Social gender and derivational morphology: a distributional study of the gendered import of learned morphology in French

Marine Wauquier Olivier Bonami Université de Paris, LLF, CNRS

1 Introduction

French suffixes *-euse* and *-rice* are clearly morphological rivals for the formation of both feminine instrument nouns (*agrafeuse* 'stapler', *excavatrice* 'excavator') and agent nouns denoting women (*danseuse* 'female dancer', *rédactrice* 'female author'). However the literature on nouns designating women gives circumstantial evidence for differences in meaning: agent nouns in *-euse* are said to denote lower-level professions, such as *coiffeuse* 'hairdresser' or *serveuse* 'waitress' (Lenoble-Pinson, 2008), or nouns with a pejorative connotation, such as *entraîneuse* 'barmaid' or *allumeuse* 'tease' (Dawes, 2003), while *-rice* is favored for more socially valued positions (*directrice* 'female manager'). This has recently been confirmed quantitatively on the basis of distributional semantics (Wauquier et al., 2020a).

While it is plausible that the two suffixes have specialized to convey classes of meanings related to gender stereotypes, previous studies have not taken into account the fact that the two suffixes also differ in their place in the French morphological system. In parallel with other suffixes such as *-ion*, *-if*, etc., *-rice* originates in learned vocabulary borrowed from Latin from Middle French on (see Rainer & Buridant 2015 for an overview). While all these suffixes then became productive in their own right, their learned origin may have an influence on the types of concepts that they are used to designate. Crucially, the *-euse/-rice* pair is paralleled by a distinction between two processes using the same suffix *-eur* to form masculine agent and instrument nouns: learned *-eur* attaches to the same learned stems as *-ion*, *-if* or *-rice* (Bonami et al.'s (2009) 'hidden stem'),¹ while nonlearned *-eur* attaches to the same ordinary, nonlearned stems as *-euse* or other nonlearned suffixes such as *-age*.

Against this background, the present study attempts to assess to what extent the observed differences between *-euse* and *-rice* follow from their status as learned vs. nonlearned formations: if the differences in meaning between *-euse* and *-rice* follow from their learned status, we expect them to be paralleled by differences between learned and nonlearned masculine nouns in *-eur*, which are otherwise morphologically parallel, modulo gender. If the effect of learned vs. nonlearned is strong enough, we might even be able to document parallel effects for other morphosemantic types such as action nouns in *-ion* vs. *-age*.

2 Data

We built three datasets of deverbal feminine agent nouns (AGF), masculine agent nouns (AGM), and action nouns (ACT), with a contrast between a learned and a nonlearned alternative in each case. Feminine agent nouns and action nouns were extracted from Lexeur (Wauquier et al., 2020b), while masculine agent nouns were borrowed from the dataset documented in Huyghe & Wauquier (2020). All agent nouns were manually filtered so as to exclude polysemy with

¹Most learned formations in *-eur* end in *-teur*, but there are exceptions in both directions: *professeur* 'professor' is learned, *acheteur* 'buyer' is not.

an instrument reading, and only nouns with a frequency of 50 or more in the FrCoW corpus (Schäfer, 2015; Schäfer & Bildhauer, 2012) were retained. The size of our final datasets are given in table 1.

	Learned	Nonlearned
Feminine agent nouns (-rice vseuse)	158	301
Masculine agent nouns	141	462
Action nouns (-ion vsage)	750	629

Table 1: Description of our dataset

To assess the semantic properties of these nouns, we used a distributional semantic model (DSM) obtained by applying the gensim (Řehůřek & Sojka, 2010) implementation of word2vec (Mikolov et al., 2013) to a tagged and lemmatized version of the FrCoW corpus.²

3 Quantitative assessment

We first assessed whether our DSM captures differences between learned and nonlearned derivatives in our three datasets. To this end, we trained classifiers to predict from the semantic representation of a lexeme whether it was formed using a learned or the corresponding nonlearned process. Specifically, we used gradient boosting (Friedman, 2001; Mason et al., 2000) applied to decision trees as our binary classification method.³ To avoid differences in accuracy due to differences in dataset size, we randomly subsampled each of the subdatasets to 141 items, the size of the smallest of our 6 subdatasets. We report the aggregated accuracy of 10-fold cross-validation. The results of this first assessment is suprisingly good: despite a small training set, each of the three classifiers reaches an accuracy between 0.77 and 0.83, well above the 0.5 baseline. This clearly indicates that there are distributional cues separating learned and nonlearned nouns. Importantly, this holds across feminine agent nouns, masculine agent nouns, and action nouns.

