

Oja's rule: Derivation, Properties *

October 2019

1 Derivation of Oja's rule

The update for each weight of the weight vector $\mathbf{w} = [w_1, \dots, w_D]^T \in \mathbb{R}^D$ for Oja's rule [1, 2] reads:

$$w_i^{n+1} = \frac{w_i^n + \beta y^n x_i^n}{\sqrt{\sum_{i=0}^{D-1} (w_i^n + \beta y^n x_i^n)^2}} \quad (1)$$

where the index n denotes the iteration number, while D is the dimension of the data vector, β is the learning rate, and i is the neuron number. In vector notation

$$\mathbf{w}^{n+1} = \frac{\mathbf{w}^n + \beta y^n \mathbf{x}^n}{\|\mathbf{w}^n + \beta y^n \mathbf{x}^n\|_2} \quad (2)$$

The gradient of expression (1) with respect to the learning rate β is given by:

$$\frac{\partial w_i^{n+1}}{\partial \beta} = \frac{y^n x_i^n \sqrt{\sum_{i=0}^{D-1} (w_i^n + \beta y^n x_i^n)^2} - \frac{1}{2} \left(\sum_{i=0}^{D-1} (w_i^n + \beta y^n x_i^n)^2 \right)^{-\frac{1}{2}} \left(\sum_{i=0}^{D-1} 2y^n x_i^n (w_i^n + \beta y^n x_i^n) \right) (w_i^n + \beta y^n x_i^n)}{\sum_{i=0}^{D-1} (w_i^n + \beta y^n x_i^n)^2} \quad (3)$$

where we used the quotient rule $(f/g)' = (f'g - g'f)/g^2$ and that $(f(g(x)))' = f'(g(x))g'(x)$. By simplifying the expression 3, we get:

$$\frac{\partial w_i^{n+1}}{\partial \beta} = \frac{y^n x_i^n}{\sqrt{\sum_{i=0}^{D-1} (w_i^n + \beta y^n x_i^n)^2}} - \frac{\left(\sum_{i=0}^{D-1} y^n x_i^n (w_i^n + \beta y^n x_i^n) \right) (w_i^n + \beta y^n x_i^n)}{\left(\sum_{i=0}^{D-1} (w_i^n + \beta y^n x_i^n) \right)^{\frac{3}{2}}} \quad (4)$$

Next, we evaluate the derivative at the linearization point $\beta^* = 0$ to get:

$$\left. \frac{\partial w_i^{n+1}}{\partial \beta} \right|_{\beta^*=0} = \frac{y^n x_i^n}{\sqrt{\sum_{i=0}^{D-1} (w_i^n)^2}} - \frac{\left(\sum_{i=0}^{D-1} y^n x_i^n (w_i^n) \right) (w_i^n)}{\left(\sum_{i=0}^{D-1} (w_i^n) \right)^{\frac{3}{2}}} \quad (5)$$

However, since we assume that the weights are normalized to one, we get $\|\mathbf{w}^n\|_w = \sum_{i=0}^{D-1} (w_i^n)^2$, that simplifies equation (5) to:

$$\left. \frac{\partial w_i^{n+1}}{\partial \beta} \right|_{\beta^*=0} = y^n x_i^n - w_i^n \sum_{i=0}^{D-1} y^n x_i^n w_i^n = y^n \left(x_i^n - w_i^n \sum_{i=0}^{D-1} x_i^n w_i^n \right) = y^n \left(x_i^n - w_i^n y^n \right) \quad (6)$$

In vector notation:

$$\left. \frac{\partial \mathbf{w}^{n+1}}{\partial \beta} \right|_{\beta^*=0} = y^n \left(\mathbf{x}^n - \mathbf{w}^n y^n \right) \quad (7)$$

This expression is the same as Oja's rule. Now, if we write down the Taylor series of \mathbf{w}^{n+1} , we get:

$$\mathbf{w}^{n+1} = f(\beta) = f(\beta^*) + \frac{f'(\beta^*)}{1!} (\beta - \beta^*) + \frac{f''(\beta^*)}{2!} (\beta - \beta^*)^2 + \dots \quad (8)$$

Since we are interested in a local approximation of f , terms in the order $\mathcal{O}(\beta^2)$ can be ignored, since:

$$(\beta - \beta^*) \approx 10^{-5} \implies (\beta - \beta^*)^2 \approx 10^{-10}. \quad (9)$$

*For typos/questions please contact: pvlachas@ethz.ch

As a consequence, we can write:

$$\mathbf{w}^{n+1} = f(\beta^*) + f'(\beta^*)\beta + \mathcal{O}(\beta^2) \quad (10)$$

with $\beta^* = 0$ and ignore the $\mathcal{O}(\beta^2)$ terms. Since $f(\beta^*) = \mathbf{w}^n$, we get:

$$\mathbf{w}^{n+1} \approx \mathbf{w}^n + \beta y^n (\mathbf{x}^n - \mathbf{w}^n y^n) \quad (11)$$

deriving Oja's rule as a linearization of the normalized Hebb's rule around $\beta^* = 0$.

