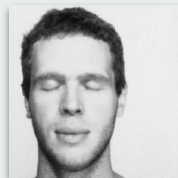




Namograph_{.antonomase.fr/}



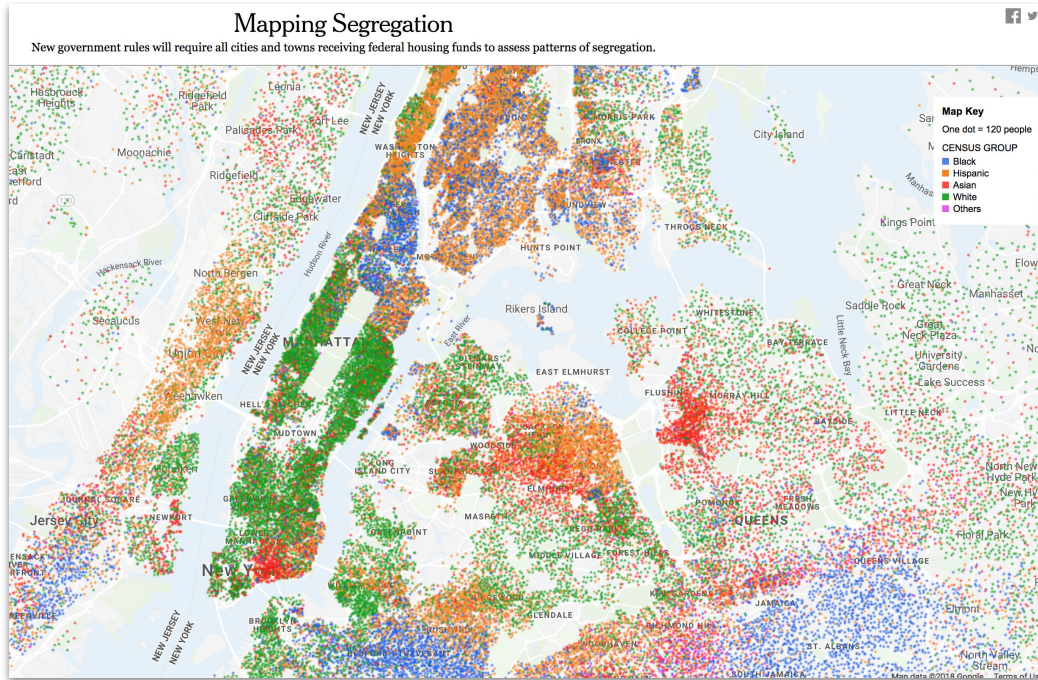
Antoine Mazieres
<https://antonomase.fr/>



