



# Text Mining a Noisy Corpus of Old Regime French Periodicals A War Story

**François Dominic Laramée**  
**PhD Candidate in History, Université de Montréal**  
**September 12th, 2017**

## THE BEST LAID PLANS OF MICE AND MEN...

**« ... no plan of operations extends with any certainty beyond the first contact with the main hostile force. »**

**— Helmuth von Moltke the Elder (1800-1891)**



## THE ENEMY

Don Joseph Zagnol, Premier Médecin du Roi, aura la direction de cet établissement.

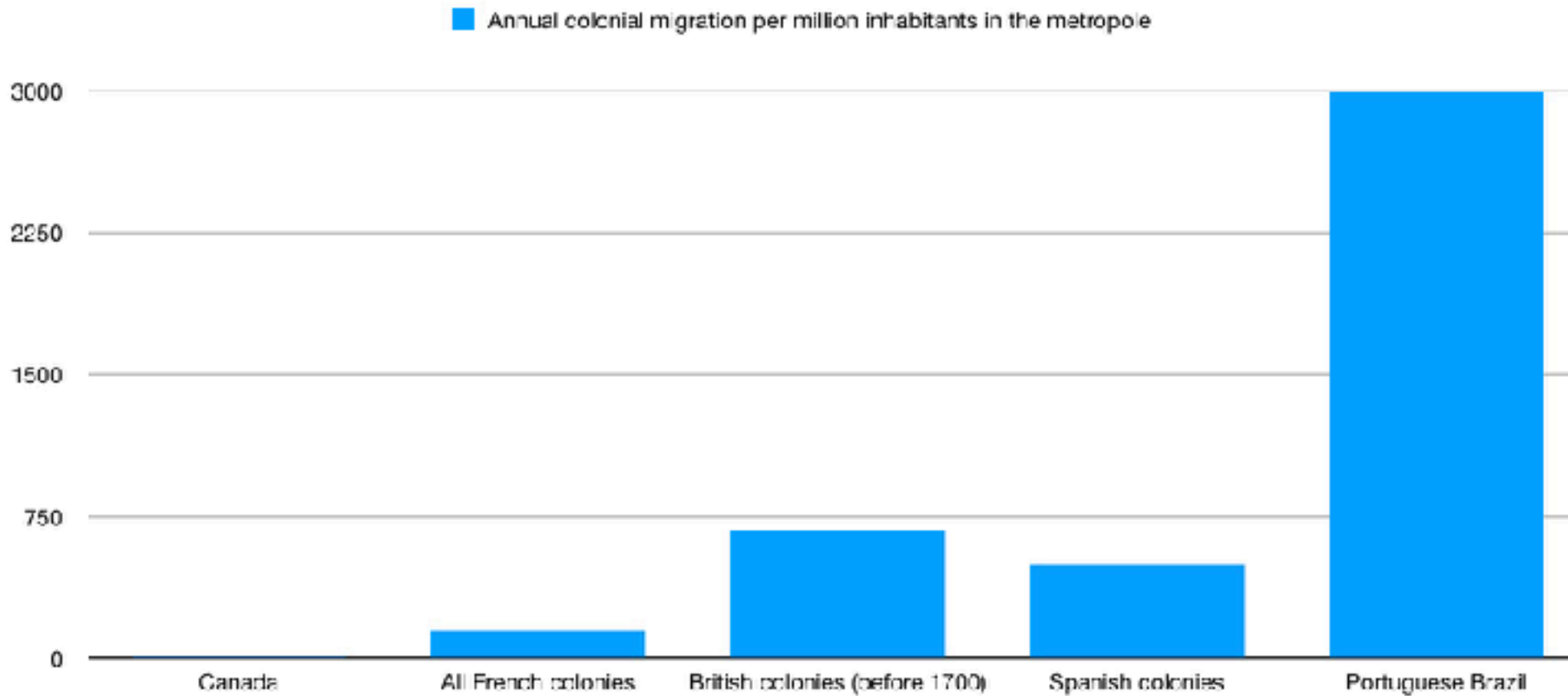
Le 1<sup>r</sup> du mois dernier, on ressentit à Maroc un affreux tremblement de terre, à la même heure qu'en Espagne. La plupart des maisons & des édifices publics de cette Ville ont été totalement renversés, & une grande multitude d'habitans a été ensevelie sous les ruines. Le nombre des personnes, qui ont péri, ne peut pas encore se fixer. A huit lieues de cette Ville, la terre s'est ouverte & a englouti une multitude de personnes.