2 Proof of Oja's rule properties

In this section, we are going to prove some interesting properties of Oja's rule. Namely, if $\lim_{n \rightarrow \infty} \mathbf{w} = \mathbf{w}^*$:

1. If C is the covariance matrix of the data $C = \mathbb{E}[\mathbf{x}\mathbf{x}^T]$, with $C \in \mathbb{R}^{D \times D}$, it holds that \mathbf{w}^* is an eigenvalue of C , i.e. $C\mathbf{w}^* = \lambda^*\mathbf{w}^*$. (this shows that \mathbf{w}^* is a PCA component).
2. $|\mathbf{w}^*|_2 = 1$
3. \mathbf{w}^* is the principal component of the data, i.e. $\lambda^* \geq \lambda_i, \forall$ other PCA components i .
4. \mathbf{w}^* maximizes the variance $\mathbb{E}[y^2]$ at the output of the perceptron over all \mathbf{w} with $|\mathbf{w}|_2 = 1$.

3 Item 1: Eigenvector of C (PCA component of data)

We are going to prove these statements, one by one. Firstly, we need to show that $|\mathbf{w}^*|_2 = 1$. Rewriting Oja's update rule (12), we have that:

$$\mathbf{w}^{n+1} = \mathbf{w}^n + \beta y^n (\mathbf{x}^n - \mathbf{w}^n y^n) \quad (12)$$

At equilibrium $\mathbb{E}[\Delta\mathbf{w}] = \mathbb{E}[\mathbf{w}^{n+1} - \mathbf{w}^n] = 0$ and $\lim_{n \rightarrow \infty} \mathbf{w}^n = \mathbf{w}^*$. As a consequence, we get:

$$\mathbb{E}[y^n (\mathbf{x}^n - \mathbf{w}^* y^n)] = 0 \implies \mathbb{E}[y^n \mathbf{x}^n] - \mathbb{E}[\mathbf{w}^* (y^n)^2] = 0 \implies \mathbb{E}[y^n \mathbf{x}^n] - \mathbf{w}^* \mathbb{E}[(y^n)^2] = 0. \quad (13)$$

Note that:

$$\mathbb{E}[y^n \mathbf{x}^n] = \mathbb{E}\left[\left(\sum_j w_j^* x_j^n\right) x_i^n\right] = \sum_j w_j^* \underbrace{\mathbb{E}[x_j^n x_i^n]}_{C_{j,i}} \implies \mathbb{E}[y^n \mathbf{x}^n] = C\mathbf{w}^*. \quad (14)$$

So equation (13) by plugging in the expression $\mathbb{E}[y^n \mathbf{x}^n] = C\mathbf{w}^*$ and denoting $\mathbb{E}[(y^n)^2] = \sigma_y^2$ becomes:

$$C\mathbf{w}^* - \mathbf{w}^* \sigma_y^2 = 0 \implies C\mathbf{w}^* = \sigma_y^2 \mathbf{w}^*, \quad (15)$$

which proves item 1. Moreover, we have proven that the eigenvalue is equal to:

$$\lambda^* = \sigma_y^2. \quad (16)$$

4 Item 2: Norm

In order to prove that the norm of \mathbf{w}^* is one, we first start by equation (15) and right-multiply with the vector \mathbf{w}^{*T} to get:

$$C\mathbf{w}^* = \sigma_y^2 \mathbf{w}^* \implies \mathbf{w}^{*T} C\mathbf{w}^* = \sigma_y^2 \mathbf{w}^{*T} \mathbf{w}^* \implies \mathbf{w}^{*T} C\mathbf{w}^* = \sigma_y^2 |\mathbf{w}^*|_2^2. \quad (17)$$

However, by definition, we have that

$$\sigma_y^2 = \mathbb{E}[(y^n)^2] = \mathbb{E}[(\mathbf{w}^{*T} \mathbf{x})(\mathbf{w}^{*T} \mathbf{x})^T] = \mathbb{E}[\mathbf{w}^{*T} (\mathbf{x}\mathbf{x}^T) \mathbf{w}^*] = \mathbf{w}^{*T} \mathbb{E}[(\mathbf{x}\mathbf{x}^T)] \mathbf{w}^* = \mathbf{w}^{*T} C\mathbf{w}^*. \quad (18)$$

By plugging (17) in equation (18), we get:

$$\sigma_y^2 = \mathbf{w}^{*T} C\mathbf{w}^* \stackrel{eq.(17)}{=} \sigma_y^2 |\mathbf{w}^*|_2^2 \implies |\mathbf{w}^*|_2^2 = 1, \quad (19)$$

proving item 2.