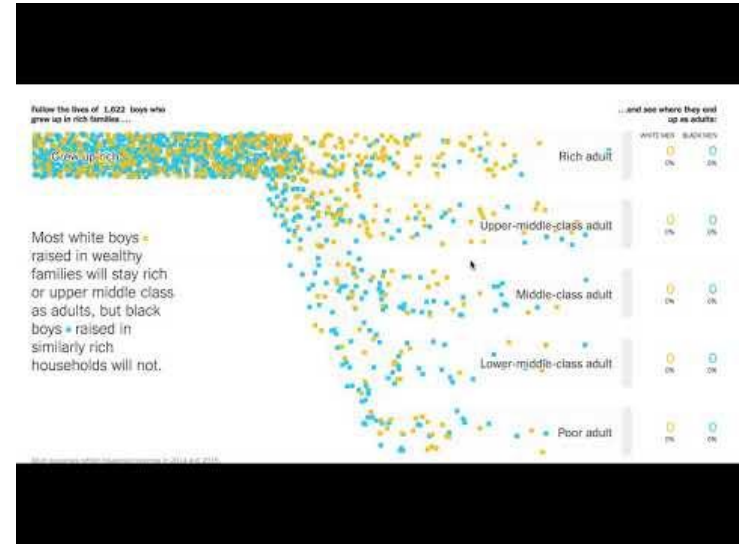
Camille Roth
<https://camilleroth.eu/>

Large-scale diversity estimation through
surname origin inference

Ethnicity statistics and discrimination studies

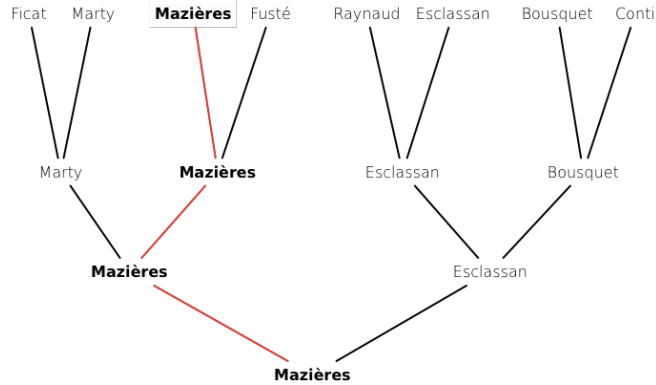


Source: <https://www.nytimes.com/interactive/2015/07/08/us/census-race-map.html>

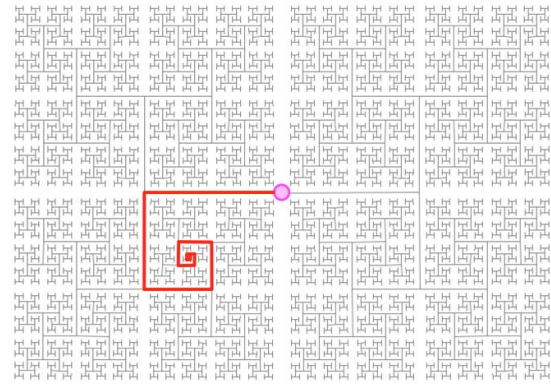


Source: <https://www.nytimes.com/interactive/2018/07/19/upshot/race-class-white-and-black-men.html>

Why surnames (still) mean something ?



My surname over 3 generations



A surname over 12 generations

Endogamy makes the name

Ethnicity *versus* Origin for onomastics



Ethnicity

- Subjective
- *Plural*
- Changing
- Identified individually
- Sometimes polemical

Origin

- **Objective**
- *Plural*
- **Stable**
- **Aggregated estimations**
- **More explicit biases**

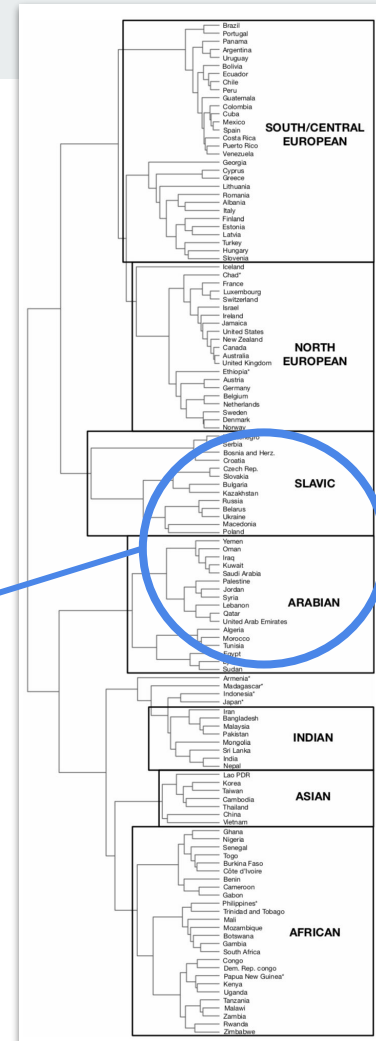
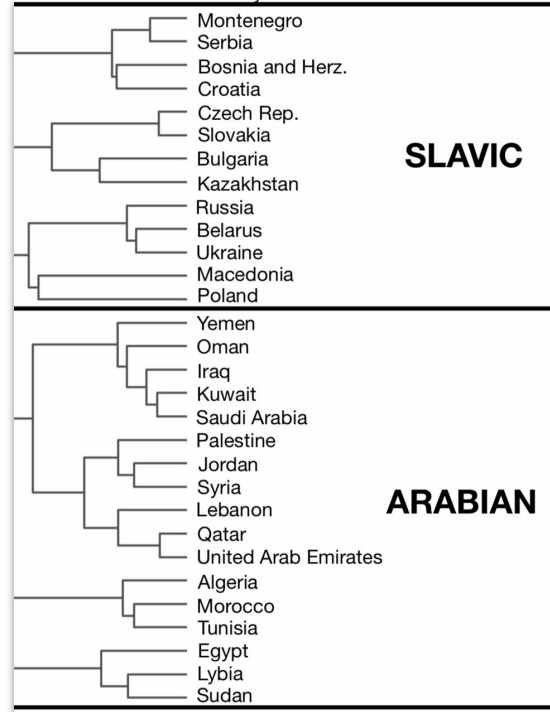
Crafting the learning data

