



GREThA

Groupe de Recherche en
Économie Théorique et Appliquée

How To Kill Inventors: Testing The Massacrator[©] Algorithm For Inventor Disambiguation

Michele PEZZONI

*University of Milano-Bicocca
KiTES-Università Bocconi
Observatoire des Sciences et des Techniques
&*

Francesco LISSONI

*KiTES-Università Bocconi
GREThA, CNRS, UMR 5113
Université de Bordeaux
&*

Gianluca TARASCONI

KiTES-Università Bocconi

Cahiers du GREThA

n° 2012-29

December

GREThA UMR CNRS 5113

Université Montesquieu Bordeaux IV
Avenue Léon Duguit - 33608 PESSAC - FRANCE
Tel : +33 (0)5.56.84.25.75 - Fax : +33 (0)5.56.84.86.47 - www.gretha.fr

Comment tuer les inventeurs: une évaluation de l'Algorithme Massacrator© pour désambiguïser les inventeurs

Résumé

La désambiguïstation de noms des inventeurs est un problème de plus en plus important pour les utilisateurs de données de brevets. Nous proposons et testons un certain nombre d'améliorations à l'algorithme Massacrator ©, proposé initialement par Lissoni et al. (2006) et maintenant appliqué à APE-INV, une base de données en accès libre soutenue par l'European Science Foundation. D'après Raffo et Lhuillery (2009), nous décrivons la désambiguïstation comme un processus en 3 étapes: nettoyage et analyse, sélection et filtrage. Par le biais d'une analyse de sensibilité, basée sur des simulations MonteCarlo, nous montrons comment divers critères de filtrage peuvent être manipulés afin d'obtenir des combinaisons optimales de précision et de recall (type I et type II des erreurs). Nous montrons aussi comment ces combinaisons différentes produisent des résultats différents, plus ou moins fiables en fonction des applications prévues (études sur la productivité, la mobilité ou les réseaux des inventeurs). Les critères de filtrage basés sur les informations sur les adresses des inventeurs sont sensibles à la qualité des données, alors que celles fondées sur l'information sur les réseaux de co-inventeurs sont toujours efficaces. Des détails sur l'accès aux données et sur la collecte des retours d'information par les utilisateurs (ayant pour but l'amélioration de la qualité des données) sont également discutés.

Mots-clés : données de brevets, inventeurs, désambiguïstation de noms

How To Kill Inventors: Testing The Massacrator© Algorithm For Inventor Disambiguation

Abstract

Inventor disambiguation is an increasingly important issue for users of patent data. We propose and test a number of refinements to the Massacrator© algorithm, originally proposed by Lissoni et al. (2006) and now applied to APE-INV, a free access database funded by the European Science Foundation. Following Raffo and Lhuillery (2009) we describe disambiguation as a 3-step process: cleaning&parsing, matching, and filtering. By means of sensitivity analysis, based on MonteCarlo simulations, we show how various filtering criteria can be manipulated in order to obtain optimal combinations of precision and recall (type I and type II errors). We also show how these different combinations generate different results for applications to studies on inventors' productivity, mobility, and networking. The filtering criteria based upon information on inventors' addresses are sensitive to data quality, while those based upon information on co-inventorship networks are always effective. Details on data access and data quality improvement via feedback collection are also discussed.

Keywords: patent data, inventors, name disambiguation

JEL: C15, C81, O34

<p>Reference to this paper: PEZZONI Michele, LISSONI Francesco, TARASCONI Gianluca (2012) How To Kill Inventors: Testing The Massacrator© Algorithm For Inventor Disambiguation, <i>Cahiers du GREThA</i>, n°2012-29. http://ideas.repec.org/p/grt/wpegrt/2012-29.html.</p>
--

1. Introduction¹

Economic studies of innovation have for long made use of patent data (Griliches,1990; Nagahoka et al., 2010). Assisted by digitalization of records and increasing computational power, economists and other social scientists have extracted increasing quantities of information from patent documents, such as the applicants' identity and location, the technological contents of the invention, or the latter's impact, as measured by citations. More recently, information on inventors has attracted a good deal of attention. Identifying inventors allows studying their mobility patterns, both in space and across companies, (Agrawal et al., 2006; Marx et al., 2009) as well as their social capital, as measured within co-inventor networks (Fleming, 2007; Breschi and Lissoni, 2009; Lissoni, 2010). Finally, it makes it possible to look for additional information at the individual level, ranging from professional identities (does the inventor appear also on a list of R&D employees? or a list of academic scientists working for a university or a public lab?) to other type archival data on knowledge-related activities (such as scientific publications; see Azoulay et al., 2009; Breschi et al., 2007; Lissoni et al., 2008).

Identifying inventors within any given set of patent data, as well as matching them to any other list of individuals, requires the elaboration of complex "disambiguation" algorithms. They are necessary to analyse in a non-trivial way the text strings containing the inventors' names, surnames, and addresses. Yet, it is only of late that users of inventor data have started discussing openly about the disambiguation techniques they employ, and examine their implications in terms of data quality and reliability of the evidence produced (Raffo and Lhuillery, 2009; Lai et al., 2011).

This paper deals with Massacrator[©] 2.0, the disambiguation algorithm we elaborated to create the APE-INV inventor database, an open-access initiative funded by Research Networking Programme of the European Science Foundation (<http://www.esf-ape-inv.eu>). The APE-INV inventor database has been conceived as a subset of the PatStat-Kites database (<http://db.kites.unibocconi.it/>), which contains all patent applications filed at EPO, as derived from the October 2011 release of the Worldwide Patent Statistical Information Database (better known as PatStat²). As such, it can be more generally described as a PatStat-compatible dataset, which addresses the needs of the increasingly large community of PatStat users.

Massacrator[©] 2.0 is a revised form of the original Massacrator[©] algorithm, which was originally conceived for the *ad hoc* purpose of identifying inventors in selected countries, and with the intent of maximizing precision (that is, minimizing type I errors, or false positives; Lissoni et al.,2006). Our revision has transformed it into a more general tool, one that users can calibrate also to maximize recall (minimize type II errors, or false negatives) or to achieve the best possible combination of recall and precision (that is, to strike a balance between different types of errors).

¹ Acknowledgements:

Financial support from the Research Networking Programme of the European Science Foundation is acknowledged (APE-INV – Academic Patenting in Europe Project). Early drafts of the paper benefitted of comments from participants to the APE-INV NameGame workshops. We are also grateful to Nicolas Carayol, Lorenzo Cassi, Stephan Lhuillery and Julio Raffo for providing us with core data for the two benchmark datasets. Monica Coffano and Ernest Miguelez provided extremely valuable research assistantship. Andrea Maurino's expertise on data quality has been extremely helpful.

² Access information for PatStat at: <http://forums.epo.org/epo-worldwide-patent-statistical-database/> - last visited: 12/13/2012

In what follows, we first describe the general workflow (cleaning & parsing → matching → filtering) of the Massacrator[©] 2.0 algorithm (section 2). Then, in section 3, we present our calibration methodology for the filtering stage, which crucially affects the algorithm's performance. In section 4 we perform a validation exercise, based on two "benchmark" datasets. In the same section, we move on to apply the validated algorithm to the entire PatStat data, in order to generate the APE-INV inventor database. Section 5 concludes.

2. An Overview of Massacrator[©] 2.0

Disambiguation of inventors consists in assigning a unique code to several inventors listed on different patents who are homonyms or quasi-homonyms, and share of a set of similar characteristics (e.g. they have the same addresses or patents with the same technological content). Inventors with same code are then treated as one individual. Following Raffo and Luhillery (2009), we describe disambiguation as three-step process:

1. *Cleaning & Parsing*: the relevant text strings (chiefly, those containing information on name, surname and address of the inventor) are purged of typographical errors, while all characters are converted to a standard character set. If necessary, any relevant string is parsed into a several substrings, according to various criteria (punctuation, blank spaces, etc.). Typically, the string containing the inventors' complete name (e.g. Duck, Prof. Donald) is parsed into name, surname and title (if any). The address is parsed, too.
2. *Matching*: the algorithm selects pairs of inventors, from different patents, who are likely candidates to be the same person, due to homonymy or similarity of names.
3. *Filtering*: the selected pairs are filtered according to additional information retrieved either from the patent documentation or external sources. Typical information from within the patent documentation are the address (e.g. quasi-homonyms sharing the same address are believed to be the same person) or some characteristics of the patent, such as the applicant's name (e.g. homonyms whose patents are owned by the same company may be presumed to be the same person) or its technological contents (as derived from the patent classification system or patent citations).