The fact that all three morphosemantic types of nouns differ in their distribution does not entail that they differ in the same fashion. Further exploration however indicates that the relevant distributional properties overlap strongly. First, the analysis of dimension importance indicates that one and the same specific dimension has markedly more predictive power for all three models. Examination of the three datasets confirms in each case a highly significant contrast in values of that dimension between learned and nonlearned exemplars, although the distributions strongly overlap. Second, we used each of the three models to conduct *extrinsic prediction* on data from another semantic type: for instance, the model trained on feminine agent nouns is used to predict the constrast between learned and nonlearned masculine agent nouns. The results are shown in Table 2, with intrinsic prediction results on the diagonal.

The striking conclusion is that intrinsic and extrinsic prediction lead to very similar accuracy. All nine 95% confidence intervals overlap, so that one may not conclude that the learned vs. nonlearned distinction in one dataset is better predicted by vectors for the corresponding morphosemantic type or another one.

²We used the cbow variant of the algorithm with the following hyperparameters: 2 training epochs, 5 negative samples, window size 5, vector size 100. We used the tagging provided with the corpus and improved the lemmatization semi-automatically to correctly have separate lemmas for nouns of different grammatical genders but proper gender neutralizations of all non-nouns.

³We used the Python implementation of gradient boosting in the scikit-learn package (Pedregosa et al., 2011), with the following hyperparameters for all models: 500 estimators, max depth of 2, deviance loss function.

	Test data		
Training data	AGF	AGM	ACT
AGF	0.80	0.77	0.79
AGM	0.77	0.77	0.82
ACT	0.76	0.79	0.83

Table 2: Accuracy of the three classifiers applied to the three datasets

The evidence thus strongly suggests that there are general distributional differences between learned and nonlearned deverbal formations in French that are not limited to feminine agent nouns. A likely cause of these contrasts is the fact that learned formations entered the language in particular sociolinguistic circumstances, and that analogical extension of their use led to a partial specialization for some type of concepts.

4 Qualitative evaluation

While we have shown that the difference between *-rice* and *-euse* follows at least in part from their respective learned vs. nonlearned character, it remains to be seen whether there is a link between this difference and the observed difference in connotations. To assess this, we built, for each of the 6 processes under consideration, the *centroid* representing the average of their respective vector representations. Intuitively this should capture what the processes have in common, neutralizing individual lexical semantics. We identified their 100 nearest neighbors in the DSM, and examined qualitatively the semantic properties of these neighbors. Because we are interested in connotations linked to social gender, we focus on agent nouns.

Two main oppositions emerge from the comparison of the four lists of neighbors. The first involves the overall axiological valence of learned and nonlearned centroids, regardless of the targeted gender. Both feminine and masculine nonlearned centroids display a much higher proportion of negatively valued neighbors than their learned equivalent, for which neighbors are at best positively valued (dirigeante 'female leader', chirurgienne 'female surgeon', avocate 'female lawyer' for the feminine; érudit 'scholar', académicien 'academician', orateur 'orator' for the masculine), at worst neutral (entrepeneure 'businesswoman', sculptrice 'female sculptor', collaboratrice 'female associate' for the feminine; exécutant 'subordinate', journaliste 'journalist', comptable 'accountant' for the masculine). The second opposition concerns the types of axiological properties displayed by the neighbors of the nonlearned centroids with regard to gender. Neighbors of the feminine nonlearned centroid involve connotation with respect to sexuality (nymphomane 'nymphomaniac', tapineuse 'prostitute', catin 'harlot') and physical characterization (laideron 'plain Jane', monstresse 'monstress', midinette 'starry-eyed girl'). On the other hand, the axiological valence of the masculine nonlearned centroid's neighbors also involves sexuality (dragueur 'womanizer', séducteur 'seducer'), but mainly builds on other domains such as criminal activities (truand 'gangster', voleur 'thief') or behavioral characterization (tâcheron 'drudge', poivrot 'drunkard').

