

## The OT online model of the acquisition of phonotactics

### Class 2: converge of cautious promotion, and invariants

**Summary** — This class introduces the goal of devising promotion/demotion OT update rules that converge for any input OT compatible tableau; introduces a new family of OT update rules, that perform promotion, but cautiously “not too much”; shows that OT-compatibility of the data entails conic independence of the update vectors used by OT online algorithms; uses this fact to prove convergence for the proposed family of cautious promotion/demotion OT update rules; introduces the use of invariants to characterize the behavior of OT online algorithms. Thus, the main point of this class is that convergent constraint promotion is possible, as far as we cautiously do not “promote too much”.

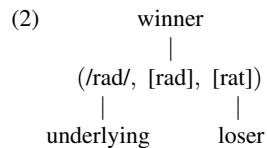
## 1 The core research question

■ **Online models of the early stage.** My modeling task is the *early stage* of the acquisition of phonotactics, characterized by the following input/output properties:

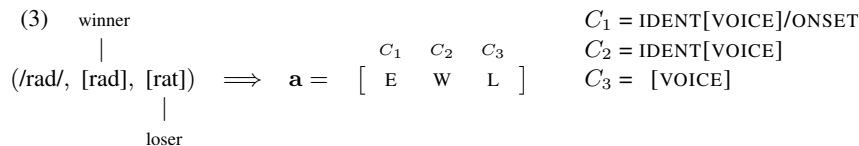
- (1) a. *Properties of the input.* Throughout this stage, the child does not yet have knowledge of morphology and is thus blind to alternations.
- b. *Properties of the output.* By the end of this stage, the child is able to distinguish legal from illegal sequences w.r.t. his target language.

My modeling tool is online algorithms within the framework of OT. Let me quickly review the main properties of OT online algorithms.

■ **First property.** At each time, OT online algorithms take a current *data triplet* consisting of an underlying form, an intended winner form and a loser form, as (2).



A data triplet can be represented as a *comparative row* **a** that specifies whether a constraint is *winner-preferring* (WPC) or *loser-preferring* (LPC) or even, as in (3).



Many data triplets correspond to many comparative rows that can be stacked one on top of the other into a *comparative tableau*.

■ **Second property.** OT online algorithms entertain a current *OT grammar* at every time, represented as a *ranking vector*  $\theta$ , namely an  $n$ -tuple of numerical *ranking values*  $\theta_1, \dots, \theta_n$ , one for every constraint.

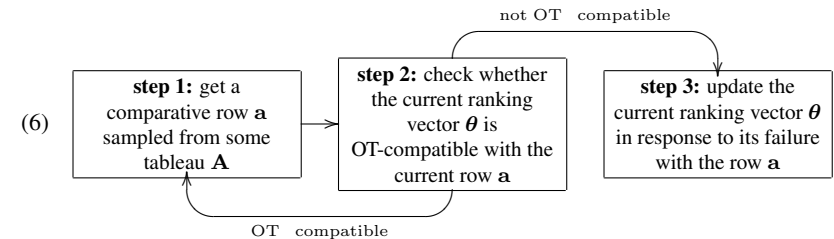
$$(4) \quad \theta = \begin{pmatrix} C_1 & \dots & C_k & \dots & C_n \\ \theta_1, & \dots & \theta_k, & \dots & \theta_n \end{pmatrix}.$$

A ranking vector *represents* a ranking  $\gg$  iff (5) holds for any pair of constraints  $C_h, C_k$ , namely  $\gg$  respects the ordering implicit in the relative size of the ranking values.

$$(5) \quad \text{The ranking value } \theta_h \text{ is larger than the ranking value } \theta_k \implies C_h \gg C_k$$

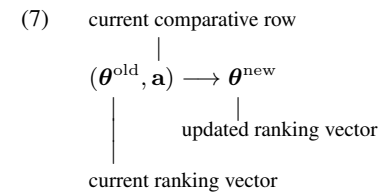
A ranking vector is *OT-compatible* with a comparative tableau iff all the rankings it represents are OT-compatible with it.

■ **Third property.** OT online algorithms maintain a current ranking vector. Initialize it to a pre-defined initial ranking vector. And keep updating it through the three steps (6).



Rows fed to the algorithm at step 1 are sampled from a fixed OT-compatible comparative tableau **A**, called the *input tableau*.

■ **Research question.** OT online algorithms differ for the OT update rule (7) they use in step 3. An *OT update rule* has the following form:



Thus, the core research question is: what is an OT update rule that yields a proper model of the early stage of the acquisition of phonotactics?

## 2 We need constraint promotion

■ **Demotion-only vs promotion/demotion.** OT update rules come in two varieties:

(8)

	demote LPCs	promote WPCs
<i>Demotion-only update rules:</i>	✓	
<i>Promotion-demotion update rules</i>	✓	✓

Class 1 presented the beautiful example (9) of demotion-only update rules, from Tesar and Smolensky (1998) and Boersma (1997):

$$(9) \theta_k^{\text{new}} = \begin{cases} \theta_k^{\text{old}} - 1 & \text{if } C_k \text{ is an undominated LPC} \\ \theta_k^{\text{old}} & \text{otherwise} \end{cases}$$

as well as example (10) of promotion/demotion update rules, from Boersma (1997):

$$(10) \theta_k^{\text{new}} = \begin{cases} \theta_k^{\text{old}} - 1 & \text{if } C_k \text{ is a LPC} \\ \theta_k^{\text{old}} + 1 & \text{if } C_k \text{ is a WPC} \\ \theta_k^{\text{old}} & \text{if } C_k \text{ is an even constraint} \end{cases}$$

■ **We need constraint promotion.** Suppose we try to model the early stage of the acquisition of phonotactics with an OT online algorithm that performs demotion only. Then:

- (11) a. during the early stage, the learner has no access to alternations;
- b. thus, the model posits fully faithful underlying forms;
- c. thus, the faithfulness constraints are never LPCs;
- d. thus, the faithfulness constraints are never re-ranked by demotion-only;

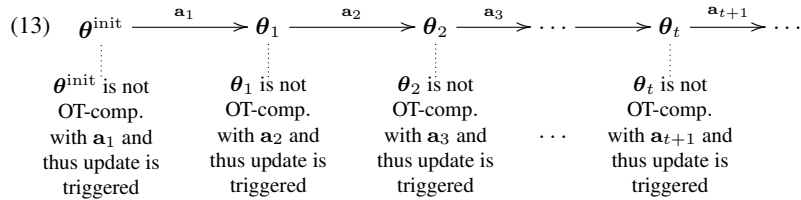
And this cannot be right, as seen in class 1. Thus:

- (12) The OT online model of the early stage cannot work with demotion only; rather, we need constraint promotion in order to re-rank the faithfulness constraints too, despite the fact that they are always WPCs.

But how should we choose a proper promotion/demotion update rule?

### 3 We want universally convergent constraint promotion

■ **Convergence.** The OT online algorithm (6) *converges* on an input comparative tableau **A** iff it is not possible to construct an infinite sequence  $\mathbf{a}_1, \mathbf{a}_2, \dots, \mathbf{a}_t, \dots$  of rows of **A** such that each row  $\mathbf{a}_t$  in the sequence triggers an update if fed to the algorithm at time  $t$ :



To illustrate, recall from class 1 that the OT online algorithm:

- (14) a. does *indeed* converge for *any* OT-compatible input tableau, in the case of the demotion-only update rule (9);

- b. does *not* converge for *every* OT-compatible input tableau, in the case of the promotion/demotion update rule (10): Pater's (2008) counterexample.

We want convergent update rules. But convergent for which tableaux? Two options:

- (15) a. require convergence for *any* OT-compatible input tableau  

we take convergence to follow from the general ranking logic of OT
- b. require convergence for a *special* subset of OT-compatible input tableaux  

we take convergence to follow from peculiar properties of this special family of tableaux

I will call (15a) *universal convergence*. Which option should we pick, (15a) or (15b)?

■ **Ranking problem.** The simplest problem in computational OT is (16).

- (16) *given:* an OT-compatible OT-comparative tableau **A**;
- find:* a ranking  $\gg$  OT-compatible with the tableau **A**.

