mathbf{W}$ , while  $DL(\mathbf{W})$  is a functional on the same space. In our situation:

$$(5) \quad \nabla L(\mathbf{W}) \in L(\mathbb{R}^D, \mathbb{R}^C)$$

$$(6) \quad DL(\mathbf{W}) \in L(L(\mathbb{R}^D, \mathbb{R}^C), \mathbb{R}) = L(\mathbb{R}^C, \mathbb{R}^C)^*$$

where the notation  $X^* = L(X, \mathbb{R})$  defines the *dual space* of the vector space  $X$ . Thus, an expression  $\mathbf{W} + \eta \nabla L(\mathbf{W})$  can be evaluated for  $\eta \in \mathbb{R}$ , but  $\mathbf{W} + \eta DL(\mathbf{W})$  makes no sense.

We recall that the notion of gradient depends on the inner product in the underlying vector space. In our case, it is the space of weight matrices  $\mathbf{W}$ , i.e.  $L(\mathbb{R}^D, \mathbb{R}^C)$ . We assume the most simple form of the inner product: the Frobenius inner product:

$$(7) \quad \langle \mathbf{U}, \mathbf{V} \rangle = \sum_{j=1}^C \sum_{k=1}^D U_{jk} V_{jk} = \text{tr } \mathbf{U}^T \mathbf{V}.$$

We then have the definition of gradient by duality:

$$(8) \quad DL(\mathbf{W})\mathbf{V} = \langle \nabla L(\mathbf{W}), \mathbf{V} \rangle.$$

for all matrices  $\mathbf{V}$ . Then, for  $N = 1$  we look for

$$-(\mathbf{t} - \mathbf{y})^\top \mathbf{V} \mathbf{x} = \text{tr } \nabla L(\mathbf{W})^\top \mathbf{V}$$

Using the identity  $\text{tr } (\mathbf{A} \mathbf{B}) = \text{tr } (\mathbf{B} \mathbf{A})$ , and for vectors  $\mathbf{a}$  and  $\mathbf{b}$ ,  $\mathbf{a}^\top \mathbf{b} = \text{tr } (\mathbf{b} \mathbf{a}^\top)$ , we obtain:

$$-\text{tr } \{ \mathbf{V} \mathbf{x} (\mathbf{t} - \mathbf{y})^\top \} = \text{tr } \{ \mathbf{V} (\nabla L(\mathbf{W}))^\top \}$$

or, more invariantly,

$$-\langle \mathbf{V}, (\mathbf{t} - \mathbf{y}) \mathbf{x}^\top \rangle = \langle \mathbf{V}, \nabla L(\mathbf{W}) \rangle.$$

Because  $\mathbf{V}$  is arbitrary:

$$\nabla L(\mathbf{W}) = -(\mathbf{t} - \mathbf{y}) \mathbf{x}^\top.$$

For arbitrary  $N$  we obtain the following gradient formula:

$$(9) \quad \nabla L(\mathbf{W}) = - \sum_{n=1}^N \left( \mathbf{t}^{(n)} - \mathbf{y}^{(n)} \right) \left( \mathbf{x}^{(n)} \right)^\top = -(\mathbf{T} - \mathbf{Y}) \mathbf{X}^\top.$$

where

$$\begin{aligned} \mathbf{T} &= [\mathbf{t}^{(1)} \quad \mathbf{t}^{(2)} \quad \dots \quad \mathbf{t}^{(N)}], \\ \mathbf{Y} &= [\mathbf{y}^{(1)} \quad \mathbf{y}^{(2)} \quad \dots \quad \mathbf{y}^{(N)}], \\ \mathbf{X} &= [\mathbf{x}^{(1)} \quad \mathbf{x}^{(2)} \quad \dots \quad \mathbf{x}^{(N)}], \end{aligned}$$

are matrices containing the elements of the training data in their columns. Thus, calculating the gradient can be expressed through simple matrix arithmetic, which leads to very efficient implementations of the gradient method.

