

Some theoretical and experimental observations on naive discriminative learning

Stefan Evert
Friedrich-Alexander University of Erlangen-Nürnberg
Email: stefan.evert@fau.de

Antti Arppe
University of Alberta
Email: arppe@ualberta.ca

Abstract—The abstract goes here.

I. INTRODUCTION

Naïve Discriminative Learning (NDL) [citation needed] performs linguistic classification on the basis of direct associations between cues and outcomes, which are learned with the Rescorla-Wagner (R-W) equations. It is well known that the R-W model is related to neural networks (through the “delta rule” for gradient-descent training of a single-layer perceptron or SLP) and to linear least-squares regression [CITE e.g. Danks 2003, Baayen 2011]. However, most authors do not seem to be aware of the true depth of this equivalence and of its implications.

[CITE Danks (2003)] argued that if the R-W equations successfully acquire the true associations between cues and outcomes, they should approximate an equilibrium state in which the expected change in associations $E[\Delta V] = 0$ if a cue-outcome event is randomly presented to the learner. This equilibrium state can be computed directly by solving a matrix equation, without carrying out many iterations of R-W updates. Note that – unless the learning rate is gradually decreased – the R-W model cannot converge to its equilibrium state (and will rather oscillate around the equilibrium).

In this paper, we show that the R-W equations are identical to gradient-descent training of a single-layer feed-forward neural network, which we refer to as a single layer perceptron (SLP)¹ here (Sec. ??). Based on this result, we present a new, simpler derivation of the equilibrium conditions of [CITE Danks (2003)] and prove that they correspond to the solution of a linear least-squares regression problem (Sec. ??). In Sec. ?? we discuss some consequences of these new insights.

II. THE RESCORLA-WAGNER EQUATIONS

This section gives a mathematically precise definition of the R-W model, following the notation of [CITE Danks (2003)]. The purpose of the R-W equations is to determine associations between a set of cues C_1, \dots, C_n and a single outcome O in a population of event tokens $(\mathbf{c}^{(t)}, o^{(t)})$ with

¹The term SLP is often reserved for a particular form of such a single-layer network using a Heavyside activation function [see e.g. its Wikipedia entry]. Here, we use it more generally to refer to any single-layer feed-forward network.

$c_i^{(t)} = \begin{cases} 1 & \text{if } C_i \text{ is present} \\ 0 & \text{otherwise} \end{cases}$ $o^{(t)} = \begin{cases} 1 & \text{if } O \text{ results} \\ 0 & \text{otherwise} \end{cases}$ and $t = 1, \dots, m$ (where $m = \infty$ can be allowed).

When presented with an event (\mathbf{c}, o) , the R-W equations update the associations V_i between cues and the outcome according to Eq. (??), which is a more formal notation of Eq. (1) from [CITE Danks (2003: 111)].

$$\Delta V_i = \begin{cases} 0 & \text{if } c_i = 0 \\ \alpha_i \beta_1 (\lambda - \sum_{j=1}^n c_j V_j) & \text{if } c_i = 1 \wedge o = 1 \\ \alpha_i \beta_2 (0 - \sum_{j=1}^n c_j V_j) & \text{if } c_i = 1 \wedge o = 0 \end{cases} \quad (1)$$

Here, α_i is a measure of the salience of cue C_i , β_1 and β_2 are learning rates for positive ($o = 1$) and negative ($o = 0$) events, and λ is the maximal activation level of the outcome O . A simplified form of the R-W equations proposed by [CITE Widrow & Hoff (1960)] assumes that $\alpha_1 = \dots = \alpha_n = 1$, $\beta_1 = \beta_2 = \beta$ and $\lambda = 1$ (known as the W-H rule).

[CITE Danks (2003)] argue that a successful R-W model should approach an equilibrium state of the association vector $\mathbf{V} = (V_1, \dots, V_n)$ where the expected update $E[\Delta V_i] = 0$ if a random event token is sampled from the population. If we make the simplifying assumption that $\beta_1 = \beta_2 = \beta$, this condition corresponds to the equality

$$\lambda \frac{1}{m} \sum_{t=1}^m c_i^{(t)} o_i^{(t)} - \sum_{j=1}^n V_j \frac{1}{m} \sum_{t=1}^m c_i^{(t)} c_j^{(t)} = 0 \quad (2)$$

In Danks’s notation, Eq. (??) can be written as

$$\lambda P(O, C_i) - \sum_{j=1}^n V_j P(C_i, C_j) = 0 \quad (3)$$

and is equal to his Eq. (11) [CITE (Danks 2003: 113)] multiplied by $P(C_i)$.

