

Naïve Discriminative Learning: Theoretical and Experimental Observations

Stefan Evert¹ & Antti Arppe²

¹Friedrich-Alexander-Universität Erlangen-Nürnberg, Germany
stefan.evert@fau.de

²University of Alberta, Edmonton, Canada
arppe@ualberta.ca

QITL-6, Tübingen, 6 Nov 2015



FRIEDRICH-ALEXANDER
UNIVERSITÄT
ERLANGEN-NÜRNBERG
PHILOSOPHISCHE FAKULTÄT
UND FACHBEREICH THEOLOGIE

Outline

- 1 Introduction
 - Naïve Discriminative Learning
 - An example
- 2 Mathematics
 - The Rescorla-Wagner equations
 - The Danks equilibrium
 - NDL vs. the Perceptron vs. least-squares regression
- 3 Insights
 - Theoretical insights
 - Empirical observations
 - Conclusion

Outline

- 1 Introduction
 - Naïve Discriminative Learning
 - An example
- 2 Mathematics
 - The Rescorla-Wagner equations
 - The Danks equilibrium
 - NDL vs. the Perceptron vs. least-squares regression
- 3 Insights
 - Theoretical insights
 - Empirical observations
 - Conclusion

Objectives

- Present the mathematic underpinnings of NDL in one place, in a systematic way
- High-light the theoretical similarities of NDL with linear/logistic regression and perceptron
- Present some empirical simulations of NDL, in light of the theory

Naive Discriminative Learning

- Baayen et al. 2011; Baayen 2011
- Rescorla-Wagner (1972) incremental learning equations
- Danks (2003) equilibrium equations
- Implementation as an R package `nd1`: Arppe et al. 2011; Shaoul et al. 2013

Danks (2003) equilibrium equations – verbally

- presume an ideal ‘adult/stable’ state where all the cue-outcome associations have been fully learnt – any more data points bring nothing ‘new’ to learn, i.e. have zero impact on the cue-outcome associations.
- make it possible to estimate the weights for a system using relatively simple matrix algebra.
- provide a convenient short-cut to calculating the consolidated cue-outcome association weights resulting from incremental learning.
- the learning parameters of the Rescorla-Wagner equations drop out of the equilibrium equations.
- circumvent the problem that a simulation of an Rescorla-Wagner learner does not converge to a single state unless the learning rate is gradually decreased.

Rescorla-Wagner equations (1972) – verbally

Represent incremental learning and subsequently on-going adjustments to an accumulating body of knowledge:

Changes in association strengths:

- If a cue is not present in the input, no change
- Increased when the cue and outcome co-occur
- Decreased when the cue occurs without the outcome
- The more cues are present simultaneously, the smaller the adjustments are

Only the results of the incremental adjustments to the cue-outcome associations are kept – no need for remembering the individual adjustments, however many there are.

Naive Discriminative Learning

- Naive: cue-outcome associations estimated separately for each outcome (this simplifying assumption of independence similar to a naive Bayesian classifier).
- Discriminative: direct associations with each outcome given a set of cues.
- Learning: based on incremental learning.

Rescorla-Wagner equations – traditional vs. linguistic applications

- traditionally: simple controlled experiments on item-by-item learning, with only a couple of cues and some perfect associations.
- natural language: full of choices among multiple possible alternatives – phones, words, or constructions – which are influenced by a large number of contextual factors, and which rather exhibit asymptotic, imperfect tendencies favoring one or more of the alternatives, instead of single, categorical, perfect choices.
- these messy, complex types of problems as a key area of interest in modeling and understanding language use.
- the application of the Rescorla-Wagner equations in the form of a Naive Discriminative Learning classifier to such linguistic phenomena of considerable utility.

Simple vs. complex settings - QITL-01 revisited

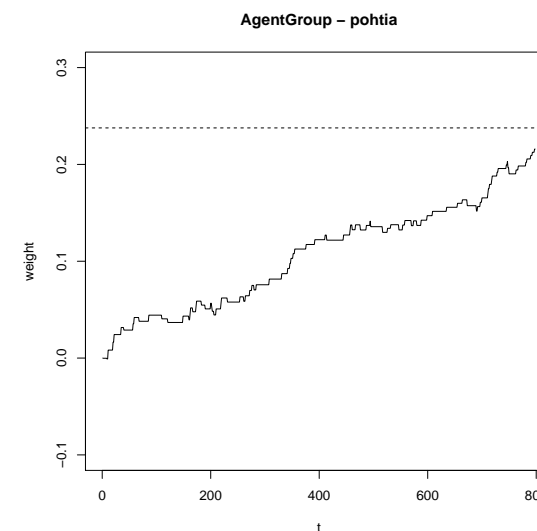
- Arppe & Järvikivi (2002, 2007)
- *Person* (FIRST PERSON SINGULAR or not) and *Countability* (COLLECTIVE or not) of AGENT/SUBJECT of Finnish verb synonym pair *mieltiä* vs. *pohtia* 'think, ponder':

Forced-choice Dispreferred	Preferred	Frequency (relative)	Unacceptable	Acceptability Acceptable
∅	mieltiä+SG1 pohtia+COLL	Frequent	∅	mieltiä+SG1 pohtia+COLL
mieltiä+COLL pohtia+SG1	∅	Rare	mieltiä+COLL	pohtia+SG1