L-e 1<sup>r</sup> du mois demih, On reRèntic à Maroc un affreux t- ViMe même beure qu'en Efpa◇ne, La plâpart des imifons & des-édificies publics de ctteVille ont été totalement renver \* une grande mul," : titude d'habitans aéreenfevelie fous les ruines. Le nom- PfS encore sebre des personnes, qui ne peut pas encore [e ftrer.

## Five parts

- Historical problem
- Corpus of Old Regime periodicals
- Methodology
- Gazette: a mostly digital analysis
- Other periodicals: a weakly digital analysis

# COLONIAL MIGRATION DURING THE OLD REGIME



*Data source: Landry (2006), « Les Français passés au Canada avant 1760 »*

# CORPUS

RECUEIL  
DES  
GAZETTES  
DE  
FRANCE.

*Année MDCCLVI.*



G. I. S. L. A.  
44

A P. A.

Du Bureau d'Adresse, aux  
vis-à-vis la rue S.

AVEC PRIVILEGE

MERCURE  
DE FRANCE  
DÉDIÉ AU ROI.

NOVEMBRE 1745.



LE  
JOURNAL  
DES  
SCAVANS,

POUR  
L'ANNÉE M. DCC. LV. I.  
JANVIER.

PARIS,  
Chez M. CAVELIER  
au Salon de la  
SOT, Quai de Conty  
de du Pont-Neuf.  
NULLY, au Palais.

XLV.

À Privilege du Roi

# THE SILENCE OF POPULAR MEDIA

ALMANACH

LITTÉRAIRE;

OU

ÉTRENNES

D'APOLLON;

CONTENANT de jolies Pièces en prose, & en vers,  
des saillies ingénieuses, des variétés intéressantes,  
& beaucoup d'autres Morceaux-curieux.

AVEC une Notice des Ouvrages nouveaux, remplie  
d'Anecdotes piquantes.

PAR

M. D'AQUIN DE CHATEAU-LYON.

Prix, trente-six sols.



Chez } MME LA VEUVE DUCHESNE, rue S. Jacques;  
DEPER DE MAISONNEUVE, rue du Foia S.

## Bibliothèque Bleue

[Database Home](#)

[The ARTFL Project](#)

### Search in Texts or Find Documents

Search for:

Display:  Context  KWIC  Similarity Search

## CORPUS-SPECIFIC OCR ISSUES

- Problematic source documents
- High error variance
- Errors do not follow usual patterns



# METHODOLOGY OVERVIEW

## HEURISTIC PROJECT DESIGN

- Exploratory methods => Salient patterns
- Salient patterns => Small keyword set
- Small keyword set => Noise neutralization

## KEYWORD-DRIVEN DATA RECONSTRUCTION

- Focused OCR correction
- Article subset selection
- Metadata extraction

# LEVENSHTEIN'S ALGORITHM

Keyword	Candidate	Levenshtein Distance	Operations
Amérique	Amérique	0	—
Amérique	Amrique	1	Delete « é »
Amérique	Cmévrique	2	A → C Insert « v »
Amérique	Musique	3	Delete « A » é → u r → s

# LEVENSHTEIN APPLIED TO THE *GAZETTE*

## SOME OF THE RECOVERED INSTANCES

Keyword	Keyword tokens	Alternate types	Alternate type occurrences	Alternate types as % of all tokens	Examples
Amérique/ d'Amérique/ l'Amérique	485	36	128	20,9 %	l'amerique (59), d'amcrique (2), ramèrique
Brésil	5	10	159	97,0 %	bresil (143), bretîl, brcfil
Canada/ Canadiens	139	3	3	2,1 %	canada*, en.canada
Colonie(s)	411	15	16	3,7 %	5lonie, coioniej
Domingue/ Saint- Domingue	88	11	13	12,9 %	jjomingue, saintdomingu e
<b>TOTAL</b>	<b>1867</b>	<b>103</b>	<b>532</b>	<b>22,2%%</b>	<b>Increase: 28,5%</b>

