**ANNOUNCING NAMSORML DEEP LEARNING FOR #GEOINT :** DECRYPTING IDENTITY IN SPACE AND TIME THROUGH PERSONAL NAMES, GEOGRAPHIC, SEMANTIC, GRAPH DATA

2018-01

Elian CARSENAT, NamSor

# Founder Bio

**Elian CARSENAT**, a computer scientist trained at ENSIIE/INRIA, started his career at JP Morgan in Paris in 1997. He later worked as consultant and managed business & IT projects in London, Paris, Moscow and Shanghai.

In 2012, Elian created **NamSor**, a piece of sociolinguistics software to mine the 'Big Data' and better understand international flows of money, ideas and people.

http://fr.linkedin.com/in/eliancarsenat/en

# Two NamSor product lines

## NamSor CORE

- Optimized for global coverage : coding names to a large multi-class taxonomy (all countries / regions /ethnicities)

- The only input is NAMES : not other information is required

- SaaS API or on-premises software licensed per CPU

- Includes 4 proprietary onomastics models (.ONO)

## NamSor ML                    NEW!

- Deep-learning capability to re-train models towards a focused research or a customized taxonomy (binary classifier, or just a few classes)

- Name information is combined with other data (geographic, behavioural, semantic …)

- On-premises software licensed per CPU
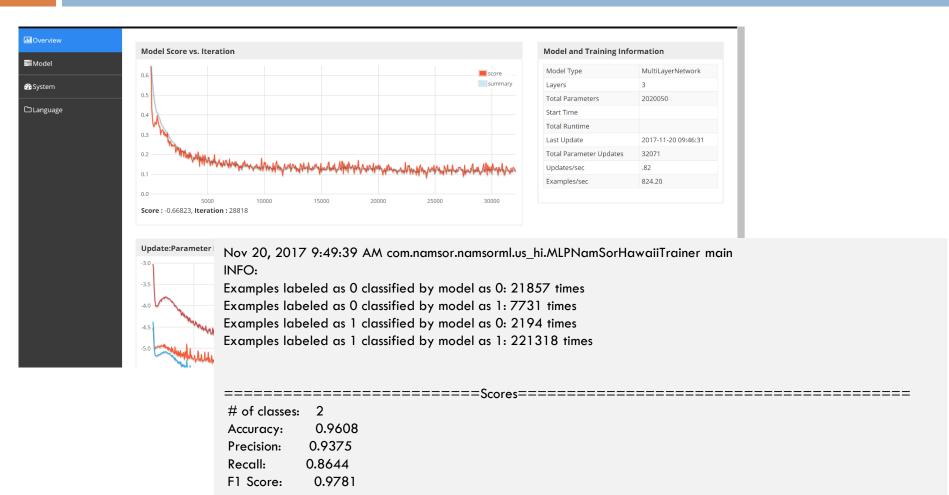
- Includes 4 word-embedding models (FastText or W2V)

# NamSor ML key functions
## custom classifiers for race / ethnicity / migration studies

- State-of-the-art **Deep Learning** (Neural Network)
- Same NamSor CORE models are available as **pre-trained models** for word embedding, both in Word2Vec (W2V) and FastText (BIN) formats:
  - COUNTRY_MELTINGPOTS_GENDER_SCRIPT
  - COUNTRY_REGION_SCRIPT
  - COUNTRY_SCRIPT
  - ETHNO_COUNTRY_SCRIPT
- Software is licensed as a SDK/Framework with samples for use and integration.
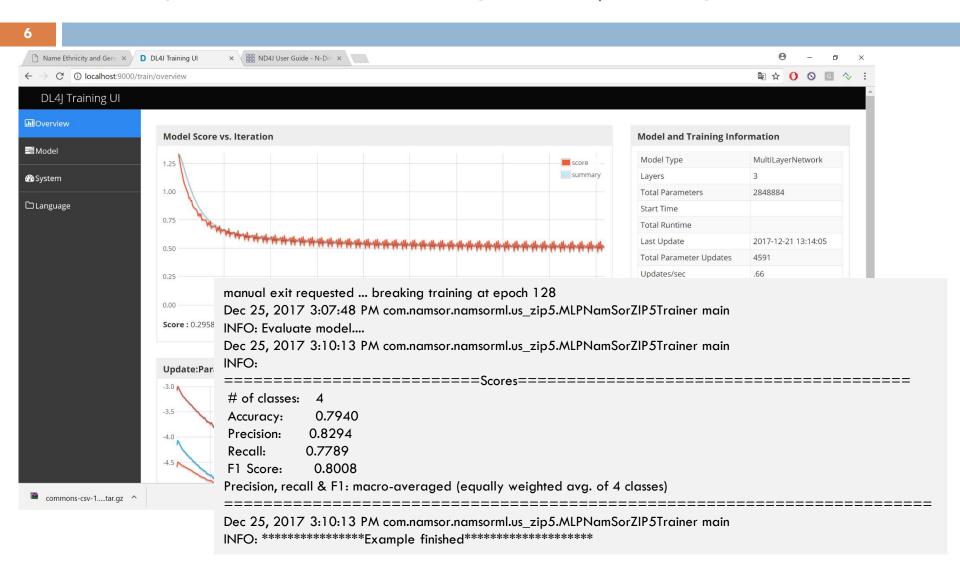
# NamSor ML :
# Native Hawaiian Names - A Binary Classifier Example



