

An incremental linear perceptron based on the BFGS algorithm

Flávio Eler De Melo

April 5, 2019

Abstract

This document briefly describes a new procedure for linear incremental learning that minimises the least-square residues recursively based on the Broyden–Fletcher–Goldfarb–Shanno (BFGS) algorithm.

1 Least squares regressor

Let X and Y be the input and output matrices with dimensions $n_x \times m$ and $n_y \times m$ respectively, where m is the number of examples. The least squares problem requires to find the weight matrix W (of dimension $n_x \times n_y$) that minimizes the loss:

$$\mathcal{L}(X, Y, W) = \frac{1}{2} \sum_{j=1}^{n_y} \sum_{k=1}^m \left(y_{jk} - \sum_{\ell=1}^{n_x} w_{\ell j} x_{\ell k} \right)^2. \quad (1.1)$$

Note that:

$$\begin{aligned} \nabla_W \mathcal{L}(X, Y, W) &= [\nabla_W \mathcal{L}(X, Y, W)]_{ij} = - \sum_{k=1}^m x_{ik} \left(y_{jk} - \sum_{\ell=1}^{n_x} w_{\ell j} x_{\ell k} \right) \\ &= -X (Y^T - X^T W) \equiv \mathbf{0}_{n_x \times n_y}. \\ XX^T W &= XY^T, \\ HW &= XY^T, \\ W &= H^{-1} XY^T, \end{aligned}$$

where $H = H^T = XX^T$ (symmetric). Therefore, a(n) (initial) solution can be calculated for the first collected examples (X_0, Y_0) by:

$$W_0 = H_0^{-1} X_0 Y_0^T, \quad (1.2)$$

where $H_0 = X_0 X_0^T$, $X_0 = [x_{0,ij}]_{i \in [1..n_x], j \in [1..m]}$ at step $k = 0$ for n_x dimensions (including the constant term) of the independent variable vector and m training examples, $Y_0 = [y_{0,ij}]_{i \in [1..n_y], j \in [1..m]}$ at step $k = 0$ for n_y dimensions of the dependent variable vector.

2 Incremental update

Now suppose that we will search for an incremental minimisation of $\mathcal{L}(X_k, Y_k, W_k)$. This can be done by finding the direction p_k of minimisation via the (quasi) Newton equation:

$$H_k p_k = -\nabla_W \mathcal{L}(X_k, Y_k, W_k) = X_k (Y_k^T - X_k^T W_k). \quad (2.1)$$

Note that for $H_0 = X_0 X_0^T$,

$$\begin{aligned} p_0 &= H_0^{-1} X_0 (Y_0^T - X_0^T W_0) \\ &= H_0^{-1} X_0 Y_0^T - H_0^{-1} X_0 X_0^T W_0 \\ &= W_0 - W_0 = \mathbf{0}, \end{aligned}$$

i.e., the loss functional is already at a minimum. By assumption, H_{k-1}^{-1} and W_{k-1} are parameters kept to induce the incremental regression. Therefore, given new examples (X_k, Y_k) , our incremental task involves inducing the procedure by finding the incremental direction via

$$p_k = -H_{k-1}^{-1} \nabla_W \mathcal{L}(X_k, Y_k, W_{k-1}) = H_{k-1}^{-1} X_k (Y_k^T - X_k^T W_{k-1}). \quad (2.2)$$

In the context of the Broyden–Fletcher–Goldfarb–Shanno (BFGS) algorithm, the problem requires line search to find a proper incremental factor as:

$$\alpha_k = \operatorname{argmin}_{\alpha} \mathcal{L}(X_k, Y_k, W_k(\alpha)), \quad (2.3)$$

where $W_k(\alpha) = W_{k-1} + \alpha p_k$, which can be computed analytically by

$$\begin{aligned} \nabla_{\alpha} \mathcal{L}(X_k, Y_k, W_k(\alpha)) &= \nabla_W^T \mathcal{L}(X_k, Y_k, W_k) \nabla_{\alpha} W_k(\alpha) \\ &= -(Y_k - W_k^T X_k) X_k^T p_k \equiv \mathbf{0}_{n_y \times n_y}, \\ \alpha_k p_k^T X_k X_k^T p_k &= (Y_k - W_{k-1}^T X_k) X_k^T p_k, \\ \alpha_k &= (Y_k - W_{k-1}^T X_k) X_k^T p_k (p_k^T X_k X_k^T p_k)^{-1}. \end{aligned} \quad (2.4)$$

