



VNiVERSiDAD

DSALAMANCA CAMPUS OF INTERNATIONAL EXCELLENCE

> Dr. Carlos Arcila Calderón University of Salamanca @carlosarcila

BIG DATA ANALYSIS FOR MIGRATION STUDIES

- Collecting and modelling <u>unstructured data</u> (e.g. social media contents, parliamentary speeches, video records, etc.).
 - Case of study: Large-scale detection of hate speech against migrants and refugees.
- Collecting and modelling <u>structured data</u> (e.g. databases, surveys, official statistics, etc.).
 - Case of study: Using machine learning and synthetic populations to predict social acceptance of asylum seekers in European regions

INTRODUCTION

- The mass arrival of forced displaced people in Europe is generating an unprecedented political and demographic impact
- In this context, **migration management has become one of the main challenges** for politicians and competent organizations
- Related to this topic, some of the main programs of the European Union and the UN Refugee Agency (UNHCR) in recent years was intended settle new refugees arriving

in Europe and to relocate refugees and asylum seekers who remained in first-reception

countries

CONTEXT

- One of the most important variables for resettlement processes to occur in a positive way is the probability of acceptance of the receiving societies
- Until now, the criteria taken into account to manage migration and refuge have been fundamentally economic and political (Jones, Teytelboym 2017), but the contexts of reception by the local population have not been taken into account



 The prediction of attitudes at the regional level would mitigate the risks of exclusion through the planning of integration policies, anticipating which populations could accept or oppose more resistance to the arrival of refugees

CONTEXT

- Previous studies indicate that the perception of migration is strongly influenced by individual sociodemographic characteristics (Finney and Peach, 2004)
- Some of these variables may be *nationality*, *age*, *education*, *income* or *political ideology*
- This gives enormous importance to sociodemographic factors to predict refugee acceptance, since they can condition the degree of acceptance or rejection of refugees (Finney and Peach, 2004; Schweitzer et al., 2005; Cameron et al., 2006; Song, 1992)



CONTEXT

 It is also known that the level of acceptance of migrants and refugees varies between countries, but also between regions (Crawley, Drinkwater and Kauser, 2013)



 However, most refugee support surveys (such as the Eurobarometer or the European Social Survey) do not have enough geo-localized data to make estimates at the regional or local level

BIG DATA SOLUTIONS



- We can estimate the probability of acceptance of refugees in each European region based on:
 - data from social media using supervised text classification to detect hate speech against migrants
 - data from existing surveys using computational approaches that combine machine learning and synthetic populations

DATA-BASED MIGRATION PLANNING

- Administrations continue to base their decisions on traditional studies and data
- However, the academy has already begun to understand the great potential of big data and computational methods to study social processes and movements, including migrant and refugee flows
- These studies based on large amounts of data should serve as a basis for making decisions and managing migration flows



DATA-BASED MIGRATION PLANNING



Some examples of these types of studies are:

<u>Hawelka et al (2014) and Zagheni et al (2014)</u>: They use geographically placed Twitter messages to identify global mobility patterns and predict migratory flows

Lamana et al (2018): They use data extracted from Twitter, such as language and spatial-temporal communication patterns of individuals, to estimate the power of social integration of different cities in the world

DATA-BASED MIGRATION PLANNING

Bansak et al (2018): These authors develop a model that uses multiple data to assign a specific place to refugees according to their profile, and thus improve integration



They use information about the basic characteristics of refugees (country of origin, language skills, gender, age, etc.), time of arrival, assigned location and <u>average employment</u> <u>success</u>

So they create a set of supervised machine learning models to predict the expected job success of refugees based on their background characteristics

SOLUTION 1

 Big data from social media using supervised text classification to detect hate speech against migrants

INTRODUCTION

- Online hate speech has increased exponentially in recent years
- This increase is especially worrying on social media like Twitter (Bartlett et al., 2014)
- Twitter has updated its Hate
 Speech Policy to counter the hateful conduct



- Hate speech implies promoting messages that encourage rejection, humiliation, harassment, discrediting and stigmatization of individuals or social groups
- Any type of expression that propagates, promotes or justifies hate and intolerance towards the otherness

(Council of Europe, Recommendation Nº 97, 1997)