**Model Score vs. Iteration**

score
summary

Score : -0.66823, **Iteration** : 28818

**Model and Training Information**

| Model Type | MultiLayerNetwork |
|---|---|
| Layers | 3 |
| Total Parameters | 2020050 |
| Start Time | |
| Total Runtime | |
| Last Update | 2017-11-20 09:46:31 |
| Total Parameter Updates | 32071 |
| Updates/sec | .82 |
| Examples/sec | 824.20 |

**Update:Parameter**

Nov 20, 2017 9:49:39 AM com.namsor.namsorml.us_hi.MLPNamSorHawaiiTrainer main
INFO:
Examples labeled as 0 classified by model as 0: 21857 times
Examples labeled as 0 classified by model as 1: 7731 times
Examples labeled as 1 classified by model as 0: 2194 times
Examples labeled as 1 classified by model as 1: 221318 times

============================Scores============================
# of classes:     2
Accuracy:        0.9608
Precision:       0.9375
Recall:          0.8644
F1 Score:        0.9781
============================================================
Nov 20, 2017 9:49:39 AM com.namsor.namsorml.us_hi.MLPNamSorHawaiiTrainer main
INFO: ***************Example finished*******************

# NamSor ML :
# US – input Name+ZIP code; output : race/ethnicity

manual exit requested ... breaking training at epoch 128
Dec 25, 2017 3:07:48 PM com.namsor.namsorml.us_zip5.MLPNamSorZIP5Trainer main
INFO: Evaluate model....
Dec 25, 2017 3:10:13 PM com.namsor.namsorml.us_zip5.MLPNamSorZIP5Trainer main
INFO:
==========================Scores==========================================
 # of classes:    4
 Accuracy:        0.7940
 Precision:       0.8294
 Recall:          0.7789
 F1 Score:        0.8008
Precision, recall & F1: macro-averaged (equally weighted avg. of 4 classes)
==========================================================================
Dec 25, 2017 3:10:13 PM com.namsor.namsorml.us_zip5.MLPNamSorZIP5Trainer main
INFO: ***************Example finished*******************

- ABOUT NAMSOR

# NamSor sorts *Names*

☐ Names reflect cultural *Identity*

- ☐ **NamSor** data mining software recognizes the **linguistic or cultural origin of names** in any alphabet / language, with fine grain and high accuracy.

☐ Names are meaningful : we use sociolinguistics to extract their semantics and deliver actionable intelligence.

