# Kaili's very official dissertation draft

For purposes of graduating in the near-ish future

## 4/5. Learning biases and optimal learning conditions

Yellow highlighting indicates updates/follow-up needed.

Notes highlighted in blue throughout.

Prior to this chapter, I will have introduced:
- The balto-finnic languages / vowel pattern typology
- Representative sample languages (patterns) for the above
- My constraint set & OT analysis of the sample languages
- Learning algorithms, in particular GLA-type learners
  - Definition of idempotence
  - Definition of θ-notation for ranking values
  - Explain my method of calculating "average frequency of correct results"
- My methods for simulations - python script etc
- Simulated learning data for the sample languages

Simulating acquisition of a grammar with an algorithmic learner involves many potential variables, parameters, and biases. In this chapter I introduce those factors that are relevant to learning Balto-Finnic vowel patterns, and discuss the impact that each has on the sample languages Finnish, North Estonian, and North Seto.

In Section 4.1 I introduce the factors that are assumed to remain constant, present learning results given these foundational assumptions, and discuss challenges to be overcome from this starting point. Sections 4.2, 4.3, and 4.4 introduce various biases, each of which contributes to solving one of the problems previously identified while simultaneously uncovering more subtle obstacles not previously apparent. Section 4.2 investigates the specific-over-general faithfulness bias; Section 4.3 focuses on options for varying the promotion rate applied at each learning update; and Section 4.4 explores the general-over-specific markedness bias. Finally, in Section 4.5, I summarize all of the learning simulations performed with various combinations of values for the biases introduced in the preceding sections and generalize a set of ideal conditions for learning these Balto-Finnic languages.

# 4.1. Learning simulations with default settings

The GLA, as specified by Boersma and Hayes (2001) and discussed in Section 3.X, describes the general procedure for this type of gradual, error-driven learning. The bare bones of the learning algorithm as described lay the foundation for additional potential parameters or biases to be included.

## 4.1.1. Learning parameters/biases assumed to remain constant

For the purposes of this project, there are a number of parameters that I considered allowing to vary, but ultimately decided to keep constant.

The first of these determines whether all constraints have the same initial ranking values or if faithfulness constraints should start lower than markedness constraints. I consistently apply a low-faithfulness bias in these simulations.[1] The bias toward low initial faithfulness is widely used in the learning literature, as it helps to ensure that the acquired grammar is as restrictive as possible; that is, it mitigates the Subset Problem (Angluin 1980, Baker 1979). Readers can find more detailed discussion in, e.g., Gnanadesikan (1995), Smolensky (1996), Hayes (2004), Prince and Tesar (2004), and Jesney and Tessier (2011). The default implementation of this bias in this project is to set the initial ranking value of faithfulness constraints to be 0, and that of markedness constraints to be 100. There are some other biases discussed in later parts of this chapter that will set initial markedness values to be different from the default; however, these will continue to preserve the overarching low-faithfulness bias.

---

[1] I also ran a small number of exploratory simulations in which faithfulness constraints experience a more persistent downward bias, being demoted at regular intervals through the learning process. However, these experiments did not produce any promising results so I set the notion of "gravity" aside and did not pursue it any further.

The second parameter that will remain constant is that of demotion eligibility; that is, whether all loser-preferrering constraints get demoted at each learning update, or just the undominated ones. In all learning simulations, I demote *all* such constraints rather than choosing to run some simulations in which only *undominated* loser-preferrers get demoted. Boersma and Hayes (2001) find that demoting only undominated losers caused the GLA to fail on their test data. On the other hand, Magri (2012) shows that doing so can prevent the learner from converging efficiently. Suffice to say that even if choosing to demote all lower-preferrers affects the learner's ability to converge *efficiently*, it will not affect whether or not the learner converges *at all*.

The fourth constant is the number of learning trials, which is fixed at 20,000 for each simulation. All simulations described herein converged well before iterating through this many trials, providing a long enough timeline to ensure that even the odd later error (caused by a particularly noisy evaluation) did not affect the overall ranking.