Source: Miró (2016)

JUSTIFICATION

 Recent studies establish correlations between hate speech spread through social media such as Twitter and hate crimes that take place in particular locations



Müller, Karsten and Schwarz, Carlo, Fanning the Flames of Hate: Social Media and Hate Crime (December Mutter, Karsten and Schwarz, Carto, Fanning the Flames of Hate: Social Media and Hate Chine (Decentoer 6, 2017). Available at SSRN: https://ssm.com/abstract=3082972 or http://dx.doi.org/10.2139/ssm.3082972 Fanning the Flames of Hate: Social Media and Hate Crime* Karsten Müller[†] and Carlo Schwarz[‡] December 6, 2017 This paper investigates the link between social media and hate crime using hand-Abstract collected data from Facebook and Twitter. We study the case of Germany, where the recently emerged right-wing party Alternative für Deutschland (AfD) has developed a major social media presence. Using a difference-in-differences design, we show that right-wing anti-refugee sentiment on Facebook predicts violent crimes against refugees in otherwise similar municipalities with higher social media usage. Consistent with social media being the driving force, the effect decreases with internet outages; increases with user network interactions; is not driven by the news cycle; and does not hold for posts unrelated to refugees. We find similar evidence for the United States, where President Trump's twitter activity strongly predicts hate crimes against the minorities targeted in his tweets, but not other minorities. We find no effect for the period before Trump's presidential campaign or measures of general anti-minority sentiment. JEL classification: D74, J15, Z10, D72, O35, N32, N34. Keywords: social media, hate crime, minorities, Germany, US, AfD, Trump

JUSTIFICATION

 Recently several prototypes of online hate speech detectors have been developed

 All of them have significant limitations





nen et al. Hum. Cent. Comput. Int. Sci., (2020) 10:1 2//dol.org/10.1186/s13673-019-0205-6 (2020) 10:1 Human-centric Computin Developing an online hate classifier and Information Sciences for multiple social media platforms Joni Salminen 1.2 Maximilian Hopr³, Shammur A. Chowdhury¹, Soon-gyo Jung¹, Hind Almerekhj⁴ Foxicity, Social media, Machine I/ Toduction sline hate, described as abusive language [1], aggression [2], cyberbullying [3, 4], hate. invare [c] investige [c] normanial attacks [c] normanian [c] racism [c] environ [tail] Inline hate, described as abusive language [1], aggression [2], cyberbuilying [3, 4], fiate-fulness [5], insuits [6], personal attacks [7], provocation [8], racism [9], sexism [10], sector (1) on traventities [10], has been identified as a major thread on radius envirol module tuness [5], insuits [6], personal attacks [7], provocation [6], faction [6], sexism [10], threats [11], or toxicity [12], has been identified as a major threat on online social media infarforme. Dow Research Control [13] resource that among data while in the limited threats [11], or toxicity [12], has been identified as a major threat on online social media platforms. Pew Research Center [13] reports that among 4248 adults in the United crasses are have been associated by a second by puttoring, rew sessarcin center [15] reports that allong 4246 adults in the clined States, 41% have personally experienced harassing behavior online, whereas 66% with sent directed towards others. Around 22% of adults have e nicence turescolar contraction of the contraction o the most prominent grounds for such teach weilay and seven transfer they ways of flagging offensive and hateful content, only 17% of all adults have i summent (0%), annung outer (1985 or narassument, 50000 interna-comment grounds for such toxic behavior. Even though they ing conversation, whereas only 12% of adults have reported so



OBJECTIVE

 Development and evaluation of a detector of online hate speech in Spanish transmitted on Twitter, and for 3 specific types of prejudice:



Racism/xenophobia
 Sexual orientation
 Politic ideology

LARGE-SCALE DETECTION STRATEGY



SUPERVISED CLASSIFICATIO N



 The predictive model is trained and generated from the given examples, manually classified for it



-6

-4

-2

0

2 4



GENERATING THE TRAINING CORPUS FOR STOP-HATE

1. To define the type of hate speech to detect, and generate a manual or codebook in which those types of hate to be identified are explained with clear examples