5 Item 3: Principal Component

In order to prove that \mathbf{w}^* is a principal component, we will perform a linear stability analysis around the convergence point \mathbf{w}^* . We make the Ansatz

$$\mathbf{w}^{*n} = \mathbf{w}^* + \boldsymbol{\epsilon}^n, \quad (20)$$

where $\boldsymbol{\epsilon}^n$ is a small perturbation. This is expected to cause an update

$$\mathbf{w}^{*(n+1)} = \text{Oja}(\mathbf{w}^{*n}) = \text{Oja}(\mathbf{w}^*) + \boldsymbol{\epsilon}^{n+1}, \quad (21)$$

where $\text{Oja}(\cdot)$ is the classical application of the Oja's rule (12). The question we ask is how will an initial perturbation propagate due to the update rule? In other words, how does $\mathbb{E}[\delta\boldsymbol{\epsilon}] = \mathbb{E}[\boldsymbol{\epsilon}^{n+1} - \boldsymbol{\epsilon}^n]$ look like?

The expected update of Oja's rule around \mathbf{w} reads:

$$\mathbb{E}[\Delta\mathbf{w}] = \beta \mathbb{E}[y\mathbf{x} - y^2\mathbf{w}] = \beta \left(C\mathbf{w} - \underbrace{(\mathbf{w}^T C\mathbf{w})}_{\sigma_y^2 = \mathbb{E}[y^2]} \mathbf{w} \right), \quad (22)$$

where we used previous results, from equations (14) and equation (18). Now, let's plug $\mathbf{w} = \mathbf{w}^* + \boldsymbol{\epsilon}$ in the update rule. As we are simplifying the expression, we are ignoring $\mathcal{O}(\boldsymbol{\epsilon}^2)$ terms:

$$\frac{\mathbb{E}[\Delta(\mathbf{w}^* + \boldsymbol{\epsilon}^n)]}{\beta} = C(\mathbf{w}^* + \boldsymbol{\epsilon}^n) - ((\mathbf{w}^* + \boldsymbol{\epsilon}^n)^T C(\mathbf{w}^* + \boldsymbol{\epsilon}^n))(\mathbf{w}^* + \boldsymbol{\epsilon}^n) = \quad (23)$$

$$= C\mathbf{w}^* + C\boldsymbol{\epsilon}^n - (\mathbf{w}^{*T} C\mathbf{w}^*)\mathbf{w}^* - 2(\boldsymbol{\epsilon}^{nT} C\mathbf{w}^*)\mathbf{w}^* - \underbrace{(\boldsymbol{\epsilon}^{nT} C\boldsymbol{\epsilon}^n)\mathbf{w}^*}_{\mathcal{O}(\boldsymbol{\epsilon}^2) \approx 0} - \underbrace{(\mathbf{w}^{*T} C\mathbf{w}^*)\boldsymbol{\epsilon}^n}_{\mathcal{O}(\boldsymbol{\epsilon}^2) \approx 0} - \underbrace{2(\boldsymbol{\epsilon}^{nT} C\mathbf{w}^*)\boldsymbol{\epsilon}^n}_{\mathcal{O}(\boldsymbol{\epsilon}^2) \approx 0} - \underbrace{(\boldsymbol{\epsilon}^{nT} C\boldsymbol{\epsilon}^n)\boldsymbol{\epsilon}^n}_{\mathcal{O}(\boldsymbol{\epsilon}^3) \approx 0}. \quad (24)$$

Now, taking into account expression (24) and separating the terms that depend on $\boldsymbol{\epsilon}$ from those that do not, we get:

$$\frac{\mathbb{E}[\Delta(\mathbf{w}^* + \boldsymbol{\epsilon}^n)]}{\beta} = \underbrace{C\mathbf{w}^* - (\mathbf{w}^{*T} C\mathbf{w}^*)\mathbf{w}^*}_{\text{Oja}(\mathbf{w}^*)/\beta} + \underbrace{C\boldsymbol{\epsilon}^n - 2(\boldsymbol{\epsilon}^{nT} C\mathbf{w}^*)\mathbf{w}^* - (\mathbf{w}^{*T} C\mathbf{w}^*)\boldsymbol{\epsilon}^n}_{\boldsymbol{\epsilon}^{n+1}/\beta}. \quad (25)$$