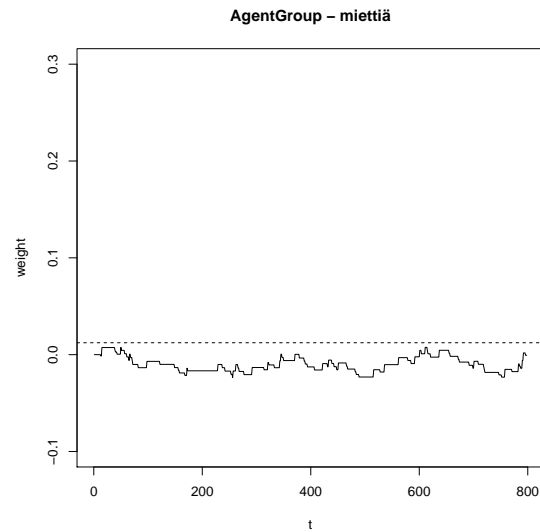
Outline

- 1 Introduction
 - Naïve Discriminative Learning
 - An example
- 2 Mathematics
 - The Rescorla-Wagner equations
 - The Danks equilibrium
 - NDL vs. the Perceptron vs. least-squares regression
- 3 Insights
 - Theoretical insights
 - Empirical observations
 - Conclusion

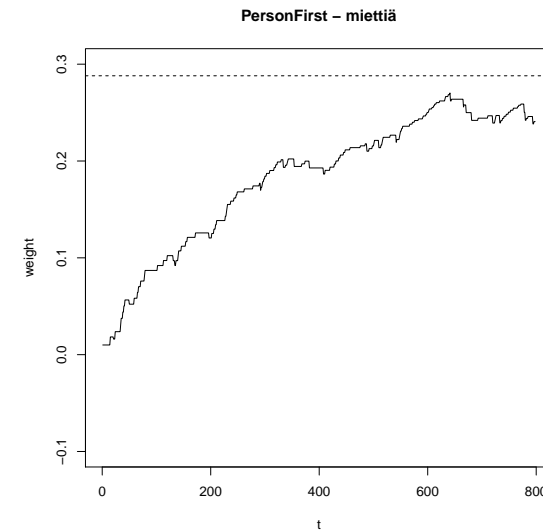
QITL-1 through the lenses of NDL: 1/4



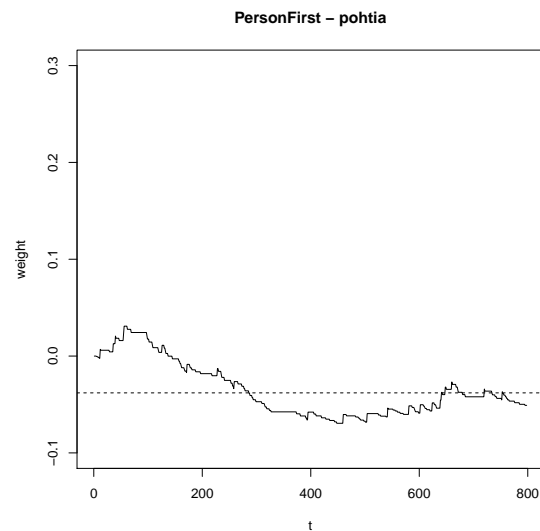
QITL-1 through the lenses of NDL: 2/4



QITL-1 through the lenses of NDL: 3/4



QITL-1 through the lenses of NDL: 4/4



QITL-01: Linguistic production vs. judgments

Forced-choice Dispreferred	Preferred	Frequency (relative)	Unacceptable	Acceptability Acceptable
∅	+	Frequent	∅	+
+	∅	Rare	+	+

$Frequency \Rightarrow Acceptability$

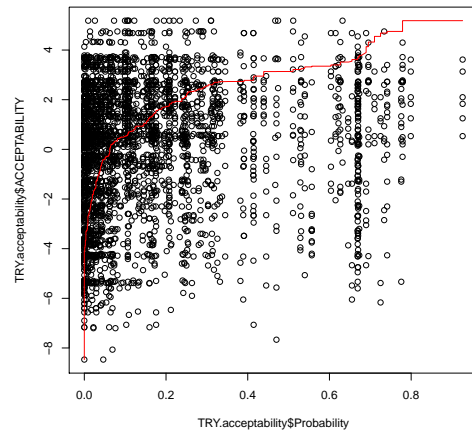
$Unacceptability \Rightarrow Rarity$

$\neg(Acceptability \Rightarrow Frequency)$

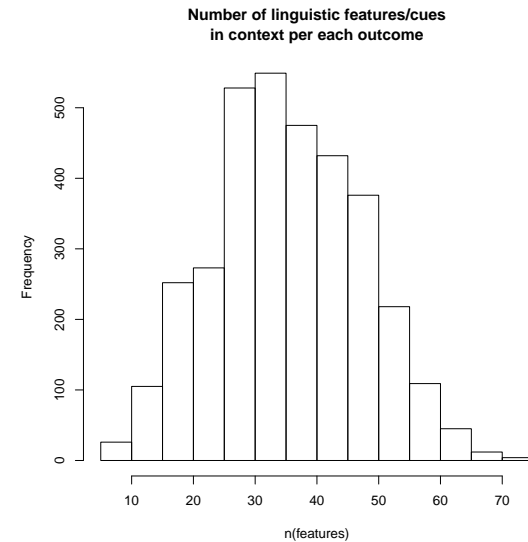
$\neg(Rarity \Rightarrow Unacceptability)$

QITL-01 through the lenses of QITL-6

(courtesy of Dagmar Divjak)



Simple vs. complex settings – QITL-2 revisited



QITL-4 revisited – comparison of NDL with statistical methods – Classification Accuracy & Recall

	$\lambda_{\text{prediction}}$	$\tau_{\text{classification}}$	Accuracy
Polytomous logistic regression (One-vs-rest)	0.447	0.516	0.646
Polytomous mixed logistic regression (Poisson reformulation)			
1—Register	0.435	0.505	0.638
1—Genre	0.433	0.504	0.637
1—Lexeme	0.428	0.499	0.634
1—Register + 1—Lexeme	0.431	0.502	0.636
Support Vector Machine	0.414	0.487	0.625
Memory-Based Learning (TiMBL)	0.287	0.376	0.543
Random Forests	0.445	0.515	0.645
Naive Discriminative Learning	0.442	0.511	0.642

Table: Classification diagnostics for five models fitted to the English data set ($n = 909$).