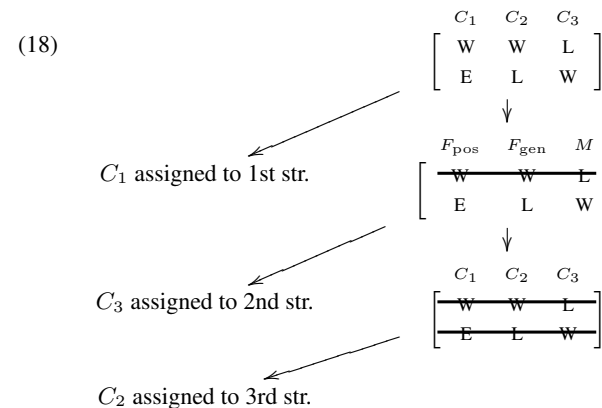
I denote by  $\text{RP}(\mathbf{A})$  the instance of the *Ranking problem* (16) for a tableau **A**.

■ **Claim 1** The Ranking problem  $\text{RP}(\mathbf{A})$  is easily solvable (in time linear in the number of columns of **A**), for *any* OT-compatible tableau **A**.

*Proof.* T&S note that the  $\text{RP}(\mathbf{A})$  can be solved by repeating the two steps (17).

- (17) a. Assign to the currently available top stratum, a constraint that has no L's in the current tableau;
- b. update the current tableau by deleting all rows where the currently ranked constraint has a w.

An example of the procedure (17) is provided in (18): diagonal arrows in (18) correspond to step (17a) and vertical arrows to step (17b).



Algorithm (17) is *Recursive Constraint Demotion*. It is a restatement from an algorithmic perspective of T&S’s characterization of OT-compatibility (claim 1, class 1). □

■ **We want universally convergent update rules.** I now reason as follows:

- (19) a. The OT online algorithm converges on an input tableau **A** iff it solves the corresponding Ranking problem  $RP(\mathbf{A})$ ;
- b. if the Ranking problem  $RP(\mathbf{A})$  cannot be solved for all tableaux **A**, the would have to “settle” for option (15b) of non-universal convergence;
- c. but claim 1 ensures that the Ranking problem  $RP(\mathbf{A})$  can indeed be solved for *any* (OT compatible) tableaux **A**;
- d. thus let’s be bold, and try to go for universal convergence (15a), investigating whether the bare logic of OT is able to ensure convergence.

This is an instance of the modeling strategy (20) outlined in class 1; we will see another example in class 4.

- (20) a. Single out the building blocks of the learning task (e.g. the Ranking problem);
- b. study their complexity (e.g. claim 1);
- c. devise algorithms that solve these problems up to their complexity class.

Thus, this class has two main goals:

- (21) a. develop a family of promotion/demotion OT update rules that are provably universally convergent;
- b. provide some preliminary support for the heuristics of requiring universal convergence, by taking a preliminary look at correctness.

## 4 Cautious constraint promotion

■ **Constraint promotion is hard to get.** T&S warn against constraint promotion:

- (22) “Couldn’t the learner just as easily promote constraints toward the correct hierarchy? The answer is no” Tesar and Smolensky (1998, pp. 244-245).

The only promotion-demotion update rule considered in the literature is Boersma’s (10). Yet, Pater (2008) shows that:

- (23) There are cases where the OT online algorithm run with Boersma’s update rule (10) performs an infinite number of updates and thus never converges.

Thus, constraint promotion has so far eluded the efforts of the computational OT literature. Is it possible to devise convergent promotion/demotion OT update rules?

■ **First step.** Suppose that the current comparative row has a unique L and a unique W:

$$(24) \quad \left[ \dots \begin{array}{c} C_h \\ W \end{array} \dots \begin{array}{c} C_\ell \\ L \end{array} \dots \right]$$

I submit the following intuition:

- (25) a. The comparative row (24) unambiguously says that the WPC  $C_h$  must be ranked above the LPC  $C_\ell$  in the end;
- b. thus, in this case we can be *bold*;
- c. namely, we can promote by the same amount we demote, say 1.

Thus, I suggest to update as follows, namely the same as Boersma’s update rule (10):

- (26) a. Demote the unique LPC by 1;
- b. promote the unique WPC by 1.

■ **Second step.** Suppose the current comparative tableau has a unique L but multiple W’s:

$$(27) \quad \left[ \dots W \quad W \quad W \quad \dots \quad L \quad \dots \right]$$

Boersma’s update rule (10) is not *cautious*, because it overall promotes too much, namely more (i.e.  $1 + 1 + 1 = 3$ ) than it demotes (i.e. 1). Indeed, recall from class 1 that the ranking values kept growing in the case of Pater’s counterexample.

$$(28) \quad \left[ \dots \begin{array}{c} W \\ \text{add 1} \end{array} \begin{array}{c} W \\ \text{add 1} \end{array} \begin{array}{c} W \\ \text{add 1} \end{array} \dots \begin{array}{c} L \\ \text{subtract 1} \end{array} \dots \right]$$

3 times

I submit the following intuition:

- (29) a. This comparative row (27) does *not* unambiguously say which one of the multiple WPCs needs in the end be ranked above the LPC;
- b. thus, we should *not* be bold and promote each WPC by the same amount we demote the unique LPC, as in the preceding case;
- c. rather, we should be *cautious* and split that amount over the multiple WPCs.

In other words, I suggest the following *cautious* update:

$$(30) \quad \left[ \dots \begin{array}{c} W \\ \text{add } \frac{1}{3} \end{array} \begin{array}{c} W \\ \text{add } \frac{1}{3} \end{array} \begin{array}{c} W \\ \text{add } \frac{1}{3} \end{array} \dots \begin{array}{c} L \\ \text{subtract 1} \end{array} \dots \right]$$

3 times

More in general, I suggest the following *cautious* update:

- (31) a. Demote the unique LPC by 1;
- b. promote each WPC by  $\frac{1}{\text{total number of WPCs}}$ .

Of course, if I multiply (31) by the total number of WPCs, nothing changes (OT update rules are invariant by positive rescaling — provided the initial ranking vector has identical components). Thus, I can describe (31) as follows:

- (32) a. Demote the unique LPC by the total number of WPCs;  
 b. Promote each WPC by 1.

■ **Third step.** Suppose now the current comparative tableau has multiple *undominated* L's and multiple W's:

$$(33) \quad \begin{array}{cccc} & C_h & C_k & C' & C'' \\ [ \dots & W & W & \dots & L & L & \dots ] \end{array}$$

This comparative row (33) is *equivalent* to (i.e. OT-compatible with the same rankings as) the two comparative rows (34): the two L's of (33) can be split over two rows.

$$(34) \quad \begin{array}{l} \text{a.} \quad \begin{array}{cccc} & C_h & C_k & C' & C'' \\ [ \dots & W & W & \dots & L & & \dots ] \\ & C_h & C_k & C' & C'' \\ \text{b.} \quad [ \dots & W & W & \dots & & L & \dots ] \end{array} \end{array}$$

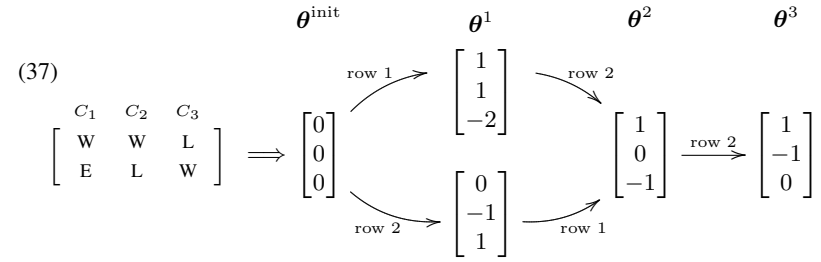
I thus submit the following intuition:

- (35) a. Update triggered by the original row (33) should be equivalent to two consecutive updates triggered by the two new rows (34a) and (34b).  
 b. Update by the first row (34a) can be performed according to (32), since it has a unique LPC: each WPC is promoted by 1 and the unique LPC is demoted by the total number of WPCs.  
 c. Subsequent update by the second row (34b) can be performed again according to (32): each WPC is promoted again by 1 and the unique LPC is demoted by the total number of WPCs.  
 d. In conclusion, each WPC gets promoted by the total number of LPCs and each LPC gets demoted by the total number of WPCs.