Massacrator 2.0 deals with 2,806,516 inventors listed on the EPO patent applications contained in the October 2011 version of PatStat, and it implements the three steps as follows:

2.1 Cleaning & Parsing:

C&P step 1: characters from an *ad hoc* list are removed, as well as punctuation and double blanks. All remaining characters are converted into plain ASCII³. As a result, a new field is created ("Inventor's name"), which contains the inventor's surname (possibly composed of several words, as it happens, for example, with Spanish surnames) and all of his/her names (including second, third or fourth names, and suffixes, such as "junior", "senior", "III" etc.). Similar steps are followed to create the following fields: "Inventor's address" (street's name and the number), "Inventor's city", "Inventor's county", "Inventor's region", and "Inventor's state" (to be intended as sub-national units, as in federal nations such as the US or Germany).

³ See: <http://rawpatentdata.blogspot.com/2010/05/converting-patstat-text-fields-into.html>

"Inventor's country" is derived directly from PatStat (ISO_3166-2 country codes).

C&P step 2: The original "Inventor's name" string from PatStat is parsed in as many substrings as the number of blanks it contains plus one. In the remainder of the paper we will refer to these substrings as "tokens". Due to EPO's conventions in reporting surnames and names, we can safely assume that the first token always contains the inventor's surname (or part of it, in case of double or triple surnames), while the last one always contains the given name (or part of it, in case of multiple names). Most cases are easy to manage since they are written in the form "surname, name" so using comma as separator we can easily parse different components.

Substrings whose contents matches a list of surname prefixes (such as "Van" or "De", respectively typical of Dutch and French/Italian surnames) are re-joined to the Surname string. Substrings whose contents matches a list of personal titles (such as "Professor" or "Prof.") are stored in a field different from the name (intitle).

2.2 Matching methodology

Massacrator 2.0 matches not only inventors with identical names, but also inventors with similar names, such as those hiding minor misspellings (ex.: "Duck, Donald" and "Duck, Donnald") as well as those resulting from the omission or inversion of words within the name or surname (ex.: "Duck, Donald D." and "Duck, Donald" or "Duck, D. Donald"), for a total of about 10 millions matches. In order to do so, it mixes the Token approach just described with an edit distance approach, in particular one based upon the 2-gram distances.

In detail, the algorithm sorts alphabetically all the tokens extracted from the original PatStat inventor's name text strings, without distinguishing between surnames and names (for a total of 444,215 tokens; tokens of 2 letters or less are discarded). It then computes the 2-gram (2G) distance between consecutive tokens (e.g. tokens appearing in row n and $n+1$ in the sorted list). The 2G can be described as the vector distance between two strings of different lengths, normalized by the total length of the strings. In our case it will be:

$$2G(t1, t2) = \frac{\sqrt{\sum_{i=aa}^{zz(N)} (G1_i - G2_i)^2}}{num(t1) + num(t2)}$$

Equation 1

where:

- $G1_i$ and $G2_i$ are the number of occurrences of the i -th bigram appears in tokens $t1$ and $t2$, respectively;
- $num(t1)$ and $num(t2)$ are the number of characters in tokens $t1$ and $t2$, respectively;
- N is the number of possible combinations of two consecutive letters (bigrams) in the alphabet of choice (in our case, plain ASCII, from which $N = 650$)⁴.

⁴ As an example, consider token "ABCABC" as $t1$ and token "ABCD" as $t2$. The bigram sets for $t1$ and $t2$ will be respectively: (AB,BC,CA,AB,BC) and (AB,BC,CD). Applying Equation 1 returns:

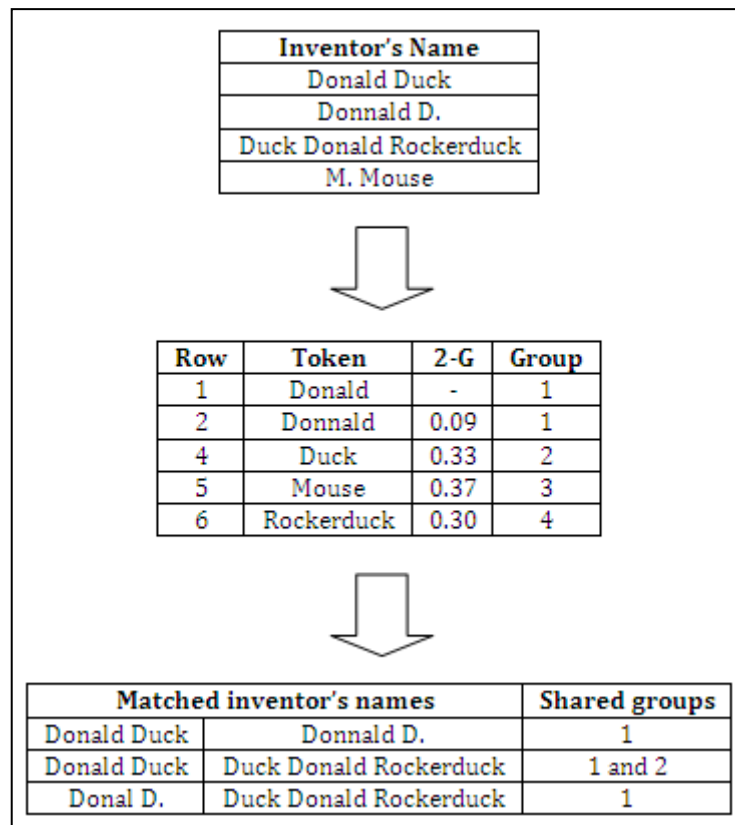
$$2G(t1, t2) = \frac{\sqrt{(2-1)_{AB}^2 + (2-1)_{BC}^2 + (1-0)_{CA}^2 + (1-0)_{CD}^2}}{5+3}$$

Once all $2G(t_1, t_2)$ distances are computed, consecutive tokens can be assigned to "groups", on the basis of their reciprocal distance, as follows:

- starting from the top of the token list, token in row 1 is assigned to group 1;
- then token in row 2 is also assigned to group 1 if its $2G$ distance from token in row 1 is less than or equal to an arbitrary threshold value δ (in the case of Massacrator 2.0: $\delta=0.1$) ; otherwise the algorithm creates a new group (group 2);
- The algorithm then proceeds in a similar fashion for all rows n and $n+1$

Once all groups are defined, the algorithm substitutes to each token the number of its corresponding group. As a result, each "Inventor's name" string is now replaced by a vector of numbers, each of which corresponds to a group of tokens. Any pair of inventors whose "Inventor's name" string contains identical group numbers (no matter in which order) are then treated as a match. In case the "Inventor's name" string are composed by a different number of tokens, the minimum common number of tokens (groups) is considered (see Figure 1 for a practical example). All matches obtained in this way are then passed on to the filtering stage.

Figure 1 - Example of Massacrator “mixed” matching rule



2.3 Filtering

For each pair of inventors in a match, Massacrator calculates a "similarity score", based upon a large set of weighted criteria. By comparing this score to a threshold value (*Threshold*), Massacrator then decides which matches to retain as valid (positive matches), and which to discard (negative matches). The criteria considered are 17,

grouped in 6 families: *network, geographical, applicant, technology, citations, and others*. A number of these criteria are derived from the original Massacrator[©] and they are quite intuitive, so we do not discuss them (see Table 1 for a short description). We discuss instead the *Approximated structural equivalence [ASE]* criterion, which is not present in the original Massacrator[©] and is rather complex.

The concept of Structural Equivalence was first introduced to social network analysis by Burt (1987). ASE adapts it to networks of patent citation and it was first proposed by Huang et al. (2011) as a method for inventor disambiguation⁵. The basic intuition is that the higher the number of citations two patents have in common, the higher the probability that any two inventors of such patents are the same person. Consider inventors I 's and J 's sets of patents:

$$P_i = \{p_{i,1}, p_{i,2}, \dots, p_{i,N_i}\} \quad P_j = \{p_{j,1}, p_{j,2}, \dots, p_{j,N_j}\}$$

Consider also P_{cit} as the set of all patents in our dataset receiving at least one citation:

$$P_{cit} = \{p_1, p_2, \dots, p_N\}$$

P_i (P_j) and P_{cit} are then used to compute matrix D_i (D_j), which has as lines patents in P_i (P_j) and as columns the cited patent in P_{cit} . If a patent in P_i (P_j) cites a patent in P_{cit} , the corresponding element in matrix D_i (D_j) takes value 1 ($D_z[p_z, p_{cit}] = 1; z = i, j$); if no citation occurs, it takes value zero ($D_z[p_z, p_{cit}] = 0; z = i, j$). Intuitively, inventors with similar matrixes are more likely to be the same person (their patents cite the same patents).

