

# Multi Layered Perceptron

Parakram Majumdar

May 28, 2018

## 1 Introduction

We are given a training set  $T$  of  $t$  points in  $\mathbb{R}^d$ , a set  $L$  of  $l$  labels, and a *classification*  $C : T \rightarrow L$ .

Machine learning is interested in learning classifiers of the form  $f : \mathbb{R}^d \rightarrow L$  that can closely approximate  $C$ .

## 2 Single Perceptron, Single Point, Single Label

Having truly one label is an uninteresting case, as basically all points would have the same label, and there is nothing to “learn”. Instead we assume two labels,  $L_1$  and  $L_2$ . We then train a classifier for  $L_1$ , and assign to  $L_2$  whatever is not in  $L_1$ .

Given a single point  $p \in \mathbb{R}^d$ , a single perceptron  $c$  with *weight* vector  $w_c \in \mathbb{R}^d$ , and a *bias*  $b_c \in \mathbb{R}$ , for classifying a single Label  $L_1$ , accepts  $p$  with the probability

$$\hat{y}_{p,L_1} = \hat{\mathbb{P}}[p \in L_1] = \text{sigmoid}(w_c \cdot p + b_c)$$

where

$$\text{sigmoid}(x) = \frac{e^x}{e^x + 1}$$

If the actual classification of the training point is given by  $y_{p,L_1} = \mathbb{I}[p \in L_1]$  then training criterion is to minimize the cross entropy:

$$\begin{aligned} \mathbb{H} &= \mathbb{H}[\mathbb{P}, \hat{\mathbb{P}}] \\ &= \mathbb{E}\left[\log \frac{1}{\hat{y}}\right] \\ &= \sum_{L_i \in L} \mathbb{H}[y_{p,L_i}, \hat{y}_{p,L_i}] \\ &= - \sum_{L_i \in L} y_{p,L_i} \log \hat{y}_{p,L_i} \\ &= -y_{p,L_1} \log(\hat{y}_{p,L_1}) - (1 - y_{p,L_1}) \log(1 - \hat{y}_{p,L_1}) \end{aligned}$$

### 3 Multiple Points

Given a set of  $t$  points represented as a 2 dimensional tensor  $T \in \mathbb{R}^{t \times d}$ , the perceptron  $c$  can then be *applied* to all the points in  $T$  to get

$$\hat{y}_{T,c} = \text{map}(T \cdot w_c + b_c \vec{u}_t, \text{sigmoid}(\cdot))$$

where

- $\text{map}(\vec{x}, f)$  is the element-wise application of  $f$  on  $\vec{x}$
- $\vec{u}_n = [1, 1, 1, \dots n \text{ times}]$

The  $i^{th}$  element of  $\hat{y}_{T,c}$  then gives us the probability of the  $i^{th}$  point being in  $L_1$ .

$$\hat{y}_{T,c}[i] = \hat{\mathbb{P}}[T[i] \in L_1] = \hat{y}_{T[i],c}$$

And therefore, if the actual classification of all training points is given by

$$y_{T,c} = \text{map}(T, \mathbb{I}[\cdot \in L_1])$$

then the training criterion is to minimize the cross entropy

$$\begin{aligned} \mathbb{H} &= \sum_{i=1}^t \mathbb{H}[y_{T,c}[i], \hat{y}_{T,c}[i]] \\ &= \text{sum}(\text{map}(y_{T,c}, \hat{y}_{T,c}, \mathbb{H}[\cdot, \cdot])) \end{aligned}$$

where

- $\text{sum}(x)$  gives the sum of all the elements in the tensor  $x$
- $\text{map}(x, y, f(\cdot, \cdot))$  is the element-wise application of the binary function  $f$  to the elements of  $x$  and  $y$ . In other words,

$$\text{map}(x, y, f(\cdot, \cdot))[i] = f(x[i], y[i])$$

### 4 Multiple Labels

Classifying a point across multiple labels may be achieved using multiple perceptrons. However, if we simply use the multiple perceptrons independently, we end up with unconstrained probabilities for a point belonging to each label. Instead, we would ideally want the probabilities to add up to 1, since a point can only belong to one label.