**Corollary 1** (Characterization of critical points). *A weight matrix  $\mathbf{W}$  is a critical point of  $L$  iff the sample uncentered covariance matrix of the errors*

$$\mathbf{e}^{(n)} = \mathbf{t}^{(n)} - \mathbf{y}^{(n)}$$

*and of the input vectors  $\mathbf{x}^{(n)}$  is zero, where by definition the sample uncentered covariance matrix is the right-hand side of (9).*

Let us address the matter of uniqueness of the critical point. If we have two critical points  $\mathbf{W}$  and  $\tilde{\mathbf{W}}$ , the corresponding matrices  $\mathbf{Y}$  and  $\tilde{\mathbf{Y}}$  satisfy

$$0 = -(\mathbf{T} - \mathbf{Y}) \mathbf{X} = -(\mathbf{T} - \tilde{\mathbf{Y}}) \mathbf{X}$$

Thus

$$0 = (\mathbf{Y} - \tilde{\mathbf{Y}}) \mathbf{X}$$

Let  $\Delta \mathbf{Y} = \mathbf{Y} - \tilde{\mathbf{Y}}$ . Thus

$$\Delta \mathbf{Y} \mathbf{X}^\top = 0.$$

Hence

$$\mathbf{X} \Delta \mathbf{Y}^\top = 0.$$

Furthermore, every row of  $\mathbf{Y}$  is in the nullspace of  $\mathbf{X}$ . The matrix  $\mathbf{X}$  has the elements of the training set laid out in columns and it has  $N$  columns. If the rank of  $\mathbf{X}$  is maximal (equal to  $D$ , assuming  $N \gg D$ ) then the nullspace is trivial. Thus  $\mathbf{Y} = \tilde{\mathbf{Y}}$ . Let  $\mathbf{y} = \mathbf{y}^{(n)}$  and  $\tilde{\mathbf{y}} = \tilde{\mathbf{y}}^{(n)}$  for some  $n$ . Then  $\mathbf{y} = \tilde{\mathbf{y}}$ . Let  $\mathbf{w}_i$  and  $\tilde{\mathbf{w}}_i$  be the  $i$ -th rows of  $\mathbf{W}$  and  $\tilde{\mathbf{W}}$ . Let  $\mathbf{x} = \mathbf{x}^{(n)}$  be the corresponding element of the training set. Then

$$\frac{\exp(\mathbf{w}_i^\top \mathbf{x})}{\exp(\tilde{\mathbf{w}}_i^\top \mathbf{x})} = \exp(\Delta \mathbf{w}_i^\top \mathbf{x})$$

is independent of  $i$ , and equals

$$\frac{\sum_{j=1}^C \exp(\tilde{\mathbf{w}}_j^\top \mathbf{x})}{\sum_{j=1}^C \exp(\mathbf{w}_j^\top \mathbf{x})}.$$

Hence,

$$\Delta \mathbf{w}_i^\top \mathbf{x} = c$$

where  $c = c_n$  is a scalar depending on the element of the training set. Therefore

$$\Delta \mathbf{W} \mathbf{x} = c \mathbb{1}.$$

Combining all training set elements,

$$\Delta \mathbf{W} \mathbf{X} = \mathbb{1} \mathbf{c}^\top$$

where  $\mathbf{c} = (c_1, c_2, \dots, c_N)$ .

Let  $\tilde{\mathbf{X}}$  be a  $D \times D$  non-singular submatrix of  $\mathbf{X}$  and let  $\tilde{\mathbf{c}}$  be the corresponding vector of coefficients  $c_n$ . Then

$$\Delta \mathbf{W} \tilde{\mathbf{X}} = \mathbb{1} \tilde{\mathbf{c}}^\top$$

and

$$\Delta \mathbf{W} = \mathbb{1} \tilde{\mathbf{c}}^\top \tilde{\mathbf{X}}^{-1} = \mathbb{1} \tilde{\tilde{\mathbf{c}}}^\top.$$

Thus, all columns of  $\Delta \mathbf{W}$  must be multiples of  $\mathbb{1}$ , i.e. the weights  $\tilde{\mathbf{W}}$  are obtained from  $\mathbf{W}$  by a shift depending only on the element of the training set. Such a shift does not change the activations  $\mathbf{y}^{(n)}$  and thus produces the same value of the cross-entropy loss function. Thus, the weights differing by such a shift are equivalent in all respects.