# NamSor is focused on classification

Data Mining/ Predictive analytics

Social Networks

Geo-demographics/ GEOINT

NamSor

1.Classification

2.Transliteration & Identification

3. Named Entity Extraction

Watch Lists/ Anti-Fraud/ Counter-Terrorism

Indexing & Text Mining

# NamSor sorts Names

- <u>Classification</u> with various taxonomies
    - Gender (female/ male / unknown)
    - Script (LATIN, ARABIC, GUJARATI,...)
    - Origine (Country ex. France vs. Inde)
    - Region (ex. Gujarat vs. Andhra Pradesh)
    - Diaspora (ex. Indian Diaspora in US vs Indian Diaspora in Mauricius)
- <u>Sorting</u> according to a numerical score, allowing combining NamSor with other algorithm (graph, semantics, predictive …)
- <u>Flexibility</u> to learn new taxonomies (machine learn.)
- <u>Ease of integration</u> (NamSor API, Java/Python SDK, ESRI, RapidMiner, NationBuilder …)

# A global coverage -142+ countries

| DIMENSION | CURRENT COVERAGE |
|---|---|
| SCRIPT (22) | LATIN, ARABIC, CYRILLIC, ARMENIAN, BENGALI, DEVANAGARI, GEORGIAN, GREEK, GUJARATI, GURMUKHI, HAN, HANGUL, HEBREW, HIRAGANA, KANNADA, KATAKANA, MALAYALAM, MYANMAR, ORIYA, TAMIL, TELUGU, THAI |
| COUNTRY (142+) | AE, AF, AL, AM, AO, AR, AT, AZ, BA, BD, BE, BF, BG, BH, BI, BJ, BN, BR, BT, BW, BY, CA, CD, CF, CG, CH, CI, CL, CM, CN, CO, CR, CV, CY, CZ, DE, DK, DZ, EE, EG, ER, ES, ET, FI, FJ, FR, GA, GB, GE, GH, GM, GN, GR, HK, HR, HT, HU, ID, IE, IL, IN, IQ, IR, IS, IT, JO, JP, KE, KG, KH, KM, KP, KR, KW, KZ, LA, LB, LK, LR, LS, LT, LU, LV, LY, MA, MD, ME, MG, MK, ML, MM, MN, MR, MU, MV, MW, MX, MY, MZ, NA, NE, NG, NL, NO, NP, OM, PE, PH, PK, PL, PS, PT, QA, RO, RS, RU, RW, SA, SD, SE, SI, SK, SN, SO, SR, SY, TD, TG, TH, TJ, TM, TN, TO, TR, TT, TW, TZ, UA, UG, US, UZ, VE, VN, YE, ZA, ZM, ZW |
| COUNTRY / REGION (15) | RU (80), IN (~30), FR (22), IT (17), LB (14), BF (13), CD (8), TR (7), ID (7), GB (4), ES (17), ML (50), GN (8), CI (34), AF(16) |
| COUNTRY / DIASPORA | US, CA, SG, GB, (EU) |

# Detailed functional coverage

A complex n-dimension matrix

☐ # According to script / country :

❑ Ex. recognizing an Armenian name in LATIN, ARABIC, CYRILLIC, or GEORGIAN ?

nikoloz doborjginidze
Sարիելիվաա Sարիելաշվիլի
ნანული ყაზარაშვილი
николоз доборджгинидзе
كريكور انتر انيك طاكسيان

|  | ARABIC | ARMENIAN | CYRILLIC | GEORGIAN | LATIN |
|---|---|---|---|---|---|
| Armenian | X | X | X | . | X |
| Georgian | . | X | X | X | X |

| | Maghrebi | | | | Egyptian | Gulf | | | | | | Hassaniya | Levantine | | | | Mesopotamian | Sudanese | Yemeni | | | Other | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | DZ | LY | MA | TN | EG | AE | BH | KW | OM | QA | SA | MR | JO | LB | PS | SY | IQ | SD | YE | SO | DJ | ID | IN | PK | MY | IR | AM | Other |
| LATIN | X | x | X | X | X | x | X | x | x | x | X | X | X | X | X | X | x | x | x | x | . | x | X | X | x | X | X | X |
| ARABIC | x | x | X | X | X | x | X | x | x | x | X | X | X | X | X | X | X | x | X | . | . | x | x | x | x | X | X | . |

☐ # According to host / origin country :

❑ Ex. recognizing an Georgian name in the US, in Russia, in Europe?