Thus, I suggest the following *cautious* update:

- (36) a. Demote all undominated LPCs by the total number of WPCs;  
 b. promote all WPCs by the total number of undominated LPCs.

■ **An example.** The behavior of the OT online algorithm (6) with the new cautious promotion-demotion update rule (36) is illustrated in the diagram (37).



Note the non-monotonic dynamics of the ranking values (they go up and down), which is the hallmark of update rules that perform promotion too.

## 5 Cautious promotion works!

■ **Claim 2** The OT online algorithm (6) run with the new cautious promotion-demotion update rule (36) converges for *any* OT-compatible input tableau.

■ **Generalization.** Let me say that a promotion-demotion update rule is *cautious* iff it does not “promote too much”, in the following sense:

- (38) For any comparative row, the sum of the amounts by which the WPCs are promoted is smaller than or at most equal to the sum of the amounts by which the undominated LPCs are demoted.

Constraint promotion is delicate. T&S are drastic: they suggest that we do not promote at all. It turns out we can be a bit less drastic: we can promote, but cautiously. In fact, convergence actually holds for any cautious update rule (provided it only demotes undominated LPCs and that all WPCs are promoted).

■ **The actual parameterization.** In class 1, we saw that OT-grammars can be parameterized equivalently by *rankings* and *ranking vectors*. The following question thus arises:

- (39) Do speakers assume that OT languages are parameterized by *combinatoric* objects such as rankings or by *continuous* objects such as ranking vectors?

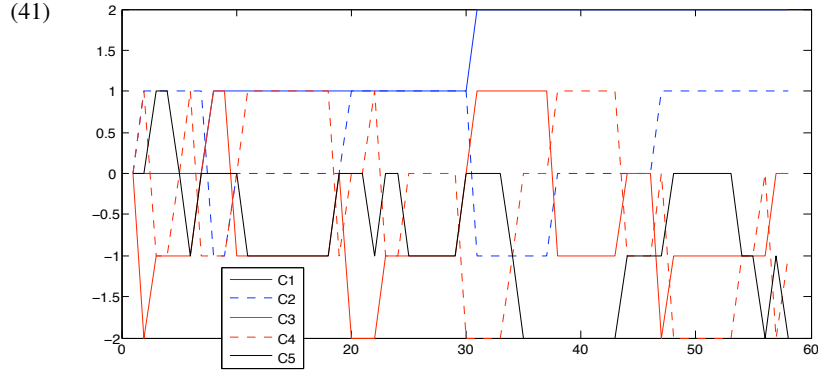
The computational perspective adopted here offers a handle on question (39) — other perspectives being modeling gradualness, gradience, etc. In fact:

- (40) a. *Demotion-only can be implemented with rankings:*  
 for example, T&S’s original demotion-only OT online algorithm represents the current grammar by means of rankings;  
 b. *Promotion/demotion can only be implemented with ranking vectors:*  
 as we have just seen, it is crucial to promote by the proper promotion amount; hence, the rule is inherently numerical and it thus requires ranking vectors.

If indeed promotion is needed to model, say, the early stage of the acquisition of phonotactics, then the answer to (39) is: *continuous ranking vectors*, not discrete rankings.

## 6 Sketch of the proof of convergence

■ **Preliminaries.** An example of the dynamics of the current ranking values entertained by the OT online algorithm (6) with the cautious promotion/demotion update rule (36) run on Pater's comparative tableau (see class 1) with the rows sampled uniformly is (41).

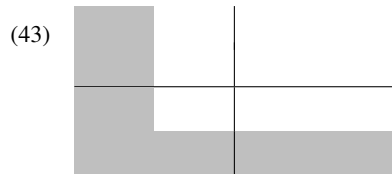


This dynamics is complicated and non-monotonic, contrary to the case of demotion-only update rules. Thus, the proof of convergence is more involved. It has four steps.

■ **First step.** The following fact holds:

(42) The ranking values entertained by OT online algorithms with any update rule (that only demotes *undominated* LPCs) cannot get smaller than a given constant.

In other words, the current ranking vector cannot live in the shaded region of (43).

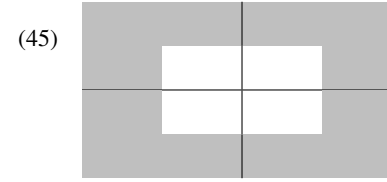


*Proof.* As seen in class 1, T&S show that (42) holds for update rules that perform demotion only (of undominated LPCs). But a close look at their proof reveals that it holds for pretty much any update rule (that demotes only undominated LPCs). □

■ **Second step.** The following fact holds:

(44) The ranking values entertained by OT online algorithms with a *cautious* promotion-demotion update rule cannot get larger than a given constant.

Together with (42), this means that the current ranking vector entertained by the algorithm cannot live in the shaded region in (45).



*Proof.* Let me focus on the specific cautious promotion-demotion update rule (36). The ranking vectors entertained by the algorithm in the run (37) are repeated in (46).

$$(46) \quad \begin{bmatrix} 1 \\ 1 \\ -2 \end{bmatrix}, \begin{bmatrix} 0 \\ 0 \\ 0 \end{bmatrix}, \begin{bmatrix} 1 \\ 0 \\ -1 \end{bmatrix}, \begin{bmatrix} 1 \\ -1 \\ 0 \end{bmatrix}, \begin{bmatrix} 0 \\ -1 \\ 1 \end{bmatrix}$$

Note that the components of each ranking vector in (46) sum up to zero. This is a general property of the new cautious promotion-demotion update rule (36):

(47) The sum of the current ranking values at any given time is zero.

In fact, suppose that the sum  $\sum_{k=1}^n \theta_k^{\text{old}}$  of the current rankings values is zero; suppose that the current comparative tableau has  $w$  WPCs and  $\ell$  undominated LPCs (ULPCs). Then, the sum  $\sum_{k=1}^n \theta_k^{\text{new}}$  of the updated ranking values is zero too:

$$(48) \quad \sum_{k=1}^n \theta_k^{\text{new}} \stackrel{(a)}{=} \sum_{k \in \text{WPC}} \theta_k^{\text{new}} + \sum_{k \in \text{ULPC}} \theta_k^{\text{new}} + \sum_{\text{remaining } k} \theta_k^{\text{new}}$$

split the set  $\{1, \dots, n\}$  that  $k$  runs over into the set WPC of winner-preferrers, the set ULPC of undominated loser-preferrers and their complement

$$\stackrel{(b)}{=} \sum_{k \in \text{WPC}} (\theta_k^{\text{old}} + \ell) + \sum_{k \in \text{ULPC}} (\theta_k^{\text{old}} - w) + \sum_{\text{remaining } k} \theta_k^{\text{old}}$$

use (36) to express the updated ranking values  $\theta_k^{\text{new}}$  in terms of the current ranking values  $\theta_k^{\text{old}}$

$$\stackrel{(c)}{=} \sum_{k=1}^n \theta_k^{\text{old}} + \sum_{k \in \text{WPC}} \ell - \sum_{k \in \text{ULPC}} w$$

by reordering the terms

$$\stackrel{(d)}{=} \sum_{k \in \text{WPC}} \ell - \sum_{k \in \text{ULPC}} w$$

by the hypothesis  $\sum_{k=1}^n \theta_k^{\text{old}} = 0$  that the current ranking values sum to zero

$$= w\ell - \ell w = 0$$

I can now conclude the reasoning as follows:

- (49) a. since the current ranking values cannot become too small, by (42),  
b. since furthermore they must always add up to zero, by (47),  
c. then they cannot become too large either, as stated in (44).

For an arbitrary initial ranking vector and an arbitrary cautious update rule, the sum of the current ranking values is always smaller than a fixed constant (instead than equal to zero); the reasoning in (49) thus carries over.  $\square$

■ **Third step.** The following fact holds:

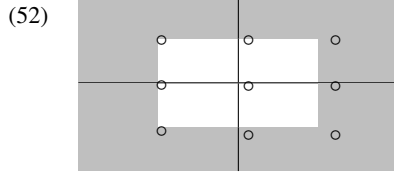
(50) The OT online algorithm with any update rule can never loop, namely it can never entertain again a current ranking vector that it had previously dismissed (as long as all the WPCs are promoted and the input tableau is OT-compatible).