Massacrator then calculates weights W_{Citing} and W_{Cited} . The former is the inverted number of citations p_i (p_j), that is the inverse of the number of citations received by patent p_{cit} element. These weights allow to give less importance to matrix elements $D_z[p_z, p_{cit}] = 1$ ($z = i, j$) corresponding to "popular" patents (that is, patents sending out and/or receiving many citations).

⁵ Huang et al.'s original formula was proposed to compare inventors with no more than one patent each. We have adapted it to the case of inventors with multiple patents.

Table 1 Description of the criteria and classification in 5 families of filtering criterion

N.	Name of the criterion [names of variables in squared brackets]	Description
	<i>Network</i>	This family of criteria bases on the intuition that two matched inventors who turn out to be socially close are more likely to be the same person. As most patents are invented by two or more inventors, we consider each patent as a social tie between the listed co-inventors (Breschi and Lissoni, 2004).
1	Common coinventor [Coinventor]	Any two inventors <i>I</i> and <i>J</i> who have both signed patents with inventor are defined as having a common coinventor.
2	3 degrees of separation [Three degrees]	Any pair of inventors <i>I</i> and <i>J</i> , are said to stand at three degrees of separation when at least one of <i>I</i> 's coinventor and one of <i>J</i> 's coinventor have collaborated on the same patent.
	<i>Geographical</i>	This family exploits the inventor's address information.
3	City [City]	Two inventors share the same city within the address field (eg. Paris, Rome, Dijon)
4	Province [Province]	Two inventors share the same province within the address field (eg. Cote-d'Or)
5	Region [Region]	Two inventors share the same region within the address field (eg. Bourgogne)
6	State [State]	Two inventors have in common the same state within the address field (eg Texas).
7	Street [Street]	Two inventors share the same street and number within the address field. (eg. Boulevard Pasteur 32)
	<i>Applicant related variables</i>	This family exploits the characteristics of the patent applicant.
8	Applicant [Applicant]	Two inventors have signed at least one patent each for the same applicant.
9	Small Applicant [Small applicant]	As with Applicant, when the applicant has less than 50 inventors affiliated. If this criterion is satisfied also Applicant is satisfied.
10	Group [Group]	two inventors have signed at least one patent each for two distinct applicants belonging to the same group
	<i>Technology classes</i>	This family of criteria bases on the IPC code that identifies the technology class of a patent. The more digits two codes defining the IPC class have in common, the less the technological distance between the patents. The three criteria in this family are strictly related. In the case inventors <i>I</i> and <i>J</i> share at least one patent each with 12 digits in common, the other two criteria will be satisfied by definition (they have also 6 and 4 digits in common)
11	IPC 12 [IPC 12]	Within the stock of patents attributed to inventor <i>I</i> there is at least one patent with 12 digits of IPC code in common with another patent belonging to the inventor's <i>J</i> stock
12	IPC 6 [IPC 6]	Within the stock of patents attributed to inventor <i>I</i> there is at least one patent with 6 digits of IPC code in common with another patent belonging to the inventor's <i>J</i> stock
13	IPC 4 [IPC 4]	Within the stock of patents attributed to inventor <i>I</i> there is at least one patent with 4 digits of IPC code in common with another patent belonging to the inventor's <i>J</i> stock
	<i>Citation</i>	This family exploits citation links between patents.
14	Citations [Citation]	When a patent belonging to the stock of patents of inventor <i>I</i> is cited by a patent belonging to the stock of patents of inventor <i>J</i> , or vice versa, the pair of inventors has in common one citation.
15	Approximated structural equivalence [ASE]	Discussed in detail by the end of section 3.3
	<i>Others</i>	This family includes two criteria that cannot be classified in all the other four families
16	Rare surname [Rare surname]	At least one among the matched inventors' surnames is uncommon within the inventor's country. We identify rare surnames according to the frequency (by country) of first token (which we know to contain surnames) from the "inventor's name" PatStat field.
17	Priority date differs for less than 3 years [Three years]	A patent's priority dates is the earliest date of application in the patent's family ⁶ . For each pair of inventors we first calculate the minimum temporal distance (that is, the distance in time between the most recent among inventor <i>I</i> 's patents and the least recent among <i>J</i> 's, or vice versa). The distribution of minimum distances is very skewed, we set a threshold value of 3 years as a filtering criterion (<i>I</i> and <i>J</i> are more likely to be the same person if temporal distance is less than three years).

⁶ For a definition of patent family, see Martinez (2011)

Finally the algorithm divides the resulting index by the sum of $num(P_i)$ and $num(P_j)$ in order to normalize for the total number of patents of each inventor in the $[I, J]$ pair (Equation 2)

$$ASE[I, J] = \frac{\sum_{p_{i,1}=p_i}^{p_{i,N_i}} \sum_{p_{j,1}=p_j}^{p_{j,N_j}} \sum_{p_{cit}=p_1}^{p_N} D[p_i, p_{cit}] * W_{Citing_{p_i}} * D[p_i, p_{cit}] * W_{Citing_{p_j}} * W_{Cited_{p_{cit}}}}{num(P_i) + num(P_j)}$$

Equation 2

The higher the index, the closer inventors I and J are to the “perfect” structural equivalence (same position in the network of citations).

Massacrator find only 291469 non-null $ASE[I, J]$ scores, out of the >10 million matches analysed. The ASE filtering criterion is then considered satisfied by all these matches, no matter the score's exact value.

All the filtering criteria reported in Table 1 are used to compute a similarity score of the matched inventors as follows:

$$\alpha_m = \sum_{i=1}^{17} x_{i,m}$$

where $x_{i,m}$ is a dummy variable that equals 1 if match m meets criterion i , 0 otherwise. The number of retained (positive) matches depends upon the value assigned to the threshold variable (*Threshold*); when the similarity score α_m is larger than *Threshold* inventors in match m are considered to be the same person. This is the most delicate aspect of the algorithm implementation because values assigned arbitrarily can affect strongly the algorithm's performance. For this reason, Massacrator 2.0 relies on a calibration methodology, based upon a MonteCarlo simulation exercise, to which we now move on.

3. Filtering Calibration

The final output of the filtering phase must consist in a list of inventor pairs:

$$[m, I, J, D_{\alpha_m}] \quad I \neq J$$

where I and J are the two inventors forming pair m . D_{α_m} is a binary variable that takes value 1 if the two inventors in pair m are believed to be the same person (positive match) and 0 otherwise (negative match), based on their similarity score α_m and the chosen *Threshold* value. Notice that the output varies according to the number of filtering criteria we decide to use, and the *Threshold* value we choose. Calibration serves the purpose of guiding our selection of filtering criteria and *Threshold* value, on the basis of the efficiency of the resulting output.

We measure efficiency in terms of precision and recall:

$$Precision = \frac{tp}{tp + fp}$$

$$Recall = \frac{tp}{tp + fn}$$

where

tp = number of true positives
tn = number of true negatives
fp = number of false positives
fn = number of false negatives

We establish whether a positive (negative) match is true (false) by comparing the algorithm's results to information contained in two benchmark databases, namely the "Noise Added French Academic" (NAFA) and the "Noise Added EPFL" (NAE). Each benchmark database consists of a certain number of inventors, all matched one with the other, plus hand-checked information on whether the match is negative or positive. NAFA contains information on 530 inventors from France, 424 of which result affiliated to a university, the others being homonyms added *ad hoc* for testing purposes (that is, they represent added false positives or "noise"). NAE contains information on 342 inventors, 312 of which are faculty members at EPFL (the Federal Polytechnic of Lausanne, Switzerland), the others being added noise.⁷

For any match in the benchmark datasets we define D_{γ_m} analogously to D_{α_m} . It follows:

$$\begin{aligned} D_{\alpha_m} = 1 \cup D_{\gamma_m} = 1 &\Rightarrow \text{true positive} \\ D_{\alpha_m} = 0 \cup D_{\gamma_m} = 0 &\Rightarrow \text{true negative} \\ D_{\alpha_m} = 1 \cup D_{\gamma_m} = 0 &\Rightarrow \text{false positive} \\ D_{\alpha_m} = 0 \cup D_{\gamma_m} = 1 &\Rightarrow \text{false negative} \end{aligned}$$

(Equation 3)

We expect to observe a trade-off between precision and recall; any identification algorithm can decrease the number of false positives only by increasing the number of false negatives and *vice versa*. The smaller the trade-off, the better the algorithm. However, to the extent that a trade-off exists, we want to calibrate the algorithm in order to:

- discard suboptimal sets of filtering criteria, namely those sets which increase recall by decreasing too much precision (and *vice versa*)
- choose among optimal sets, according to the research objectives (some of which may require precision to be sacrificed to recall, or *vice versa*).