GENERATING THE TRAINING CORPUS FOR STOP-HATE

2. To **generate filters**, identifying and selecting terms and combinations of words that could indicate and represent hate motivated by the 3 predefined types of prejudice:

- Racism/xenophobia, - sexual orientation and - political ideology

LINGUISTIC FILTERS

Racism / Xenophobia panchito/a/o/os/as guachupin/o/a/os/as tiraflechas sudaca machupichu chino/a/os/as AND mandarin/o/a/os/as sin AND papeles, cierre AND frontera/s invasion AND silenciosa efecto AND llamada migracion/inmigracion AND masiva cerrar/cierre AND frontera/s migrante/s, inmigrante/s, refugiado/a/os/as, desplazado/a/os/as, africano/a/os/as, asiatico/a/os/as, chino/a/os/as, negrato/a/os/as, moreno/a/os/as, mulato/a/os/as, moro/a/os/as, moreno/a/os/as, mulato/a/os/as, arabe/s, islamista/s, islamico/a/os/as	Hate reason	0. Fixed terms or combinations	1. Primary words (descriptive within the hate group)
	Racism / Xenophobia	panchito/a/o/os/as guachupin/o/a/os/as tiraflechas sudaca machupichu chino/a/os/as AND mandarin/o/a/os/as sin AND papeles, cierre AND frontera/s invasion AND silenciosa efecto AND llamada migracion/inmigracion AND masiva cerrar/cierre AND frontera/s	migrante/s, inmigrante/s, refugiado/a/os/as, desplazado/a/os/as, africano/a/os/as, asiatico/a/os/as, chino/a/os/as, latino/a/os/as, negro/a/os/as, negrato/a/os/as, moreno/a/os/as, mulato/a/os/as, moro/a/os/as, musulman/a/es/as, arabe/s, islamista/s, islamico/a/os/as

2. Secondary words (accompanying and expressing hatred)

carga, amenaza, inseguridad, peste, mierda, lacra, gentuza, basura, escoria, fuera, largo, culo, puto/a/os/as, jodido/a/os/as, maldito/a/os/as, sucio/a/os/as, peligro, peligroso/a/os/as, radical/es, terror, terrorista/s, asco/asqueroso/a/os/as, ilegal/es, indocumentado/a/os/as, invasion, invasor/a/es/as, invadir, invade, invaden, criminal/es, robo/s, robar, roba, roban, delincuencia, delincuente/s, delinquir, delinque, delinquen, vago/a/os/as, zangano/a/os/as, violencia, violento/a/os/as, violacion, violador/es, violar, viola, violan, trafico, traficar, trafica, trafican, ocupacion, ocupar, ocupa, ocupan, saqueo/s, saquear, saquean, asalto, asalta, asaltan, paga/paguita/a, subvencion/es, ayuda/s, carterista/s, hurto/s, deportacion/es, deportar, deportaria, deportaba, deporte/deporten, deportarlo/a/os/as, repatriacion/es, repatriar, repatriaria, repatrie, repatrien, repatriaba, reventar, revienta, revienta, revienta, revienta, reventaba, reventarlo/a/os/as, matanza, matar, mato, mata, mataria, mataba, matarlo/a/os/as

Generating the training corpus for STOP-HATE

3. To **download the filtered tweets** from the Twitter API using the pre-set combinations as

word = sryingeres 'cierre\nfrontera', 'cierre\nfronteras', 'anti\ninmigracion', u'anti\ninmigración', u'invasión\neuropea', 'invasion\neuropea', u'invasión\nsilenciosa', 'invasion\nsilenciosa', 'efecto\nllamada', u'inmigración\nmasiva', u'migración\nmasiva', u'inmigracion\nmasiva', u'inmigración\nmasiva', 'machupicchu', 'machupicchu', 'guachupin', 'guachupino', 'guachupina', 'sudaca', u'tráfico\nmigrante', u'tráfico\ninmigrante', u'tráfico\ninmigrantes', 'repatriar', 'repatriarlo', 'repatriarlos', u'repatriación', 'repatriación', 'deportar', 'deportarlos', u'deportación', 'deportacion', 'deportación', 'depor