- Need for **(Surname -> Origin)** tuples
- 25 millions (Author's surname - Affiliation's Country)
- Normalize and seeking "**Core Surnames**" with frequency
and *Herfindahl-Hirschmann Index*: **650k** core names
- Features with *n-grams*
- How to regroup countries ?



A data-driven typology of origins

- Regrouping *n-grams* by countries
- Hierarchical clustering (Ward) yields intuitive representation
- Hand-made tree cutting.
- A few irregularities



Learning a model to infer surname origin

- 650k ("core surnames" -> Origin)

(15% for testing)

- Multinomial Naive-Bayesian

Classifier w/ smoothing

- Heterogeneous results

Cluster	Core names		Class. Perf.	
	<i>Total</i>	<i>Evaluation</i>	<i>Precision</i>	<i>Recall</i>
African	30 748	4 529	0.43	0.61
Arabian	31 272	4 596	0.52	0.72
Asian	44 658	6 754	0.61	0.77
CS-European	189 624	28 668	0.81	0.71
Indian	68 145	10 067	0.63	0.72
N-European	216 465	32 469	0.78	0.62
Slavic	65 259	9 843	0.64	0.84
<i>Total</i>	<i>646 171</i>	<i>96 926</i>		

Measuring diversity

- Why **diversity** instead of **discrimination** ?
- Simple method: **ratio** between target and reference