We proceed in three steps. First, by means of a MonteCarlo simulation exercise, the algorithm generates a large number of observations, each of which consists of a random set of weights assigned to the filtering criteria, a *Threshold* value, and the corresponding results in terms of precision and recall (*Data generation* step).

Second, the simulation results are split into two sets (*dominant vs dominated*), with the dominant results further split into three regions of interest, each of which is characterized by a different mix of precision and recall (*Mapping* step).

⁷ More precisely, NAFA and NAE contain matches between an inventor and one of his/her patents, and another inventor and one of his/her patents, plus information on whether the two inventors are the same person, according to information collected manually. Having been hand-checked, the matches in the benchmark databases are expected to contain neither false positives nor false negatives. Notice that both NAFA and NAE are based upon the PatStat October 2009 release. A detailed description is available online (Lissoni et al., 2010)

Finally, weights are assigned to the filtering criteria, according to the desired results in terms of precision and recall (*Weight calibration*). Notice that weights are binary values (0,1), which amounts to say that our weight calibration consists in including some filtering criteria (1) and excluding others (0). However, extensions of Massacrator[©] may be conceived, which make use of continuous weights (comprised between 0 and 1).

Sections 3.1, 3.2 and 3.3 describe in details the three steps.

Table 2 Satisfied criteria in benchmark datasets, x^k

	NAE (EPFL)	NAFA (French academics)
City	0.15	0.24
Province	0.01	0.3
Region	0.02	0.42
State	0.02	0
Street	0.04	0.02
IPC 4	0.32	0.31
IPC 6	0.2	0.19
IPC 12	0.1	0.07
Three Years	0.49	0.44
Applicant	0.22	0.25
Small Applicant	0.06	0.03
Group	0.01	0.02
Coinventor	0.09	0.1
Three Degrees	0.13	0.12
Citations	0.08	0.08
Rare Surname	0.07	0.05
ASE	0.07	0.06

Table 3 Correlation between criteria k1 and k2 corr(x^{k1} , x^{k2})

		1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17
1.	City	1	0.45	0.3	-0.02	0.18	0.21	0.17	0.12	0	0.28	0.05	0.11	0.15	0.17	0.12	0.05	0.09
2.	Province	0.45	1	0.79	0	0.02	0.21	0.17	0.2	-0.01	0.26	0.02	0.11	0.28	0.28	0.16	0.09	0.12
3.	Region	0.3	0.79	1	0	-0.01	0.15	0.13	0.17	-0.07	0.24	-0.01	0.09	0.25	0.23	0.15	0.09	0.1
4.	State	-0.02	0	0	1	0.02	0.02	0.01	0.03	0.02	0.06	-0.02	0.01	0.02	0.03	0.01	0.06	0
5.	Street	0.18	0.02	-0.01	0.02	1	0.08	0.11	0.05	0.03	0.11	0.08	-0.01	0.03	0.06	0.07	0	0.07
6.	IPC 4	0.21	0.21	0.15	0.02	0.08	1	0.71	0.43	-0.01	0.35	0.16	0.06	0.27	0.27	0.29	0.07	0.26
7.	IPC 6	0.17	0.17	0.13	0.01	0.11	0.71	1	0.6	0.08	0.37	0.19	0.03	0.29	0.29	0.33	0.07	0.31
8.	IPC 12	0.12	0.2	0.17	0.03	0.05	0.43	0.6	1	0.15	0.3	0.16	0.02	0.32	0.3	0.33	0.05	0.3
9.	3 Years	0	-0.01	-0.07	0.02	0.03	-0.01	0.08	0.15	1	0.13	0.08	0.03	0.16	0.16	0.08	-0.06	0.13
10.	Applicant	0.28	0.26	0.24	0.06	0.11	0.35	0.37	0.3	0.13	1	0.34	0.1	0.34	0.37	0.31	0.07	0.3
11.	Small Applicant	0.05	0.02	-0.01	-0.02	0.08	0.16	0.19	0.16	0.08	0.34	1	-0.01	0.16	0.16	0.19	0.02	0.19
12.	Group	0.11	0.11	0.09	0.01	-0.01	0.06	0.03	0.02	0.03	0.1	-0.01	1	0.05	0.05	0.04	-0.01	0.04
13.	Coinventor	0.15	0.28	0.25	0.02	0.03	0.27	0.29	0.32	0.16	0.34	0.16	0.05	1	0.84	0.28	0.13	0.3
14.	3 Degrees	0.17	0.28	0.23	0.03	0.06	0.27	0.29	0.3	0.16	0.37	0.16	0.05	0.84	1	0.29	0.11	0.32
15.	Citations	0.12	0.16	0.15	0.01	0.07	0.29	0.33	0.33	0.08	0.31	0.19	0.04	0.28	0.29	1	0.06	0.47
16.	Rare Surname	0.05	0.09	0.09	0.06	0	0.07	0.07	0.05	-0.06	0.07	0.02	-0.01	0.13	0.11	0.06	1	0.04
17.	ASE	0.09	0.12	0.1	0	0.07	0.26	0.31	0.3	0.13	0.3	0.19	0.04	0.3	0.32	0.47	0.04	1

3.1. Data generation

We generate data for calibration as follows:

1. *Vectors of criteria*: for each pair of inventors m , a set of k dummy variables x_m^k ($k=1\dots 17$) is generated, each of them corresponding to one of the 17 filtering criteria described in section 2.3. x_m^k takes value 1 if the filtering criterion is satisfied at least once by the inventors' pair, zero otherwise. Tables 2 and 3 report the percentage of pairs satisfying each criterion and the resulting correlation matrix.
2. *Vectors of weights and computation of similarity scores*: We draw randomly W vectors of weights from a uniform Bernoulli multivariate distribution, where W is set to 2000. The dimensions of the multivariate distribution are as many as the number of variables in vector x (i.e. $K=17$). Each draw generates a different vector of weights ω_w , where each k -th weight (ω_w^k) can take value one or zero (i.e. binomial weight). Each pair of matched inventors from NAFA and NAE benchmark databases is then weighted as follows:

$$\alpha_{m,w} = x_m X \omega_w$$

where: $w = 1 \dots 2000$; $m = \{m_{NAFA}, m_{NAE}\}$; $m_{NAFA} = 1 \dots 2817$; $m_{NAE} = 1 \dots 1011$; and sizes of the matrixes are: $x_m [1X17]$; $\omega_w [17X1]$; $\alpha_{m,w} [1X1]$.

Binomial weights can be interpreted as a way to exclude/include randomly the k -th filtering criterion in the x_m set. The product of two vectors x_m and ω_w returns in the $\alpha_{m,w}$ similarity score of match m , for a specific set w of weights.

3. *Threshold value* : In order to determine whether a match is positive or negative the algorithm compares each similarity score $\alpha_{m_{NAFA},w}$ and $\alpha_{m_{NAE},w}$ to a *Threshold* value. We treat the latter as a parameter subject to calibration, too. Therefore, we add to each vector of weights ω_w , a random threshold value, extracted from a uniform distribution with upper bound 4 and lower bound zero:

$$Threshold_w = U(0,4)$$

4. *Observations*: Each vector of weights w generates 2817 $\alpha_{m,w}$ values in case of NAFA and 1011 $\alpha_{m,w}$ values in case of NAE, one for each inventor pair in the dataset. They come along with a threshold value ($Threshold_w$), which allows us to define $D_{\alpha_{m,w}}$ as follows

$$\begin{aligned} D_{\alpha_{m,w}} &= 1 \text{ if } \alpha_{m,w} \geq Threshold_w \\ D_{\alpha_{m,w}} &= 0 \text{ if } \alpha_{m,w} < Threshold_w \\ m &= \{m_{NAE}, m_{NAFA}\} \end{aligned}$$

By comparing $D_{\alpha_{m,w}}$ and D_{γ_m} as in Equation 3, we then compute the number of true (false) positives (negatives) obtained by applying different sets of weights and threshold values $[\omega_w, Threshold_w]$. That is, we generate 4000 records (2000 for NAE and 2000 for NAFA), to be used in our calibration exercise, each record being characterized by a different combination of precision rate, recall rate, vectors of weights and threshold value.