However, instead we will define a learning rate α_k proportional to a value of reference α_r (a preset learning rate) that balances the effectiveness of negatives and positives for the regression. For instance, given $m = p + n$ examples, with p positives and n negatives, $\alpha_{p,k} = \alpha_r n/m$ and $\alpha_{n,k} = \alpha_r p/m$ so that $\alpha_{p,k} + \alpha_{n,k} = \alpha_r$. Note that for $p > n$ the negative examples will have a higher factor to compensate for their rarer (and so weaker) contribution whereas for $n > p$ the positive examples will have a higher factor.

Now, for each example, compute the weights increment and the loss gradient increment by:

$$\delta W_k = \alpha_k p_k, \quad (2.5)$$

$$\begin{aligned} z_k &= \nabla_W \mathcal{L}(X_k, Y_k, W_{k-1} + \delta W_k) - \nabla_W \mathcal{L}(X_k, Y_k, W_{k-1}) \\ &= -X_k (Y_k^T - X_k^T W_{k-1} - X_k^T \alpha_k p_k) + X_k (Y_k^T - X_k^T W_{k-1}) \\ &= X_k X_k^T \alpha_k p_k = X_k X_k^T \delta W_k. \end{aligned} \quad (2.6)$$

And finally, update the weights and the inverse Hessian matrix by the BFGS algorithm and the Sherman-Morrison formula by:

$$\mathbf{W}_k = \mathbf{W}_{k-1} + \delta \mathbf{W}_k, \quad (2.7)$$

$$\mathbf{H}_k^{-1} = \left(\mathbb{I} - \frac{\delta \mathbf{W}_k \mathbf{z}_k^T}{\mathbf{z}_k^T \delta \mathbf{W}_k} \right) \mathbf{H}_{k-1}^{-1} \left(\mathbb{I} - \frac{\mathbf{z}_k \delta \mathbf{W}_k^T}{\mathbf{z}_k^T \delta \mathbf{W}_k} \right) + \frac{\delta \mathbf{W}_k \delta \mathbf{W}_k^T}{\mathbf{z}_k^T \delta \mathbf{W}_k}, \quad (2.8)$$

where:

$$\delta \mathbf{W}_k \mathbf{z}_k^T = \delta \mathbf{W}_k (\mathbf{X}_k \mathbf{X}_k^T \delta \mathbf{W}_k)^T = \delta \mathbf{W}_k (\mathbf{X}_k^T \delta \mathbf{W}_k)^T \mathbf{X}_k^T = \delta \mathbf{W}_k \delta \mathbf{W}_k^T \mathbf{X}_k \mathbf{X}_k^T, \quad (2.9)$$

$$\mathbf{z}_k \delta \mathbf{W}_k^T = \mathbf{X}_k \mathbf{X}_k^T \delta \mathbf{W}_k \delta \mathbf{W}_k^T = (\delta \mathbf{W}_k \mathbf{z}_k^T)^T, \quad (2.10)$$

$$\mathbf{z}_k^T \delta \mathbf{W}_k = (\mathbf{X}_k \mathbf{X}_k^T \delta \mathbf{W}_k)^T \delta \mathbf{W}_k = \delta \mathbf{W}_k^T \mathbf{X}_k \mathbf{X}_k^T \delta \mathbf{W}_k. \quad (2.11)$$

The procedure is repeated until all examples have been taken into account.