Suppose we have  $l$  labels, and corresponding  $l$  perceptrons. Each perceptron produces an output  $x_i \in \mathbb{R}^t$ , where  $t$  is the number of points. Then the probability of point  $T_j$  belonging to label  $L_k$  is

$$\hat{y}_{T_j, L_k} = \frac{e^{x_k[j]}}{\sum_{m=1}^l e^{x_m[j]}}$$

To simplify this notation, we introduce the softmax function:

$$\begin{aligned}\text{softmax}(\vec{z}) &= \text{softmax}(z_1, z_2, \dots, z_n) \\ &= \left[ \frac{e^{z_1}}{\sum_{i=1}^n e^{z_i}}, \frac{e^{z_2}}{\sum_{i=1}^n e^{z_i}}, \dots, \frac{e^{z_n}}{\sum_{i=1}^n e^{z_i}} \right] \\ \text{softmax}(\vec{z}; j) &= \frac{e^{z_j}}{\sum_{i=1}^n e^{z_i}}\end{aligned}$$

And concisely represent all the probabilities as:

$$\hat{y}_{T,L} = \text{map}(x_1, x_2, \dots, x_l, \text{softmax})$$

Therefore, the training criteria is to minimize the cross-entropy:

$$\mathbb{H}(\hat{y}_{T,L}, y_{T,L}) = \sum_{T_i \in T, L_j \in L} \mathbb{H}(\hat{y}_{T_i, L_j}, y_{T_i, L_j})$$

## 5 Multiple Layers of Perceptrons

Multiple layers of perceptrons, with non-linear activation functions, are necessary for classifying data that is not linearly separable. An activation function is applied to the output of a neuron before feeding it to the next layer.

## 6 Sample 2-Layer Perceptron

Let the input data be organized into a 2-dimensional structure  $T \in \mathbb{R}^{t \times d}$ , where  $t$  is the number of points and  $d$  is the dimensionality of each point.

We feed these points to the first layer consisting of, say,  $m$  perceptrons. Thus, we have  $m$  weight vectors:

$$w_{1,1}, w_{1,2}, \dots, w_{1,m} \in \mathbb{R}^d$$

and  $m$  biases:

$$b_{1,1}, b_{1,2}, \dots, b_{1,m} \in \mathbb{R}$$

to give the unactivated output  $p_{1,i} \in \mathbb{R}^t$  of the  $i^{th}$  perceptron:

$$p_{1,i} = T \cdot w_{1,i} + b_{1,i} \vec{u}_t$$

These unactivated outputs  $p_{1,\cdot}$  are passed through a sigmoid activation function to get the final activated outputs  $z_{1,\cdot} \in \mathbb{R}^t$  of this layer:

$$z_{1,i} = \text{map}(p_{1,i}, \text{sigmoid}(\cdot))$$

These outputs  $z_{1,\cdot}$  then become the inputs for the next and final layer. For convenience, we define  $z_1 \in \mathbb{R}^{t \times m}$  as the output vector of the first layer, created by concatenating all the individual outputs.

Since we have  $l$  labels for classification, the final layer should have  $l$  perceptrons. Thus, we have  $l$  weight vectors:

$$w_{2,1}, w_{2,2}, \dots, w_{2,l} \in R^m$$

and  $l$  biases:

$$b_{2,1}, b_{2,2}, \dots, b_{2,l} \in R$$

to give the unactivated outputs  $p_{2,\cdot} \in \mathbb{R}^t$ :

$$p_{2,i} = z_1 \cdot w_{2,i} + b_{2,i} \vec{u}_t$$

We leave the outputs unactivated for this final layer, and concatenate the outputs to get  $p_2 \in R^{t \times l}$ . Thus, the probabilities  $\hat{y} \in \mathbb{R}^{t \times l}$  of the points belonging to classes is given by:

$$\hat{y} = \text{map}(p_2, \text{softmax}(\cdot))$$

If the real probabilities are given by  $y \in R^{t \times l}$  then the training criterion is to minimize the cross entropy:

$$\mathbb{H} = \mathbb{H}[\hat{y}, y] = \sum_{i \in T, j \in L} \mathbb{H}[\hat{y}_{i,j}, y_{i,j}]$$