In other words, the OT online algorithm can never walk through a sequence of ranking vectors as (51), whereby the algorithm loops back to  $\theta$  after having dismissed it for  $\theta^\theta$ .

$$(51) \quad \dots \longrightarrow \theta \longrightarrow \dots \longrightarrow \theta^\theta \longrightarrow \dots \longrightarrow \theta \longrightarrow \dots$$

*Proof.* This is the non-trivial step; I will come back to this in the next section.  $\square$

■ **Fourth step.** Since the ranking vectors entertained by the algorithm must live on a lattice and must furthermore live within a bounded region by (42) and (44), then the search space of the algorithm is finite, as illustrated in (52).



Since the algorithm cannot loop by (50), finiteness of the search space entails finite time convergence.

## 7 The third step of the proof, in some detail

■ **Conic combinations.** Consider  $\ell$  vectors  $\bar{\mathbf{a}}_1, \dots, \bar{\mathbf{a}}_\ell$  with the same number  $n$  of components. Consider a combination of these vectors with some coefficients  $\alpha_1, \dots, \alpha_\ell$ :

$$(53) \quad \alpha_1 \bar{\mathbf{a}}_1 + \dots + \alpha_\ell \bar{\mathbf{a}}_\ell = \alpha_1 \begin{bmatrix} \bar{a}_1^1 \\ \vdots \\ \bar{a}_n^1 \end{bmatrix} + \dots + \alpha_\ell \begin{bmatrix} \bar{a}_1^\ell \\ \vdots \\ \bar{a}_n^\ell \end{bmatrix}$$

Here is a little example:

$$(54) \quad 2 \begin{bmatrix} \bar{a}_1 \\ 1 \\ -2 \\ 1 \end{bmatrix} + 4 \begin{bmatrix} \bar{a}_2 \\ 1 \\ 0 \\ 1 \end{bmatrix} + 3 \begin{bmatrix} \bar{a}_3 \\ 0 \\ 1 \\ -1 \end{bmatrix} = \begin{bmatrix} 6 \\ -1 \\ 3 \end{bmatrix}$$

If the coefficients  $\alpha_1, \dots, \alpha_\ell$  are all non-negative, then (53) is a *conic combination*.

■ **Conic independence.** The following two properties (55) and (56) are equivalent.

(55) A conic combination of  $\bar{\mathbf{a}}_1, \dots, \bar{\mathbf{a}}_\ell$  adds up to the null vector iff the coefficients  $\alpha_1, \dots, \alpha_\ell$  are all null:

$$\left. \begin{array}{l} \alpha_1 \geq 0 \\ \vdots \\ \alpha_\ell \geq 0 \\ \alpha_1 \bar{\mathbf{a}}_1 + \dots + \alpha_\ell \bar{\mathbf{a}}_\ell = \mathbf{0} \end{array} \right\} \implies \begin{cases} \alpha_1 = 0 \\ \vdots \\ \alpha_\ell = 0 \end{cases}$$

(56) Two conic combinations of  $\bar{\mathbf{a}}_1, \dots, \bar{\mathbf{a}}_\ell$  add up to the same vector despite the fact that the coefficients  $\alpha_1, \dots, \alpha_\ell$  of one of the two are not smaller than the coefficients  $\alpha_1^\theta, \dots, \alpha_\ell^\theta$  of the latter iff the coefficients are all identical:

$$\left. \begin{array}{l} \alpha_1 \geq \alpha_1^\theta \\ \vdots \\ \alpha_\ell \geq \alpha_\ell^\theta \\ \alpha_1 \bar{\mathbf{a}}_1 + \dots + \alpha_\ell \bar{\mathbf{a}}_\ell = \alpha_1^\theta \bar{\mathbf{a}}_1 + \dots + \alpha_\ell^\theta \bar{\mathbf{a}}_\ell \end{array} \right\} \implies \begin{cases} \alpha_1 = \alpha_1^\theta \\ \vdots \\ \alpha_\ell = \alpha_\ell^\theta \end{cases}$$

Let me show for instance that property (55) entails (56):

$$(57) \quad \alpha_1 \bar{\mathbf{a}}_1 + \dots + \alpha_\ell \bar{\mathbf{a}}_\ell = \alpha_1^\theta \bar{\mathbf{a}}_1 + \dots + \alpha_\ell^\theta \bar{\mathbf{a}}_\ell \implies$$

$$\stackrel{(a)}{\implies} (\alpha_1 - \alpha_1^\theta) \bar{\mathbf{a}}_1 + \dots + (\alpha_\ell - \alpha_\ell^\theta) \bar{\mathbf{a}}_\ell = \mathbf{0}$$

by reordering

$$\stackrel{(b)}{\implies} \begin{cases} (\alpha_1 - \alpha_1^\theta) = 0 \\ \vdots \\ (\alpha_\ell - \alpha_\ell^\theta) = 0 \end{cases}$$

by property (55)

$$\stackrel{(c)}{\implies} \begin{cases} \alpha_1 = \alpha_1^\theta \\ \vdots \\ \alpha_\ell = \alpha_\ell^\theta \end{cases}$$

If either (55) or (56) holds, then  $\bar{\mathbf{a}}_1, \dots, \bar{\mathbf{a}}_\ell$  are called *conically independent*.

■ **OT compatibility entails conic independence.** To warm up, consider the following simple OT-compatible comparative tableau:

$$(58) \quad \begin{array}{c} \text{row 1} \\ \text{row 2} \\ \text{row 3} \end{array} \begin{array}{ccc} C_1 & C_2 & C_3 \\ \left[ \begin{array}{ccc} W & L & W \\ W & & L \\ & W & L \end{array} \right] \end{array}$$

Update triggered by the three rows of (58) according to (36) can be described as follows:

(59)

1st row:	2nd row:	3rd row:
$\begin{bmatrix} W & L & W \end{bmatrix}$	$\begin{bmatrix} W & & L \end{bmatrix}$	$\begin{bmatrix} & W & L \end{bmatrix}$
add 1	add 1	add 1
subtract 2	subtract 1	subtract 1

Equivalently, update triggered by the three rows of (58) according to (36) can be described as adding to the current ranking vector one of the vectors in (68), called the corresponding *update vector*:

(60)

1st row:	2nd row:	3rd row:
$\bar{\mathbf{a}}_1 = \begin{bmatrix} 1 \\ -2 \\ 1 \end{bmatrix}$	$\bar{\mathbf{a}}_2 = \begin{bmatrix} 1 \\ 0 \\ 1 \end{bmatrix}$	$\bar{\mathbf{a}}_3 = \begin{bmatrix} 0 \\ 1 \\ -1 \end{bmatrix}$

Consider a combination of the vectors in (68) that adds up to the null vector. Since the first entry in each vector is positive or null, then the coefficients  $\alpha_1$  and  $\alpha_2$  corresponding to vectors with non-null first entries needs to be 0. And

As noted above, update triggered by one of the three rows of (67) according to (36) consists of adding to the current ranking vector the corresponding update vector in (68).

$$(68) \quad \begin{array}{ccc} \text{1st row:} & \text{2nd row:} & \text{3rd row:} \\ \bar{\mathbf{a}}_1 = \begin{bmatrix} 1 \\ -2 \\ 1 \end{bmatrix} & \bar{\mathbf{a}}_2 = \begin{bmatrix} 1 \\ 0 \\ 1 \end{bmatrix} & \bar{\mathbf{a}}_3 = \begin{bmatrix} 0 \\ 1 \\ -1 \end{bmatrix} \end{array}$$

Furthermore, consider the three coefficients  $\alpha_1^t, \alpha_2^t, \alpha_3^t$  as follows:

$$(69) \quad \begin{array}{ccc} \alpha_1^t & \alpha_2^t & \alpha_3^t \\ \parallel & \parallel & \parallel \\ \# \text{ of updates up to} & \# \text{ of updates up to} & \# \text{ of updates up to} \\ \text{time } t \text{ triggered by the} & \text{time } t \text{ triggered by the} & \text{time } t \text{ triggered by the} \\ \text{1st row of tableau (67)} & \text{2nd row of tableau (67)} & \text{3rd row of tableau (67)} \end{array}$$