'desplazado\nputo', 'desplazado\njodido', 'desplazado\nmaldito', 'desplazado\nsucio', 'desplazado\nasqueroso', 'desplazado\nasco', 'desplazado\nmierda', 'desplazados\nputos', 'desplazados\njodidos', 'desplazados\nmalditos', 'desplazados\nsucios', 'desplazados\nasquerosos', 'desplazados\nasco', 'desplazados\nmierda', 'refugiado\nputo', 'refugiado\njodido', 'refugiado\nmaldito', 'refugiado\nsucio', 'refugiado\nasqueroso', 'refugiado\nmierda', 'refugiados\nmalditos', 'refugiados\nmald

'moro\nputo', 'moro\nsucio', 'moro\nasco', 'moro\nmierda', 'moro\njodido', 'moro\nmaldito', 'moro\nasqueroso', 'moros\nputos', 'moros\nputos', 'moros\nputos', 'moros\nputos', 'moros\nsucios', 'moros\nsucios', 'moros\nsucios', 'moros\nsucios', 'moros\nasco', 'moras\nputa', 'mora\njodida', 'mora\njodida', 'mora\nasquerosa', 'moras\nmierda', 'moras\nputas', 'moras\nputas', 'moras\nputas', 'moras\nputas', 'moras\nputas', 'moras\nsucios', 'moras\nsucios', 'moras\nsucios', 'moras\nputas', 'moras\nsucios', 'moras\nsuc

'chino\mandarin', u'chino\mandarín', 'chino\nmandarino', 'china\nmandarina', 'chino\namarillo', 'china\namarilla', 'chino\nputo', 'china\nputa', 'chino\njodido', 'chino\nmierda', 'chino\nasqueroso', 'chino\nsucio', 'chino\nasco', 'chino\nmaldito', 'chinos\nputos', 'chinos\njodidos', 'chinos\nmalditos', 'chinos\nsucios', 'chinos\nsucios', 'chinos\nasquerosos', 'chinos\nasco', 'chinos\nasco', 'chinos\nmalditos', 'chinos\nmalditos', 'chinos\nsucios', 'chinos\nasquerosos', 'chinos\nasco', 'chinos\nmalditos', 'chinos\nsucios', 'chinos\nasquerosos', 'chinos\nmalditas', 'chinas\njodidas', 'chinas\nmalditas', 'chinas\nmalditas', 'chinas\nsucias', 'chinas\nsucias', 'chinas\nsucias', 'chinas\nasco', 'chinas\nsucias', 'chinas\nasco', 'chinas\nsucias', 'chinas\nsucias',

'chinas\nsucias', 'chinas\nasquerosas', 'chinas\nasco', 'chinas\nasco', 'chinas\nasco', 'chinas\nasco', 'u'asiático\nsucio', u'asiático\nsucio', u'asiático\nasqueroso', u'asiático\nputo', u'asiático\nputo', 'u'asiático\nmaldito', u'asiático\nmaldito', u'asiático\nsucio', u'asiático\nasqueroso', u'asiático\nputo', u'asiático\nputo', 'u'asiático\nmaldito', u'asiático\nsucio', u'asiático\nasqueroso', u'asiático\nputo', 'u'asiático\nputo', 'u'asiático\nmaldito', u'asiático\nsucio', u'asiático\nsucio', u'asiático\nasqueroso', u'asiático\nputo', 'u'asiático\nputo', 'u'asiático\nputo', u'asiático\nputo', 'u'asiático\nputo', 'u'asiático\nputo', u'asiático\nputo', u'asiático\nputo', u'asiático\nputo', u'asiáticos\nputo', u

'latino\nputo', 'latino\njodido', 'latino\nmaldito', 'latino\nsucio', 'latino\nasqueroso', 'latino\nasco', 'latino\nmierda', 'latinos\nputos', 'latinos\nputos', 'latinos\nsucios', 'latinos\nasquerosos', 'latinos\njodidos', 'latinos\nmalditos', 'latinos\nsucios', 'latinos\nasquerosos', 'latinos\nasco', 'latinos\nmierda', 'latinos\nsucias', 'latinos\nsucios', 'latinos\nasquerosos', 'latinos\nasco', 'latinos\nsucias', 'latinos\nsucios', 'latinos\nsucias', 'latinos\nsu

Generating the training corpus for STOP-HATE

Compilation of a total set of 72,000 filtered tweets:

- 24,000 tweets from the racism / xenophobia group
- 24,000 tweets from the political ideology group
- 24,000 tweets from the sexual orientation group
- 4. To **annotate and classify the downloaded filtered tweets** with 2 coders per message on the **Doccano** platform:
 - All the downloaded tweets from each group were classified by a project researcher and various secondary coders

ANNOTATION ON THE DOCCANO PLATFORM

← → C â stophatespanish.herokuapp.com/projects/37/	#	2	3
Aplicaciones		Projects Logout	
Search document	3000 /3000		
About 12 results (page 2 of 3)	DESCARTES S-z NO ODIO z ODIO x		
 Gracias @madrid_vox por que estamos hartos de moros y mierda 	ODIO ×		
 ver la tele con mi tío es esto; ponemos quién quiere ser mil @McAvenalta @pavisnovis Moro de mierda! 	Gracias @madrid_vox por que estamos hartos de r España	noros y mierdas en	
 @yaniravargasf Necesitaríamos un Sergio Morodifícil ! Ya v 			
 @Santi_ABASCAL MIERDA DE MOROS DESAGRADECIDOS 	< € 6 / 3000	> >	

Generating the training corpus for STOP-HATE

- 5. To check the inter-coder reliability, crossing all classified tweets to reach TOTAL inter-coder agreement
- 6. To clean the resultant dataset, discarding the previously rejected tweets and tweets without agreement

Distribution of frequencies and percentages of tweets coded for each of the hate reasons



* The label discards corresponds to tweets without agreement plus those discarded in the classification process



Generating the predictive models



RESULTS

- We generated 6 predictive models with shallow algorithms, one based on the votes of those models; and a model generated with Recurrent Neural Networks
- In summary, we developed 8 models for each of the 3 groups, a total of 24 classification models

Shallow Algorithms	Accuracy	F-Score	AUC- ROC
Original Naive Bayes	.67	.73	.65
Naïve Bayes for Multinomial Models	.75	.83	.64
Naïve Bayes for Bernoulli's multivariate Models	.68	.81	.50
Logistic Regression	.78	<mark>.</mark> 84	.71
Lineal classifiers with SGD training	.75	.82	.69
Support Vector Machines	.76	.83	.71
Classifier based on the votes of the other models	.76	.82	.68
Deep modelling	Accuracy	F-Score	AUC- ROC
Recurrent Neural Networks	.86	.78	.92

Evaluation measures of the clasification algorithms

Evaluation measures of the models generated with each of the classification algorithms

1



NEXT STEP

PHARM project aims to develop a detector of hate speech towards migrants and refugees in Spain, Italy and Greece, in more online sources, and with the added implementation of early warning systems that will allow predicting possible increases in geo-located hate crimes



PREVENTING HATE AGAINST REFUGEES AND MIGRANTS

SOLUTION 11

 Big data from existing surveys using computational approaches that combine machine learning and synthetic population METHOD 1. Predicting individual acceptance



- In the first stage, data from the available Eurobarometers that included the objective question "–Our country– should help refugees?" (objective characteristic) were extracted, along with 8 other variables that could explain this attitude (descriptive characteristics)
- In total, data were collected from 165,089 people distributed in 5 periods between 2015 and 2017
- Of the 47 descriptive variables, only sociodemographic characteristics that could later be identified in a census were used:

Country, origin, gender, age, occupation, marital status, education type of community and household composition



METHOD 1. Predicting individual acceptance

- To facilitate the application of the algorithms, all measurements were standardized, turning them into dichotomous variables
- After eliminating the remaining values, a final matrix of 112,837 cases x 48 characteristics was obtained

With this data, supervised machine learning models were built to predict the individual support of Europeans to refugees



