# An analysis of a French as a Foreign Language Corpus for Readability Assessment

Thomas François
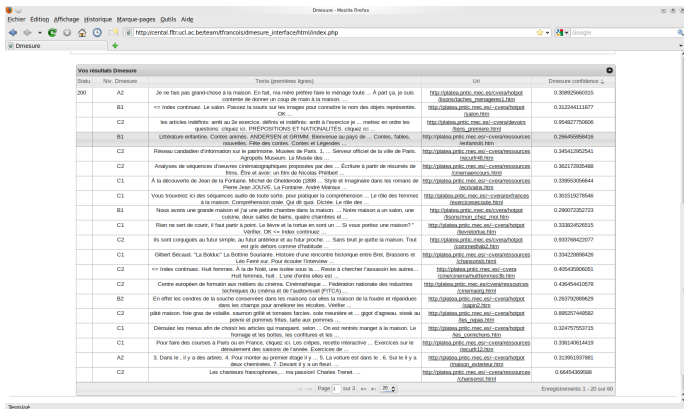
NLP4CALL, Uppsala

13-11-2014

# Context

Readability models have roles to play in iCALL:
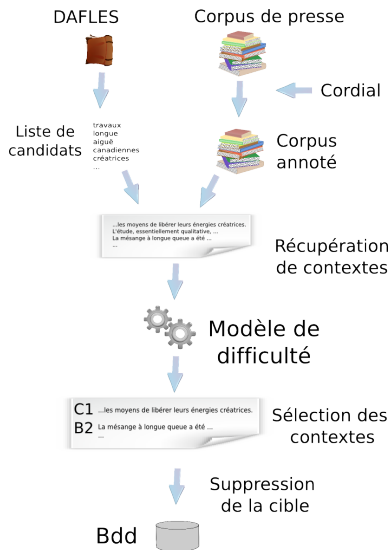
1. Find educational material of a given level, on the web

# Context

DAFLES

Corpus de presse

Cordial

Liste de candidats

travaux
longue
aiguë
canadiennes
créatrices
...

Corpus annoté

...les moyens de libérer leurs énergies créatrices.
L'étude, essentiellement qualitative, ...
La mésange à longue queue a été ...
...

Récupération de contextes

Modèle de difficulté

C1 ...les moyens de libérer leurs énergies créatrices.
B2 La mésange à longue queue a été ...
...

Sélection des contextes

Suppression de la cible

Bdd

2 Enhance systems for automatic exercise generation (proposal based on ALFALEX) [Verlinde et al., 2003]

## Problematic

- Many formulas "available" for L1, especially for English
  [Collins-Thompson and Callan, 2005, Feng et al., 2010, Vajjala and Meurers, 2012]
- Limited amount of models "available" for L2
  [Heilman et al., 2007, Schwarm and Ostendorf, 2005]
- AND only two formulas use the CEFR scale (current standard)
  [François and Fairon, 2012, Pilán et al., 2014]

- Problem: Large amount of efforts required to collect the annotated data to train a readability model!

# Objectives of the paper

1. Gathering a corpus of texts annotated for difficulty, as large as possible
   - Critical review of 5 common annotation approaches in readability
   - We opt for the extraction of texts from textbook series (with some conditions) and giving the textbook level to all texts extracted from it.
   - We report the collect process of about 2,000 texts for French as a Foreign Language (FFL)

2. Investigate the shortcomings of this criterion (which is mainstream in the field)
   - Produce unbalanced corpus
   - Homogeneity of annotations across textbook series is questionable

# 2 experiments

## Class imbalanced effect

- We compared the results of models either based on a balanced or an unbalanced corpus
- The ordinal logistic models used includes two well-acknowledged features (mean number of words/sentence and of letters/word)
- It confirmed that majority class may deform the prediction space

## Homogeneity of the annotations

- We compared the lexical and syntactic difficulty of textbook series among them (for each levels)
- We used ANOVA and MANOVA and noticed significant homogeneity issues in the corpus
- A relation between the type of pedagogical approach and homogeneity issues also appeared.

Take-home message: Test the homogeneity of annotations when using texts from educational sources (just as we would do for human annotators).

# References I

Collins-Thompson, K. and Callan, J. (2005).
Predicting reading difficulty with statistical language models.
*Journal of the American Society for Information Science and Technology*, 56(13):1448–1462.

Feng, L., Jansche, M., Huenerfauth, M., and Elhadad, N. (2010).
A Comparison of Features for Automatic Readability Assessment.
In *COLING 2010: Poster Volume*, pages 276–284.

François, T. and Fairon, C. (2012).
An "AI readability" formula for French as a foreign language.
In *Proceedings of the 2012 Conference on Empirical Methods in Natural Language Processing (EMNLP 2012)*, pages 466–477.

Heilman, M., Collins-Thompson, K., Callan, J., and Eskenazi, M. (2007).
Combining lexical and grammatical features to improve readability measures for first and second language texts.
In *Proceedings of NAACL HLT*, pages 460–467.

Pilán, I., Volodina, E., and Johansson, R. (2014).
Rule-based and machine learning approaches for second language sentence-level readability.
In *Proceedings of the Ninth Workshop on Innovative Use of NLP for Building Educational Applications*, pages 174–184.

# References II

Schwarm, S. and Ostendorf, M. (2005).
Reading level assessment using support vector machines and statistical language models.
*Proceedings of the 43rd Annual Meeting on Association for Computational Linguistics*, pages 523–530.

Vajjala, S. and Meurers, D. (2012).
On improving the accuracy of readability classification using insights from second language acquisition.
In *Proceedings of the Seventh Workshop on Building Educational Applications Using NLP*, pages 163–173.

Verlinde, S., Selva, T., and Binon, J. (2003).
Alfalex: un environnement d'apprentisage du vocabulaire français en ligne, interactif et automatisé.
*Romaneske*, 28(1):42–62.