

MAXENT IS AS TRANSPARENT AS OT AND HG

- Tesar (2014) extends the distinction between transparency and opacity from rule-based to arbitrary categorical phonological grammars through the notion of output-drivenness. Anttila and Magri (2018) extend entailments such as the one used to define output-drivenness from categorical to probabilistic phonological grammars. Building on these two developments, this paper investigates transparency/opacity for probabilistic MaxEnt grammars. The main result is that the same condition on the faithfulness constraints which is responsible for the transparency of categorical OT and HG also suffices to ensure the transparency of MaxEnt. This result is surprising in light of Anttila and Magri’s finding that ME breaks many of the entailments predicted by OT and HG.
- Let $\langle \mathbf{x}, \mathbb{A}, \mathbb{B} \rangle = \mathbf{y}$ denote the SPE derivation which starts with the UR \mathbf{x} and applies rule \mathbb{A} followed by rule \mathbb{B} to yield the SR \mathbf{y} . Hence $\langle \mathbf{x}, \mathbb{A}, \mathbb{B} \rangle = \langle \langle \mathbf{x}, \mathbb{A} \rangle, \mathbb{B} \rangle$. Recall that rule \mathbb{B} **counter-feeds** rule \mathbb{A} in this derivation iff \mathbb{A} applies vacuously to \mathbf{x} (1a) but not to the result of applying \mathbb{B} to \mathbf{x} (1b).

- (1) a. $\langle \mathbf{x}, \mathbb{A} \rangle = \mathbf{x}$ My first contribution is that counter-feeding entails **chain shifts**.
 b. $\langle \mathbf{x}, \mathbb{B}, \mathbb{A} \rangle \neq \langle \mathbf{x}, \mathbb{B} \rangle$ More precisely, consider the form $\mathbf{z} = \langle \mathbf{x}, \mathbb{B} \rangle$. Suppose that $\langle \mathbf{z}, \mathbb{A}, \mathbb{B} \rangle$ is not a **Duke of York** (DOY) derivation. This assumption is mild because DOYs are rare. Under this mild assumption, if $\mathbf{y} = \langle \mathbf{x}, \mathbb{A}, \mathbb{B} \rangle$ is a counter-feeding derivation, the SPE phonology corresponding to the ordered rules $\langle \mathbb{A}, \mathbb{B} \rangle$ yields the chain shift $\mathbf{x} \rightarrow \mathbf{y} \not\rightarrow \mathbf{y}$.

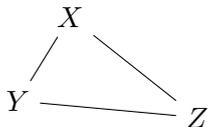
- Here is the proof. Recall that the derivation $\langle \mathbf{z}, \mathbb{A}, \mathbb{B} \rangle$ would be a DOY if \mathbb{A} modifies \mathbf{z} (2a) but \mathbb{B} subsequently gets back the initial form \mathbf{z} (2b). The claim follows by applying the SPE phonology $\langle \mathbb{A}, \mathbb{B} \rangle$ to the form $\mathbf{y} = \langle \mathbf{x}, \mathbb{A}, \mathbb{B} \rangle$ as in (3), obtaining a form different from \mathbf{y} , whereby the chain shift. In step (3a), I have used the definition of the form \mathbf{y} as the result of applying $\langle \mathbb{A}, \mathbb{B} \rangle$ to \mathbf{x} . In step (3b), I have used the counter-feeding condition (1a), whereby \mathbb{A} applies vacuously to \mathbf{x} . Since $\mathbf{z} = \langle \mathbf{x}, \mathbb{B} \rangle$, the counter-feeding condition (1b) can be restated as $\langle \mathbf{z}, \mathbb{A} \rangle \neq \mathbf{z}$. In other words, the DOY condition (2a) holds. The other DOY condition (2b) must therefore fail. Since $\mathbf{z} = \langle \mathbf{x}, \mathbb{B} \rangle$, failure of (2b) means $\langle \langle \mathbf{x}, \mathbb{B} \rangle, \mathbb{A}, \mathbb{B} \rangle \neq \langle \mathbf{x}, \mathbb{B} \rangle$, whereby (3c). In step (3d), I have used again the counter-feeding condition (1a). Finally in step (3e), I have used again the definition $\mathbf{y} = \langle \mathbf{x}, \mathbb{A}, \mathbb{B} \rangle$.