We are interested in the error term $\boldsymbol{\epsilon}^{n+1}$ that shows how an initial perturbation to \mathbf{w}^* evolves. By rewriting the error and simplifying, we get:

$$\boldsymbol{\epsilon}^{n+1}/\beta = C\boldsymbol{\epsilon}^n - 2(\boldsymbol{\epsilon}^{nT} C\mathbf{w}^*)\mathbf{w}^* - (\mathbf{w}^{*T} C\mathbf{w}^*)\boldsymbol{\epsilon}^n = C\boldsymbol{\epsilon}^n - 2\lambda^*(\boldsymbol{\epsilon}^{nT} \mathbf{w}^*)\mathbf{w}^* - \lambda^*(\mathbf{w}^{*T} \boldsymbol{\epsilon}^n)\boldsymbol{\epsilon}^n \implies \quad (26)$$

$$\boldsymbol{\epsilon}^{n+1}/\beta = C\boldsymbol{\epsilon}^n - 2\lambda^*(\boldsymbol{\epsilon}^{nT} \mathbf{w}^*)\mathbf{w}^* - \lambda^*\boldsymbol{\epsilon}^n \quad (27)$$

where we used the fact that \mathbf{w}^* is an eigenvector with eigenvalue λ^* . Now let's say \mathbf{w}^* is the eigenvector \mathbf{v}^* . We need to prove that $\mathbf{v}^* = \mathbf{v}_1$. Equation (27) becomes:

$$\boldsymbol{\epsilon}^{n+1}/\beta = C\boldsymbol{\epsilon}^n - 2\lambda^* \underbrace{(\boldsymbol{\epsilon}^{nT} \mathbf{v}^*)}_{\in \mathbb{R}} \mathbf{v}^* - \lambda^*\boldsymbol{\epsilon}^n. \quad (28)$$

If we left multiply with any other eigenvector \mathbf{v}_j^T with $C\mathbf{v}_j = \lambda_j\mathbf{v}_j$ and $\lambda_j \neq \lambda^*$, we get:

$$\frac{\mathbf{v}_j^T \boldsymbol{\epsilon}^{n+1}}{\beta} = \mathbf{v}_j^T C\boldsymbol{\epsilon}^n - 2\lambda^* \underbrace{(\boldsymbol{\epsilon}^{nT} \mathbf{v}^*)}_{\in \mathbb{R}} \underbrace{\mathbf{v}_j^T \mathbf{v}^*}_{=0} - \lambda^* \mathbf{v}_j^T \boldsymbol{\epsilon}^n, \quad (29)$$

where $\mathbf{v}_j^T \mathbf{v}^* = 0$ due to orthogonality of the eigenvectors. Further, we know that

$$C\mathbf{v}_j = \lambda_j\mathbf{v}_j \implies (C\mathbf{v}_j)^T = \lambda_j\mathbf{v}_j^T \implies \mathbf{v}_j^T C^T = \lambda_j\mathbf{v}_j^T \implies \mathbf{v}_j^T C = \lambda_j\mathbf{v}_j^T \quad (30)$$

where we used that C is a symmetric real matrix (covariance matrix of independent variables). Then equation (29) becomes:

$$\frac{\mathbf{v}_j^T \boldsymbol{\epsilon}^{n+1}}{\beta} = \lambda_j \mathbf{v}_j^T \boldsymbol{\epsilon}^n - \lambda^* \mathbf{v}_j^T \boldsymbol{\epsilon}^n = (\lambda_j - \lambda^*) \mathbf{v}_j^T \boldsymbol{\epsilon}^n. \quad (31)$$

Expression (31) shows that convergence to \mathbf{w}^* can only be stable if $\lambda_j - \lambda^* < 0$, so that the expected increase towards any other principal direction \mathbf{v}_j is suppressed as $\frac{\mathbf{v}_j^T \boldsymbol{\epsilon}^{n+1}}{\beta} < 0$. If there is a principal component with $\lambda_j > \lambda^*$, then $\frac{\mathbf{v}_j^T \boldsymbol{\epsilon}^{n+1}}{\beta} > 0$ will be positive and \mathbf{w}^* is not a stable convergence point of Oja's rule, and the updates will diverge towards the principal component. This shows that if Oja's rule converges, $\mathbf{w}^* = \mathbf{v}_1$ is the principal component, with $\lambda^* > \lambda_j$. Cases where the geometric multiplicity of the principal eigenvalue/eigenvector is greater than one (e.g. $\lambda_1 = \lambda_2$) are not considered in this analysis.