Then, the ranking vector  $\theta^t$  entertained by the algorithm at time  $t$  can be described as follows (we already saw an analogous description, in the explanation of Pater's counterexample in class 1).

$$(70) \quad \begin{array}{ccc} \# \text{ of updates in response to the} & \# \text{ of updates in response to the} & \\ \text{1st row of the tableau (67)} & \text{3rd row of the tableau (67)} & \\ \theta^t = \begin{bmatrix} \theta_1^t \\ \theta_2^t \\ \theta_3^t \end{bmatrix} = \alpha_1^t \begin{bmatrix} 1 \\ -2 \\ 1 \end{bmatrix} + \alpha_2^t \begin{bmatrix} 1 \\ 0 \\ -1 \end{bmatrix} + \alpha_3^t \begin{bmatrix} 0 \\ 1 \\ -1 \end{bmatrix} & & \\ \text{update vector that} & \text{update vector that} & \\ \text{corresponds to the 1st} & \text{corresponds to the 3rd} & \\ \text{row of the tableau (67)} & \text{row of the tableau (67)} & \\ \text{according to rule (36)} & \text{according to rule (36)} & \end{array}$$

Suppose now that the algorithm entertains the same ranking vector  $\theta^{t'} = \theta^{t''}$  at time  $t^0$  and at some later time  $t^{00}$ . This means that:

$$(71) \quad \underbrace{\alpha_1^{t'} \begin{bmatrix} 1 \\ -2 \\ 1 \end{bmatrix} + \alpha_2^{t'} \begin{bmatrix} 1 \\ 0 \\ -1 \end{bmatrix} + \alpha_3^{t'} \begin{bmatrix} 0 \\ 1 \\ -1 \end{bmatrix}}_{\theta^{t'}} = \underbrace{\alpha_1^{t''} \begin{bmatrix} 1 \\ -2 \\ 1 \end{bmatrix} + \alpha_2^{t''} \begin{bmatrix} 1 \\ 0 \\ -1 \end{bmatrix} + \alpha_3^{t''} \begin{bmatrix} 0 \\ 1 \\ -1 \end{bmatrix}}_{\theta^{t''}}$$

Recall that the update vectors are conically independent. Since  $\alpha_i^{t'} \leq \alpha_i^{t''}$  (as time  $t^{00}$  follows  $t^0$ ), characterization (56) of conic independence thus entails (72).

$$(72) \quad \alpha_1^{t'} = \alpha_1^{t''}, \quad \alpha_2^{t'} = \alpha_2^{t''}, \quad \alpha_3^{t'} = \alpha_3^{t''}$$

These identities (72) say that not a single update has happened between times  $t^0$  and  $t^{00}$ . Hence, the algorithm could not have entertained a different ranking vector in between these two times. The reasoning straightforwardly extends to the general case.

■ **Remark.** No matter the details, I think there is something exciting in the reasoning just presented, that can be brought out as follows:

- (73) a. No matter the framework, many online algorithms work by adding to a current vector a certain update vector, see Cesa-Bianchi and Lugosi (2006);
- b. thus, it is a general property of online algorithms that the current vector is a conic combination of a finite set of update vectors;
- c. thus, it is important to study properties of conic combinations of the update vectors, of their conic geometry;
- d. among all possible conic properties, one of the most important ones is of course conic independence (55);
- e. by optimal design, we might thus want to assume a framework which ensures conic independence of the update vectors;
- f. Voilá: OT is such a model!

This reasoning contributes to the recent debate on the comparison of OT with alternative models; see Bane and Riggle (2010), Tesar (2007) and Pater (2009) a.o. We will come back to this in class 3.

## 8 An invariant for promotion/demotion OT online algorithms

■ **Invariants for demotion-only.** As seen in class 1, T&S's analysis of the OT online algorithm with the beautiful demotion-only update rule (9) includes the following *invariant* on the current ranking values:

- (74) If the input tableau is OT compatible with the ranking  $C_1 \gg C_2 \gg \dots \gg C_n$ , then the current ranking value of constraint  $C_k$  can never drop below  $-(k-1)$ .

In general, invariants for demotion-only OT algorithms bound the ranking value  $\theta_k$  of the  $k$ th constraint in a given run of the OT online algorithm, i.e. have the following form:

- (75) If the input tableau is such and such, then  $\theta_k \geq bound$ .

■ **Invariants for promotion-demotion.** Of course, we want invariants also for promotion/demotion OT online algorithms. Yet:

- (76) a. in the promotion/demotion case, the dynamics of the ranking values  $\theta_k$  is very complicated because non-monotonic, as in picture (41);
- b. thus, we bound instead the number of times  $\alpha_i$  that the  $i$ th row triggers an update, whose dynamics is very simple ( $\alpha_i$  can only increase);
- c. and then we use the relationship (70) between the ranking values  $\theta_k$  and the coefficients  $\alpha_i$  to derive invariants on the current ranking values.

Thus, I want to bound the number of updates  $\alpha_i$  triggered by the  $i$ th row  $\mathbf{a}_i$  in a given run of the OT online algorithm. To simplify the discussion, I will assume that:

- (77) a. the initial ranking vector has all identical components;



- b. the input tableau has a unique L per row.

By (77a), at least one update by  $\mathbf{a}_i$  is necessary. Our question thus becomes (78): how many further updates will  $\mathbf{a}_i$  trigger in the worst-case?

- (78) If the input tableau is such and such, then  $\alpha_i - 1 \leq ??$

I will concentrate on the cautious promotion-demotion update rule (36), repeated in (79);

- (79) a. Promote every WPC by 1;  
b. demote the unique LPC by the total number of WPCs.

Yet, the reasoning extends to any promotion-demotion update rule and does not require assumptions (77), made only to simplify the presentation.

■ **General shape of the bound.** Intuitively, we might reason as follows:

- (80) a. If the input tableau only consisted of row  $\mathbf{a}_i$ , one update by row  $\mathbf{a}_i$  would be enough to ensure that the current ranking vector is OT compatible with  $\mathbf{a}_i$ ;  
b. in order for  $\mathbf{a}_i$  to trigger another update, that ranking configuration needs to be disrupted by update by some *disruptor row*  $\mathbf{a}_j$  of the input tableau;  
c. it might then take a certain number of further updates by row  $\mathbf{a}_i$  in order to re-establish a ranking configuration OT-compatible with  $\mathbf{a}_i$ ;  
d. that number of further updates by row  $\mathbf{a}_i$  will depend on the *amount of disruption* caused by the disruptor row  $\mathbf{a}_j$ .

Thus, let:

- (81) a.  $disruptor(\mathbf{a}_i)$  = the set of rows of the input tableau that trigger an update that disrupts the ranking configuration necessary for OT compatibility with row  $\mathbf{a}_i$ ;  
b.  $disruption_j$  = the amount of disruption caused by a disruptor row  $\mathbf{a}_j$  to the ranking configuration necessary for OT compatibility with row  $\mathbf{a}_i$ .

By (80), we expect the number  $\alpha_i - 1$  of further updates triggered by row  $\mathbf{a}_i$  to be bound by the sum over the numbers  $\alpha_j$  of updates triggered by disruptor rows  $\mathbf{a}_j$ , each multiplied by the corresponding amount of disruption:

$$(82) \quad \alpha_i - 1 \leq \sum_{j \in disruptor(\mathbf{a}_i)} disruption_j \cdot \alpha_j$$

To make (82) explicit, we need to make (81) explicit, namely to define disruptor rows and their amount of disruption.