Outline

- 1 Introduction
 - Naïve Discriminative Learning
 - An example
- 2 Mathematics
 - The Rescorla-Wagner equations
 - The Danks equilibrium
 - NDL vs. the Perceptron vs. least-squares regression
- 3 Insights
 - Theoretical insights
 - Empirical observations
 - Conclusion

The Rescorla-Wagner equations

- Goal of naïve discriminative learner: predict an **outcome** O based on presence or absence of a set of **cues** C_1, \dots, C_n
- An **event** (\mathbf{c}, o) is formally described by indicator variables

$$c_i = \begin{cases} 1 & \text{if } C_i \text{ is present} \\ 0 & \text{otherwise} \end{cases} \quad o = \begin{cases} 1 & \text{if } O \text{ results} \\ 0 & \text{otherwise} \end{cases}$$

- Given cue-outcome **associations** $\mathbf{v} = (V_1, \dots, V_n)$ of learner, the **activation level** of the outcome O is

$$\sum_{j=1}^n c_j V_j$$

- Associations $\mathbf{v}^{(t)}$ as well as cue and outcome indicators $(\mathbf{c}^{(t)}, o^{(t)})$ depend on time step t

The Rescorla-Wagner equations

- Rescorla and Wagner (1972) proposed the **R-W equations** for the change in associations given an event (\mathbf{c}, o) :

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

with parameters

- $\lambda > 0$ target activation level for outcome O
- $\alpha_i > 0$ salience of cue C_i
- $\beta_1 > 0$ learning rate for positive ovents ($o = 1$)
- $\beta_2 > 0$ learning rate for negative ovents ($o = 0$)

The Widrow-Hoff rule

- The **W-H rule** (Widrow and Hoff 1960) is a widely-used simplification of the R-W equations:

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

$$= c_i \beta (o - \sum_{j=1}^n c_j V_j)$$

with parameters

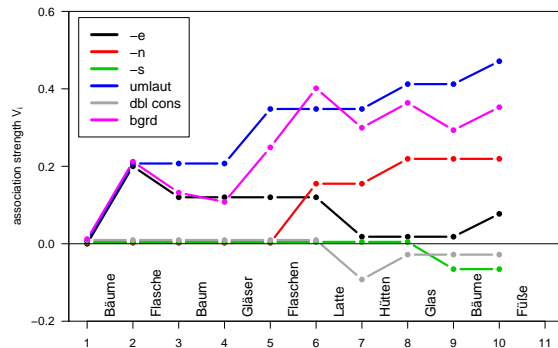
- $\lambda = 1$ target activation level for outcome O
- $\alpha_i = 1$ salience of cue C_i
- $\beta_1 = \beta_2$ global learning rate for positive and negative events
- $= \beta > 0$

A simple example: German noun plurals

t	word	o pl?	c_1 -e	c_2 -n	c_3 -s	c_4 umlaut	c_5 dbl cons	c_6 bgrd
1	Bäume	1	1	0	0	1	0	1
2	Flasche	0	1	0	0	0	0	1
3	Baum	0	0	0	0	0	0	1
4	Gläser	1	0	0	0	1	0	1
5	Flaschen	1	0	1	0	0	0	1
6	Latte	0	1	0	0	0	1	1
7	Hütten	1	0	1	0	1	1	1
8	Glas	0	0	0	1	0	0	1
9	Bäume	1	1	0	0	1	0	1
10	Füße	1	1	0	0	1	0	1

A simple example: German noun plurals

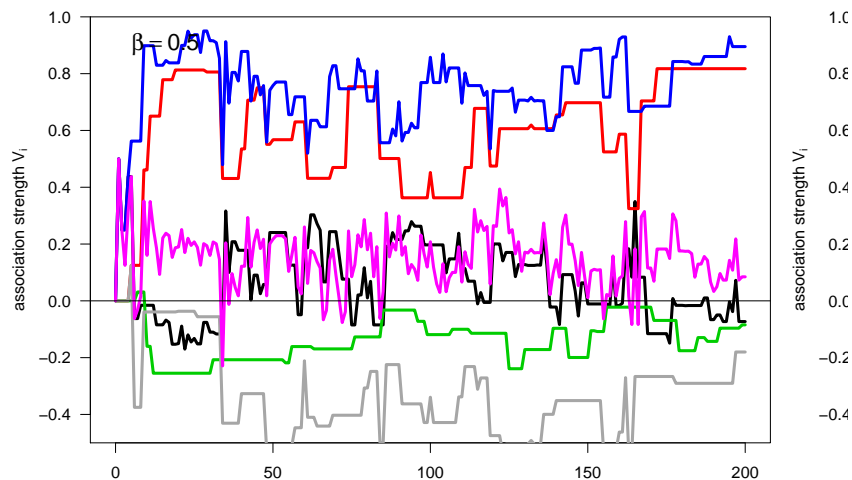
t	$\sum c_j V_j$	V_1	V_2	V_3	V_4	V_5	V_6
10	.882	.077	.217	-.070	.464	-.038	.340
Füße	1	1	0	0	1	0	1
	o	c_1	c_2	c_3	c_4	c_5	c_6



A stochastic NDL learner

- A specific event sequence $(\mathbf{c}^{(t)}, o^{(t)})$ will only be encountered in controlled experiments
- For applications in corpus linguistics, it is more plausible to assume that events are randomly sampled from a population of event tokens $(\mathbf{c}^{(k)}, o^{(k)})$ for $k = 1, \dots, m$
 - ☞ event types listed repeatedly proportional to their frequency
- I.i.d. random variables $\mathbf{c}^{(t)} \sim \mathbf{c}$ and $o^{(t)} \sim o$
 - ☞ distributions of \mathbf{c} and o determined by population
- NDL can now be trained for arbitrary number of time steps, even if population is small (as in our example)
 - ▶ study asymptotic behaviour of learners
 - ▶ convergence → stable “adult” state of associations