One way to choose the matrix  $\mathbf{W}$  is to shift every column by subtracting the mean, which will have the effect of choosing the equivalent weight matrix with minimum Frobenius norm:

$$\sum_{j=1}^C \sum_{k=1}^D w_{jk}^2 = \min.$$

A practical application of the above considerations is that we can subtract the mean from columns as part of the iteration process, when implementing an optimization method, such as gradient descent, which should help with stability of the method.

#### 4. THE SECOND DERIVATIVE AND THE HESSIAN

In this section we consider the second (Fréchet) derivative of the loss function which is closely related to the Hessian. We recall the expression:

$$L = Q_{\mathbf{t}} \circ \rho \circ P_{\mathbf{x}}.$$

In view of linearity of the  $Q_{\mathbf{t}}$  and  $P_{\mathbf{x}}$ , we have:

$$D^2 L(\mathbf{W})(\mathbf{U}, \mathbf{V}) = Q_{\mathbf{t}} D^2 \rho(\mathbf{a})(P_{\mathbf{x}} \mathbf{U}, P_{\mathbf{x}} \mathbf{V})$$

where  $\mathbf{a} = P_{\mathbf{x}}(\mathbf{W}) = \mathbf{W} \mathbf{x}$ . We recall that the second derivative  $D^2 L(\mathbf{W})$  is a bi-linear, symmetric function on pairs of vectors  $\mathbf{U}$ ,  $\mathbf{W}$ , which in our case are also  $C \times D$  matrices.

We also have the formula for  $D\rho$ :

$$D\rho(\mathbf{u}) = -I + \mathbb{1} \sigma(\mathbf{u})^\top.$$

Differentiating again, we obtain

$$D^2\rho(\mathbf{a})(\mathbf{g}, \mathbf{h}) = \mathbb{1} (D\sigma(\mathbf{a}) \mathbf{g})^\top \mathbf{h} = \mathbb{1} \mathbf{g}^\top D\sigma(\mathbf{a})^\top \mathbf{h}.$$

Therefore,

$$Q_{\mathbf{t}} D^2\rho(\mathbf{a})(\mathbf{g}, \mathbf{h}) = (\mathbf{t}^\top \mathbb{1}) \mathbf{g}^\top D\sigma(\mathbf{a})^\top \mathbf{h} = \mathbf{g}^\top D\sigma(\mathbf{a})^\top \mathbf{h}$$

Hence,

$$\begin{aligned} D^2L(\mathbf{W})(\mathbf{U}, \mathbf{V}) &= (P_{\mathbf{x}} \mathbf{U})^\top D\sigma(\mathbf{a})^\top (P_{\mathbf{x}} \mathbf{V}) \\ &= (\mathbf{Ux})^\top D\sigma(\mathbf{a})^\top (\mathbf{Vx}) \\ &= \mathbf{x}^\top \mathbf{U}^\top D\sigma(\mathbf{a})^\top \mathbf{Vx}. \end{aligned}$$

We also find

$$\begin{aligned} D\sigma(\mathbf{a}) &= D \exp(-\rho(\mathbf{a})) \\ &= -\text{diag}(\exp(-\rho(\mathbf{a}))) D\rho(\mathbf{a}) \\ &= -\text{diag}(\sigma(\mathbf{a}))(-I + \mathbb{1} \sigma(\mathbf{a})^\top) \\ &= \text{diag}(\sigma(\mathbf{a})) - \text{diag}(\sigma(\mathbf{a})) \mathbb{1} \sigma(\mathbf{a})^\top. \end{aligned}$$

Hence,

$$D\sigma(\mathbf{a})^\top = \text{diag}(\sigma(\mathbf{a})) - \sigma(\mathbf{a}) \mathbb{1}^\top \text{diag}(\sigma(\mathbf{a})).$$