$$(3) \quad \langle \mathbf{y}, \mathbb{A}, \mathbb{B} \rangle \stackrel{(a)}{=} \langle \langle \mathbf{x}, \mathbb{A}, \mathbb{B} \rangle, \mathbb{A}, \mathbb{B} \rangle \stackrel{(b)}{=} \langle \langle \mathbf{x}, \mathbb{B} \rangle, \mathbb{A}, \mathbb{B} \rangle \stackrel{(c)}{\neq} \langle \mathbf{x}, \mathbb{B} \rangle \stackrel{(d)}{=} \langle \langle \mathbf{x}, \mathbb{A} \rangle, \mathbb{B} \rangle \stackrel{(e)}{=} \mathbf{y}$$

- A **categorical** grammar G (construed as a mapping from URs to SRs) is called **idempotent** provided it yields no chain shifts. Equivalently, it satisfies the implication (4): whenever G maps a UR \mathbf{x} to a SR \mathbf{y} , it also maps \mathbf{y} (construed as a UR) to itself. Since counter-feeding entails chain-shifts as shown above, non-idempotency generalizes counter-feeding opacity beyond rule-based phonology.

$$(4) \quad G(\mathbf{x}) = \mathbf{y} \implies G(\mathbf{y}) = \mathbf{y}$$

$$(5) \quad F(\mathbf{x}, \mathbf{z}) \leq F(\mathbf{x}, \mathbf{y}) + F(\mathbf{y}, \mathbf{z})$$



Magri (2018) shows that OT and HG grammars are idempotent provided each faithfulness constraint F satisfies the **triangle inequality** (TI) in (5): the violations $F(\mathbf{x}, \mathbf{z})$ assigned to the mapping of a UR \mathbf{x} to a SR \mathbf{z} are never larger than the violations $F(\mathbf{x}, \mathbf{y})$ assigned to the mapping of the UR \mathbf{x} to some other SR \mathbf{y} plus the violations $F(\mathbf{y}, \mathbf{z})$ assigned to the mapping of \mathbf{y} (construed as a UR) to the SR \mathbf{z} . The intuition is that F measures the phonological distance between URs and SRs and it thus complies with an axiom on distances which says that any side XZ of a triangle is shorter than the sum of the other two sides XY and YZ . The well-known resistance of OT and HG to chain shifts is then explained by the fact that most faithfulness constraints in the literature indeed satisfy the TI.

- This paper extends these results to **MaxEnt** (ME). A ME grammar G^{ME} maps a UR \mathbf{x} to a probability distribution $G^{\text{ME}}(\cdot | \mathbf{x})$ over the candidate set of \mathbf{x} . Following Anttila and Magri (2018),

$$(6) \quad G^{\text{ME}}(\mathbf{y} | \mathbf{x}) \leq G^{\text{ME}}(\mathbf{y} | \mathbf{y})$$

I generalize the idempotency entailment (4) to this probabilistic setting as the condition (6) that the probability $G^{\text{ME}}(\mathbf{y} | \mathbf{y})$ assigned to the consequent faithful mapping of \mathbf{y} to itself is not smaller than the probability $G^{\text{ME}}(\mathbf{y} | \mathbf{x})$ of the antecedent mapping of the UR \mathbf{x} to this SR \mathbf{y} . Condition (4) is indeed a special case of (6) for probability distributions which are categorical, namely concentrated on a single candidate. My second contribution is that the faithfulness TI (5) also suffices to ensure that ME satisfies (6) and is thus idempotent (whenever the URs \mathbf{x} and \mathbf{y} share the same candidate set).

- Here is the proof. Using the analytical expression of ME, the idempotency inequality (6) is made explicit in (7). Here, $H_{\mathbf{w}}(\text{UR}, \text{SR})$ is the HG **harmony** of the mapping (UR, SR) relative to the **weight vector** \mathbf{w} . The sums in the denominators run over all candidates \mathbf{z} in the shared candidate set which are different from \mathbf{y} . I multiply the numerator and the denominator of the left hand side of

$$(7) \quad \frac{e^{H_{\mathbf{w}}(\mathbf{x}, \mathbf{y})}}{e^{H_{\mathbf{w}}(\mathbf{x}, \mathbf{y})} + \sum_{\mathbf{z} \neq \mathbf{y}} e^{H_{\mathbf{w}}(\mathbf{x}, \mathbf{z})}} \leq \frac{e^{H_{\mathbf{w}}(\mathbf{y}, \mathbf{y})}}{e^{H_{\mathbf{w}}(\mathbf{y}, \mathbf{y})} + \sum_{\mathbf{z} \neq \mathbf{y}} e^{H_{\mathbf{w}}(\mathbf{y}, \mathbf{z})}}$$