A stochastic NDL learner

Effect of the learning rate β 

Outline

- 1 Introduction
 - Naïve Discriminative Learning
 - An example
- 2 Mathematics
 - The Rescorla-Wagner equations
 - The Danks equilibrium
 - NDL vs. the Perceptron vs. least-squares regression
- 3 Insights
 - Theoretical insights
 - Empirical observations
 - Conclusion

Expected activation levels

- Since we are interested in the general behaviour of a stochastic NDL, it makes sense to average over many individual learners to obtain **expected associations** $E[V_j^{(t)}]$

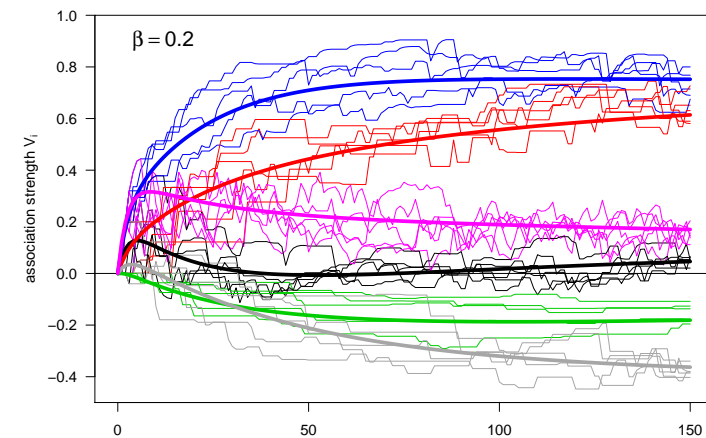
$$E[V_{j+1}^{(t)}] = E[V_j^{(t)}] + E[\Delta V_j^{(t)}]$$

$$\begin{aligned} E[\Delta V_j^{(t)}] &= E\left[c_i \beta (o - \sum_{j=1}^n c_j V_j^{(t)})\right] \\ &= \beta \cdot \left(\Pr(C_i, O) - \sum_{j=1}^n \Pr(C_i, C_j) E[V_j^{(t)}]\right) \end{aligned}$$

- c_i and c_j are independent from $V_j^{(t)}$
- indicator variables: $E[c_i o] = \Pr(C_i, O)$; $E[c_i c_j] = \Pr(C_i, C_j)$

Expected activation levels

$$E[\Delta V_j^{(t)}] = \beta \cdot (\Pr(C_i, O) - \sum_{j=1}^n \Pr(C_i, C_j) E[V_j^{(t)}])$$



The Danks equilibrium

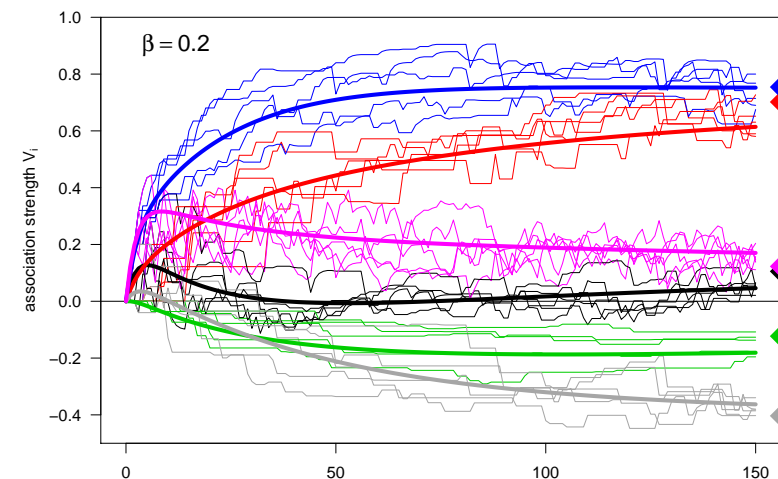
- If $E[V_i^{(t)}]$ converges, the asymptote $V_i^* = \lim_{t \rightarrow \infty} E[V_i^{(t)}]$ must satisfy the **Danks equilibrium** conditions $E[\Delta V_i^*] = 0$, i.e.

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

(Danks 2003, p. 113)

- Now there is a clear interpretation of the Danks equilibrium as the stable average associations reached by a community of stochastic learners with input from the same population
- allows us to compute the “adult” state of NDL without carrying out a simulation of the learning process

The Danks equilibrium



Matrix notation

$$\mathbf{X} = \begin{bmatrix} c_1^{(1)} & \dots & c_n^{(1)} \\ c_1^{(2)} & \dots & c_n^{(2)} \\ \vdots & & \vdots \\ c_1^{(m)} & \dots & c_n^{(m)} \end{bmatrix} \quad \mathbf{z} = \begin{bmatrix} o^{(1)} \\ o^{(2)} \\ \vdots \\ o^{(m)} \end{bmatrix} \quad \mathbf{w} = \begin{bmatrix} V^{(1)} \\ \vdots \\ V^{(n)} \end{bmatrix}$$

$$\begin{bmatrix} \Pr(C_1, O) \\ \vdots \\ \Pr(C_n, O) \end{bmatrix} = \frac{1}{m} \mathbf{X}^T \mathbf{z} \quad \begin{bmatrix} \Pr(C_1, C_1) & \dots & \Pr(C_1, C_n) \\ \vdots & & \vdots \\ \Pr(C_n, C_1) & \dots & \Pr(C_n, C_n) \end{bmatrix} = \frac{1}{m} \mathbf{X}^T \mathbf{X}$$