Therefore, taking into account that  $\text{diag}(\mathbf{s}) \mathbb{1} = \mathbf{s}$  and  $\mathbb{1}^\top \text{diag}(\mathbf{s}) = \mathbf{s}^\top$ , we obtain:

$$\begin{aligned} D\sigma(\mathbf{a})^\top &= \text{diag}(\sigma(\mathbf{a})) - \sigma(\mathbf{a}) \sigma(\mathbf{a})^\top. \\ D^2L(\mathbf{W})(\mathbf{U}, \mathbf{V}) &= \mathbf{x}^\top \mathbf{U}^\top \text{diag}(\sigma(\mathbf{a})) \mathbf{Vx} - \mathbf{x}^\top \mathbf{U}^\top \sigma(\mathbf{a}) \sigma(\mathbf{a})^\top \mathbf{Vx} \\ &= \mathbf{x}^\top \mathbf{U}^\top (\text{diag}(\sigma(\mathbf{a})) - \sigma(\mathbf{a}) \sigma(\mathbf{a})^\top) \mathbf{Vx} \\ &= \mathbf{x}^\top \mathbf{U}^\top (\text{diag}(\mathbf{y}) - \mathbf{y} \mathbf{y}^\top) \mathbf{Vx}. \end{aligned}$$

This bi-linear form is positive definite, as

$$D^2L(\mathbf{W})(\mathbf{U}, \mathbf{U}) = (\mathbf{Ux})^\top (\text{diag}(\mathbf{y}) - \mathbf{y} \mathbf{y}^\top) (\mathbf{Ux}) \geq 0.$$

Indeed, the matrix  $\text{diag}(\mathbf{y}) - \mathbf{y} \mathbf{y}^\top$  is symmetric and thus has real spectrum. It suffices to show that the spectrum is non-negative.

**Proposition 1.** *Let  $\mathbf{y} \in \mathbb{R}^C$  be a vector satisfying*

- (1)  $y_i > 0$ ;
- (2)  $\sum_{i=1}^C y_i = 1$ .

*Then the eigenvalues of the symmetric matrix*

$$\text{diag}(\mathbf{y}) - \mathbf{y} \mathbf{y}^\top$$

*are all non-negative. Furthermore, the only eigenvector with eigenvalue 0 is  $\mathbb{1}$  up to a scalar factor.*

*Proof. Method 1:* Let  $\lambda$  be an eigenvalue and  $\mathbf{z}$  be an eigenvector for eigenvalue  $\lambda$ ; thus

$$\text{diag}(\mathbf{y}) \mathbf{z} - \mathbf{y} (\mathbf{y}^\top \mathbf{z}) = \lambda \mathbf{z}.$$

Therefore,

$$\begin{aligned} y_i z_i - y_i \langle \mathbf{y}, \mathbf{z} \rangle &= \lambda z_i. \\ (y_i - \lambda) z_i &= y_i \langle \mathbf{y}, \mathbf{z} \rangle \end{aligned}$$

We know that  $0 < y_i < 1$ . Therefore, unless  $0 < \lambda < 1$ ,  $y_i \neq \lambda$ , and

$$z_i = \frac{y_i}{y_i - \lambda} \langle \mathbf{y}, \mathbf{z} \rangle.$$

Since  $\mathbf{z} \neq 0$ , we assume WLOG that  $\langle \mathbf{y}, \mathbf{z} \rangle \neq 0$ . Hence

$$\langle \mathbf{y}, \mathbf{z} \rangle = \sum_{i=1}^C \frac{y_i^2}{y_i - \lambda} \langle \mathbf{y}, \mathbf{z} \rangle.$$

and

$$\sum_{i=1}^C \frac{y_i^2}{y_i - \lambda} = 1.$$

We need to show that all roots  $\lambda$  of this equation are non-negative. Indeed, if  $\lambda < 0$  then

$$\sum_{i=1}^C \frac{y_i^2}{y_i - \lambda} < \sum_{i=1}^C \frac{y_i^2}{y_i} = \sum_{i=1}^C y_i = 1$$

which is a contradiction. Thus  $\lambda \geq 0$ .