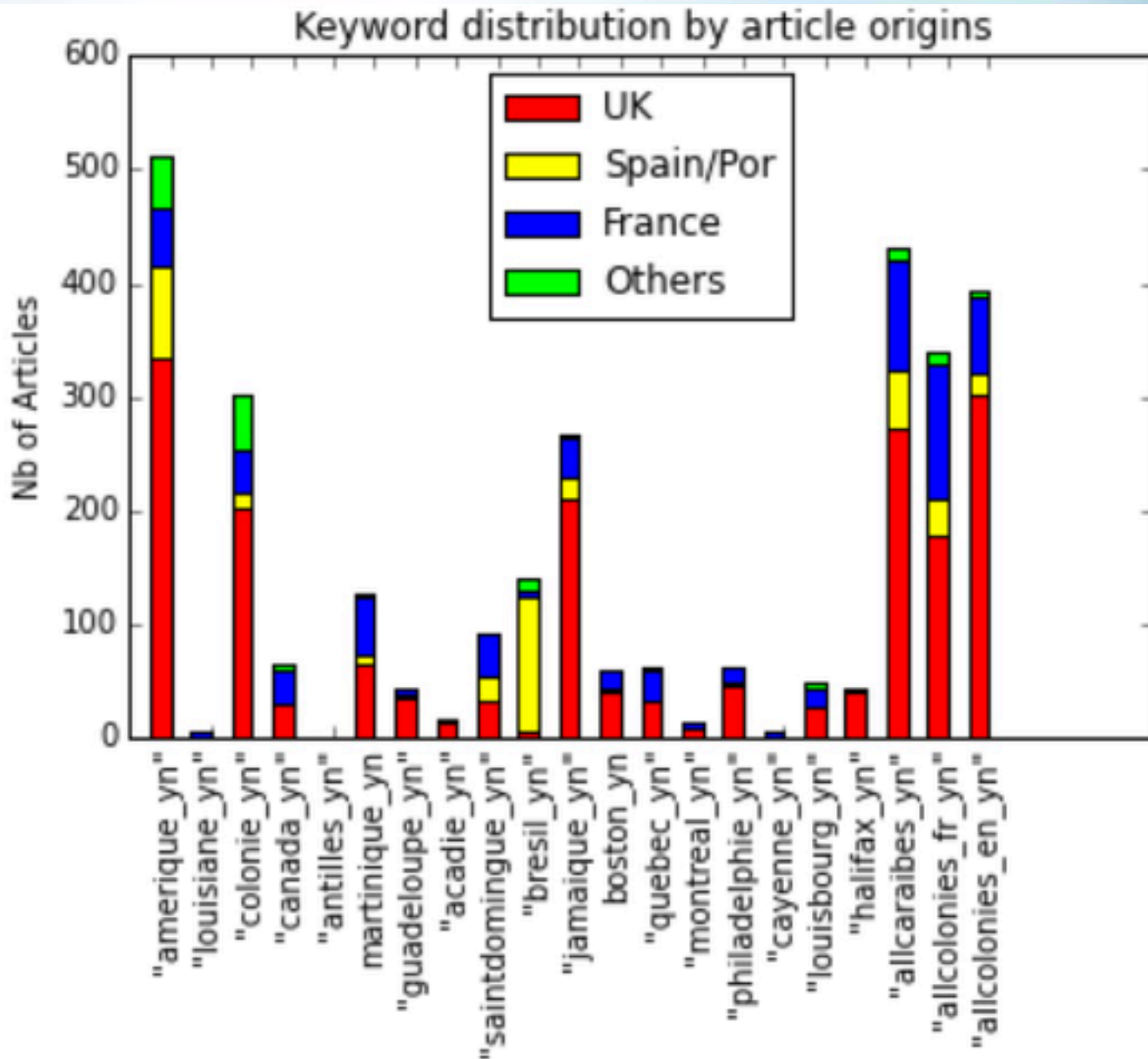
# GAZETTE CORPUS METADATA FILE (sample)