Danks equilibrium: $\mathbf{X}^T \mathbf{z} = \mathbf{X}^T \mathbf{X} \mathbf{w}^*$

Matrix notation: German noun plurals

$$\mathbf{X} = \begin{bmatrix} 1 & 0 & 0 & 1 & 0 & 1 \\ 1 & 0 & 0 & 0 & 0 & 1 \\ 0 & 0 & 0 & 0 & 0 & 1 \\ 0 & 0 & 0 & 1 & 0 & 1 \\ 0 & 1 & 0 & 0 & 0 & 1 \\ 1 & 0 & 0 & 0 & 1 & 1 \\ 0 & 1 & 0 & 1 & 1 & 1 \\ 0 & 0 & 1 & 0 & 0 & 1 \\ 1 & 0 & 0 & 1 & 0 & 1 \\ 1 & 0 & 0 & 1 & 0 & 1 \end{bmatrix} \quad \mathbf{z} = \begin{bmatrix} 1 \\ 0 \\ 0 \\ 1 \\ 1 \\ 0 \\ 1 \\ 0 \\ 1 \\ 1 \end{bmatrix} \quad \mathbf{w} = \begin{bmatrix} V^{(1)} \\ \vdots \\ V^{(n)} \end{bmatrix}$$

$$\begin{bmatrix} .3 \\ .2 \\ .0 \\ .5 \\ .1 \\ .6 \end{bmatrix} = \frac{1}{m} \mathbf{X}^T \mathbf{z} \quad \begin{bmatrix} .5 & .0 & .0 & .3 & .1 & .5 \\ .0 & .2 & .0 & .1 & .1 & .2 \\ .0 & .0 & .1 & .0 & .0 & .1 \\ .3 & .1 & .0 & .5 & .1 & .5 \\ .1 & .1 & .0 & .1 & .2 & .2 \\ .5 & .2 & .1 & .5 & .2 & 1 \end{bmatrix} = \frac{1}{m} \mathbf{X}^T \mathbf{X}$$

Outline

- 1 Introduction
 - Naïve Discriminative Learning
 - An example
- 2 Mathematics
 - The Rescorla-Wagner equations
 - The Danks equilibrium
 - NDL vs. the Perceptron vs. least-squares regression
- 3 Insights
 - Theoretical insights
 - Empirical observations
 - Conclusion

The single-layer perceptron (SLP)

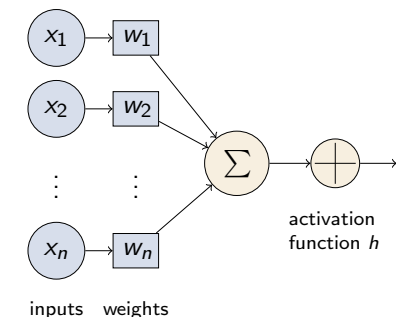
SLP (Rosenblatt 1958) is most basic feed-forward **neural network**

- numeric inputs x_1, \dots, x_n
- output activation $h(y)$ based on weighted sum of inputs

$$y = \sum_{j=1}^n w_j x_j$$

- h = Heaviside step function in traditional SLP
- even simpler model: $h(y) = y$
- cost wrt. target output z :

$$E(\mathbf{w}, \mathbf{x}, z) = \left(z - \sum_{j=1}^n w_j x_j \right)^2$$



SLP training: the delta rule

- SLP weights are learned by **gradient descent** training: for a single training item (\mathbf{x}, z) and learning rate $\delta > 0$

$$\begin{aligned}\Delta w_i &= -\delta \frac{\partial E(\mathbf{w}, \mathbf{x}, z)}{\partial w_i} \\ &= 2\delta x_i \left(z - \sum_{j=1}^n x_j w_j \right) \\ &= \beta c_i (o - \sum_{j=1}^n c_j V_j)\end{aligned}$$

- Perfect **correspondence to W-H rule** with

$$V_i = w_i \quad c_i = x_i \quad o = z \quad \beta = 2\delta$$

Batch training

- Neural networks often use **batch training**, where all training data are considered at once instead of one item at a time
- The corresponding batch training cost is

$$\begin{aligned}E(\mathbf{w}) &= \frac{1}{m} \sum_{k=1}^m E(\mathbf{w}, \mathbf{x}^{(k)}, z^{(k)}) \\ &= \frac{1}{m} \sum_{k=1}^m \left(z^{(k)} - \sum_{j=1}^n w_j x_j^{(k)} \right)^2\end{aligned}$$

- Similar to stochastic NDL, batch training computes the expected weights $E[\mathbf{w}^{(t)}]$ for SLP with stochastic input
- Minimization of $E(\mathbf{w})$ = linear **least-squares regression**

Linear least-squares regression

- Matrix formulation of the linear least-squares problem:

$$\begin{aligned}E(\mathbf{w}) &= \frac{1}{m} \sum_{k=1}^m \left(z^{(k)} - \sum_{j=1}^n w_j x_j^{(k)} \right)^2 \\ &= \frac{1}{m} (\mathbf{z} - \mathbf{X}\mathbf{w})^T (\mathbf{z} - \mathbf{X}\mathbf{w})\end{aligned}$$

- Minimum of $E(\mathbf{w})$, the L_2 solution, must satisfy $\nabla E(\mathbf{w}^*) = \mathbf{0}$, which leads to the **normal equations**

$$\mathbf{X}^T \mathbf{z} = \mathbf{X}^T \mathbf{X} \mathbf{w}^*$$

- Normal equations = Danks equilibrium conditions
- Regression theory shows that batch training / stochastic NLP converges to the unique* solution of the L_2 problem

What have we learned?