It remains to see that the eigenvector  $\mathbf{z}$  for the eigenvalue  $\lambda = 0$  is  $\mathbf{z} = \mathbb{1}$  up to a multiplicative constant. Indeed, if  $\lambda = 0$  then for  $i = 1, 2, \dots, C$ :

$$y_i z_i = y_i \langle \mathbf{y}, \mathbf{z} \rangle.$$

Dividing by  $y_i > 0$ , we obtain:

$$z_i = \langle \mathbf{y}, \mathbf{z} \rangle.$$

Hence  $z_i$  is independent of  $i$ , i.e. proportional to  $\mathbb{1}$ .

**Method 2:** (Gershgorin Circle Theorem) The matrix in question has diagonal entry  $y_i - y_i^2$  at position  $i$ , and the off-diagonal entries in row  $i$  are  $y_i y_j$ ,  $j = 1, 2, \dots, C$ ,  $j \neq i$ . Hence, for every eigenvalue  $\lambda$  there exists  $i$  such that:

$$|\lambda - (y_i - y_i^2)| \leq \sum_{\substack{j=1 \\ j \neq i}}^C y_i y_j.$$

In particular

$$\lambda \geq (y_i - y_i^2) - \sum_{\substack{j=1 \\ j \neq i}}^C y_i y_j = y_i - y_i \sum_{j=1}^C y_j = y_i - y_i = 0.$$

**Method 3:** Let  $\mathbf{M} = \text{diag}(\sqrt{\mathbf{y}})$ . Obviously this is a symmetric diagonal matrix. Then we have the following factorization:

$$\text{diag}(\mathbf{y}) - \mathbf{y} \mathbf{y}^\top = \mathbf{M}^\top (I - \mathbf{u} \mathbf{u}^\top) \mathbf{M}$$

where  $\mathbf{u} = M^{-1} \mathbf{y} = \sqrt{\mathbf{y}}$ . Clearly,  $\sum_{i=1}^C u_i^2 = \sum_{i=1}^C y_i = 1$ . Thus  $\mathbf{u}$  is a unit vector. The matrix  $I - \mathbf{u} \mathbf{u}^\top$  is non-negative definite as it is an orthogonal projection on the hyperplane normal to  $\mathbf{u}$ , and therefore its eigenvalues are 0 (multiplicity 1) and 1 (multiplicity  $C - 1$ ). Moreover, the eigenvector with eigenvalue 0 is  $\mathbf{u}$ . Hence, the eigenvector  $\mathbf{z}$  with eigenvalue 0 for  $\text{diag}(\mathbf{y}) - \mathbf{y} \mathbf{y}^\top$  satisfies (up to a scalar multiple):

$$\mathbf{M} \mathbf{z} = \mathbf{u}.$$

and  $\mathbf{z} = \mathbf{M}^{-1} \mathbf{u} = \mathbb{1}$ . □



Let us consider arbitrary  $N$  (the size of the training set). We then have

$$D^2L(\mathbf{W})(\mathbf{U}, \mathbf{U}) = \sum_{n=1}^N \left( \mathbf{U} \mathbf{x}^{(n)} \right)^\top \left( \text{diag}(\mathbf{y}^{(n)} - \mathbf{y}^{(n)}) \left( \mathbf{y}^{(n)} \right)^\top \right) \left( \mathbf{U} \mathbf{x}^{(n)} \right) \geq 0.$$

The only way to get a value of 0 is for every summand to be 0, i.e. to have for  $n = 1, 2, \dots, N$ :

$$\mathbf{U} \mathbf{x}^{(n)} = c_n \mathbf{1}$$

for some scalars  $c_n$ ,  $n = 1, 2, \dots, N$ . These equations can also be written as a single matrix identity:

$$(10) \quad \mathbf{U} \mathbf{X} = \mathbf{1} \mathbf{c}^\top$$

where  $\mathbf{c} = (c_1, c_2, \dots, c_N)$ .