$$(8) \quad \sum_{\mathbf{z} \neq \mathbf{y}} e^{H_{\mathbf{w}}(\mathbf{x}, \mathbf{z}) + H_{\mathbf{w}}(\mathbf{y}, \mathbf{y})} \geq \sum_{\mathbf{z} \neq \mathbf{y}} e^{H_{\mathbf{w}}(\mathbf{x}, \mathbf{y}) + H_{\mathbf{w}}(\mathbf{y}, \mathbf{z})}$$

(7) by $e^{H_{\mathbf{w}}(\mathbf{y}, \mathbf{y})}$. I also multiply the numerator and the denominator of the right hand side by $e^{H_{\mathbf{w}}(\mathbf{x}, \mathbf{y})}$. After some simplifications, I can thus rewrite (7) as in (8) (with the inequality reversed). Since a faithfulness constraint F assigns zero violations to the faithful mapping (\mathbf{y}, \mathbf{y}) , I can add $F(\mathbf{y}, \mathbf{y})$ to the left hand side of the TI (5), obtaining (9a). Furthermore, since a markedness constraint M is oblivious to the URs, the identity (9b) holds. Since the harmony $H_{\mathbf{w}}(\cdot, \cdot)$ is the weighted sum of the constraint violations multiplied by -1 , the two

$$(9) \quad \begin{array}{ll} \text{a. } F(\mathbf{x}, \mathbf{z}) + F(\mathbf{y}, \mathbf{y}) \leq F(\mathbf{x}, \mathbf{y}) + F(\mathbf{y}, \mathbf{z}) & \text{inequalities (9) yield the harmony inequality} \\ \text{b. } M(\mathbf{x}, \mathbf{z}) + M(\mathbf{y}, \mathbf{y}) = M(\mathbf{x}, \mathbf{y}) + M(\mathbf{y}, \mathbf{z}) & H_{\mathbf{w}}(\mathbf{x}, \mathbf{z}) + H_{\mathbf{w}}(\mathbf{y}, \mathbf{y}) \geq H_{\mathbf{w}}(\mathbf{y}, \mathbf{z}) + H_{\mathbf{w}}(\mathbf{x}, \mathbf{y}) \end{array}$$

(the inequality is reversed because of -1), whereby (8) (because the exponential is increasing).

- I now turn to counter-bleeding opacity. Recall that rule \mathbb{B} **counter-bleeds** rule \mathbb{A} in the derivation $\langle \mathbf{x}, \mathbb{A}, \mathbb{B} \rangle$ provided \mathbb{A} applies non-vacuously to \mathbf{x} (10a) but vacuously to the result of applying \mathbb{B} to \mathbf{x}

$$(10) \quad \begin{array}{ll} \text{a. } \langle \mathbf{x}, \mathbb{A} \rangle \neq \mathbf{x} & (10\text{b}). \text{ My } \boxed{\text{third contribution}} \text{ is that counter-bleeding entails} \\ \text{b. } \langle \mathbf{x}, \mathbb{B}, \mathbb{A} \rangle = \langle \mathbf{x}, \mathbb{B} \rangle & \text{saltations. More precisely, suppose that } \mathbb{A} \text{ does not bleed } \mathbb{B} \text{ in the} \\ & \text{derivation } \langle \mathbf{x}, \mathbb{A}, \mathbb{B} \rangle. \text{ This assumption is mild because the deriva-} \end{array}$$

tion $\langle \mathbf{x}, \mathbb{B}, \mathbb{A} \rangle$ would otherwise be a case of **mutual bleeding**, which is uncommon (Bakovic 2011). Under this mild assumption, if $\mathbf{z} = \langle \mathbf{x}, \mathbb{A}, \mathbb{B} \rangle$ is a counter-bleeding derivation, the SPE phonology corresponding to the ordered rules $\langle \mathbb{A}, \mathbb{B} \rangle$ yields the saltation $\mathbf{x} \xrightarrow{\mathbf{y}} \mathbf{z}$: the form $\mathbf{y} = \langle \mathbf{x}, \mathbb{B} \rangle$ is mapped to itself but \mathbf{x} is mapped to \mathbf{z} , despite the fact \mathbf{x} is closer to \mathbf{y} than to \mathbf{z} , as \mathbf{y} is obtained from \mathbf{x} through rule \mathbb{B} only while \mathbf{z} is obtained from \mathbf{x} through both non-vacuous rules \mathbb{A} and \mathbb{B} . The proof is omitted for space, but it is similar in spirit to the proof of my first contribution.

