

From Performance Errors to Optimal Competence Learnability of OT and HG with Simulated Annealing

TAMÁS BIRÓ

University of Amsterdam, Netherlands

t.s.biro@uva.nl



COMPETENCE as linguistic optimization

Generative linguistics as an optimization problem: how to map underlying form U onto surface form $SF(U)$?

$$SF(U) = \arg \text{opt}_{w \in \text{Gen}(U)} H(w)$$

Target function ("Harmony") $H(w)$ derived from elementary functions $C_i(w)$ ("constraints" – a misnomer!):

1. Hard constraints: $H(w) = C_1(w) + C_2(w) + \dots + C_N(w) \rightarrow$ P&P
2. Weighted sum: $H(w) = g_N \cdot C_N(w) + g_{N-1} \cdot C_{N-1}(w) + \dots + g_1 \cdot C_1(w) \rightarrow$ HG
3. Exponential weights: $H(w) = -C_N(w) \cdot q^N - C_{N-1}(w) \cdot q^{N-1} - \dots - C_1(w) \cdot q \rightarrow$ q -HG
4. OT tableau row: $H(w) = \begin{bmatrix} C_N(w) & C_{N-1}(w) & \dots & C_1(w) \end{bmatrix} \rightarrow$ OT

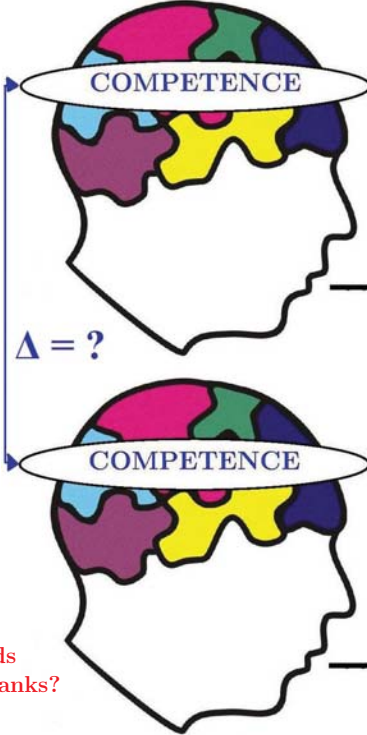
Rank r_i for each constraint C_i .

We focus on OT and q -HG.

Grammar = set of ranks.

Sort constraints by rank:

- OT hierarchy: $C_i \gg C_j$ iff $r_i > r_j$.
- HG weights: $g_i = q^{K_i}$ such that $K_i > K_j$ iff $r_i > r_j$, and $K_i \in \{1, \dots, N\}$.



Can you look into the brain of the teacher?

Online learning from teacher's performance.

Any need to learn until convergence on the grammar? Until the learner finds the teacher's set of ranks?

PERFORMANCE or production as implementation

1. <i>Competence</i> : the static knowledge	grammatical forms	(explained by) grammar
2. <i>Mental computation</i> in the brain	produced forms	implementation of grammar
3. <i>Performance</i> in its outmost sense	produced forms	phonetics, pragmatics, etc.

Cf. Biró (2006:43); Smolensky and Legendre (2006:vol. 1. p. 228). Ways to implement HG and OT:

- **Grammatical**: return the most harmonic candidate (exhaustive search; FS-OT; dynamic programming).
- **Simulated annealing**: return local optima, depending on *cooling schedule* (t_{step} : step by which temperature is decreased in each iteration, "speed").
 - HG: sa converges to gr (frequency of global optimum converges to 1) if $t_{\text{step}} \rightarrow 0$ (more iterations).
 - OT: grammatical forms, irregular forms and fast speech forms are returned (Biró 2007):
 - * Grammatical form: globally optimal.
 - * Fast speech form: globally not optimal; its frequency converges to 0 if $t_{\text{step}} \rightarrow 0$.
 - * Irregular form: globally not optimal; its frequency converges to some positive value if $t_{\text{step}} \rightarrow 0$.

Simulated Annealing

Originating in physics, *simulated annealing* (Boltzmann Machines or stochastic gradient ascent) is a widespread heuristic technique for combinatorial optimization. A random walk is performed on the search space until being trapped in the global or in another local optimum. If target function is real-valued, as in HG, then the slower the speed of the algorithm, the closer to 1 the probability of finding the global optimum. Biró (2006) demonstrates how to apply simulated annealing in the non-real-valued OT case. In this SA-OT, irregular forms also emerge, which persist even at slow speed.