stochastic	=	batch	=	L_2 regression
NDL	=	SLP		

- These equivalences also hold for the general R-W equations with arbitrary values of α_i , β_1 , β_2 and λ (see paper)

Outline

- 1 Introduction
 - Naïve Discriminative Learning
 - An example
- 2 Mathematics
 - The Rescorla-Wagner equations
 - The Danks equilibrium
 - NDL vs. the Perceptron vs. least-squares regression
- 3 Insights
 - Theoretical insights
 - Empirical observations
 - Conclusion

Effects of R-W parameters

- $\beta > 0$: learning rate \rightarrow convergence of individual learners
- $\lambda \neq 1$: linear scaling of association / activation (obvious)
- $\alpha_i \neq 1$: salience of cue C_i determines how fast associations are learned, but does not affect the final stable associations (same L_2 regression problem)
- $\beta_1 \neq \beta_2$: different positive/negative learning rates *do* affect the stable associations; closely related to prevalence of positive and negative events in the population

What about logistic regression?

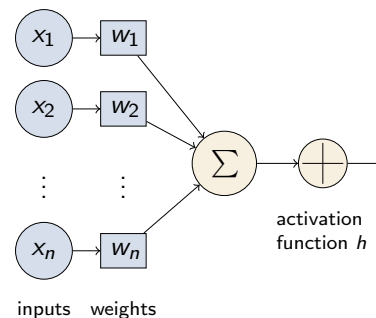
Logistic regression is the standard tool for predicting a categorical response from binary features

- can be expressed as SLP with probabilistic interpretation
- uses logistic activation function

$$h(y) = \frac{1}{1 + e^{-y}}$$

- and Bernoulli cost

$$E(\mathbf{w}, \mathbf{x}, z) = \begin{cases} -\log h(y) & \text{if } z = 1 \\ -\log(1 - h(y)) & \text{if } z = 0 \end{cases}$$



What about logistic regression?

- Gradient descent training leads to delta rule that corresponds to a modified version of the R-W equations

$$\Delta V_i = \begin{cases} 0 & \text{if } c_i = 0 \\ \beta \left(1 - h(\sum_{j=1}^n c_j V_j) \right) & \text{if } c_i = 1 \wedge o = 1 \\ \beta \left(0 - h(\sum_{j=1}^n c_j V_j) \right) & \text{if } c_i = 1 \wedge o = 0 \end{cases}$$

- Same as original R-W, except that activation level is now transformed into probability $h(y)$
- But no easy way to analyze stochastic learning process (batch training \neq expected value of single-item training)
- Less robust for highly predictable outcomes $\rightarrow \mathbf{w}$ diverges

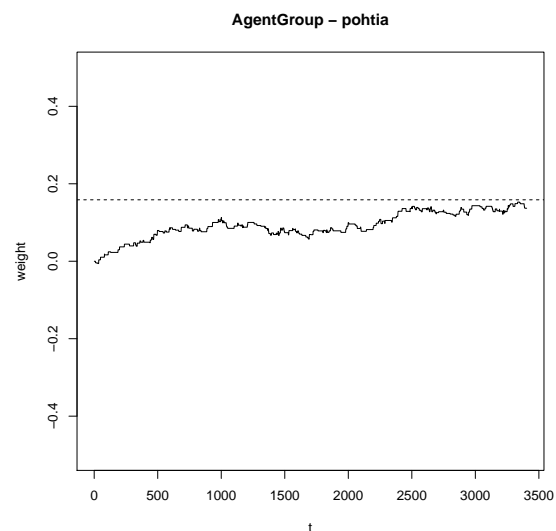
Outline

- 1 Introduction
 - Naïve Discriminative Learning
 - An example
- 2 Mathematics
 - The Rescorla-Wagner equations
 - The Danks equilibrium
 - NDL vs. the Perceptron vs. least-squares regression
- 3 Insights
 - Theoretical insights
 - Empirical observations
 - Conclusion

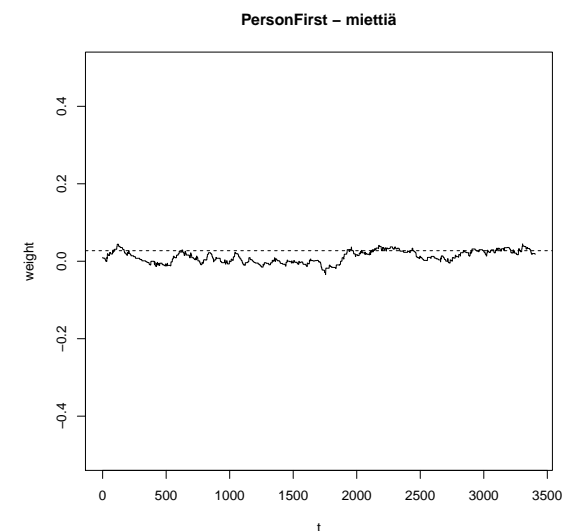
Empirical questions

- How much data is needed for R-W learning convergence with the Danks equilibria
- Are there cases where we observe non-convergence between the R-W learning associations and Danks equilibria – if yes, why?
- Does NDL accuracy always improve with increasing cues? If not, why?