# NamSor CORE key functions
to quickly re-calibrate a model

- Machine learning
- Hierarchical taxonomies
  - Ex. Country / Region
- Automatic clustering
  - Name – Surname, Father/Mother Name
  - Real Name – Pseudonym
- Multiple iterations, using Score to sort data
+ a proven methodology to qualify data sources

# NamSor CORE Technical features

- Software as a Service (SaaS),
  - Online processing (Excel ou UTF-8 txt)
  - NamSor API (REST)
  - NamSor SDK (Java, Python, Scala)
  - NamSor add-on for RapidMiner, NationBuilder, …
  - 3 taxonomies : Gender, Origin, Diaspora
  - 99.9% availability
  - Throughput: ~1000 names per second (~100 millions names per day)

- On-site software deployment,
  - Hardware: 2x Servers (4c/8t, 3GHz+, 64GB, 2x2TB)
  - Linux/Java or Scala/Spark

# NamSor CORE SaaS architecture

# Confidentiality, data protection

☐ Software Licence + Confidentiality Agreement

☐ As Input : names only (no other personal data)

☐ Secure Data Center and 'Cloud', accessed via SSH

☐ NamSor API over HTTPS

Technical Details

**Connection Encrypted (TLS_ECDHE_RSA_WITH_AES_128_CBC_SHA, 128 bit keys, TLS 1.2)**

The page you are viewing was encrypted before being transmitted over the Internet.

Encryption makes it difficult for unauthorized people to view information traveling between computers. It is therefore unlikely that anyone read this page as it traveled across the network.

☐ Two logging options :

  ☐ In clear (for machine learning)

  ☐ As SHA-256 (for confidentiality) ex. Jean Durieux becomes ef61a579c907bbed674c0dbcbcf7f7af8f851538eef7b8e58c5bee0b8cfdac4a

- NamSor use cases

# Tunisian Diaspora

## La BIAT lance un road show pour les tunisiens en île de France

Le 17 mars 2016, BIAT France lance le « BIAT France Tour » et part à la rencontre des Tunisiens résidents à Paris et en région parisienne afin de leur présenter les produits et services qui leur sont destinés.

G f 181 ✉ + 1



C'est aux 720 000 Tunisiens vivant en France que la BIAT s'adresse avec sa filiale BIAT France, sous la signature « Ici pour vous », un service de transfert d'argent leur est proposé à des prix très compétitifs et dans des conditions de rapidité et de sécurité exemplaires.

# Mining 3M twitter names to map *Diasporas*
## *Who are they, where are they and what are they doing?*

Source: Twitter
Visualization : CartoDB
Data Mining: NamSor

# Flow view – who travels where?

| Source | Target | Type | Id | Onoma | Weight |
|---|---|---|---|---|---|
| **United Kingdom** | **France** | **Directed** | **16** | **Great Britain** | **37** |
| Spain | France | Directed | 55 | Spain | 14 |
| United States | France | Directed | 75 | Great Britain | 12 |
| Turkey | France | Directed | 79 | Turkey | 11 |
| Brazil | France | Directed | 87 | Portugal | 10 |
| United Kingdom | France | Directed | 112 | Ireland | 9 |
| Italy | France | Directed | 152 | Italy | 7 |
| Switzerland | France | Directed | 226 | France | 5 |
| Belgium | France | Directed | 247 | France | 5 |
| **United Kingdom** | **France** | **Directed** | **258** | **France** | **5** |
| Mexico | France | Directed | 287 | Spain | 4 |
| Ireland | France | Directed | 317 | Great Britain | 4 |
| United Kingdom | France | Directed | 333 | Italy | 4 |
| United States | France | Directed | 375 | France | 4 |



Source: Twitter
Visualization : Gephi
Data Mining: NamSor

# Mapping Talents in Cancer Research
## (in collaboration with French INSERM)

**Who's in Cancer Research - an onomastic mille-feuille**
abscysse: country of affiliation; ordinate: likely country of origin
Source: Thomson WoS

NamSor™ 2015

**Thomson Reuters WebOfScience (6 countries, 250k scientists, 50k papers)**

"Analysts uncovered amazing patterns in the way scientists' names correlate with whom they publish, and who they cite in their papers - not just in case of a particular country, but globally. Tania Vichnevskaia of the French National Institute for Health (INSERM) presented the paper 'Applying onomastics to scientometrics' at IREG International symposium 2015 organised by University of Maribor and Shanghai Jiao Tong University. The paper was prepared jointly with NamSor, a private start-up company specialized in mapping international Diasporas."