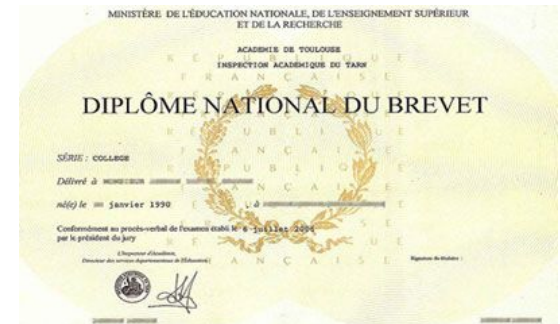
dataset

- **Brevet des collèges**

as reference

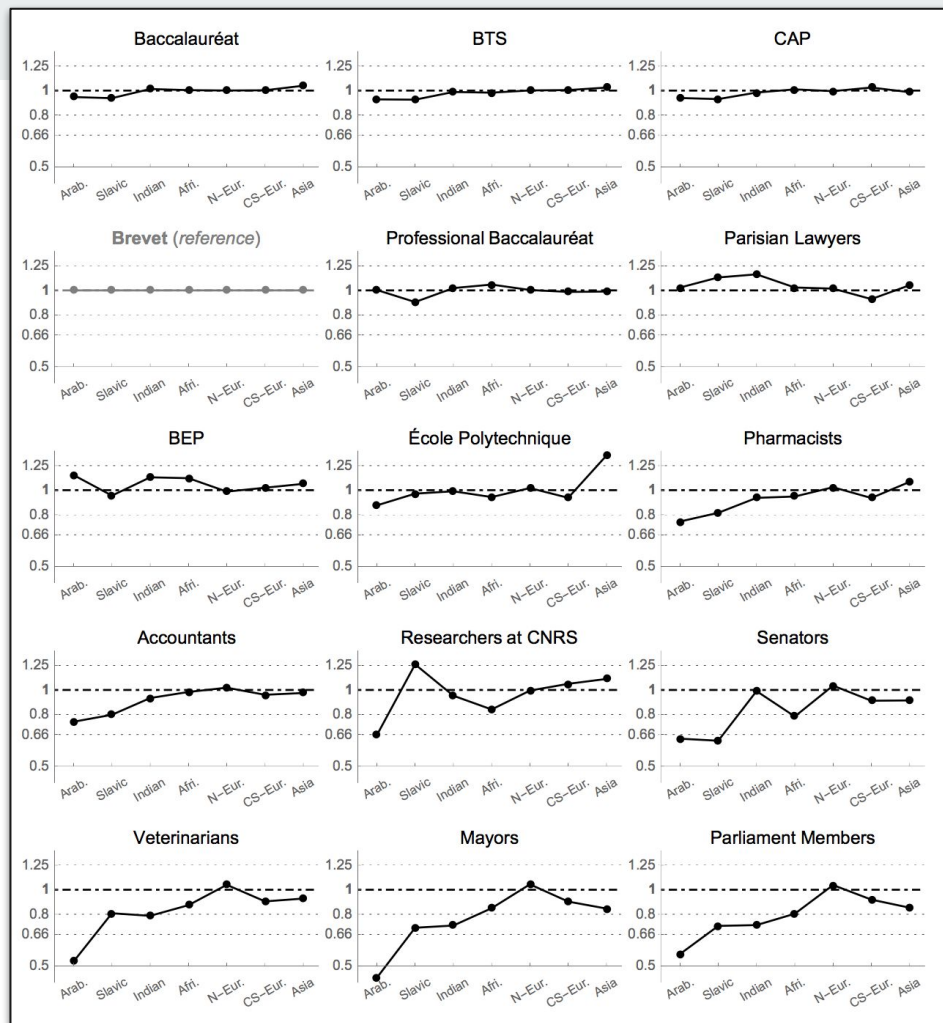
- Targets:

Name	List of surnames of all ...	nb. obs.
<i>Brevet</i>	Candidates to <i>Diplôme National du Brevet</i> in 2008 ⁵	562,952
<i>Baccalauréat</i>	Candidates to the nationwide <i>Baccalauréat (Général and Technologique)</i> in 2008	435,645
<i>BEP</i>	Candidates to <i>Brevets d'Études Professionnelles</i> in 2008	116,814
<i>CAP</i>	Candidates to <i>Certificats d'Aptitude Professionnelle</i> in 2008	98,364
<i>BTS</i>	Candidates to <i>Brevets de Technicien Supérieur</i> in 2008	87,917
<i>Professional Baccalauréat</i>	Candidates to <i>Baccalauréats Professionnels</i> in 2008	80,672
<i>Pharmacists</i>	Pharmacists registered in their <i>Ordre Professionnel</i> in 2017 ⁶	73,422
<i>Mayors</i>	Mayors of French cities ("communes") in 2014 ⁷	36,628
<i>Parisian Lawyers</i>	Lawyers registered in the Parisian Bar Association in 2017 ⁸	32,021
<i>École Polytechnique</i>	Students at <i>École Polytechnique</i> (1958-2016) ⁹	23,058
<i>Accountants</i>	Accountants registered in their <i>Ordre Professionnel</i> in 2017 ¹⁰	20,946
<i>Veterinarians</i>	Veterinary physicians registered in their <i>Ordre Professionnel</i> in 2017 ¹¹	15,710
<i>Researchers</i>	Researchers at <i>Centre National de la Recherche Scientifique</i> in 2017 ¹²	12,657
<i>Parliament Members</i>	Parliament Members of <i>Assemblée Nationale</i> (1958-2016) ¹³	8,326



Results

- Plots and origins ordered by proximity
- Clusters of types of occupational groups (political functions, state exams) and particular profiles
- Certain comparison with discrimination studies results



Questions



Advertisement

An advertisement graphic. On the left is a small photo of a man with glasses. To his right is a large speech bubble containing text.

I'm looking for
a **post-doc position**
for another
great project !
<3

Project site: <https://namograph.antonomase.fr/>

Get in touch: {mazieres,roth}@cmb.hu-berlin.de