Figure 2 is a scatter plot for the precision and recall rates, where dots correspond to observations and dot colors indicate the relative threshold value. The figure shows the extent of the trade-off between precision and recall. It also shows how the trade-off

depends on the threshold value: higher precision and lower recall for higher thresholds, and *vice versa*. Yet, we observe that for different threshold values we can obtain a similar combination of precision and recall, depending on the values assigned to weights w^k (overlapping regions of dots). Also figure 3 is a scatterplot for precision and recall rates; in this case the dots are grouped according to the benchmark databases they refer to, NAFA and NAE. We notice that NAFA dots tend to exhibit higher precision rates, given the recall rate, and *vice versa*; this suggests that our algorithm fares better when applied to NAFA than to NAE, that is, it is sensitive to the benchmark chosen for calibration.

Figure 2 - Precision and Recall values according to different threshold (t) values (4000 sets of weights)

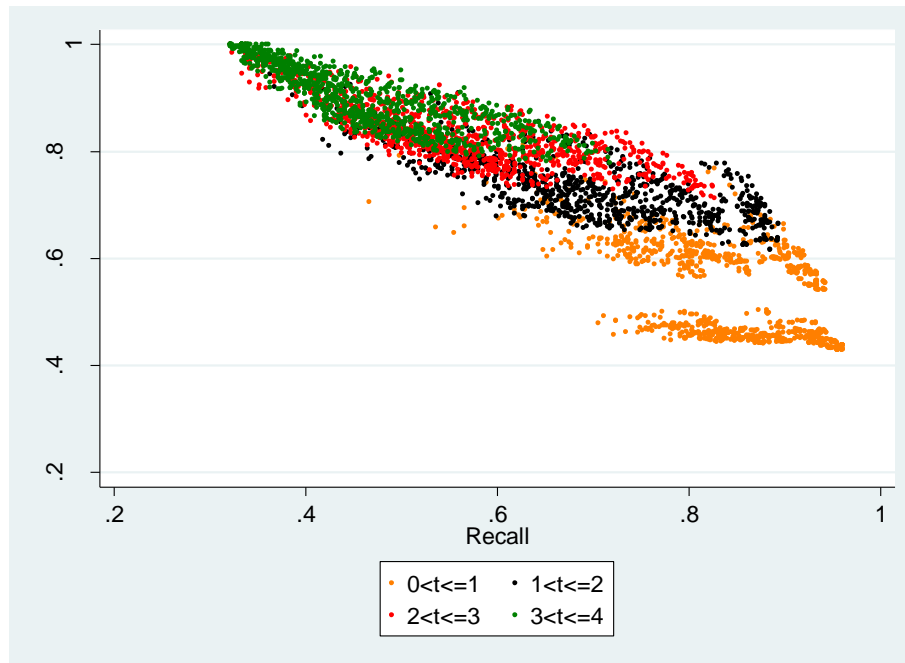
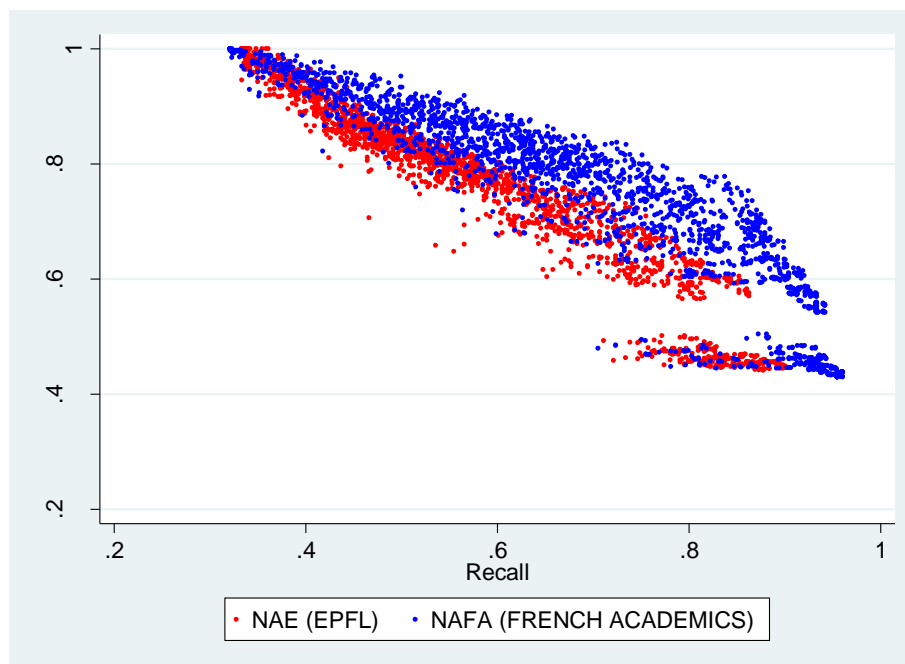


Figure 3 - Precision and Recall values according to NAFA and NAE datasets (4000 sets of weights)



3.2. Mapping

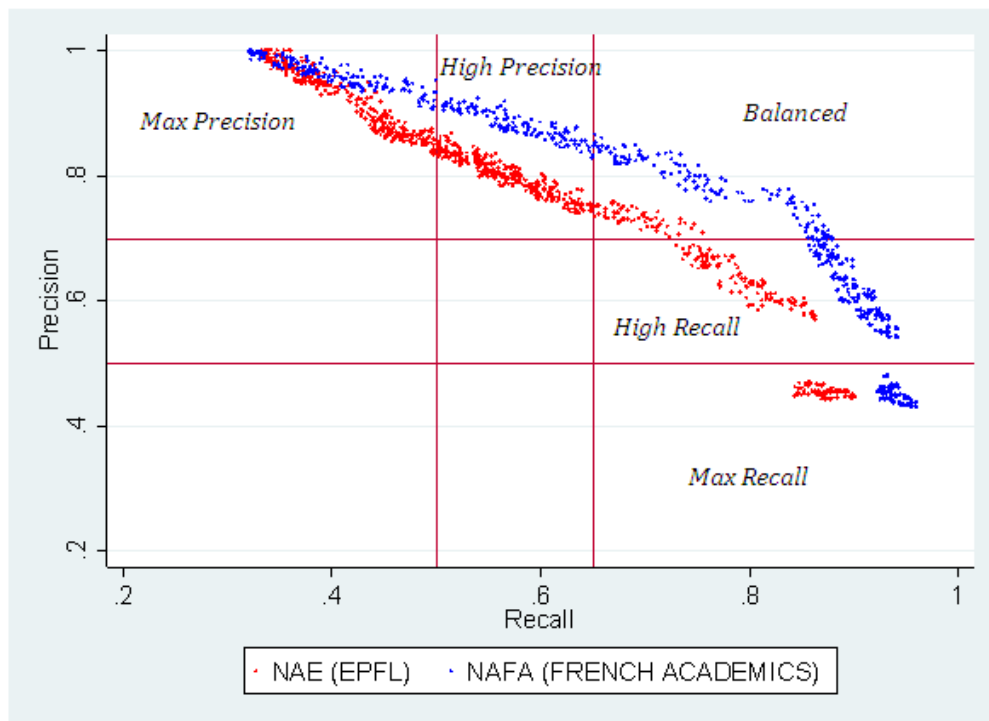
This second step identifies the most efficient combinations of weights with respect to pre-defined objective regions: *high precision*, *high recall* and *balanced* mix of recall and precision levels. Outcomes (observations) in each region are first split in two groups: *dominant* and *dominated*. An outcome is dominated whenever another outcome exists which has both higher precision and higher recall; it is dominant whenever no such other outcome exists.

$$\text{Dominant outcomes } (\bar{o}) \rightarrow \{ \bar{o} : \nexists \underline{o} \text{ Precision}(\underline{o}) > \text{Precision}(\bar{o}) \cap \text{Recall}(\underline{o}) > \text{Recall}(\bar{o}) \}$$

$$\text{Dominated outcomes } (\underline{o}) \rightarrow \{ \underline{o} : \exists \bar{o} \text{ Precision}(\underline{o}) < \text{Precision}(\bar{o}) \cap \text{Recall}(\underline{o}) < \text{Recall}(\bar{o}) \}$$

Dominant outcomes can be seen in figure 3 as dots at the upper frontier of the cloud of observations. If we consider separately the clouds for NAFA and NAE results, we would obtain two distinct sets of dominant outcomes, one for each benchmark dataset, as in figure 4⁸. Vertical lines in the figure identify nine areas, three of which include outcomes corresponding to our objectives of high precision, high recall and balanced results. In particular, the high precision area includes all dominant outcomes with precision rate higher than 0.7, and recall rate between 0.5 and 0.65; the high recall region includes all dominant outcomes with a recall rate higher than 0.65, and precision rate between 0.5 and 0.7; the balanced results region includes all dominant outcomes with a recall higher than 0.65 and precision higher than 0.7.