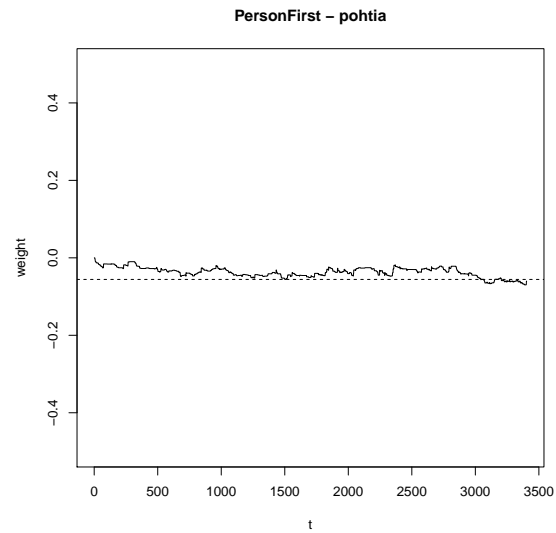
Non-equivocal positive assoc.: convergence



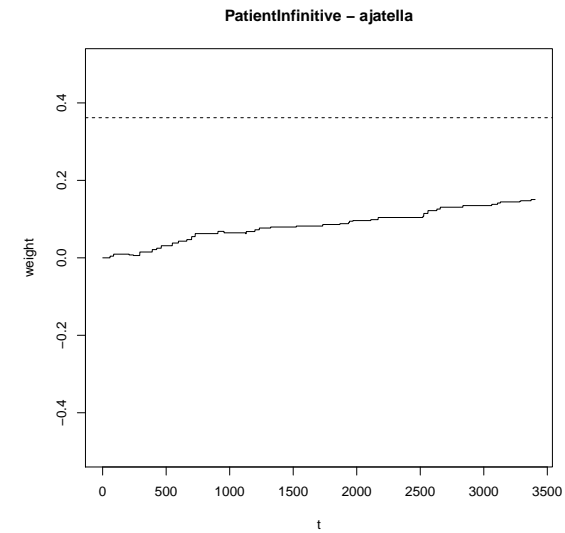
Non-equivocal positive assoc.: convergence with 1x data



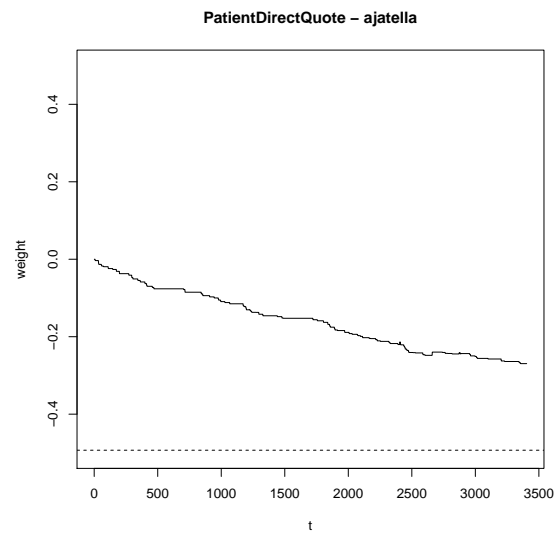
Non-equivocal negative assoc.: convergence with 1x data



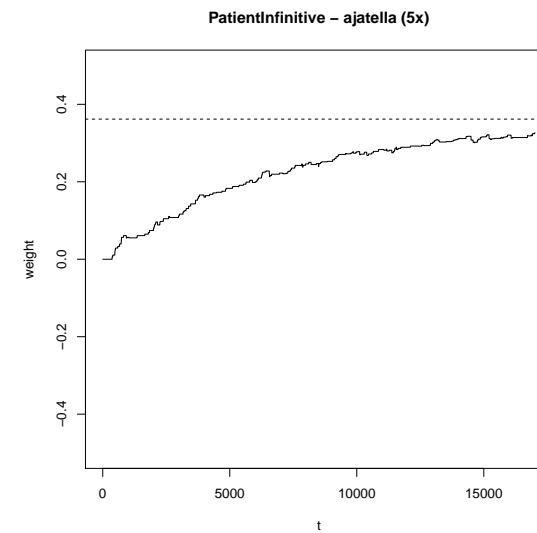
Near-perfect positive assoc.: non-convergence with 1x data



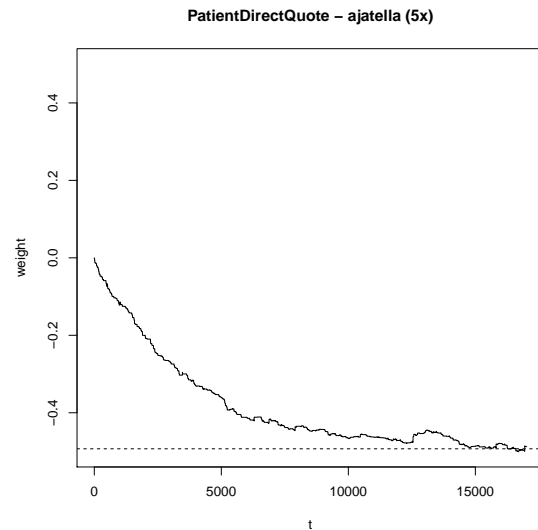
Near-perfect neg. assoc.: non-convergence with 1x data



Near-perfect positive assoc.: convergence with 5x data



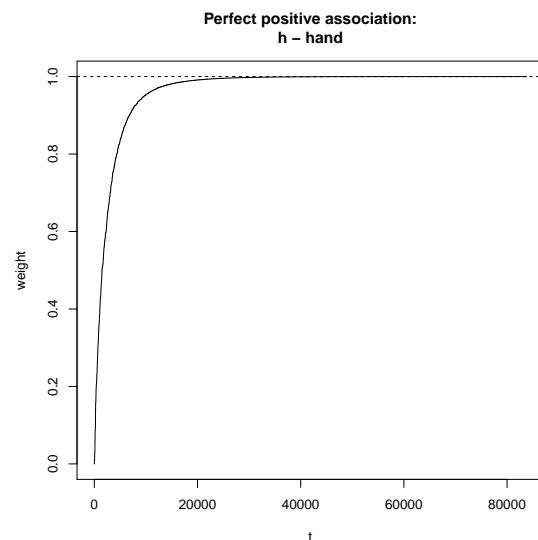
Near-perfect negative assoc.: convergence with 5x data



