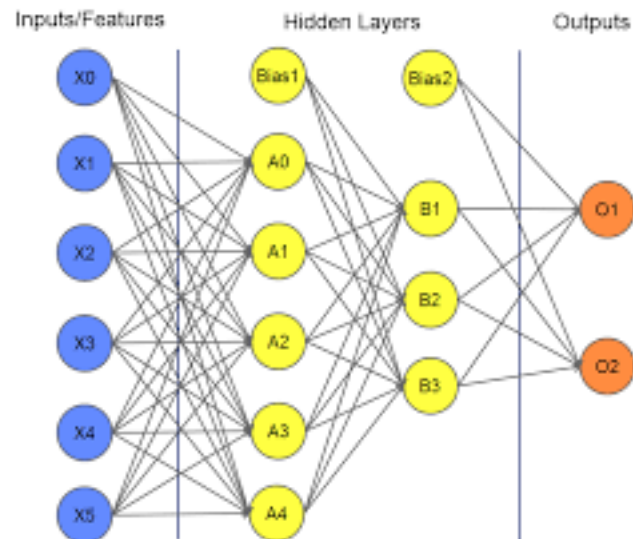


Perceptron

Multi Layer Perceptron



Learning with Back Propagation Algo :

Notations

w_{ij}^k : weight for node j in layer l_k for incoming node i

b_i^k : bias for node i in layer l_k

a_i^k : product sum plus bias (activation) for node i in layer l_k

o_i^k : output for node i in layer l_k

r_k : number of nodes in layer l_k

g : activation function for the hidden layer nodes

g_o : activation function for the output layer nodes

Activation function

$$a_i^k = b_i^k + \sum_{j=1}^{r_{k-1}} w_{ji}^k o_j^{k-1} = \sum_{j=0}^{r_{k-1}} w_{ji}^k o_j^{k-1}$$

Backpropagation attempts to minimize

$$E(X, \theta) = \frac{1}{2N} \sum_{i=1}^N (\hat{y}_i - y_i)^2$$

with respect to the neural network's weights:

→ for each weight w_{ij}^k , evaluate $\frac{\partial E}{\partial w_{ij}^k}$. By decomposing into a sum over individual error terms for each individual input-output pair

$$\frac{\partial E(X, \theta)}{\partial w_{ij}^k} = \frac{1}{N} \sum_{d=1}^N \frac{\partial}{\partial w_{ij}^k} \left(\frac{1}{2} (\hat{y}_d - y_d)^2 \right) = \frac{1}{N} \sum_{d=1}^N \frac{\partial E_d}{\partial w_{ij}^k}.$$

Error function derivatives (index d is omitted below)

Apply the chain rule to the error function partial derivative

$$\frac{\partial E}{\partial w_{ij}^k} = \frac{\partial E}{\partial a_j^k} \frac{\partial a_j^k}{\partial w_{ij}^k},$$

$$\delta_j^k \equiv \frac{\partial E}{\partial a_j^k}.$$

Then

$$\frac{\partial E}{\partial w_{ij}^k} = \delta_j^k o_i^{k-1}.$$

$$\frac{\partial a_j^k}{\partial w_{ij}^k} = \frac{\partial}{\partial w_{ij}^k} \left(\sum_{l=0}^{r_k-1} w_{lj}^k o_l^{k-1} \right) = o_i^{k-1}.$$

Output Layer (MLP with $M + 1$ layers)

Assume a one-output neural network, so there is only one output node $j = 1$). Expressing E in terms of the value a_1^m (since δ_1^m is a partial derivative with respect to a_1^m) gives

$$E = \frac{1}{2} (\hat{y} - y)^2 = \frac{1}{2} (g_o(a_1^m) - y)^2,$$

where $g_o(x)$ is the activation function for the output layer.

thus,

$$\delta_1^m = (g_o(a_1^m) - y) g'_o(a_1^m) = (\hat{y} - y) g'_o(a_1^m).$$

and finally

$$\frac{\partial E}{\partial w_{i1}^m} = \delta_1^m o_i^{m-1} = (\hat{y} - y) g'_o(a_1^m) o_i^{m-1}.$$

The Hidden Layers

One has for the error term δ_j^k in layer $1 \leq k < m$:

$$\delta_j^k = \frac{\partial E}{\partial a_j^k} = \sum_{l=1}^{r^{k+1}} \frac{\partial E}{\partial a_l^{k+1}} \frac{\partial a_l^{k+1}}{\partial a_j^k},$$

where l ranges from 1 to r^{k+1}

(The bias input o_0^k corresponds to w_{0j}^{k+1} is fixed, does not depend on the outputs of previous layers, thus l does not take on the value 0.)

Plugging in the error term :

$$\delta_j^k = \sum_{l=1}^{r^{k+1}} \delta_l^{k+1} \frac{\partial a_l^{k+1}}{\partial a_j^k}.$$

From the definition of a_l^{k+1}

$$a_l^{k+1} = \sum_{j=1}^{r^k} w_{jl}^{k+1} g(a_j^k),$$

where $g(x)$ is the activation function for the hidden layers,

$$\frac{\partial a_l^{k+1}}{\partial a_j^k} = w_{jl}^{k+1} g'(a_j^k).$$

Plugging this into the latter equation yields a final equation for the error term δ_j^k in the hidden layers, called the *backpropagation* formula:

$$\delta_j^k = \sum_{l=1}^{r^{k+1}} \delta_l^{k+1} w_{jl}^{k+1} g'(a_j^k) = g'(a_j^k) \sum_{l=1}^{r^{k+1}} w_{jl}^{k+1} \delta_l^{k+1}.$$