Figure 4 - Dominant Solutions for NAE and NAFA benchmarks



Notice that two other areas of potential interest are “maximum precision” and “maximum recall” (see figure 4). However, these are not reasonable objectives to

⁸ The NAFA and NAE frontiers in figure 6, include not only the most extreme points, but are extended to include all outcomes with precision and recall values higher than $\text{Precision}(\bar{o})-0.02$ and $\text{Recall}(\bar{o})-0.02$ for any \bar{o} . This will turn out useful for the ensuing statistical exercise.

pursue, as they come at too high a cost in terms of recall and precision, respectively (e.g. to achieve max precision we should stand a recall rate of less than 0.5, which is worse than the result of just guessing).

We also calculate the number of positive weights characterizing the sets of weights within the region of interest (*AVG nr filtering criteria* positively weighted). In the *Balanced* region for the NAFA benchmark of Figure 4, we have 132 vectors of weights, one for each algorithm run falling within the region, having on average 8.77 filtering criteria positively weighted. However, we count the number of positive weights assigned to criteria with integer numbers, then we can conclude that the 132 observations are on average characterized by nine positive weights..

3.3. Weight Assignment and threshold selection

Once defined the three regions of interest, we assess which of the filtering criteria are over-represented (or under-represented) within each region, and consequently we select them for inclusion in the vector of weights representing the calibrated parametrization of the algorithm. Criterion k is over-represented (under-represented) if the expected value of its weight $E[\omega^k]$ in the region of interest is significantly higher (lower) than 0.5⁹.

We test the over-representation (under-representation) hypothesis by means of one-tail t -tests, with 95% significance, as follows:

- ✓ Over-representation test for criterion k $H_0: E[\omega^k]=0.5$ $H_1: E[\omega^k]>0.5$
- ✓ Under-representation test for criterion k $H_0: E[\omega^k]=0.5$ $H_1: E[\omega^k]<0.5$

We then proceed by including in the algorithm all over-represented criteria (that is, we assign them weight $\omega^k=1$), and excluding the under-represented ones (assign $\omega^k=0$), depending on the objective region.

Table 4a and 4b report for each filtering criterion, the sample mean of its weight and the p-values of the one-tail t -tests. Separate tests are run for NAFA and NAE benchmark datasets and for the three regions of interest.

For illustration consider *City* and *State* criteria from table 4a (NAFA dataset) in the *Balanced* precision-recall region¹⁰.

We observe a sample mean equal to 0.42 for *City* criterion, which translates into a rejection of the null hypothesis in the under-representation test (p-value=0.03), but not in the over-representation test (p-value=0.97). As the *City* criterion is significantly under-represented in the observations characterized by *Balanced* objective, then we exclude it by assigning to *city* criterion a zero-weight ($\omega^{City}=0$).

⁹ Remember that ω_w^k is a random variable with expected value equal to 0.5. By definition, any sample with a different mean cannot be randomly drawn, and must be considered either over- or under-represented by comparison to a random distribution.

¹⁰ Regression analysis can be applied to the same set of results in order to estimate the marginal impact of each filtering criterion and the *Threshold* on either precision and recall, other things being equal. In general, we expect all filters to bear a negative influence on recall (in that they increase the number of negative matches, both true and false), and a positive influence on precision (they eliminate false positives). In case the estimated impact of a criterion is not significantly different than zero for recall, but positive for precision, then it is desirable to include it in any parametrization, as it increases precision at no cost in terms of recall. Conversely, any filter with zero impact on precision, but significantly negative for recall, ought to be excluded from any parametrization, as it bears a cost in the terms of the latter, and no gains in terms of precision. We have conducted this type of analysis, and found it helpful to understand the relative importance of the different filtering criteria. We do not report it for reasons of space, but it is available on request.

On the contrary, for the *State* criterion, the null hypothesis cannot be rejected either in the under-representation test nor in the over-representation test (p-values being respectively 0.64 and 0.36). This means that *t-tests* do not give a clear (and statistically significant) evidence to help us deciding whether to include or exclude the *State* criterion. In this case we give a positive weight to the *State* criterion only if it contributes to reach, in the calibrated parametrization, the average number of positively weighted filtering criteria characterizing the observations in the objective region. (that is, if the positively weighted criteria selected in the calibrated parametrization are less than the nine observed on average in the *Balanced* case, *State* is included).

Results of the tests provide us with a guide for choosing the filtering criteria to include (assign positive weight) in calibrated parametrization of the algorithm, according to the precision and recall objectives we aim at. For sake of simplicity we identify the positively weighted criteria with an asterisk in table 4a and 4b.

In the case of NAFA benchmark, whatever the objective region, *network* criteria are always assigned a positive weight. In case of NAE only *Three degrees* is assigned a positive weight, for all the three regions.

Table 4a - Averages [$H_0: E[\omega^k]=0.5$ $H_1: E[\omega^k]>0.5$, $H_0: E[\omega^k]=0.5$ $H_0: E[\omega^k]<0.5$]

NAFA (French Academics)	Balanced			High Precision			High Recall		
	Mean	Pvalue	Pvalue	Mean	Pvalue	Pvalue	Mean	Pvalue	Pvalue
		$H_0: E[\omega^k]=0.5$	$H_1: E[\omega^k]>0.5$		$H_0: E[\omega^k]=0.5$	$H_1: E[\omega^k]>0.5$		$H_0: E[\omega^k]=0.5$	$H_1: E[\omega^k]>0.5$
<i>Network variables</i>									
Coinventor	0.67*	0	1	0.75*	0	1	0.52*	0.29	0.71
Three Degrees	0.48*	0.7	0.3	0.67*	0	1	0.57*	0.02	0.98
<i>Geographical variables</i>									
City	0.42	0.97	0.03	0.29	1	0	0.47	0.83	0.17
Province	0.93*	0	1	0.71*	0	1	0.72*	0	1
Region	0.93*	0	1	0.77*	0	1	0.81*	0	1
State	0.48*	0.64	0.36	0.51*	0.4	0.6	0.46	0.89	0.11
Street	0.33	1	0	0.38	1	0	0.49	0.61	0.39
<i>Applicant related variables</i>									
Applicant	0.49*	0.57	0.43	0.53*	0.22	0.78	0.51*	0.34	0.66
Small Applicant	0.52*	0.36	0.64	0.41	0.98	0.02	0.56*	0.05	0.95
Group	0.52*	0.36	0.64	0.5*	0.53	0.47	0.5	0.5	0.5
<i>Technology classes</i>									
IPC 4	0.37	1	0	0.4	0.99	0.01	0.6*	0	1
IPC 6	0.3	1	0	0.22	1	0	0.52*	0.25	0.75
IPC 12	0.45	0.89	0.11	0.48*	0.67	0.33	0.51*	0.39	0.61
<i>Citation related variables</i>									
Citations	0.45	0.89	0.11	0.46	0.83	0.17	0.49	0.66	0.34
ASE	0.45	0.89	0.11	0.44	0.91	0.09	0.47	0.79	0.21
<i>Other filtering criteria</i>									
Rare Surname	0.6*	0.01	0.99	0.6*	0.01	0.99	0.5	0.5	0.5
Three Years	0.39	0.99	0.01	0.43	0.93	0.07	0.16	1	0
<i>Nr of filtering criteria and threshold</i>									
AVG nr filtering criteria	8.77			8.57			8.86		
AVG threshold	2.22			3.16			0.76		
Observations	132			129			214		

The family of *geographic* criteria plays an important role in the NAFA benchmark, but not in the NAE benchmark. This is not surprising given the low quality of geographical information for Swiss inventors available on PatStat data (see Lissoni et al., 2010). *Applicant* and *Technology* families show a mixed evidence, the choice of weights being specific to any combination of benchmark dataset and objective regions. The *Citation* family does not play any role in NAFA dataset, while it has to be weighted positively in NAE dataset. Among the remaining criteria (*others* family), having a *rare surname* has to be included in NAFA database when objective regions are *Balanced* and *High precision*, as well as *Three years* in case of NAE benchmark database.

Once defined the vector of weights for the calibrated parametrization of the algorithm, a threshold value is needed, which we calculate as the average threshold

value within each region. For instance, in the *Balanced* region of the NAFA benchmark, the average threshold value for the 132 outcomes (dots) is 2.22. It means that the similarity score $\alpha_{m,w}$ must be equal to or higher than 2.22 for any match to be considered positive. As expected, the average threshold value is highest in the high precision region and lowest in the high recall one (see tables 4a and 4b).