III. R-W AND THE SINGLE LAYER PERCEPTRON

We will now formulate a single-layer feed-forward neural network (SLP) whose learning behaviour – with gradient-descent training, which corresponds to the backprop algorithm for a SLP and is also known as the “delta rule” in this case – is identical to the R-W equations with equal positive and negative learning rates $\beta_1 = \beta_2$ (but

no other restrictions). The SLP requires a slightly different representations of events as pairs $(\mathbf{x}^{(t)}, z^{(t)})$ with $\mathbf{x}_i^{(t)} = \{a_i \text{ if } C_i \text{ is present}$
 0 otherwise $z^{(t)} = \{\lambda \text{ if } O \text{ results}$
 0 otherwise Here, $a_i > 0$ is a (different) measure of the salience of cue C_i and $\lambda > 0$ the maximum activation of outcome O . Note that the event representation (\mathbf{x}, z) is connected to the representation (\mathbf{c}, o) through the equivalences $x_i = a_i c_i$ and $z = \lambda o$. In the W-H case, the two representations are identical.

The SLP uses a linear activation function $h(y) = y$ and Euclidean cost for the difference between y and the desired activation z . It computes the activation of the outcome as a linear combination $y = \sum_{i=1}^n w_i x_i$, where $\mathbf{w} = (w_1, \dots, w_n)$ is the weight vector of the SLP. The cost associated with a given event token (\mathbf{x}, z) is

$$E(\mathbf{w}, \mathbf{x}, z) = (z - y)^2 = \left(z - \sum_{i=1}^n w_i x_i\right)^2 \quad (4)$$

For batch updates based on the full population of event tokens, the corresponding batch cost is

$$E(\mathbf{w}) = \sum_t E(\mathbf{w}, \mathbf{x}^{(t)}, z^{(t)}) \quad (5)$$

If smaller batches are used, the sum j ranges over a subset of the population for each update step.

Given an event token (\mathbf{x}, z) , gradient-descent training of this SLP updates the weight vector by

$$\Delta w_i = -\delta \frac{\partial E(\mathbf{w}, \mathbf{x}, z)}{\partial w_i} \quad (6)$$

where $\delta > 0$ is the learning rate and the gradient $\partial E / \partial w_i$ is given by

$$\frac{\partial E(\mathbf{w}, \mathbf{x}, z)}{\partial w_i} = 2(z - y)(-x_i) = -2\left(z - \sum_{j=1}^n w_j x_j\right)x_i \quad (7)$$

Inserting the definition of \mathbf{x} and z from Eq. (??), we obtain

$$\Delta w_i = \{0 \text{ if } c_i = 0, 2\delta a_i(\lambda - \sum_{j=1}^n c_j a_j w_j) \text{ if } c_i = 1 \wedge o = 1, 2\delta a_i(0 - \sum_{j=1}^n c_j a_j w_j) \text{ if } c_i = 1 \wedge o = 0\} \quad (8)$$

Comparing this with Eq. (??), we can set $V_j = a_j w_j$, i.e. we interpret the weight vector \mathbf{w} of the SLP as salience-adjusted cue-outcome associations. With $\Delta V_i = a_i \Delta w_i$, we obtain

$$\Delta V_i = \{0 \text{ if } c_i = 0, 2\delta a_i^2(\lambda - \sum_{j=1}^n c_j V_j) \text{ if } c_i = 1 \wedge o = 1, 2\delta a_i^2(0 - \sum_{j=1}^n c_j V_j) \text{ if } c_i = 1 \wedge o = 0\} \quad (9)$$

which is equal to the R-W equations for $\beta = 2\delta$ and $\alpha_i = a_i^2$.

The assumption $\beta_1 = \beta_2$ can be relaxed if we change the representation of events to

$$x_i^{(t)} = \{a_i \text{ if } c_i = 1 \wedge o = 1, a_i \sqrt{\frac{\beta_2}{\beta_1}} \text{ if } c_i = 1 \wedge o = 0, 0 \text{ otherwise}\} \quad (10)$$

i.e. if we allow the salience of cues to differ between positive ($o = 1$) and negative ($o = 0$) events; the scaling factor β_2/β_1 is the same for all cues C_i . We do not pursue this extension further here because it affects the equilibrium state in an unpredictable way. As [CITE Danks (2003)] has already observed, the cue saliences α_i have no impact at all and the maximum activation level λ merely results in a linear scaling of the equilibrium state.

IV. R-W AND LEAST-SQUARES REGRESSION

We have shown in Sec. ?? that the R-W equations describe the gradient-descent training of a SLP for the linear regression problem

$$\min_{\mathbf{w}} E(\mathbf{w}) = \min_{\mathbf{w}} \sum_t E(\mathbf{w}, \mathbf{x}^{(t)}, z^{(t)}) \quad (11)$$

This equivalence holds generally, not only in the case of the simplified W-H rule. Thus, both R-W and our SLP aim to solve the same regression problem.