	Variable	Description
0	Support for refugees (target variable)	1 = Not Support, 0 = Support
1	Country	Dummy (0-1): 'country_BALGARIJA', 'country_BELGIQUE', 'country_CESKA REPUBLIKA', 'country_DANMARK', 'country_DEUTSCHLAND OST', 'country_DEUTSCHLAND WEST', 'country_EESTI', 'country_ELLADA', 'country_ESPANA', 'country_FRANCE', 'country_GREAT BRITAIN', 'country_HRVATSKA', 'country_IRELAND', 'country_ITALIA', 'country_KYPROS', 'country_LATVIA', 'country_LIETUVA', 'country_LUXEMBOURG', 'country_MAGYARORSZAG', 'country_MALTA', 'country_NEDERLAND', 'country_POLSKA', 'country_PORTUGAL', 'country_ROMANIA', 'country_SLOVENIJA', 'country_SLOVENSKA REPUBLIC', 'country_SLIOMI' 'country_SVERIGE' 'country_ÖSTERPEICH'
2	Origin	0= Foreigner, 1= Native
3	Marital status	Dummy (0-1): 'Married', 'Partnership', 'Single', 'Divorced/Separated', 'Widow'
4	Education	0= No formal education <10, 1= ISCED Level 1. Primary education 10-12, 2= ISCED Level 2. Lower secondary education 13-15, 3= ISCED Level 3. Upper secondary education 16-18, 4= ISCED Level 4. Post secondary non-tertiary education, ISCED Level 5, First stage of tertiary education, ISCED Level
5	Gender	Dummy (0-1): 'Woman', 'Man'
6	Age	0= under 15 years, 1= 15 to 29 years, 2= 30 to 49 years, 3= 50 to 64 years, 4= 65 to 84 years, 5= 85 years and over
7	Occupation	Dummy (0-1): 'Employed', 'Not active', 'Unemployed'
8	Type of community	Dummy (0-1): 'Small/middle town', 'Rural area or village', 'Large town'
9	Household composition	0= 1 person, 1= 2 persons, 2= 3 to 5 persons, 3= 6 and more persons

Variables taken from the Eurobarometer to model the level of support for refugees of a European citizen

METHOD 1. Predicting individual acceptance



- Social acceptance of refugees was modeled using logistic regression (LR), decision trees (DT), random forests (RF), support vector machines (SVM) and nearest k-neighbors (K-NN) on the data Longitudinal Eurobarometer, obtaining good prediction accuracy in all algorithms
- With these models the probability (0-1) of refugee support of any European citizen could be predicted
- Some of the factors with the most significant coefficients were the country and education in logistic regression and age and education in random forests

	Algorithm	Accuracy
Accuracy of the models generated with each of the algorithms	Logistic Regression (LR)	0.73
	Decision Tree (DT)	0.70
	Random Forest (RF)	0.69
	Support Vector Machine (SVM)	0.72
	SVM with kernel and gamma	0.65
	K Nearest Neighbour (K-NN)	0.64



Weights of the variables in the random forest algorithm

V



METHOD 2. Predicting social acceptance

To overcome the limitations of individual predictions and extend predictive models to large populations, we simulate the sociodemographic characteristics of 2,710,000 European citizens, corresponding to 10,000 in each of the 271 geographic regions (NUTS2) Specifically, we generate synthetic populations for each region using

distributions of the sociodemographic characteristics obtained in the Eurostat census data