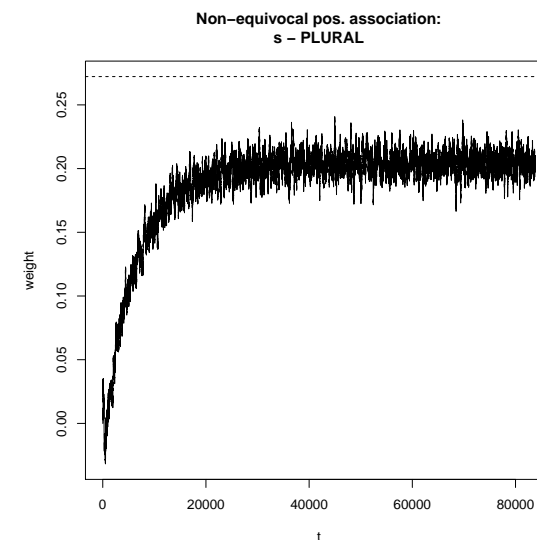
Convergence vs. non-convergence – artificial data: plurals

WordForm	Frequency	Outcomes	Cues
hand	10	hand_NIL	h_a_n_d
hands	20	hand_PLURAL	h_a_n_d_s
land	8	land_NIL	l_a_n_d
lands	3	land_PLURAL	l_a_n_d_s
and	35	and_NIL	a_n_d
sad	18	sad_NIL	s_a_d
as	35	as_NIL	a_s
lad	102	lad_NIL	l_a_d
lad	54	lad_PLURAL	l_a_d
lass	134	lass_NIL	l_a_s_s

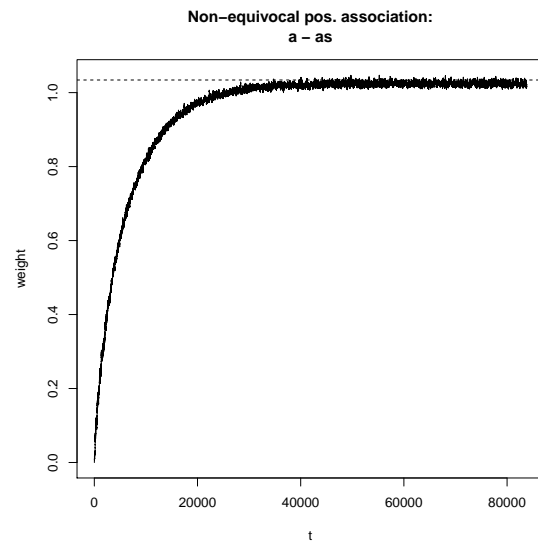
Perfect association – convergence



Non-equivocal positive association – non-convergence



Non-equivocal positive association – convergence



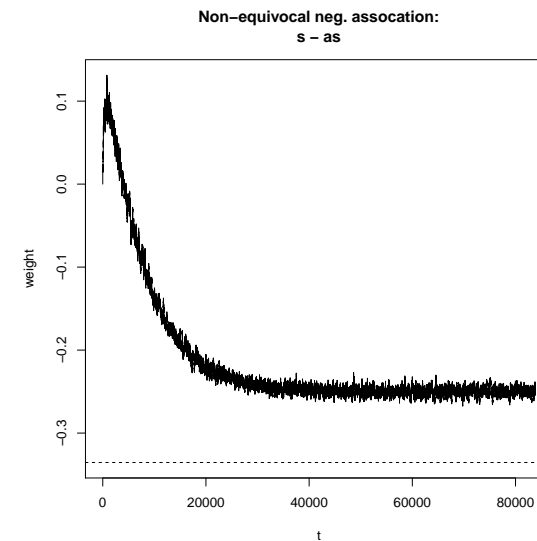
S. Evert & A. Arppe

NDL: Theory & Experiments

Tübingen, 6 Nov 2015

57 / 63

Non-equivocal negative association – non-convergence



S. Evert & A. Arppe

NDL: Theory & Experiments

Tübingen, 6 Nov 2015

58 / 63

Outline

- 1 Introduction
 - Naïve Discriminative Learning
 - An example
- 2 Mathematics
 - The Rescorla-Wagner equations
 - The Danks equilibrium
 - NDL vs. the Perceptron vs. least-squares regression
- 3 Insights
 - Theoretical insights
 - Empirical observations
 - Conclusion

S. Evert & A. Arppe

NDL: Theory & Experiments

Tübingen, 6 Nov 2015

59 / 63

Summary

stochastic	=	batch	=	L_2 regression
NDL	=	SLP		

S. Evert & A. Arppe

NDL: Theory & Experiments

Tübingen, 6 Nov 2015

60 / 63

Acknowledgements 1/2



The mathematical analysis was fuelled by large amounts of coffee and cinnamon rolls at Cinnabon, Harajuku, Tokyo

Follow me on Twitter: @RattiTheRat

Acknowledgements 2/2



The empirical analyses were conducted in the natural environment of Ninase, Saaremaa, Estonia.

References I

- Danks, David (2003). Equilibria of the Rescorla-Wagner model. *Journal of Mathematical Psychology*, **47**, 109–121.
- Rescorla, Robert A. and Wagner, Allen R. (1972). A theory of Pavlovian conditioning: Variations in the effectiveness of reinforcement and nonreinforcement. In A. H. Black and W. F. Prokasy (eds.), *Classical Conditioning II: Current Research and Theory*, chapter 3, pages 64–99. Appleton-Century-Crofts, New York.
- Rosenblatt, Frank (1958). The perceptron: A probabilistic model for information storage and organization in the brain. *Psychological Review*, **65**(6), 386–408.
- Widrow, Bernard and Hoff, Marcian E. (1960). Adaptive switching circuits. In *IRE WESCON Convention Record*, pages 96–104, New York. IRE.