■ **Disruptors.** Suppose the target row  $\mathbf{a}_i$  has an L corresponding to  $C_\ell$  and a W corresponding to some  $C_h$ , as in (83). A row  $\mathbf{a}_j$  can be a disruptor of the ranking configuration necessary for OT-compatibility with row  $\mathbf{a}_i$  in one of three ways.

$$(83) \quad \mathbf{a}_i \quad \left[ \begin{array}{cccc} \dots & C_h & \dots & C_\ell & \dots \\ \dots & W & \dots & L & \dots \end{array} \right] \quad \Longrightarrow \quad \left\{ \begin{array}{l} \text{pushes up } C_h \\ \text{pushes down } C_\ell \end{array} \right.$$

*First*, row  $\mathbf{a}_j$  can be an *L-disruptor*, namely has an L corresponding to constraint  $C_h$ , so that  $\mathbf{a}_j$  pushes down the constraint  $C_h$  that  $\mathbf{a}_i$  pushes up.

$$(84) \quad \mathbf{a}_j \quad \left[ \begin{array}{cccc} \dots & C_h & \dots & C_\ell & \dots \\ \dots & L & \dots & \dots & \dots \end{array} \right] \quad \Longrightarrow \quad \left\{ \begin{array}{l} \text{pushes down } C_h \end{array} \right.$$

*Second*, row  $\mathbf{a}_j$  can be a *W-disruptor*, namely has a W corresponding to constraint  $C_\ell$ , so that  $\mathbf{a}_j$  pushes up the constraint  $C_\ell$  that  $\mathbf{a}_i$  pushes down.

$$(85) \quad \mathbf{a}_j \quad \left[ \begin{array}{cccc} \dots & C_h & \dots & C_\ell & \dots \\ \dots & \dots & \dots & W & \dots \end{array} \right] \quad \Longrightarrow \quad \left\{ \begin{array}{l} \text{pushes up } C_\ell \end{array} \right.$$

*Third*, row  $\mathbf{a}_j$  can be an *LW-disruptor*, namely has an L corresponding to  $C_h$  and a W corresponding to  $C_\ell$ , so that  $\mathbf{a}_j$  pushes down the constraint  $C_h$  that  $\mathbf{a}_i$  pushes up and pushes up the constraint  $C_\ell$  that  $\mathbf{a}_i$  pushes down.

$$(86) \quad \mathbf{a}_j \quad \left[ \begin{array}{cccc} \dots & C_h & \dots & C_\ell & \dots \\ \dots & L & \dots & W & \dots \end{array} \right] \quad \Longrightarrow \quad \left\{ \begin{array}{l} \text{pushes down } C_h \\ \text{pushes up } C_\ell \end{array} \right.$$

Thus, the set of disruptors  $disruptor_h(\mathbf{a}_i)$  of row  $\mathbf{a}_i$  (relative to its WPC  $C_h$ ) is the collection of all rows of the input tableau of type (84), (85) and (86).

■ **Disruption.** If row  $\mathbf{a}_j$  is an L-disruptor of the target row  $\mathbf{a}_i$  as in (84), then:

- (87) a. update by row  $\mathbf{a}_j$  according to (79) demotes  $C_h$  by the total number of WPCs of row  $\mathbf{a}_j$ , despite the fact that  $\mathbf{a}_i$  would have promoted it;  
b. thus,  $\mathbf{a}_j$  disrupts the ranking configuration necessary for compatibility with  $\mathbf{a}_i$  by an amount  $disruption_j$  proportional to the number of WPCs of  $\mathbf{a}_j$ .

If row  $\mathbf{a}_j$  is a W-disruptor of the target row  $\mathbf{a}_i$  as in (85), then:

- (88) a. update by row  $\mathbf{a}_j$  according to (79) promotes  $C_\ell$  by 1, despite the fact that  $\mathbf{a}_i$  would have demoted it;  
b. thus,  $\mathbf{a}_j$  disrupts the ranking configuration necessary for compatibility with  $\mathbf{a}_i$  by an amount  $disruption_j$  proportional to 1.

If row  $\mathbf{a}_j$  is an LW-disruptor of the target row  $\mathbf{a}_i$  as in (86), then:

- (89) a. update by  $\mathbf{a}_j$  according to (79) promotes  $C_\ell$  by 1 and demotes  $C_h$  by the number of WPCs of  $\mathbf{a}_j$ , despite the fact that  $\mathbf{a}_i$  would have done the reverse;  
b. thus,  $\mathbf{a}_j$  disrupts the ranking config. necessary for compatibility with  $\mathbf{a}_i$  by an amount  $disruption_j$  proportional to 1 plus the number of WPCs of  $\mathbf{a}_j$ .

Let me summarize (87)-(89) as follows:

$$(90) \quad \text{disruption}_j \sim \begin{cases} w(\mathbf{a}_j) & \text{if row } \mathbf{a}_j \text{ is an L-disruptor, as in (84)} \\ 1 & \text{if row } \mathbf{a}_j \text{ is a W-disruptor, as in (85)} \\ 1 + w(\mathbf{a}_j) & \text{if row } \mathbf{a}_j \text{ is an LW-disruptor, as in (86)} \end{cases}$$

where  $w(\mathbf{a}_j)$  is the total number of WPCs in row  $\mathbf{a}_j$ .

■ **More on disruption.** Given the target row  $\mathbf{a}_i$  in (83), we expect the amount of disruption  $\text{disruption}_j$  caused by a disruptor row  $\mathbf{a}_j$  to depend also on the target row  $\mathbf{a}_i$ :

- (91) a. Update by the target row  $\mathbf{a}_i$  in (83) according to (79) promotes the WPC  $C_h$  by 1 and demotes the unique LPC  $C_\ell$  by the number of WPCs of row  $\mathbf{a}_i$ ;
- b. thus, the *separation* between the ranking value of the LPC  $C_\ell$  and the ranking value of the WPCs will increase by 1 plus the number of WPCs of  $\mathbf{a}_i$ ;
- c. if row  $\mathbf{a}_i$  has many WPCs, then the separation between WPCs and the LPC after one update by  $\mathbf{a}_i$  will increase a lot and will not be easily disrupted;
- d. if row  $\mathbf{a}_i$  has few WPCs, then the separation between WPCs and the LPC after one update by  $\mathbf{a}_i$  will increase only a little and will be easily disrupted.

These considerations suggest that the amount of disruption  $\text{disruption}_j$  should also be inversely proportional to 1 plus the total number of WPCs of the target row  $\mathbf{a}_i$ :

$$(92) \quad \text{disruption}_i \sim \frac{1}{w(\mathbf{a}_i) + 1}$$

where  $w(\mathbf{a}_i)$  is the total number of WPCs in row  $\mathbf{a}_i$ .

## 9 The invariant actually holds!

■ **Claim 3** Consider a row  $\mathbf{a}_i$  of an OT-compatible input tableau that has a unique L per row. Let  $C_h$  be a WPC of row  $\mathbf{a}_i$ . The number of updates  $\alpha_i^t$  triggered by  $\mathbf{a}_i$  at time  $t$  in a run of the OT online algorithm with the promotion/demotion update rule (79) starting from the initial null ranking vector can be bound as follows:

$$(93) \quad \alpha_i^t - 1 \leq \sum_{\mathbf{a}_j \text{ 2disruptor}_{C_h}(\mathbf{a}_i)} \text{disruption}_j \cdot \alpha_j^t$$

where the disruption amounts  $\text{disruption}_j$  are defined as follows:

$$(94) \quad \text{disruption}_j = \begin{cases} \frac{w(\mathbf{a}_j)}{w(\mathbf{a}_i) + 1} & \text{if row } \mathbf{a}_j \text{ is an L-disruptor, as in (84)} \\ \frac{1}{w(\mathbf{a}_i) + 1} & \text{if row } \mathbf{a}_j \text{ is an W-disruptor, as in (85)} \\ \frac{w(\mathbf{a}_j) + 1}{w(\mathbf{a}_i) + 1} & \text{if row } \mathbf{a}_j \text{ is an LW-disruptor, as in (86)} \end{cases}$$

*Proof.* The claim obviously holds if  $\alpha_i^t = 1$ . Thus, assume that  $\alpha_i^t \geq 2$ . Let  $s$  be the time preceding  $t$  when the row  $\mathbf{a}_i$  has triggered the last update:

- (95) a.  $\alpha_i^t = \alpha_i^s + 1$ ;
- b. the ranking vector  $\theta^s = (\theta_1^s, \dots, \theta_n^s)$  at time  $s$  is not OT-compatible with  $\mathbf{a}_i$ .

Let  $C_\ell$  be the unique LPC of row  $\mathbf{a}_i$ ; let  $C_h$  be any WPC of row  $\mathbf{a}_i$ . Condition (95b) thus entails in particular (96b).