Source: WoS; Data Mining: INSERM with NamSor

# Cancer Research in Poland and Slovenia

## Examining the 'brain drain'

### The Polish "brain drain"



- USA
- GBritain
- Germany
- Sweden
- Canada
- Russia
- Blgium
- Austria
- Danemark

In the Polish Corpus, we look at co-authors with Polish names, affiliated abroad.

Top countries:
1. **US,**
2. **Great-Britain,**
3. **Germany.**

### The Slovenien "brain drain"



- USA
- GBritain
- Germany
- Bosnia and Herzegovina
- Danemark
- Italia

In the Slovenian Corpus, we look at co-authors with Slovenian names, affiliated abroad.

Top countries:
1. **Great-Britain,**
2. **US,**
3. **Germany.**

Source: WoS; Data Mining: INSERM with NamSor

# "Incredible India" – 1.2 BN People
## Indian onomastics by State/Union Territory

Names in LATIN, BENGALI, DEVANAGARI, GUJARATI, GURMUKHI, KANNADA, MALAYALAM, ORIYA, TAMIL, TELUGU, ARABIC

# GUJARAT : mapping onomastics by district



Source: Voters List; Visualization : Google Fusion Tables; Data Mining: NamSor

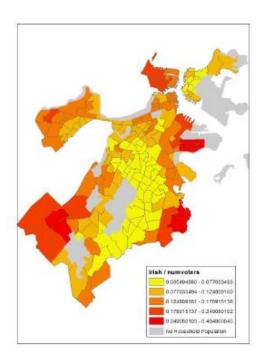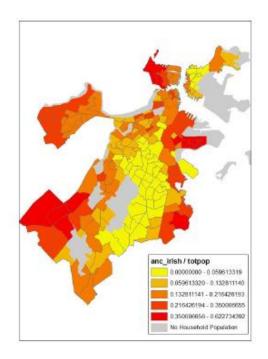# ASSAM: Karbi Anglong, within district
## Inter-caste marriages ?

| output | | Input | | | Input | |
|---|---|---|---|---|---|---|
| clusterId | clusterParentId | Firstname | LastName | parent is | FirstParen | LastParent |
| L25354:25: | L64958:2797 | Aạ\ ¹ | ¹}[ššã | husband | ¤àiꞈW | [W>ã |
| L47490:15! | L64958:2797 | ¤àꞘ¹ | [W>ã | father | ¤àiꞈW | [W>ã |
| L28582:12( | L47490:1593 | [³>à | Òꞇt̃šã | husband | ¤àꞘ¹ | [W>ã |
| L23643:66! | L35593:510 | ꞇƒ}à | [W>ãšã | father | ¤àiꞈW | [W>ã |
| L23643:66! | L35593:510 | ³à-àÒù | [W>ãšã | father | ¤àiꞈW | [W>ã |
| L47490:15! | L35593:510 | Vãꞇ=¢ | [W>ã | father | Wꞈ¢ | [W>ã |
| L23643:66! | L35593:510 | Aạ¹ | tã¹ïšã | husband | Wꞈ¢ | [W>ã |
| L35593:51( | L47490:1593 | [ƒºš | [W>ã | father | Vãiꞈ¢ | [W>ã |
| L23643:66! | L47490:1593 | [¹>à | [W>ãšã | father | Vãiꞈ¢ | [W>ã |