Table 4b Averages [$H_0: E[\omega^k]=0.5$ $H_1: E[\omega^k]>0.5$, $H_0: E[\omega^k]=0.5$ $H_0: E[\omega^k]<0.5$]

NAE (EPFL Scientists)	Balanced			High Precision			High Recall		
	Mean	Pvalue		Mean	Pvalue		Mean	Pvalue	
		$H_0: E[\omega^k]=0.5$ $H_1: E[\omega^k]>0.5$	$H_0: E[\omega^k]=0.5$ $H_1: E[\omega^k]<0.5$		$H_0: E[\omega^k]=0.5$ $H_1: E[\omega^k]>0.5$	$H_0: E[\omega^k]=0.5$ $H_1: E[\omega^k]<0.5$		$H_0: E[\omega^k]=0.5$ $H_1: E[\omega^k]>0.5$	$H_0: E[\omega^k]=0.5$ $H_1: E[\omega^k]<0.5$
<i>Network variable</i>									
Coinventor	0.56*	0.17	0.83	0.51	0.33	0.67	0.54	0.17	0.83
Three Degrees	0.58*	0.08	0.92	0.54*	0.1	0.9	0.59*	0.02	0.98
<i>Geographical variables</i>									
City	0.6*	0.05	0.95	0.55*	0.04	0.96	0.67*	0	1
Province	0.42	0.92	0.08	0.41	1	0	0.42	0.96	0.04
Region	0.51	0.41	0.59	0.47	0.84	0.16	0.42	0.96	0.04
State	0.46	0.76	0.24	0.54	0.08	0.92	0.56*	0.1	0.9
Street	0.42	0.92	0.08	0.38	1	0	0.56*	0.1	0.9
<i>Applicant related variables</i>									
Applicant	0.94*	0	1	0.75*	0	1	0.93*	0	1
Small Applicant	0.63*	0.02	0.98	0.49	0.67	0.33	0.64*	0	1
Group	0.43	0.88	0.12	0.46	0.92	0.08	0.53	0.27	0.73
<i>Technology classes</i>									
IPC 4	0.38	0.98	0.02	0.62*	0	1	0.75*	0	1
IPC 6	0.38	0.98	0.02	0.52*	0.26	0.74	0.56*	0.1	0.9
IPC 12	0.69*	0	1	0.65*	0	1	0.56	0.07	0.93
<i>Citation related variables</i>									
Citations	0.54*	0.24	0.76	0.59*	0	1	0.52	0.33	0.67
ASE	0.53*	0.32	0.68	0.52*	0.26	0.74	0.56*	0.1	0.9
<i>Other filtering criteria</i>									
Rare Surname	0.28	1	0	0.5	0.46	0.54	0.41	0.98	0.02
Three Years	0.83*	0	1	0.65*	0	1	0.23	1	0
<i>Nr of filtering criteria and threshold</i>									
AVG nr filtering chrriteria	9.17			9.15			9.43		
AVG threshold	1.42			2.42			0.75		
Observations	72			336			135		

4. Validation and Application to PatStat data

Following our calibration exercise, we produced three versions (parametrizations) of Massacrator[©], one for each precision-recall objective, with weights and threshold calculated accordingly. We then checked to what extent each of these parametrizations is satisfying in terms of the precision and recall rates it produces, conditional on its objective. Precision and recall rates are measured, once again, against the NAFA and NAE benchmarks.

We run each version of Massacrator[©], once for each combination for each benchmark, for a total of 6 six runs, with the following results:

- NAFA dataset - precision oriented parametrization -> Precision:92% Recall:54%
- NAFA dataset - recall oriented parametrization -> Precision:56% Recall:93%
- NAFA dataset - balanced parametrization -> Precision:88% Recall:68%
- NAE dataset - precision oriented parametrization -> Precision:79% Recall:62%
- NAE dataset - recall oriented parametrization -> Precision:59% Recall:85%
- NAE dataset - balanced parametrization -> Precision:74% Recall:70%

Notice that by calibrating our filtering step on either NAE or NAFA we obtain different results. This is because each dataset has a number of semantic peculiarities (variety of names and; quality of information contained in the addresses; variety in the technological classes and citations of patents), which are mirrored by differences in the number and type of criteria selected at the calibration stage.

This forced us to choose one and only benchmark dataset to perform our final calibration, the one leading to the production of the APE-INV dataset. Our choice fell on NAFA, which contains higher quality information for addresses, and more name variety. For the three alternative parametrizations of Massacrator algorithm we then obtain the following disambiguation results¹¹:

- NAFA calibrated, precision-oriented algorithm: from 2,806,516 inventors in the original PatStat database (for EPO patents) we obtain 2,520,338 disambiguated inventors (unique codes) in APE-INV, that is -10%
- NAFA calibrated, recall-oriented algorithm: from 2,806,516 inventors to 1,697,976 unique codes, that is -39%
- NAFA calibrated, balanced algorithm: from 2,806,516 inventors to 2,366,520 unique codes, that is -16%

As expected the largest reduction in the number of inventors is obtained with the recall-oriented algorithm, the smallest with the precision-oriented one. More importantly, when applying data disambiguated with different precision-recall objectives to classic problems in the economics of innovation or science and technology studies, we will get different results. As an illustration, consider three classical topics: inventors' productivity, mobility, and social networking (on the latter topic, see Borgatti et al., 2009 for technical vocabulary and basic concepts). Table 5 reports descriptive statistics for each topic, as resulting from datasets built by using different parametrizations of Massacrator, namely:

- *Avg. Patent per inventor*: it is the average number of patents per inventor in the whole dataset
- *Star inventors' productivity*: it is the share of patents belonging to the 1000 most prolific inventors in the database.
- *International mobility index*: It is the share of inventors with at least two different country addresses, over the total number of inventors with at least two patents (inventors with only one patent are not considered, as they can have only one address, by definition).
- *Connectedness*: it is the percentage of connected nodes over the total number of nodes in the network of inventors active between 2000 and 2005 in the fields of chemistry and pharmaceuticals (from now on: Net2000/05)¹². Isolated nodes represent individuals with no co-inventorship relationships over the period considered.
- *Centralization-degree*: it is a degree-based measure of graph centrality for Net2000/05, as defined in Freeman (1979). It measures the extent at which the

¹¹ The figures presented here are the result of further adjustments we introduced in order to solve transitivity problems. Transitivity problems may emerge for any triplet of inventors (such as I, J, and Z) whenever two distinct pairs are recognized to be same person (e.g, I & J and J and Z), but the same does not apply to the remaining pair (I & Z are not matched, or are considered negative matches). In this case we need to decide whether to revise the status of I & Z (and consider the two inventors as the same person as J) or the status of the other pairs (and consider either I or Z as different persons than J). When confronting this problem, we always opted for considering the two inventors the same person, then I,J and Z are the same individual according to Massacrator..

¹² Fields of chemistry and pharmaceuticals are defined as in Schmoch (2008). We consider only these fields, and years from 2000 and 2005, for ease of computation. Co-inventorship is intended as a connection between two inventors having (at least) one patent in common.

graph structure is organized around focal points, and it reaches a maximum value for a star graph.

- *Density*: it is the number of observed ties in Net2000/05, over the maximum number of possible ties (i.e. the number of ties in a fully connected network with the same number of nodes). It measures the intensity of connections between inventors.

Table 5 Descriptive statistics of inventorship: Massacrator runs with different parametrizations on the whole PatStat dataset

	(1)	(2)	(3)	(4)	(5)	(6)
Parametrization	Avg Patent per inventor	Star inventors' productivity	International mobility index	Connected nodes %	Centralization -degree %	Density %
Balanced	2.1705	3.56%	1.02%	95%	0.156	0.0034
Recall-oriented	3.0244	8.19%	4.92%	95.21%	0.515	0.0053
Precision-oriented	2.0381	3.48%	0.56%	95%	0.149	0.0032

As expected the productivity index in column (1) is higher for the recall-oriented parametrization of the algorithm, on the basis of which we treat a larger number of inventors as the same individual. The opposite happens with precision-oriented parametrization. As similar consideration is valid also for statistics on star inventors, which are assigned a maximum of 8% of patents when using a recall-oriented parametrization and only 3.5% with a precision-oriented parametrization. As for international mobility, its index ranges from 0.56% to 4.92% according to the parametrization choice. While productivity measures do not change much when moving from the *Precision*-oriented parametrization to the *Balanced* one, the same cannot be said for mobility measures: in this case, even the modest reduction in precision (increase in recall) introduced by changing algorithm changes considerably the value of the indicator.