Gradual online LEARNING needs production (that is, performance)

Repeated error-driven updates of the constraint ranks r_i , until convergence:

- Initially: fix random target grammar, fix underlying form, initial random grammar for learner.
- **Error-driven**: "winner" produced by target grammar vs. "loser" produced by learner's current grammar.
- **Update rule**: update the rank r_i of every constraint C_i , depending on whether C_i prefers the winner or the loser. Two approaches ($\epsilon = 0.1$, while ranks are initially random numbers between 0 and $N = 15$):
 - Boersma (1997): increase rank by ϵ if winner-prefering; decrease rank by ϵ if loser-prefering constraint.
 - Magri (2009): increase rank of all winner-prefering constraints by ϵ ; decrease rank of highest ranked loser-prefering constraint by $W \cdot \epsilon$, where W is the number of winner-prefering constraints.

Converging on performance: Δ is Jensen-Shannon divergence

- **Convergence criterion**: JSD between sample produced by target grammar and sample produced by learner's current grammar \leq average JSD of two samples produced by target grammar. (Sample size = 100).

A measure of the "distance" of two distributions: $JSD(P||Q) = \frac{D(P||M) + D(Q||M)}{2}$ where $D(P||Q) = \sum_x P(x) \log \frac{P(x)}{Q(x)}$ (relative entropy, Kullback-Leibler divergence), and $M(x) = \frac{P(x)+Q(x)}{2}$.

- Symmetric: $JSD(P||Q) = JSD(Q||P)$. Finite and non-negative: $0 \leq JSD(P||Q) \leq 1$.
- $JSD(P||Q) = 0$ if and only if $P(x) = Q(x)$, $\forall x$. $JSD(P||Q) = 1$ if and only if $P(x) \cdot Q(x) = 0$, $\forall x$.
- Same language: $JSD(L_t||L_t) = 0$. Not a single overlap: $JSD(L_t||L_t) = 1$.

String Grammar: "toy grammar" imitating typical OT phonology:

- *Candidates*: $\text{Gen}(U) = \{0, 1, \dots, P-1\}^L$. We have used $L = P = 4$: 0000, 0001, 0120, 0123, ..., 3333.
- *Neighborhood structure* on this candidate set: w and w' neighbors iff one basic step transforms w to w' . *Basic step*: change exactly one character $\pm 1 \pmod{P}$ (cyclicity). Each neighbor with equal probability. Example: neighbors of 0123 are exactly 1123, 3123, 0023, 0223, 0113, 0133, 0122 and 0120.
- *Constraints* (for all $n \in \{0, 1, \dots, P-1\}$):
 - No- n (number of character n in string): $*n(w) := \sum_{i=0}^{L-1} (w_i = n)$.
 - No-initial- n : $*\text{INITIAL}n(w) := (w_0 = n)$. No-final- n : $*\text{FINAL}n(w) := (w_{L-1} = n)$.
 - Assimilation (number of different adjacent character pairs): $\text{ASSIM}(w) := \sum_{i=0}^{L-2} (w_i \neq w_{i+1})$.
 - Dissimilation (number of identical adjacent character pairs): $\text{DISSIM}(w) := \sum_{i=0}^{L-2} (w_i = w_{i+1})$.
 - Faithfulness to underlying form U (using pointwise distance modulo P): $\text{FAITH}(w) := \sum_{i=0}^{L-1} d(U_i, w_i)$ where $d(a, b) = \min(|(a-b) \bmod P|, |(b-a) \bmod P|)$.

OBSERVATIONS: very long tail (significance based on Wilcoxon rank-sum test)

- Performance errors make grammars more difficult to learn: $\text{gramm} < 0.1\text{-sa} < 1\text{-sa}$.
- But reversed for HG and SA (either significant, or not significant tendency). Why?
- Magri's update rule (M) quicker than Boersma's (B) (extremely significant). Due to larger update steps?
- Grammar type (OT vs. q -HG): only minor influence ("hardly any" and "small, but very significant"). OT much easier to learn than 1.5-HG (significant difference for sa cases). NB: also quicker to produce.

Does learning speed depend on initial grammar? On order of learning data?

New experiment: Run two learners learning the same target grammar:

- With same initial grammar: strong correlation in their nr. of learning steps. Learning data not the same: slightly decreased correlation.
- With different initial grammars: correlation (almost) lost, large differences in learning time.

Cf. long tail: children must start with same initial grammar, but need not receive same (correct or erroneous) data (if algorithm is correct). SLI is being born with initial grammar causing 10-20 times higher learning steps?