ASSAM: Karbi Anlong district names clustered

parent is    husband

| Count of serial Row Labels | Column Labels L47490:1593 | L116370:3612 | L54332:2031 | L184096:2297 | L35593:510 | L168871:1819 | L135664:4438 | L51271:837 |
|---|---|---|---|---|---|---|---|---|
| L23643:669 | 6931 | 84 | 5099 | 15 | 2069 | 28 | 791 | 1924 |
| L151415:3559 | 18 | 212 | 11 | 6446 | 19 | 1217 | 55 | 6 |
| L28582:1209 | 5132 | 68 | 3565 | 10 | 1494 | 17 | 592 | 1323 |
| L116370:3612 | 66 | 10283 | 38 | 72 | 40 | 321 | 137 | 29 |
| L9839:442 | 2491 | 60 | 1851 | 9 | 774 | 11 | 321 | 660 |
| L168871:1819 | 7 | 263 | 6 | 361 | 8 | 2730 | 24 | 4 |
| L23642:141 | 1198 | 8 | 822 | 2 | 375 | 4 | 156 | 332 |
| L25354:253 | 1181 | 12 | 932 | | 375 | 7 | 100 | 323 |
| L135664:4438 | 20 | 154 | 5 | 22 | 19 | 44 | 2212 | 3 |
| L87032:1210 | 11 | 315 | 13 | 51 | 14 | 141 | 37 | 9 |
| L90333:3644 | 3 | 204 | 2 | 31 | | 190 | 5 | |
| L184096:2297 | | 13 | | 1735 | 3 | 84 | 11 | 1 |
| L87031:697 | 4 | 136 | 4 | 12 | 3 | 137 | 4 | 5 |
| L14495:131 | 614 | 10 | 432 | | 167 | 4 | 68 | 163 |
| L63724:1422 | 17 | 83 | 10 | 34 | 34 | 28 | 96 | 6 |
| L98994:891 | 31 | 161 | 46 | 21 | 19 | 59 | 21 | 5 |

Source: Voters List; Data Mining: NamSor

# Boston geo-demographics 1/2

BOSTON
REDEVELOPMENT
AUTHORITY

Irish Share, namsor

Irish Share, 2010-2014 ACS

Source: Boston Voters List ;  Visualization : ESRI ; Data Mining: NamSor[6]
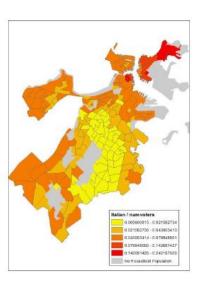
# Boston geo-demographics 2/2



Hispanic/Latino Share, namsor

Black/African-American Share, nams

Italian Share, namsor

March 7, 2016
Presentation Title

March 7, 2016
Presentation Title

March 7, 2016
Presentation Title

Source: Boston Voters List
Visualization : ESRI
Data Mining: NamSor

# Breaking down 'White' and 'Asian' into Portuguese, Spanish, Italian, India, Pakistan, China, ...



Source: Boston Voters List
Visualization : ESRI
Data Mining: NamSor

# PATENT DATABASES

**ITIF | INFORMATION TECHNOLOGY & INNOVATION FOUNDATION @10**

Get our newsletter

Search

| ISSUES | PUBLICATIONS | EVENTS | NEWS ROOM | MULTIMEDIA | ABOUT |

## The Demographics of Innovation in the United States

Adams Nager, David M. Hart, Stephen Ezell, and Robert D. Atkinson    February 24, 2016

A groundbreaking ITIF survey shows why the country needs to broaden and deepen its pool of potential innovators with better STEM immigration and education policies.

View Report    View Executive Summary    Event

Groundbreaking @ITIFdc survey shows why US needs to broaden and deepen pool of potential innovators

.@ITIFdc releases groundbreaking survey on who innovates in the United States and where and how it occurs