The fifth parameter that will remain constant is the permission of negative ranking values. Since I am working with ranked (classic OT) rather than weighted (e.g. Harmonic Grammar) constraints, there is no particular concern associated with negative ranking values; all ranking values are converted to relative ordinal rankings at evaluation so the actual numerical values themselves are irrelevant. For example, the values $\{\theta_{C1} = 100, \theta_{C2} = 50\}$ produce the exact same ranking as the values $\{\theta_{C1} = -25, \theta_{C2} = -75\}$. Given this fact, the default OTSoft (Hayes et al, 2013a) approach for GLA learning is used; that is, to permit demotion of constraints even when the resulting ranking value is negative. [I was sure that I'd seen someone else make the argument that negative values are ok when using GLA with classic OT, but I can't for the life of me find it. Does this ring a bell for anyone else? I'm happy to explain it myself but it seems silly to do so if it already exists elsewhere!]

The last few parameters that are held constant across simulations are the organization of learning trials into stages, evaluation noise, and the plasticity function. None of the results I discuss appear to depend on changes to these settings, so the default OTSoft (Hayes et al, 2013a) assumptions for GLA learning are used; they are summarized in Table 1.

| Parameter | Stage 1 | Stage 2 | Stage 3 | Stage 4 |
|---|---|---|---|---|
| Number of learning trials | 5000 | 5000 | 5000 | 5000 |
| Evaluation noise | 2 | 2 | 2 | 2 |
| Plasticity | 2 | 0.2 | 0.02 | 0.002 |

Table 1. Invariant settings for GLA learning.

## 4.1.2. Simulation results

Initial learning simulations use all of the basic parameters (those described in Section 4.1.1) at their default settings, with no additional biases or parameters introduced. Learner A is defined with the settings in Table 2.

| Parameter | Setting |
|---|---|
| All basic parameters | Default |

Table 2. Parameter settings for Learner A.

Under these default conditions, learning simulations for all three sample languages fail to acquire the target grammars, producing instead fully-faithful grammars. Such fully faithful grammars succeed 100% of the time on the input forms, since we assume idempotence. However, they are very poorly equipped to deal with illicit test forms. Test results are summarized in Table 3.

| Language | Average frequency of correct outputs |
|---|---|
| Finnish | 0.2654 |
| North Estonian | 0.2455 |
| North Seto | 0.2982 |
| Overall | 0.2697 |

Table 3. Summary of results from simulations with Learner A.

Table 4 shows the final ranking values for a selection of crucial constraints, after learning from simulated Finnish data. Id(Bk) has risen to the top of the rankings, and even though its distance from $*B_2$, the relevant no-disagreement constraints, and others with values near 100 is small enough that evaluation noise might cause it to swap rankings with one of its neighbours, the grammar does not align with the crucial rankings proposed in Section 3.X:

$$*B_2, *F_3...\underline{B}_5, *F_3\underline{B}_5, *\underline{B}_5...F_3, *\underline{B}_5F_3, Id(Bk)\sigma_1 >> Id(Bk)$$

| Constraint | Final ranking value |
|---|---|
| Id(Bk) | 116 |
| $*B_2$ | 110 |
| $*F_3...\underline{B}_5$ | 106 |
| $*\underline{B}_5...F_3$ | 106 |
| … | … |
| $*F_3\underline{B}_5$ | 104 |
| $*\underline{B}_5F_3$ | 104 |
| … | … |
| $Id(Bk)\sigma_1$ | 80 |

Table 4. Excerpt of final ranking values for Finnish after simulation with Learner A.

The final ranking values for a selection of crucial constraints, after learning from simulated North Estonian inputs, are shown in Table 5. Id(Bk) has risen to the top of the rankings, and even

though its distance from $*B_1$ (and other constraints with values near 100) is small enough that evaluation noise might cause it to swap rankings with one of its neighbours, the crucial rankings Id(Bk)$\sigma_1$ >> $*F_3$, $*B_2$ >> Id(Bk), proposed in Section 3.<mark>X</mark>, are certainly not achieved.

| Constraint | Final ranking value |
|---|---|
| Id(Bk) | 116 |
| $*B_1$ | 104 |
| $*\underline{B}_5\ldots F_3$ | 104 |
| $*\underline{B}_5F_3$ | 104 |
| … | … |
| $*F_3$ | 102 |
| $*B_2$ | 102 |
| … | … |
| $*F_5\ldots\underline{B}_2$ | 100 |
| $*F_5\underline{B}_2$ | 100 |
| … | … |
| Id(Bk)$\sigma_1$ | 92 |
| … | … |

Table <mark>5</mark>. Excerpt of final ranking values for North Estonian after simulation with Learner A. A selection of no-disagreement constraints are included here for the purpose of comparison with results of simulations discussed in subsequent sections.