The degeneracy condition (10) is rather hard to fulfill for a randomly chosen training set. If the coordinates of  $\mathbf{x}^{(n)}$  are not linearly dependent as random variables then  $\mathbf{X}$  has maximum rank  $D$  with probability 1 if the variables are continuous, or with probability asymptotically converging to 1 if the variables are discrete (or mixed). Let  $\tilde{\mathbf{X}}$  be a  $D \times D$  non-singular submatrix of  $\mathbf{X}$ . Then

$$\mathbf{U} = \mathbf{1} \mathbf{c}^\top \tilde{\mathbf{X}}^{-1} = \mathbf{1} \tilde{\mathbf{c}}^\top.$$

Thus, all columns of  $\mathbf{U}$  must be multiples of  $\mathbf{1}$ . We are already familiar with this condition. If we restrict  $L$  to only matrices  $\mathbf{W}$  for which the mean of every column is 0 then the second derivative  $D^2L(\mathbf{W})$  is strictly-positive definite, thus ensuring convergence of optimization methods, such as gradient descent.

**Corollary 2** (On convergence and regularization). *The loss function  $L$  for the multi-class logistic regression is strictly convex on the set of weight matrices  $\mathbf{W}$  with column mean zero, as long as the training set spans the vector space  $\mathbb{R}^D$ , or, equivalently, when  $\text{rank } \mathbf{X} = D$ . Therefore, it is possible to find the optimal weight without using a regularizer.*

## 5. A SAMPLE TRAINING ALGORITHM

Algorithm 1 contains the basic training loop for the logistic regression network we described in previous sections. Many modifications are possible. For instance, one can implement variable training rate, utilizing, for example, the Barzilai-Bowdoin update rule. Also, we can break out of the loop if no progress is made. We could evaluate the loss function and see if it decays, as a test of progress, etc.

It should be noted that

$$\mathbf{1}^\top \mathbf{T} = \mathbf{1}^\top \mathbf{Y} = 1$$

and therefore

$$\mathbf{1}^\top (\mathbf{T} - \mathbf{Y})^\top = 0$$

Hence, in view of  $\nabla L(\mathbf{W}) = (\mathbf{T} - \mathbf{Y}) \mathbf{X}^\top$ ,

$$\mathbf{1}^\top \nabla L(\mathbf{W}) = 0$$

and the condition

$$\mathbf{1}^\top \mathbf{W} = 0$$

is maintained during the iteration process, at least if round-off error is ignored. However, we may want to subtract the mean from each column of  $\mathbf{W}$  every so often to prevent round-off error acting as diffusion in the weight space, which clearly prevents convergence.

---

**Algorithm 1** The algorithm implements training of the logistic regression network by gradient descent method.

---

**Require:**

$\mathbf{X}$  is an  $D \times N$  matrix of rank  $D$ , containing the  $N$  training vectors as columns;

$\mathbf{W}$  is a  $C \times D$  weight matrix initialized at random;

$\eta \in \mathbb{R}$  is the learning rate;

$NumEpochs$  is the number of epochs;

**Ensure:**

$\mathbf{W}$  approximates the optimal weights minimizing the cross-entropy loss function.

**for**  $epoch = 1, 2, \dots, NumEpochs$  **do**

$\mathbf{A} \leftarrow \mathbf{W} \cdot \mathbf{X}$

    ▷ Compute activations.

$\mathbf{Y} \leftarrow \text{SOFTMAX}(\mathbf{A})$

    ▷ Compute softmax activity.

$\mathbf{E} \leftarrow \mathbf{T} - \mathbf{Y}$

    ▷ Compute errors.

$\nabla L \leftarrow -\mathbf{E} \cdot \mathbf{X}^\top$

    ▷ Find the gradient.

$\mathbf{W} \leftarrow \mathbf{W} - \eta \cdot \nabla L$

    ▷ Update weights.

**end for**

---

## REFERENCES

- [1] Christopher M. Bishop. *Neural networks for pattern recognition*. Oxford University Press, Inc. New York, NY, USA, 1996.
- [2] Logistic regression. [https://en.wikipedia.org/wiki/Logistic\\_regression](https://en.wikipedia.org/wiki/Logistic_regression).
- [3] Multinomial logistic regression. [https://en.wikipedia.org/wiki/Multinomial\\_logistic\\_regression](https://en.wikipedia.org/wiki/Multinomial_logistic_regression).