id	"ville"	"annee_pub"	"amerique_yn"	"louisiane_yn"	"colonie_yn"	"canada_yn"
am1740_1	Londres	1740	1	0	0	0
am1740_10	Madrid	1740	1	0	0	0
am1740_11	Londres	1740	1	0	0	0
am1740_12	Londres	1740	1	0	0	0
am1740_13	Dresde	1740	1	0	0	0
am1740_14	Madrid	1740	1	0	0	0
am1740_15	Londres	1740	1	0	1	0
am1740_16	Londres	1740	1	0	0	0
am1740_17	Londres	1740	1	0	0	0
am1740_18	Madrid	1740	1	0	0	0
am1740_19	Londres	1740	0	0	1	0
am1740_2	Madrid	1740	1	0	1	0

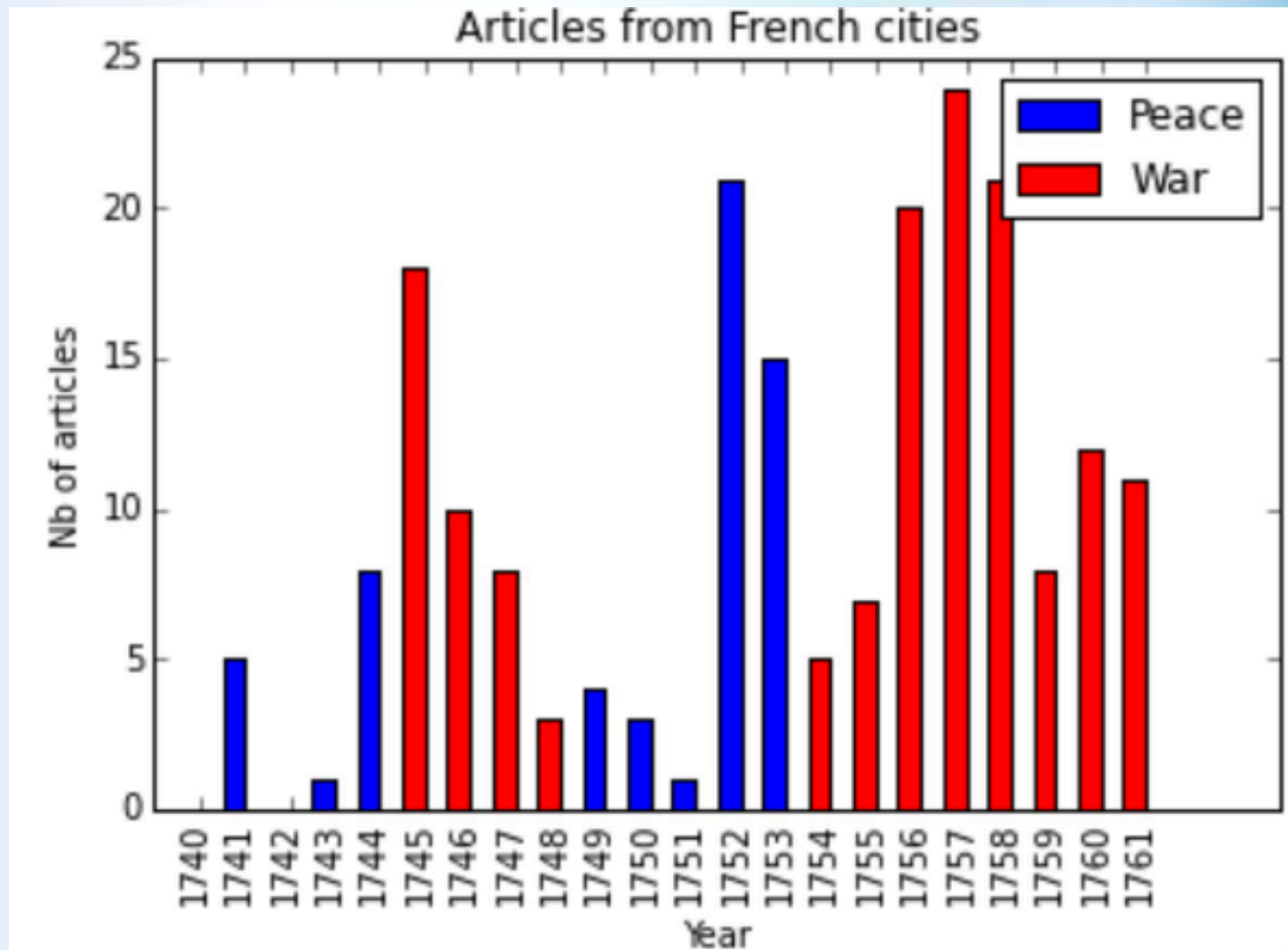
# GAZETTE CORPUS — ARTICLE ORIGINS



# GAZETTE CORPUS – GEOGRAPHIC KEYWORD DISTRIBUTION

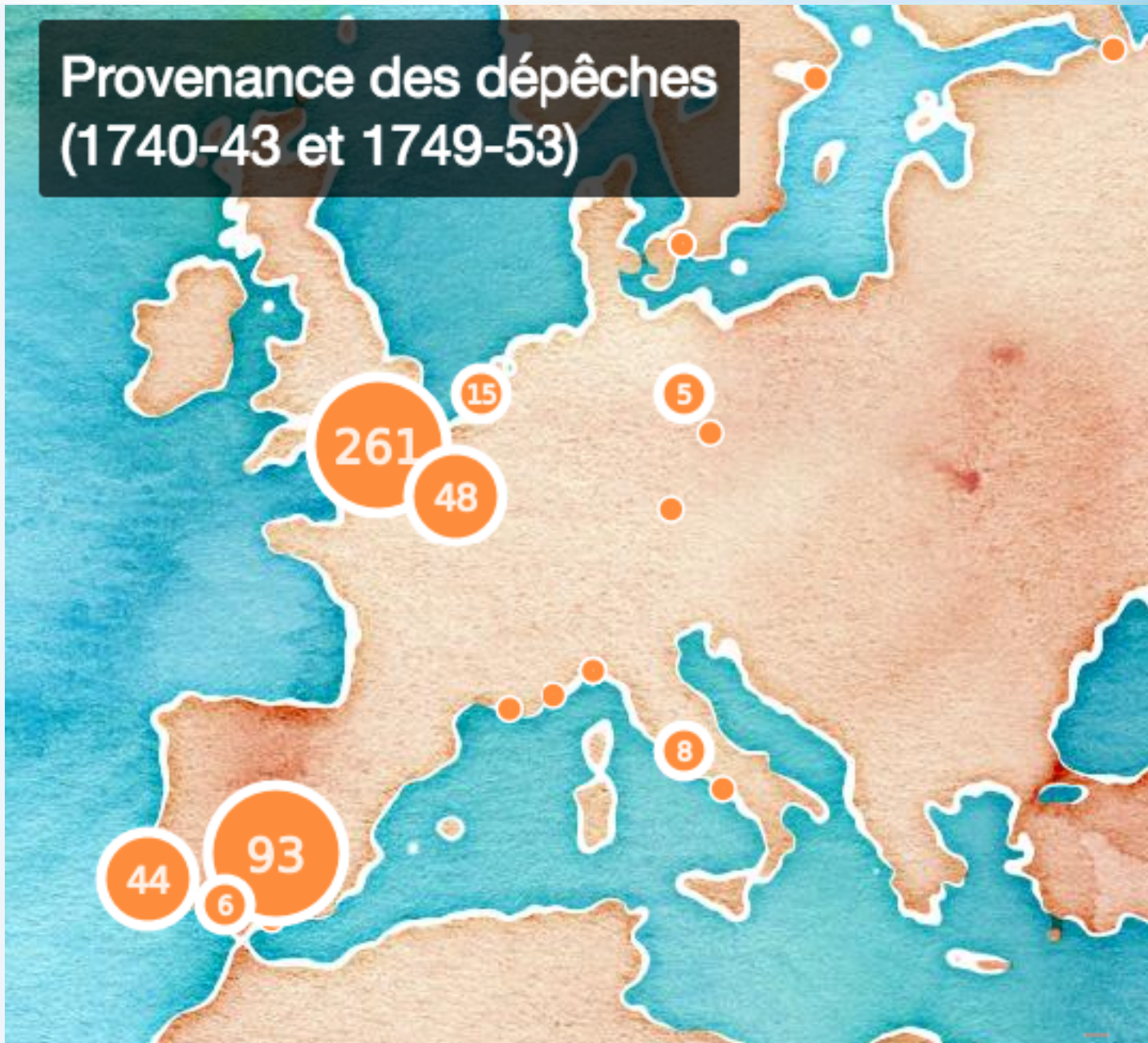


# GAZETTE CORPUS – FRENCH ARTICLES AS TIME SERIES



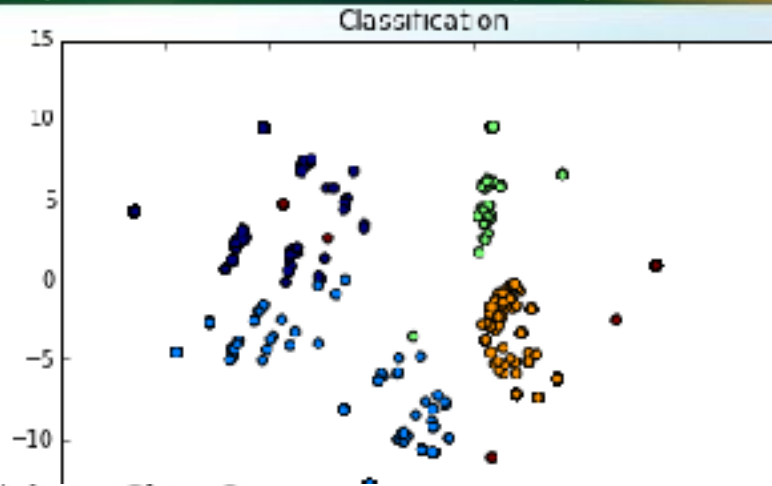
# GAZETTE CORPUS — PEACETIME GEOGRAPHIC DISTRIBUTION

Provenance des dépêches  
(1740-43 et 1749-53)