Jupyter synthetic_cities_10-10 (autosaved)		ę	Logout
File Edit View Insert Cell Kernel Widgets Help		Py Py	hon 3 O
\bullet			
<pre>country_SLOVENSKA_REPUBLIC = np.array([1,0]).cumsum().searchsorted(np.random.sample(Total_Population)) country_SUOMI = np.array([1,0]).cumsum().searchsorted(np.random.sample(Total_Population)) country_SVERIGE = np.array([1,0]).cumsum().searchsorted(np.random.sample(Total_Population)) country_ÖSTERREICH = np.array([1,0]).cumsum().searchsorted(np.random.sample(Total_Population))</pre>	on))		
<pre>#type_community #0= NO, 1= YES</pre>			
<pre>type_community_Large_town = np.array([1,0]).cumsum().searchsorted(np.random.sample(Total_Population type_community_Rural_area_village = np.array([1,0]).cumsum().searchsorted(np.random.sample(Total_Pop type_community_Small_middle_town = np.array([0,1]).cumsum().searchsorted(np.random.sample(Total_Pop type_community_Small_middle_town = np.ar</pre>	n)) opulation)) pulation))		
<pre>#occupation #0= NO, 1= YES occupation_Employed = np.array([.50,.50]).cumsum().searchsorted(np.random.sample(Total_Population)) occupation_Not_active = np.array([.52,.48]).cumsum().searchsorted(np.random.sample(Total_Population occupation_Unemployed = np.array([.98,.02]).cumsum().searchsorted(np.random.sample(Total_Population))</pre>) n)) n))		
<pre>#GENDER #0= NO, 1= YES gender_Man = np.array([.51,.49]).cumsum().searchsorted(np.random.sample(Total_Population)) gender Woman = np.array([.49,.51]).cumsum().searchsorted(np.random.sample(Total Population))</pre>			
#marital_status			
<pre>#0= NO, 1= YES marital_status_Divorced_Separated = np.array([.94,.06]).cumsum().searchsorted(np.random.sample(Total_ marital_status_Married = np.array([.53,.47]).cumsum().searchsorted(np.random.sample(Total_Population marital_status_Partnership = np.array([1,0]).cumsum().searchsorted(np.random.sample(Total_Population marital_status_Single = np.array([.60,.40]).cumsum().searchsorted(np.random.sample(Total_Population marital_status_Widow = np.array([.93,.07]).cumsum().searchsorted(np.random.sample(Total_Population)</pre>	al_Populatio on)) on)) n))))	on))	
<pre>#generation of synthetic population features = [native, educational, age, household_composition, country_BALGARIJA, country_BELGIQUE, c df = pd.DataFrame(features) df = df.T</pre>	country_CESK	A_REPU	JBLIK.
<pre>population = df.as_matrix()</pre>			

Jupyter notebook used to generate synthetic populations

METHOD 2. Predicting social acceptance



- Extrapolation at the regional level was carried out by calculating the probability of acceptance in each individual using the 6 predictive models, and for each of the 5 periods
- The average probability was then estimated for all 271 regions, and an average variable was created to summarize the 6 predictive models
- In summary, a database was generated with 271 cases (NUTS2 in Europe) and 42 (7 x 6) estimated probabilities of acceptance of refugees



 In all the maps we clearly observe that the country is a strong predictor of social support and acceptance of refugees



Comparative visualization of the 6 predictive models for social acceptance of refugees in NUTS2 regions of Europe

RESULTS

- To summarize the models, the average of the 6 estimates for each region was calculated, as well as for all periods analysed
- According to estimates, the countries most likely to socially integrate refugees are Norway, Spain, United Kingdom or Germany, and those with lower probabilities are the Czech Republic, Hungary or Romania

These estimates fully coincide with the comparative European surveys



Summary display with the average acceptance probability for the 6 models and all periods analyzed

https://afly.co/ggw2

RESULTS

 If the estimate is focused on each country, more interesting internal differences are perceived



- For example, in the case of Italy, Greece and Spain it can be seen that the most populated regions (Athens, Thessaloniki, Rome, Milan, Madrid or Catalonia) are more likely to accept refugees
- We create an interactive online map (available for reproducibility) to explore each country, with each model at each time point



Estimation of the probability of acceptance of refugees in the regions of Spain, Italy and Greece

RESULTS

- Regarding the longitudinal changes, the statistical analysis revealed that 2016 was the worst year for social integration
- There is a significant decrease in the probability of accepting refugees by Europeans in May (M = 0.66, SD = 0.13) and November 2016 (M = 0.67, SD = 0.14), compared with November 2015 (M = 0.69, SD = 0.12) or May (M = 0.68, SD = 0.14) and November 2017 (M = 0.69, SD = 0.12)



REMARKS

 In this work, the first articulated data that estimate the possible social integration of refugees in all regions of Europe have been generated

> The use of synthetic populations is an innovation that provides validated computational methods from other fields to the social and communication sciences, where they have rarely been implemented

REMARKS

- This predictive analytical approach based on data can contribute to generating new relocation and resettlement schemes
- The prediction of the probability of acceptance of refugees at the regional level will allow mitigating the risks of exclusion of these groups, through the planning of integration policies, anticipating which populations could accept or oppose more resistance to the arrival of refugees



Machine learning approaches to monitor integration of migrants

Carlos Arcila Calderón

THANK YOU



VNIVERSIDAD DSALAMANCA