- Consider two mappings (\mathbf{x}, \mathbf{z}) and (\mathbf{y}, \mathbf{z}) which share the same SR \mathbf{z} . Suppose that the UR \mathbf{y} is more similar to the shared SR \mathbf{z} than the UR \mathbf{x} is. In other words, that the mapping (\mathbf{y}, \mathbf{z}) has more “internal similarity” than the other mapping (\mathbf{x}, \mathbf{z}) . This assumption is formalized as $(\mathbf{x}, \mathbf{z}) \leq_{\text{sim}} (\mathbf{y}, \mathbf{z})$, where \leq_{sim} is a **similarity order**, namely an ordering of mappings based on their internal similarity. Tesar (2014) calls **output-driven** a categorical grammar G which satisfies the implication (11): whenever G maps the less similar UR \mathbf{x} to \mathbf{z} , it also maps the

$$(11) \quad G(\mathbf{x}) = \mathbf{z} \implies G(\mathbf{y}) = \mathbf{z} \quad \begin{array}{l} \text{more similar UR } \mathbf{y} \text{ to } \mathbf{z}. \text{ In other words, an output-driven} \\ \text{grammar yields no saltations } \mathbf{x} \xrightarrow{\mathbf{y}} \mathbf{z}. \end{array}$$

Furthermore, output-drivenness entails idempotency, as (4) is a special case of (11) where the UR \mathbf{y} more similar to \mathbf{z} is the most similar one, namely \mathbf{y} itself. Hence, an output-driven grammar also yields no chain shifts. Since opacity means counter-feeding or counter-bleeding which in turn mean chain shifts or saltations, output-drivenness generalizes transparency beyond rule-based phonology.

- In order to make progress on the theory of output-drivenness, we need to define the similarity order \leq_{sim} . Magri (2018b) shows that Tesar’s (2014) concrete definition of \leq_{sim} can be axiomatized as follows: $(\mathbf{x}, \mathbf{z}) \leq_{\text{sim}} (\mathbf{y}, \mathbf{z})$ holds provided each faithfulness constraint F satisfies (12). This inequality says that the less similar mapping (\mathbf{x}, \mathbf{z}) “makes up” not only for every faithfulness

$$(12) \quad F(\mathbf{x}, \mathbf{z}) \geq F(\mathbf{y}, \mathbf{z}) + F(\mathbf{x}, \mathbf{y}) \quad \begin{array}{l} \text{violation of the more similar mapping } (\mathbf{y}, \mathbf{z}) \\ \text{discrepancy between the two forms } \mathbf{x} \text{ and } \mathbf{y} \text{ which play the} \end{array}$$

role of the two URs. With this definition of \leq_{sim} , Magri (2018) shows that OT and HG grammars are output-driven provided again each faithfulness constraint F satisfies the TI (5).

- Again following Anttila and Magri (2018), I generalize the categorical output-drivenness entailment (11) to the probabilistic ME setting through the condition (13) that the probability $G^{\text{ME}}(\mathbf{z} | \mathbf{y})$ of the consequent more similar mapping (\mathbf{y}, \mathbf{z}) is not smaller than the probability $G^{\text{ME}}(\mathbf{z} | \mathbf{x})$ of the antecedent less similar mapping (\mathbf{x}, \mathbf{z}) . My $\boxed{\text{fourth contribution}}$

$$(13) \quad G^{\text{ME}}(\mathbf{z} | \mathbf{x}) \leq G^{\text{ME}}(\mathbf{z} | \mathbf{y}) \quad \begin{array}{l} \text{is that the faithfulness TI (5) also suffices to ensure that ME} \\ \text{satisfies (13) and is thus output-driven. The proof is omitted for space, but similar to the proof of} \end{array}$$

my second contribution. In conclusion, OT, HG, and ME are transparent under the same condition: that faithfulness constraints measure phonological distance in compliance with the metric TI (5).