6 Item 4: Maximum variance

The variance σ_y^2 is (given in equation (17)):

$$\sigma_y^2 = \mathbf{w}^{*T} C \mathbf{w}^*. \quad (32)$$

Assume that the eigenvalue decomposition of C reads $C = V \Lambda V^{-1}$, where Λ is the diagonal matrix with the eigenvalues sorted in descending order ($\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_D$) in its diagonal:

$$\Lambda = \begin{bmatrix} \lambda_1 & & \\ & \ddots & \\ & & \lambda_D \end{bmatrix}, \quad (33)$$

and $V = [\mathbf{v}_1, \dots, \mathbf{v}_D]$ is the orthogonal matrix containing the eigenvectors of C (at the same time principal components of data matrix $X \in \mathbb{R}^{N \times D}$). We have proven that $\mathbf{w}^* = \mathbf{v}_1$ is the principal component. Moreover, we know that the eigenvectors are orthogonal to each other, i.e. $\mathbf{v}_j^T \mathbf{v}_i = 0, \forall i \neq j$. **In the following we slightly abuse the notation for clarity and we denote the variance with $\sigma_y^2 = \sigma_{y, \mathbf{w}^*}^2$** to illustrate that it is the variance at the output explained by the principal component. Expression (32) becomes:

$$\sigma_y^2 = \sigma_{y, \mathbf{w}^*}^2 = \mathbf{w}^{*T} V \Lambda V^{-1} \mathbf{w}^* = \mathbf{v}_1^T V \Lambda V^{-1} \mathbf{v}_1 = \mathbf{v}_1^T [\mathbf{v}_1, \dots, \mathbf{v}_D] \Lambda [\mathbf{v}_1, \dots, \mathbf{v}_D]^T \mathbf{v}_1 = [1, 0, \dots, 0] \Lambda [1, 0, \dots, 0]^T \quad (34)$$

This leads to:

$$\sigma_{y, \mathbf{w}^*}^2 = [1, 0, \dots, 0] \begin{bmatrix} \lambda_1 & & \\ & \ddots & \\ & & \lambda_D \end{bmatrix} \begin{bmatrix} 1 \\ 0 \\ \vdots \\ 0 \end{bmatrix} \Rightarrow \sigma_{y, \mathbf{w}^*}^2 = \lambda_1. \quad (35)$$

Using the fact that $\sigma_{y, \mathbf{w}^*}^2 = \lambda^*$ given from equation (16), we end up with:

$$\sigma_{y, \mathbf{w}^*}^2 = \lambda^* = \lambda_1, \quad (36)$$

which shows that the variance "explained" in the output of the perceptron with weights \mathbf{w}^* is equal to the maximum eigenvalue of C . Is this the maximum variance that can be retained? Assume any other weight \mathbf{w} , which is **not** the principal component, and $|\mathbf{w}|_2 = 1$. V is orthonormal (C is symmetric so $V^T = V^{-1}$) and $\mathbf{z} = V^T \mathbf{w}$ also satisfies $|\mathbf{z}|_2 = 1$. The variance explained by any other weight \mathbf{w} is given by

$$\sigma_{y, \mathbf{w}}^2 = \mathbf{w}^T V \Lambda V^{-1} \mathbf{w} = \mathbf{z}^T \Lambda \mathbf{z} = [z_1, z_2, \dots, z_D] \begin{bmatrix} \lambda_1 & & \\ & \ddots & \\ & & \lambda_D \end{bmatrix} \begin{bmatrix} z_1 \\ z_2 \\ \vdots \\ z_D \end{bmatrix} = \sum_{i=1}^D \lambda_i z_i^2 \leq \lambda_1 \sum_{i=1}^D z_i^2 = \lambda_1 = \sigma_{y, \mathbf{w}^*}^2. \quad (37)$$

As a consequence, the variance explained in the output by the principal component \mathbf{w}^* is the maximum.

References

- [1] Erkki Oja. Simplified neuron model as a principal component analyzer. *Journal of Mathematical Biology*, 15(3):267–273, Nov 1982.
- [2] Erkki Oja. Neural networks, principal components, and subspaces. *Int. J. Neural Syst.*, 1:61–68, 1989.