- (96) a.  $\alpha_i^t = \alpha_i^s + 1$
- b.  $\theta_h^s \leq \theta_k^s$

By reasoning as in (70), the ranking value  $\theta_h^s$  of  $C_h$  at time  $s$  can be expressed as follows:

$$(97) \quad \theta_h^s = \underbrace{\alpha_i^s}_{\text{contribution of row } \mathbf{a}_i} + \underbrace{\sum_{f_j \neq i, j C_h \text{ 2WPC}(\mathbf{a}_j)g} \alpha_j^s}_{\text{contribution of rows } \mathbf{a}_j \text{ different from } \mathbf{a}_i \text{ where } C_h \text{ is a WPC}} - \underbrace{\sum_{f_j \neq i, j C_h \text{ 2LPC}(\mathbf{a}_j)g} w(\mathbf{a}_j) \alpha_j^s}_{\text{contribution of rows } \mathbf{a}_j \text{ different from } \mathbf{a}_i \text{ where } C_h \text{ is a LPC}}$$

And analogously for the ranking value  $\theta_\ell^s$  of the LPC  $C_\ell$  at time  $s$ :

$$(98) \quad \theta_\ell^s = -w(\mathbf{a}_i) \alpha_i^s + \sum_{f_j \neq i, j C_\ell \text{ 2WPC}(\mathbf{a}_j)g} \alpha_j^s - \sum_{f_j \neq i, j C_\ell \text{ 2LPC}(\mathbf{a}_j)g} w(\mathbf{a}_j) \alpha_j^s$$

Using expressions (97)-(98), inequality (96b) becomes:

$$(99) \quad \alpha_i^s + \sum_{f_j \neq i, j C_h \text{ 2WPC}(\mathbf{a}_j)g} \alpha_j^s - \sum_{f_j \neq i, j C_h \text{ 2LPC}(\mathbf{a}_j)g} w(\mathbf{a}_j) \alpha_j^s \leq \\ \leq -w(\mathbf{a}_i) \alpha_i^s + \sum_{f_j \neq i, j C_\ell \text{ 2WPC}(\mathbf{a}_j)g} \alpha_j^s - \sum_{f_j \neq i, j C_\ell \text{ 2LPC}(\mathbf{a}_j)g} w(\mathbf{a}_j) \alpha_j^s$$

I make the left hand side of (99) smaller by dropping positive terms and the right hand side larger by dropping negative terms:

$$(100) \quad \alpha_i^s - \sum_{f_j \neq i, j C_h \text{ 2LPC}(\mathbf{a}_j)g} w(\mathbf{a}_j) \alpha_j^s \leq -w(\mathbf{a}_i) \alpha_i^s + \sum_{f_j \neq i, j C_\ell \text{ 2WPC}(\mathbf{a}_j)g} \alpha_j^s$$

By trivially rearranging the terms of the inequality (100), I get:

$$(101) \quad \alpha_i^s \leq \sum_{f_j \neq i, j C_\ell \text{ 2WPC}(\mathbf{a}_j)g} \frac{1}{1 + w(\mathbf{a}_i)} \alpha_j^s + \sum_{f_j \neq i, j C_h \text{ 2LPC}(\mathbf{a}_j)g} \frac{w(\mathbf{a}_j)}{1 + w(\mathbf{a}_i)} \alpha_j^s$$

The hypothesis that  $s < t$  entails that  $\alpha_i^s \leq \alpha_i^t$ . I can therefore weaken (101) by replacing  $\alpha_j^s$  with  $\alpha_j^t$  on the right hand side:<sup>1</sup>

$$(102) \quad \alpha_i^s \leq \sum_{f_j \neq i, j C_\ell \text{ 2WPC}(\mathbf{a}_j)g} \frac{1}{1 + w(\mathbf{a}_i)} \alpha_j^t + \sum_{f_j \neq i, j C_h \text{ 2LPC}(\mathbf{a}_j)g} \frac{w(\mathbf{a}_j)}{1 + w(\mathbf{a}_i)} \alpha_j^t$$

The claim follows by replacing  $\alpha_i^s = \alpha_i^t - 1$  from (96) in the left hand side of (102).  $\square$

<sup>1</sup>This step is licit only because all the  $\alpha$ 's in the right hand side of (101) are multiplied by a positive coefficient, as a result of the simplification from (99) to (100).

## 10 A little application

■ **A test case.** Prince and Tesar (2004) collect benchmark test cases for models of the early stage of the acquisition of phonotactics. One of their test cases is:

$$(103) \begin{array}{c} \text{row 1} \\ \text{row 2} \end{array} \begin{array}{cccc} F_1 & F_2 & M_1 & M_2 \\ \left[ \begin{array}{cccc} & W & & L \\ & & W & W & L \end{array} \right] \end{array}$$

Given as input the tableau in (103), the algorithm needs to recognize that  $F_2$  is useless to OT-compatibility and thus should be ranked at the bottom as in (104), in order to ensure restrictiveness — we will extensively come back to this issue in class 4.

$$(104) F_1 \gg M_1 \gg M_2 \gg F_2$$

How does the OT online algorithm perform on this test case? Assume a biased initial ranking vector  $\theta^{\text{init}}$  as (105), that ranks the two markedness constraints well above the two faithfulness constraints.

$$(105) \begin{array}{l} \theta_{F_1}^{\text{ini}} = \theta_{F_2}^{\text{ini}} = 0 \\ \theta_{M_1}^{\text{ini}} = \theta_{M_2}^{\text{ini}} = \theta^{\text{init}} = \text{a large positive constant} \end{array}$$

We need constraint promotion to pull apart  $F_1$  and  $F_2$ . Thus, consider a promotion/demotion update rule (106). With  $\lambda = 1$ , this is Boersma’s update rule (10); with  $\lambda = \#$  of WPCs, this is the cautious update rule (36), as rows 1 and 2 have a unique L.

- (106) a. Promote all WPCs by 1;  
b. demote the LPCs by a demotion amounts  $\lambda$   
( $\lambda$  can depend on the row considered).

Does the OT online algorithm run on input tableau (103) starting from the initial ranking vector (105) with the update rule (106) converge to the desired final ranking (104)?

■ **Intuition.** The analysis has two steps. The *first step* is (107), corresponding to claim 4.

- (107) The algorithm converges to the desired final ranking (104) provided the total number  $\alpha_2^{\text{fn}}$  of updates triggered by row 2 is “small”.

The intuition behind (107) is straightforward:

- (108)  $F_2$  is only promoted by row 2 and thus the algorithm will keep  $F_2$  low provided row 2 triggers few updates.

The *second step* of the analysis is (109).

- (109) The total number  $\alpha_2^{\text{fn}}$  of updates triggered by row 2 is “small” (independently of the frequency with which the row is fed to the algorithm) provided the demotion amount  $\lambda_2$  in the case of update triggered by row 2 is “large”.

The intuition behind (109) is as follows:

- (110) a. if  $\lambda_2$  is large, then update by row 2 will separate the W of  $M_1$  from the L of  $M_2$  by a large amount  $\lambda_2 + 1$ ;  
b. as that separation is large, it will take a while for it to be eroded by updates triggered by row 1;  
c. thus, it will take a while for the current ranking vector to become again not OT-incompatible with row 2;  
d. as row 2 can only trigger an update when not OT compatible with the current ranking vector, then row 2 won’t have many chances of triggering an update.