If the training procedure is successful, the weight vector \mathbf{w} should approach the least-squares solution. With single updates (corresponding to the R-W model), convergence cannot be achieved unless the learning rate is gradually reduced. With batch updates treating the entire population as a single batch, the cost $E(\mathbf{w})$ is a convex function of \mathbf{w} and the gradient descent procedure converges to its unique minimum after a sufficient number of iterations.²

In order to express Eq. (??) more concisely, we define an $m \times n$ matrix $\mathbf{X} = (x_i^{(t)}) = (x_{ti})$ of cues for all event tokens in the population. The rows of this matrix correspond to event tokens t , the columns to cues i ; i.e. row number t contains the input vector $\mathbf{x}^{(t)}$. We also define the column vector $\mathbf{z} = (z^{(1)}, \dots, z^{(m)})$ of outcomes and recall that $\mathbf{w} = (w_1, \dots, w_n)$ is a column vector of SLP weights. The batch cost can now be written as

$$E(\mathbf{w}) = (\mathbf{z} - \mathbf{X}\mathbf{w})^2 \quad (12)$$

The least-squares solution must satisfy the condition $\nabla E(\mathbf{w}) = 0$, which leads to the so-called normal equations

$$\mathbf{X}^T \mathbf{X} \mathbf{w} = \mathbf{X}^T \mathbf{z} \quad (13)$$

In the case of the W-H rule, \mathbf{X} is a coincidence matrix between cues and events, with $x_{ti} \in \{0, 1\}$. A straightforward calculation shows that $\mathbf{X}^T \mathbf{X}$ is a square co-occurrence matrix with entries $f(C_i, C_j)$, and $\mathbf{X}^T \mathbf{z}$ is a vector of co-occurrence counts $f(C_i, O)$ between the cues and the outcome O . Dividing Eq. (??) by m , we obtain Eq. (??) with $\lambda = 1$, i.e.

$$P(O, C_i) - \sum_{j=1}^n V_j P(C_i, C_j) = 0$$

²In fact, the minimum of $E(\mathbf{w})$ might not be unique under certain circumstances, viz. if the correlation matrix of the cues is not positive definite [CITE cf. Danks 2003: 115–116, who considers only the special case of “coextensive” cues]. In order to keep the discussion straightforward, we assume the general case of a unique minimum in the present paper.

which is the same as Eq. (3) of [CITE Danks (2003)] with rows multiplied by $P(C_i)$. Since linear regression is invariant wrt. the salience factors a_i (the weights are simply adjusted by reciprocal factors $1/a_i$ to achieve the same regression values) and scales linearly with λ , equivalence to the equilibrium conditions [CITE (Danks 2003: 112–114)] also holds for arbitrary values of a_i and λ .

In the general case where the regression problem has a unique solution, $\mathbf{X}^T \mathbf{X}$ is symmetric and positive definite. It can therefore be inverted and the least-squares solution is given by

$$\mathbf{w}^* = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{z} \quad (14)$$

Standard software such as R [citation needed] can be used to compute \mathbf{w}^* reliably and efficiently. It is not necessary to carry out the iterative training procedure of the R-W model or the neural network, and there is no need to worry about convergence of the iterative training.

V. CONSEQUENCES

[Just some thoughts – can you please turn the bullet points into a nice text?]

- We have shown that R-W association learning, a linear SLP neural network and linear regression are fully equivalent and should ideally lead to the same least-squares solution. As long as a researcher is only interested in the result of association learning, not in the iterative process, it is sufficient to calculate the least-squares solution directly from Eq. (??).
- The R-W salience factors α_i have no effect on the learning result – because linear regression is not sensitive to such a scaling of the input variables – but only on the learning process: associations for cues with high salience α_i are learnt faster than for other cues. The parameter λ leads to a (trivial) linear scaling of the learning result, but has not effect on the learning process. Only different learning rates $\beta_1 \neq \beta_2$ affect the learning result, because they modify $\mathbf{X}^T \mathbf{X}$ in a complex way.
- If R-W association learning or SLP training does not approximate the least-squares solution, it can arguably be considered to have failed. The only research question of interest that requires R-W iteration or application of the delta rule is thus: Under which circumstances and for which parameter settings does the R-W iteration converge or at least approximate the linear regression? This is particularly relevant for single-event updates (as specified for the R-W model), which are much less robust and lead to larger fluctuations than batch updates. We plan to work on these issues with the help of simulation experiments.
- Having established NDL as linear regression with its well-known drawbacks (e.g. a propensity for overfitting the training data, especially if there is a large number n of cues), it will be interesting to contrast it with more state-of-the-art machine learning

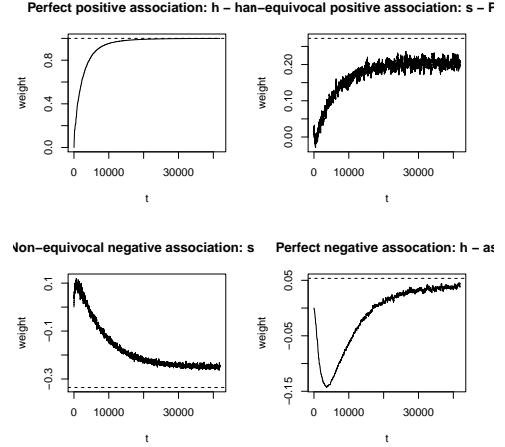


Fig. 1. Simulation results for the network.

techniques. We plan to carry out a mathematical analysis and empirical study of (i) logistic regression, which is more appropriate for dichotomous data than linear least-squares regression, and (ii) regularization techniques, which control overfitting and encourage sparse solutions.

REFERENCES

- [1] H. Kopka and P. W. Daly, *A Guide to L^AT_EX*, 3rd ed. Harlow, England: Addison-Wesley, 1999.