These results indicate that the contrast between *-rice* and *-euse* in terms of connotations exists over and above the fact that they contrast in terms of learnedness. We submit that the basic contrast between learned vs. nonlearned formations is recruited to different purposes depending on the morphosemantic type. For action nouns, it implements a contrast between intellectual and technical domains of reference (Wauquier et al., 2020b). For agent nouns, it readily encodes gendered axiological judgements, which are different in the masculine and in the feminine as a consequence of gender stereotypes.

References

- Bonami, Olivier, Gilles Boyé & Françoise Kerleroux. 2009. L'allomorphie radicale et la relation flexion-construction. In Bernard Fradin, Françoise Kerleroux & Marc Plénat (eds.), *Aperçus de morphologie du français*, 103–125. Saint-Denis: Presses de l'Université de Vincennes.
- Dawes, Elizabeth. 2003. La féminisation des titres et fonctions dans la francophonie: de la morphologie à l'idéologie. *Ethnologies* 25(2). 195–213.
- Friedman, Jerome H. 2001. Greedy function approximation: a gradient boosting machine. *Annals of statistics* 1189–1232.
- Huyghe, Richard & Marine Wauquier. 2020. What's in an agent? a distributional semantics approach to agent nouns in french. *Morphology* 30. 185–218.
- Lenoble-Pinson, Michèle. 2008. Mettre au féminin les noms de métier: résistances culturelles et sociolinguistiques. *Le français aujourd'hui* (4). 73–79.
- Mason, Llew, Jonathan Baxter, Peter L Bartlett & Marcus R Frean. 2000. Boosting algorithms as gradient descent. In *Advances in neural information processing systems*, 512–518.
- Mikolov, Tomas, Kai Chen, Greg Corrado & Jeffrey Dean. 2013. Efficient estimation of word representations in vector space. *CoRR* abs/1301.3781. http://arxiv.org/abs/1301.3781.
- Pedregosa, F., G. Varoquaux, A. Gramfort, V. Michel, B. Thirion, O. Grisel, M. Blondel, P. Prettenhofer, R. Weiss, V. Dubourg, J. Vanderplas, A. Passos, D. Cournapeau, M. Brucher, M. Perrot & E. Duchesnay. 2011. Scikit-learn: Machine learning in Python. *Journal of Machine Learning Research* 12. 2825–2830.
- Rainer, Franz & Claude Buridant. 2015. From old french to modern french. In P.O. Muller, I. Ohnheiser, S. Olsen & F. Rainer (eds.), *Word-formation: an international handbook of the languages of europe*, vol. 3, 1975–2000. Berlin/Boston: De Gruyter Mouton.
- Řehůřek, Radim & Petr Sojka. 2010. Software Framework for Topic Modelling with Large Corpora. In Proceedings of the LREC 2010 Workshop on New Challenges for NLP Frameworks, 45–50. Valletta, Malta: ELRA. http://is.muni.cz/publication/884893/en.
- Schäfer, Roland. 2015. Processing and querying large web corpora with the COW14 architecture. In *Proceedings of challenges in the management of large corpora*, 28–34.
- Schäfer, Roland & Felix Bildhauer. 2012. Building large corpora from the web using a new efficient tool chain. In *Proceedings of the eighth international conference on language resources and evaluation*, 486–493.
- Wauquier, Marine, Nabil Hathout & Cécile Fabre. 2020a. Contributions of distributional semantics to the semantic study of French morphologically derived agent nouns. In J. Audring, N. Koutsoukos & C. Manouilidou (eds.), *Rules, patterns, schemas and analogy, mmm12 online proceedings*, vol. 12, 111–121.
- Wauquier, Marine, Nabil Hathout & Cécile Fabre. 2020b. Semantic discrimination of technicality in French nominalizations. *Zeitschrift für Wortbildung / Journal of Word Formation* 4(2). 100–119.