As noted above, Boersma’s update rule (10) and the cautious update rule (36) fit into the scheme (106) but differ for  $\lambda_2$ :

$$(111) \begin{array}{cc} \text{Boersma's update rule (10)} & \text{cautious update rule (36)} \\ \downarrow & \downarrow \\ \text{"small"} \lambda_2 = 1 & \text{"large"} \lambda_2 = 2 \end{array}$$

It turns out that the demotion amount  $\lambda_2 = 1$  of Boersma’s update rule (10) is too small for the algorithm to converge to the desired final ranking (104), as shown by counterexample (112) for  $\theta^{\text{init}} = 7$  (but counterexamples can be constructed for every  $\theta^{\text{init}}$ ).

$$(112) \begin{array}{cccccccccccc} \begin{bmatrix} 0 \\ 0 \\ 7 \end{bmatrix} & \xrightarrow{\mathbf{a}_2} & \begin{bmatrix} 0 \\ 1 \\ 8 \\ 6 \end{bmatrix} & \xrightarrow{\mathbf{a}_1} & \begin{bmatrix} 1 \\ 1 \\ 7 \\ 6 \end{bmatrix} & \xrightarrow{\mathbf{a}_1} & \begin{bmatrix} 2 \\ 1 \\ 6 \\ 6 \end{bmatrix} & \xrightarrow{\mathbf{a}_1} & \begin{bmatrix} 3 \\ 1 \\ 5 \\ 6 \end{bmatrix} & \xrightarrow{\mathbf{a}_1} & \begin{bmatrix} 4 \\ 1 \\ 4 \\ 6 \end{bmatrix} & \xrightarrow{\mathbf{a}_2} & \begin{bmatrix} 4 \\ 2 \\ 3 \\ 6 \end{bmatrix} & \xrightarrow{\mathbf{a}_2} & \begin{bmatrix} 4 \\ 3 \\ 3 \\ 6 \end{bmatrix} & \xrightarrow{\mathbf{a}_1} & \begin{bmatrix} 4 \\ 3 \\ 3 \\ 5 \end{bmatrix} & \xrightarrow{\mathbf{a}_1} & \begin{bmatrix} 6 \\ 3 \\ 4 \\ 4 \end{bmatrix} & \xrightarrow{\mathbf{a}_2} & \begin{bmatrix} 6 \\ 4 \\ 5 \\ 3 \end{bmatrix} \end{array}$$

It turns out instead that the demotion amount  $\lambda_2 = 2$  of the cautious update rule (36) is large enough for the algorithm to converge to the desired final ranking (104), as proven in claim 5, using the invariant obtained in the preceding section. To illustrate, note that the algorithm with this promotion-demotion update rule and  $\theta^{\text{init}} = 7$  admits only 13 learning paths, listed in (113): in all of them, row 2 triggers only two updates.

$$(113) \begin{array}{cccccccccccccccc} \begin{bmatrix} 0 \\ 0 \\ 7 \end{bmatrix} & \xrightarrow{\mathbf{a}_1} & \begin{bmatrix} 1 \\ 0 \\ 6 \\ 7 \end{bmatrix} & \xrightarrow{\mathbf{a}_1} & \begin{bmatrix} 2 \\ 0 \\ 5 \\ 7 \end{bmatrix} & \xrightarrow{\mathbf{a}_1} & \begin{bmatrix} 3 \\ 0 \\ 4 \\ 7 \end{bmatrix} & \xrightarrow{\mathbf{a}_1} & \begin{bmatrix} 4 \\ 0 \\ 3 \\ 7 \end{bmatrix} & \xrightarrow{\mathbf{a}_2} & \begin{bmatrix} 4 \\ 1 \\ 4 \\ 5 \end{bmatrix} & \xrightarrow{\mathbf{a}_2} & \begin{bmatrix} 4 \\ 2 \\ 5 \\ 3 \end{bmatrix} & \xrightarrow{\mathbf{a}_1} & \begin{bmatrix} 5 \\ 2 \\ 4 \\ 3 \end{bmatrix} & \xrightarrow{\mathbf{a}_2} & \begin{bmatrix} 3 \\ 1 \\ 5 \\ 5 \end{bmatrix} & \xrightarrow{\mathbf{a}_2} & \begin{bmatrix} 3 \\ 2 \\ 6 \\ 3 \end{bmatrix} & \xrightarrow{\mathbf{a}_2} & \begin{bmatrix} 4 \\ 2 \\ 5 \\ 3 \end{bmatrix} & \xrightarrow{\mathbf{a}_1} & \begin{bmatrix} 5 \\ 2 \\ 4 \\ 3 \end{bmatrix} & \xrightarrow{\mathbf{a}_1} & \begin{bmatrix} 4 \\ 1 \\ 4 \\ 5 \end{bmatrix} & \xrightarrow{\mathbf{a}_1} & \begin{bmatrix} 5 \\ 1 \\ 3 \\ 5 \end{bmatrix} & \xrightarrow{\mathbf{a}_2} & \begin{bmatrix} 5 \\ 2 \\ 4 \\ 3 \end{bmatrix} \end{array}$$



Analogously, if  $\alpha_2^t = 1$ , then (121b) trivially holds. Otherwise, the proof of (121b) consists of (123) and (124).

$$\begin{aligned}
(123) \quad \alpha_2^t \geq 1 &\implies \exists s < t \text{ s.t. } \alpha_2^s = \alpha_2^t - 1 \text{ and } \theta^s \text{ not OT-compatible with } \mathbf{a}_2 \\
&\implies \exists s < t \text{ s.t. } \alpha_2^s = \alpha_2^t - 1 \text{ and } \theta_{F_2}^s \leq \theta_{M_2}^s \\
&\implies \exists s < t \text{ s.t. } \alpha_2^s = \alpha_2^t - 1 \text{ and } \alpha_2^s \leq -2\alpha_2 + \theta^{\text{init}} \\
&\implies \exists s < t \text{ s.t. } \alpha_2^s = \alpha_2^t - 1 \text{ and } 3\alpha_2^s \leq \theta^{\text{init}} \\
&\implies 3(\alpha_2^t - 1) \leq \theta^{\text{init}} \\
&\implies \alpha_2^t \leq 1 + \frac{\theta^{\text{init}}}{3} \\
(124) \quad \alpha_2^t \geq 1 &\implies \exists s < t \text{ s.t. } \alpha_2^s = \alpha_2^t - 1 \text{ and } \theta^s \text{ not OT-compatible with } \mathbf{a}_2 \\
&\implies \exists s < t \text{ s.t. } \alpha_2^s = \alpha_2^t - 1 \text{ and } \theta_{M_1}^s \leq \theta_{M_2}^s \\
&\implies \exists s < t \text{ s.t. } \alpha_2^s = \alpha_2^t - 1 \text{ and } -\alpha_1 + \alpha_2 + \theta^{\text{init}} \leq -2\alpha_2 + \theta^{\text{init}} \\
&\implies \exists s < t \text{ s.t. } \alpha_2^s = \alpha_2^t - 1 \text{ and } 3\alpha_2^s \leq \alpha_1^s \\
&\implies 3(\alpha_2^t - 1) \leq \alpha_1^t \\
&\implies \alpha_2^t \leq 1 + \frac{1}{3}\alpha_1^t
\end{aligned}$$

In order to finally derive (114) from (121), let me distinguish two cases:

$$\begin{aligned}
(125) \quad \text{a. } &\alpha_1^t < \theta^{\text{init}} - 3; \\
&\text{b. } &\alpha_1^t \geq \theta^{\text{init}} - 3.
\end{aligned}$$

In case (125a), (119) always holds, as shown by (126).

$$(126) \quad \alpha_2^t \stackrel{(a)}{\leq} 1 + \frac{1}{3}\alpha_1^t \stackrel{(b)}{<} 1 + \frac{1}{3}(\theta^{\text{init}} - 3) = \frac{1}{3}\theta^{\text{init}}$$

It is easy to see using (121b) that case (125b) can never arise.  $\square$

---

## References

- Bane, Max, and Jason Riggle. 2010. “The typological consequences of weighted constraints”. In *Proceedings of CLS 45*.
- Boersma, Paul. 1997. “How We Learn Variation, Optionality and Probability”. In *IFA Proceedings 21*, 43–58. University of Amsterdam: Institute for Phonetic Sciences.
- Cesa-Bianchi, Nicolò, and Gábor Lugosi. 2006. *Prediction, Learning, and Games*. Cambridge University Press.
- Pater, Joe. 2008. “Gradual Learning and Convergence”. *Linguistic Inquiry* 39.2:334–345.
- Pater, Joe. 2009. “Weighted Constraints in Generative Linguistics”. *Cognitive Science* 33:999–1035.

Prince, Alan, and Bruce Tesar. 2004. “Learning Phonotactic Distributions”. In *Constraints in Phonological Acquisition*, ed. R. Kager, J. Pater, and W. Zonneveld, 245–291. Cambridge University Press.

Tesar, Bruce. 2007. “A Comparison of Lexicographic and Linear Numeric Optimization Using Violation Difference Ratios”. Ms., Rutgers University; ROA-939.

Tesar, Bruce, and Paul Smolensky. 1998. “Learnability in Optimality Theory”. *Linguistic Inquiry* 29:229–268.