Learning from simulated North Seto data results in final ranking values for a selection of crucial constraints shown in Table <mark>6</mark>. Id(Bk) has risen to the top of the rankings, and even though its distance from $*B_1$, the relevant no-disagreement constraints, and others with values near 100 is small enough that evaluation noise might cause it to swap rankings with one of its neighbours, the grammar does not meet the crucial target rankings proposed in Section 3.<mark>X</mark>:

$*F_4\ldots\underline{B}_5$, $*F_4\underline{B}_5$, $*\underline{B}_5\ldots F_4$, $*\underline{B}_5F_4$, Id(Bk)syl1 >> $*B1$ >> Id(Bk)

| Constraint | Final ranking value |
|---|---|
| Id(Bk) | 116 |
| $*F_4\ldots\underline{B}_5$ | 106 |
| $*\underline{B}_5\ldots F_4$ | 106 |
| $*F_4\underline{B}_5$ | 106 |
| … | … |
| $*B_1$ | 102 |
| … | … |
| $*\underline{B}_5F_4$ | 100 |
| … | … |
| Id(Bk)$\sigma_1$ | 80 |

Table <mark>6</mark>. Excerpt of final ranking values for North Seto after simulation with Learner A.

## 4.1.3. Discussion and challenges

There are several obstacles that must be addressed on the way to acquiring better – even excellent – final grammars. However, in the results shown above in Section 4.1.2, not all of the challenges are apparent; some only become clear as the initial problems are resolved. In this section I discuss those that are immediately identifiable, and leave the others to be discussed and addressed in subsequent sections of this chapter.

With respect to the results presented in Section 4.1.2, the most glaring problem is that Id(Bk) is highest ranked in all three. This means that during the learning process, Id(Bk) rises all the way from its initial value of 0, past all of the markedness constraints starting at 100, to the very top of the rankings. Such grammars are fully faithful and therefore overgenerate to the point of excluding nothing.

The reason for Id(Bk)'s rise all the way to the top of the rankings is due to a lack of *a priori* relative ranking of the specific and general faithfulness constraints (specifically, an obligation for Id(Bk)$\sigma_1$ to outrank Id(Bk)) combined with the assumption of idempotence. Each time the learner encounters an error, Id(Bk) is always a winner-preferring constraint, since the underlying form is assumed to be identical to the heard surface form. The ERC matrix in Table 7 shows that for a surface form of /o..ɑ/ (grammatical in all three of the sample languages), Id(Bk) is a winner-preferrer for *any* error and Id(Bk)$\sigma_1$ is a winner-preferrer for only those errors involving the first syllable. Thus an error in the first syllable will result in promotion of both faithfulness constraints, but any error past the first syllable will result in promotion of only the general one. Since a crucial element of all three sample languages' target grammars is for Id(Bk)$\sigma_1$ to outrank Id(Bk), this ranking will never be achieved and the learner will only stop making errors once Id(Bk) has been promoted all the way to the top of the rankings.

| input | candidates | markedness constraints | Id(Bk) | Id(Bk)$\sigma_1$ |
|-------|-----------|------------------------|--------|------------------|
| /o..ɑ/ | o..ɑ ~ o..æ | … | W | |
| /o..ɑ/ | o..ɑ ~ ø..ɑ | … | W | W |
| /o..ɑ/ | o..ɑ ~ ø..æ | … | W | W |

Table 7. ERC matrix demonstrating that under assumption of idempotence, all learning errors have Id(Bk) as a winner-preferring constraint.

Addressing the relative ranking of specific vs general faithfulness constraints is not the only obstacle to successful learning of grammars for the sample languages. However, as it is the only one apparent under the learning conditions presented in Section 4.1, it must be addressed before any others can be revealed. Section 4.2 presents a solution for this problem.