And putting all equations together :

$$\frac{\partial E}{\partial w_{ij}^k} = \delta_j^k o_i^{k-1} = g'(a_j^k) o_i^{k-1} \sum_{l=1}^{r^{k+1}} w_{jl}^{k+1} \delta_l^{k+1}.$$

BACKPROPAGATION ALGO

Main equations

For the partial derivatives,

$$\frac{\partial E_d}{\partial w_{ij}^k} = \delta_j^k o_i^{k-1}.$$

For the final layer's error term,

$$\delta_1^m = g'_o(a_1^m) (\hat{y}_d - y_d).$$

For the hidden layers' error terms,

$$\delta_j^k = g'(a_j^k) \sum_{l=1}^{r^{k+1}} w_{jl}^{k+1} \delta_l^{k+1}.$$

For combining the partial derivatives for each input-output pair,

$$\frac{\partial E(X, \theta)}{\partial w_{ij}^k} = \frac{1}{N} \sum_{d=1}^N \frac{\partial}{\partial w_{ij}^k} \left(\frac{1}{2} (\hat{y}_d - y_d)^2 \right) = \frac{1}{N} \sum_{d=1}^N \frac{\partial E_d}{\partial w_{ij}^k}.$$

For updating the weights,

$$\Delta w_{ij}^k = -\alpha \frac{\partial E(X, \theta)}{\partial w_{ij}^k}.$$

The General algorithm

- 1) Calculate the **forward phase** for each (\vec{x}_d, y_d) ; store the results \hat{y}_d , a_j^k , and o_j^k for each node j in layer k by proceeding from layer 0, to layer m , the output layer.
- 2) Calculate the **backward phase** for each (\vec{x}_d, y_d) ; store $\frac{\partial E_d}{\partial w_{ij}^k}$ for each weight w_{ij}^k . Proceedi from output layer m , to layer 1, the input layer.
 - a) Evaluate δ_1^m
 - b) Backpropagate the error terms for the hidden layers δ_j^k , working backwards
 - c) Evaluate the partial derivatives of the individual error E_d with respect to w_{ij}^k
- 3) Combine the individual gradients for each input-output pair to get the total gradient for the entire set $X = \{(\vec{x}_1, y_1), \dots, (\vec{x}_N, y_N)\}$ (a simple average of the individual gradients).
- 4) Update the weights according to the learning rate α and total gradient $\frac{\partial E(X, \theta)}{\partial w_{ij}^k}$

In the classic formulation,
for hidden nodes ($g(x) = \sigma(x)$) (sigmoidal function)
and the output activation function is ($g_o(x) = x$)

$$g'(x) = \frac{\partial \sigma(x)}{\partial x} = \sigma(x)(1 - \sigma(x)).$$

$$g'_o(x) = \frac{\partial g_o(x)}{\partial x} = \frac{\partial x}{\partial x} = 1.$$

→ No need to remember the activation values a_1^m and a_j^k in addition to the output values o_1^m and o_j^k , greatly reducing the memory footprint of the algorithm.

BUT gradient descent algorithm may be infeasible when the training data size is huge.
Thus, a stochastic version of the algorithm is often used instead.

Stochastic Gradient method

At each iteration, rather than computing

$$\nabla_{\theta} E(\mathbf{X}) = \nabla_{\theta} \left(\sum_{d=1}^N E(x_d) \right) = \sum_{d=1}^N \nabla_{\theta} E(x_d)$$

stochastic gradient descent randomly samples d at uniform and computes $\nabla_{\theta} E(x_d)$ instead :

SGD uses $\nabla E(x_d)$ as an unbiased estimator of $\nabla E(\mathbf{X})$;

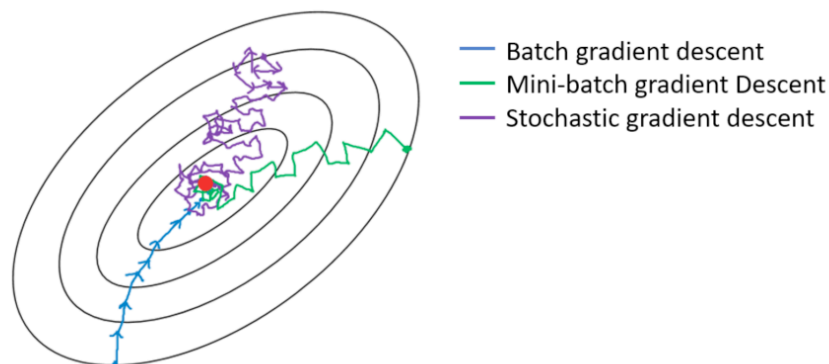
$$\mathbb{E} [\nabla_{\theta} E(x_d)] = \mathbb{E} \left[\frac{1}{k} \sum_{i=1}^k \nabla E(x_k) \right] = \nabla E(\mathbf{X}).$$

In a generalized case, at each iteration a mini-batch \mathcal{B} that consists of indices for training data instances may be sampled at uniform with replacement.

$$\nabla f E_{\mathcal{B}}(\mathbf{X}) = \frac{1}{|\mathcal{B}|} \sum_{d \in \mathcal{B}} \nabla E(x_d)$$

update θ as

$$\theta := \theta - \eta \nabla E_{\mathcal{B}}(\mathbf{X})$$



The per-iteration computational cost is $\mathcal{O}(|\mathcal{B}|)$. Thus, when the mini-batch size is small, the computational cost at each iteration is light.

If the training data set has many redundant data instances, stochastic gradients may be so close to the true gradient $\nabla_{\theta} E(\mathbf{X})$ that a small number of iterations will find useful solutions to the optimization problem

Stochastic gradient descent can be considered as offering a regularization effect especially when the mini-batch size is small due to the randomness and noise in the mini-batch sampling

Certain hardware processes mini-batches of specific sizes more efficiently.