# US AID PROJECT

| | | Explanations |
|---|---|---|
| seedname | https://ca.linkedin.com/in/yuriy-diakunchak-6a3bb6[...]_-1035308383.linkedin | |
| datetime | 2016-02-07T12:49+0000 | |
| dateextract | 07/02/2016 | |
| linkedinurl | https://ca.linkedin.com/in/yuriy-diakunchak-6a3bb61 | LinkedIn URL and Profile data |
| parentseed | - | |
| fullname | Yuriy Diakunchak | - |
| titlename | Director of Marketing at Ukrainian Credit Union Limited | - |
| orgname | Ukrainian Credit Union Limited | - |
| locality | Toronto, Canada Area | - |
| languages | English,Ukrainian, | - |
| connections | 500+ connections | - |
| descriptor | Toronto, Canada Area | - |
| profilesummary | Summary Experienced marketing manager. Focused on branding and brand management, executing programs and campaigns targeted at building br | - |
| experience | Experience Director of Marketing Ukrainian Credit Union Limited August 2010 – Present (5 years 7 months) Heading up the marketing department. R | - |
| organizations | Organizations KUMF Art Gallery Director Starting March 2013 Member of Board of Directors of KUMF Art Gallery (Ukrainian Canadian Art Foundation) | - |
| volunteering | Volunteer Experience & Causes Causes Yuriy cares about: Arts and Culture Civil Rights and Social Action Politics | - |
| projects | | |
| education | Education M.M. Robinson McMaster University BComm, Marketing, Finance Pragmatic Marketing , Product Marketing Course Ryerson University Bac | - |
| awards | | |
| skills | Skills Product Management Marketing Automation Brand Management Product Marketing Demand Generation Multi-channel Marketing Marketing S | - |
| recommandations | Recommendations A preview of what LinkedIn members have to say about Yuriy: I had the pleasure of working with Yuriy since 2006 on various vide | - |
| groups | Groups Markian Silecky Real Estate - Babiak Team Royal LePage Real Estate Services Ltd, Brokerage Cloud Computing Taxonomy & Information Mana | - |
| publications | | - |
| score | 2.368793324 | < how much Ukrainian the nams sounds |
| count | 7 | < how many Ukrainian connections (max :10) |
| titlename2 | 0 | |
| orgname3 | 0 | |
| locality2 | 0 | |
| languages3 | 1 | < Ukrainian appears in Languages |
| connections4 | 0 | |
| descriptor5 | 0 | |
| profilesummary6 | 0 | |
| experience7 | 1 | |
| organizations8 | 0 | |
| volunteering9 | 0 | |
| projects10 | 1 | |
| education11 | 0 | |
| awards12 | 0 | |
| skills13 | 0 | |
| recommandations14 | 0 | |
| groups15 | 1 | < Ukraine appears in Goups |
| publications16 | 0 | |
| TOT | 13.98341639 | < Overall score taking into account factors from L2 |

# RUSSIAN MEDIA ANALYSIS



EVOLUTION AND COMPARISON OF THE EXPOSURE OF ARAB PERSONALITIES IN RUSSIAN MEDIA
(Jan-2004 to Mar-2013, Source: Integrum WW/RussoScopie)

Source: Integrum/@RussoScopie

# Indian diaspora names - a global airline use case

*'For 93% of our customers, when NamSor recognizes an Indian name, the client has travelled to India in the past.'*

At state level : ~50%

### Analysis of NamSor's **First** Choice Country Compared to Historic Travel on ~~Emirates~~

| NamSor's First Choice Country | Count of Unique Individuals who Have Travelled to NamSor's First Choice Country | | | % of Unique Individuals who Have Travelled to NamSor's First Choice Country | | |
|---|---|---|---|---|---|---|
| | No | Yes | Grand Total | No | Yes | Grand Total |
| India | 1,633 | 20,315 | 21,948 | 7% | 93% | 100% |
| Italy | 281 | 869 | 1,150 | 24% | 76% | 100% |
| Bangladesh | 524 | 1,456 | 1,980 | 26% | 74% | 100% |
| Ethiopia | 3 | 8 | 11 | 27% | 73% | 100% |
| Iran | 701 | 1,657 | 2,358 | 30% | 70% | 100% |
| Saudi Arabia | 679 | 771 | 1,450 | 47% | 53% | 100% |
| Afghanistan | 21 | 23 | 44 | 48% | 52% | 100% |
| Pakistan | 2,171 | 2,309 | 4,480 | 48% | 52% | 100% |
| Jordan | 148 | 124 | 272 | 54% | 46% | 100% |
| Kuwait | 51 | 37 | 88 | 58% | 42% | 100% |
| Qatar | 3 | 2 | 5 | 60% | 40% | 100% |

### Analysis of NamSor's **Region Rounded** Score Compared to Historic Travel on ~~Emirates~~ for India

| Customer has Flown to NamSor's Region | Count of Unique Individuals who Have Travelled to NamSor's Region | | | % of Unique Individuals who Have Travelled to NamSor's Region | | |
|---|---|---|---|---|---|---|
| | No | Yes | Grand Total | No | Yes | Grand Total |
| 5 | | 3 | 3 | 0% | 100% | 100% |
| 4 | 569 | 696 | 1,265 | 45% | 55% | 100% |
| 3 | 2,202 | 2,774 | 4,976 | 44% | 56% | 100% |
| 2 | 2,861 | 3,226 | 6,087 | 47% | 53% | 100% |
| 1 | 2,442 | 2,523 | 4,965 | 49% | 51% | 100% |
| 0 | 1,686 | 1,423 | 3,109 | 54% | 46% | 100% |
| -1 | 702 | 519 | 1,221 | 57% | 43% | 100% |
| -2 | 180 | 109 | 289 | 62% | 38% | 100% |
| -3 | 23 | 9 | 32 | 72% | 28% | 100% |
| -4 | 1 | | 1 | 100% | 0% | 100% |
| Grand Total | 10,666 | 11,282 | 21,948 | 49% | 51% | 100% |