Stability of the algorithms: if learner reached target, will she stay there? Much improvement needed.

Experiment: Measuring number of learning steps

2000 times learning (rnd target, rnd underlying form) per grammar type, production method and learning method. Distribution of the number of learning steps until convergence: 1st quartile; median; 3rd quartile; 90th percentile

	OT	10-HG	4-HG	1.5-HG
gramm. M	13; 27; 45; 67	13; 28; 46; 70	12; 27; 48; 69	15; 30; 47; 67
B	23; 43; 65; 102	22; 41; 64; 107	22; 42; 64; 107	23; 40; 60; 90
sa, M	53; 109; 233; 497	63; 140; 328; 1081	60; 148; 366; 1517	83; 199; 508; 1302
$t_{\text{step}} = 0.1$ B	80; 171; 462; 1543	92; 240; 772; 2512	92; 239; 785; 2633	117; 290; 694; 2006
sa, M	64; 131; 305; 1022	62; 134; 304; 1127	63; 137; 329; 1278	72; 163; 437; 229
$t_{\text{step}} = 1$ B	90; 212; 560; 1906	92; 233; 572; 3116	84; 212; 646; 3005	101; 242; 616; 2001

References

- T. Biró (2006). *Finding the Right Words*. PhD thesis, University of Groningen. ROA-896.
 T. Biró (2007). 'The benefits of errors'. In *Proc. Workshop Cognitive Aspects of Comp. Lang. Acquisition*, ACL. ROA-929.
 P. Boersma (1997). 'How we learn variation, optionality, and probability'. IFA Proceedings 21: 43-58.
 G. Magri (2009). 'New update rules for on-line algorithms for the Ranking problem in OT'. Handout. LMA workshop, DGIS 31.
 P. Smolensky and G. Legendre (eds.) (2006). *The Harmonic Mind*. MIT Press, Cambridge.

From Performance Errors to Optimal Competence Learnability of OT and HG with Simulated Annealing

Tamás Biró

ACLC, University of Amsterdam, t.s.biro@uva.nl

A self-evident, and yet too often ignored fact about (child) language acquisition is that the learner acquiring her *linguistic competence* is exposed to the teacher's *linguistic performance* – hence, also to performance errors, fast speech forms, or other variations. The *performance pattern*, which may be more complex than “simple random noise”, could in theory render the learning problem extremely difficult, but a clever learning algorithm could also make use of the errors, thereby enriching the allegedly poor stimulus.

The computational approach employed in this paper has a threefold structure. Linguistic competence (both of the teacher, and of the learner) is modelled either by standard Optimality Theory (Prince and Smolensky 1993), or by a symbolic Harmonic Grammar with exponential weights (as discussed, for instance, in Biró 2009a). Performance patterns are produced either by always taking the most harmonic form, or by symbolic simulated annealing (Bíró 2006), an algorithm introducing performance errors as a function of the “speech rate”. Finally, online learning employs either Paul Boersma's update rule (Boersma 1997), or Giorgio Magri's (2009) one.

The grammar (“phenomenon”) studied is the abstract string grammar proposed by Biró (2007), arguably mimicking a simple but typical phonological grammar. As several constraint rankings or weight families correspond to the same language, the learner is not expected to converge to the teacher's competence (grammar), but to his performance (distribution of forms). In particular, the learner's distance from the teacher is measured by the *Jensen-Shannon divergence* between a sample of the teacher's performance pattern and a sample of the learner's performance pattern. The learner is said to have learnt the target language if this distance is smaller than the divergence of two random samples of the same size produced by the teacher.

The table on the poster summarizes the results of an initial experiment. Magri's approach is significantly faster than Boersma's. If performance errors are present, then learning OT is faster than learning HG. Yet, we do not want to draw far-reaching conclusions from this toy grammar. We focus more on methodological issues of this novel approach, such as the “stability” of the learning process, and its dependence on the initial conditions and the order of the learning data.

Work implemented on:

The logo for OTKit, featuring the letters 'ot' stacked above 'kit' in a stylized, lowercase font.

Tools for Optimality Theory

<http://www.birot.hu/OTKit/>

A free, Java based user interface and a Java library

for teaching or linguists wishing to implement their grammars.

Work supported by:



UNIVERSITY OF AMSTERDAM



AMSTERDAM CENTER
FOR LANGUAGE AND
COMMUNICATION



The logo of the Netherlands Organisation for Scientific Research (NWO), featuring the letters 'NWO' in a stylized, red, sans-serif font.
Netherlands Organisation for Scientific Research