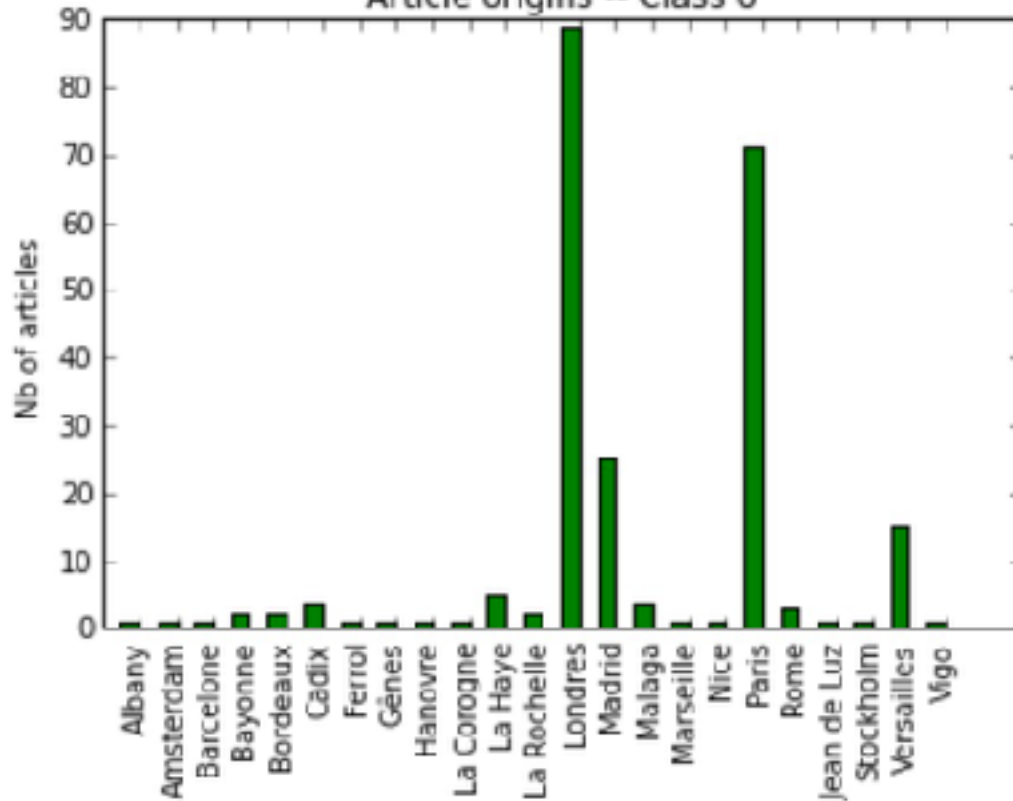




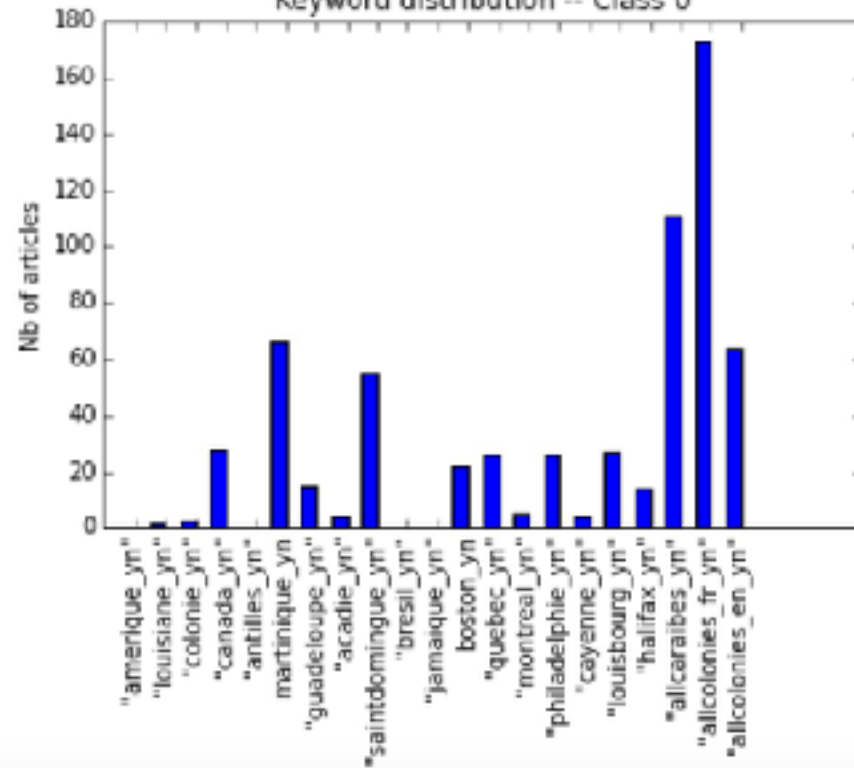
# GAZETTE CORPUS — K-MEANS CLUSTERING



Article origins -- Class 0



Keyword distribution -- Class 0



## **WHEN METADATA FAILS TOO?**

- Keyword distribution ==> Informed Close reading
- KWIC and Collocates as last remaining digital tools

## **THAT IS ENOUGH TO GET INTERESTING RESULTS**

- Land of opportunity...
- ... but also of great risk ...
- ... and irrevocably alien ...
- ... and less welcoming for the French than for others

# CONCLUSIONS

## THREE TIME SCALES, THREE IMAGINED AMERICAS

- Fast-moving current events: a dangerous, foreign land
- Slow-moving science: opportunities and risk
- Stable collective consciousness: imbalance of rewards

## DIGITAL CONTRIBUTION TO THE ANALYSIS?

- More limited than expected
- Nothing groundbreaking
- But *sufficient to support a historical argument*

# THANK YOU!

François Dominic Laramée

[fdl@francoisdominiclaramee.com](mailto:fdl@francoisdominiclaramee.com)

[www.francoisdominiclaramee.com](http://www.francoisdominiclaramee.com)

Twitter : @fdlaramee

**Work funded by a doctoral grant from the  
Fonds de recherche du Québec - Société et Culture**