# Indian diaspora in Mauritius

Registered Medical Practitioners

The annual list for the reg
category;

Register of General Practitio

**Onomastics of ~2200 GPs in Mauritius**
Using NamSor v0.0.27

544
872
35
109
148
491

- MU
- IN
- PK
- FR
- HK
- Other

c2015 namsor.com
Source: medicalcouncilmu.org
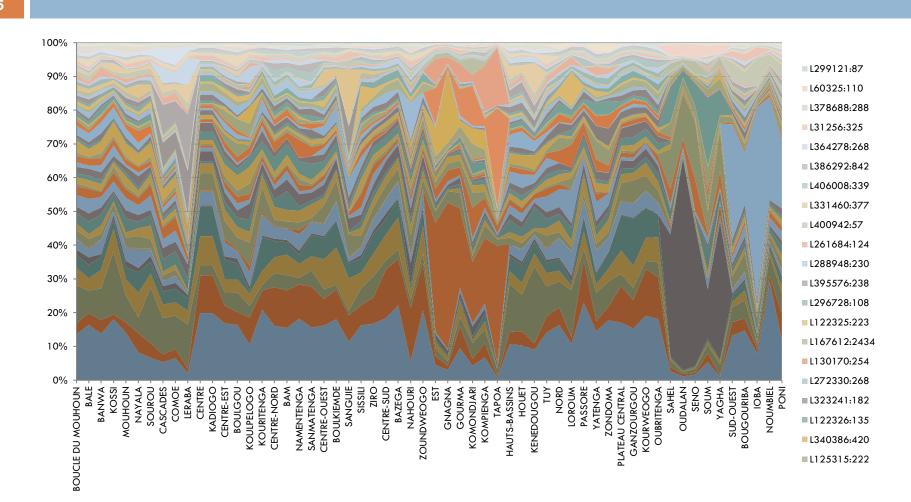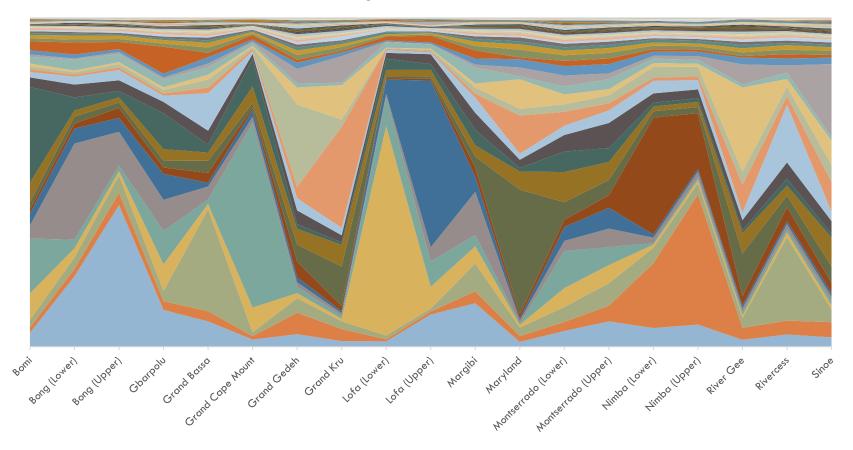
# Africa: complex identities (Congo RDC)

# Africa: complex identities (Burkina Faso)

# Africa: complex identities (Liberia)

Liberia - a regional onomastics 'mille-feuille'

# Using NamSor API

Option 1/
Online

Option 2/
RapidMiner Extension

Option 3/
NamSor API & SDK

# Thank you !

Elian CARSENAT,

elian.carsenat@namsor.com

Phone : +33 6 52 77 99 07