As for network measures, *Connectedness* is not very sensitive to the algorithm parametrization. The same cannot be said for *Centralization* and *Density*, both of which increase considerably along with recall and decline with precision.

5. Conclusions and further research

In this paper we have presented a general methodology for inventor disambiguation, with an application to EPO patent data. We have argued that producing high quality data requires calibrating the choice of weights by means of simulation analysis. Calibration is necessary to:

1. identify "frontier" results, that is the set of efficient weights that maximise the precision rate, conditional on a given recall rate (or, vice versa, recall conditional on precisions); in this way, one excludes inefficient sets of weights and make less arbitrary choices;
2. allow the researcher to choose between precision-oriented, recall-oriented or balanced"algorithms, or to combine them.

Choosing one algorithm over the others may be desirable when the research purposes require minimization of either errors of type I or errors of type II (respectively, false positives and false negatives). For example, early research on academic patenting by Lissoni et al. (2008) was aimed at proving that official estimates of the number of academic patents (namely, patents signed by at academic scientists) in Europe were wrong by defect, and thus needed to minimize errors of type I. A more recent study on the same topic, on the contrary, has produced a longitudinal database of academic

patenting in Italy, with the primary objective of detecting trends (Lissoni et al. 2012). With that objective in mind, there is no reason to prefer minimization of either errors of type I or errors of type II, so the authors make use of the APE-INV inventor database described in this paper, with balanced parametrization.

More generally, we have shown how different calibrations lead to different results for the fundamental indicators of studies on inventors' productivity, mobility, and networking. This means that the results of several studies recently published on these topics, which do not provide details on the disambiguation methods they followed, could turn out not be robust to the calibration choices presented here (this include some work by one us, such as: Balconi et al., 2004 on networks; or Breschi and Lissoni, 2009, on mobility). For sure, future research results in these area will have to be screened more closely, and disambiguation methods made explicit.

Besides conducting robustness test for different disambiguation parametrizations, authors may also pursue the road of combining the results of different calibrations in order to increase data quality. This can be done by comparing the results, for each pair of records, of the different calibrations and choose on the basis of whether a majority of the algorithms suggest the same results. The combination principles can be extended not only to different calibration algorithms, but to altogether different algorithms, as discussed by Maurino et al. (2012).

One last strategy for further data quality improvements can consist in sharing more openly inventor data and collecting feedbacks from other users. This is an integral part of the APE-INV project, for which the inventor database described in this paper was produced and made available online (<http://www.ape-inv.disco.unimib.it/>). Users who check manually the inventor data they download, or match them to other sources of information on individuals (such as lists of academics or authors of scientific papers) do inevitably find a number of false negatives or false positives. The same holds if their research requires contacting the inventors for interview or survey purposes. This user-generated information is extremely valuable, and we believe it is worth investing in finding ways to collect it. To this end we have set up the APE-INV User's Feedback project, which invites users to come back to the APE-INV data website and upload either their proposed corrections to the APE-INV inventor dataset, or the results of their own disambiguation exercises based on the same set of data (for a full description of the project, see Den Besten et al., 2012).

Collective use and quality improvement of inventor data would also serve the purpose of minimizing social cost of research, by minimizing the duplicative data disambiguation efforts pursued by the many researchers in the area. We are confident that the APE-INV database will contribute decisively in this direction.

References

- Agrawal A., Cockburn I., & McHale J. (2006). Gone but not forgotten: knowledge flows, labor mobility, and enduring social relationships. *Journal of Economic Geography*, 6(5):571.
- Azoulay, P., Ding, W. & Stuart, T. (2009). The impact of academic patenting on the rate, quality and direction of (public) research output. *The Journal of Industrial Economics*, 57, 637-676.
- Balconi M., Breschi S., & Lissoni F. (2004). Networks of inventors and the role of academia: an exploration of Italian patent data* 1. *Research Policy*, 33(1):127-145.
- Borgatti, S. P., Mehra, A., Brass, D. J., & Labianca, G. (2009). Network analysis in the social sciences. *science*, 323(5916), 892-895.
- Breschi S. & Lissoni F. (2009). Mobility of skilled workers and co-invention networks: an anatomy of localized knowledge flows. *Journal of Economic Geography*.
- Breschi S., Lissoni F., & Montobbio F. (2008). University patenting and scientific productivity: a quantitative study of Italian academic inventors. *European Management Review*, 5(2) 91-109.
- Burt, R. S. (1987). Social contagion and innovation: Cohesion versus structural equivalence. *American journal of Sociology*, 1287-1335.
- Carayol N. & Cassi L. (2009). Who's Who in Patents. A Bayesian approach. *Cahiers du GREThA*, 2009-07.
- Den Besten M., Lissoni F., Maurino A., Pezzoni M., Tarasconi G. (2012). Ape-Inv Data Dissemination And Users' Feedback Project", mimeo (<http://www.academicpatentig.eu>)
- Fleming, L., King, C., & Juda, A. I. (2007). Small Worlds and Regional Innovation. *Organization Science*, 18, 938-954.
- Freeman L. C. (1979). Centrality in social networks conceptual clarification. *Social networks*. *Social Networks* 1(3) 215-239.
- Hall B.H., Jaffe A.B., & Trajtenberg M. (2001). The NBER patent citation data _le: Lessons, insights and methodological tools. *NBER working paper*.
- Huang H. & Whalsh J.P. (2010). *A new Name-Matching Approach for searching Patent Inventors*. mimeo.
- Kim J. & Marschke G. (2005). The influence of university research on industrial innovation. *NBER working paper*.
- Lai R., D'Amour A., & Fleming L. (2009). The careers and co-authorship networks of US patent-holders, since 1975. *Unpublished Working Paper*, Harvard University.
- Lai, R., D'Amour, A., Yu, A., Sun, Y., Torvik, V., & Fleming, L. (2011). Disambiguation and co-authorship networks of the US Patent Inventor Database. *Harvard Institute for Quantitative Social Science, Cambridge, MA, 2138*.
- Lissoni F., Coffano M., Maurino A., Pezzoni M., & Tarasconi G. (2010). *APE-INV's Name Game Algorithm Challenge: A Guideline for Benchmark Data Analysis & Reporting*. mimeo.
- Lissoni F., Llerena P., McKelvey M., & Sanditov B. (2008). Academic patenting in Europe: new evidence from the KEINS database. *Research Evaluation*, 17(2):87-102.

- Lissoni, F., Llerena, P. & Sanditov, B. (2011). *Small Worlds In Networks Of Inventors And The Role Of Science: An Analysis Of France*. Bureau D'economie Théorique Et Appliquée, Uds, Strasbourg
- Lissoni F., Sanditov B., & Tarasconi G. (2006). The Keins database on academic inventors: methodology and contents. *WP cespri*, 181.
- Magerman T., Van Looy B., & Song X. (2006). Data production methods for harmonized patent statistics: Patentee name harmonization. *KUL Working Paper* No. MSI 0605.
- Marx M., Strumsky D., & Fleming L. (2009). Mobility, skills, and the Michigan non-compete experiment. *Management Science*, 55(6):875-889.
- Maurino A., Li P. (2012). *Deduplication of large personal database*. Mimeo, 2012
- Raffo J. and Lhuillery S. (2009). How to play the name game: Patent retrieval comparing different heuristics. *Research Policy*, 38(10):1617-1627.
- Schmoch U. (2008). *Concept of a technology classification for country comparisons. Final report to the World Intellectual Property Organization (WIPO)*, Fraunhofer Institute for Systems and Innovation Research, Karlsruhe.
- Tang L. & Walsh J.P. (2010). Bibliometric fingerprints: name disambiguation based on approximate structure equivalence of cognitive maps. *Scientometrics*, 84(3):763-784.
- Thoma G. and Torrisi S. (2007). *Creating powerful indicators for innovation studies with approximate matching algorithms. A test based on PATSTAT and Amadeus databases*. In Conference on Patent Statistics for Policy Decision Making, pages 2-3.
- Thoma G., Torrisi S., Gambardella A., Guellec D., Hall B.H., & Harhoff D. (2010). Harmonizing and Combining Large Datasets An Application to Firm-Level Patent and accounting Data. *NBER Working Paper*.
- Trajtenberg M., Shiff G., & Melamed R. (2006). The names game: Harnessing inventors' patent data for economic research. *NBER working paper*.

Cahiers du GREThA **Working papers of GREThA**

GREThA UMR CNRS 5113

Université Montesquieu Bordeaux IV
Avenue Léon Duguit
33608 PESSAC - FRANCE
Tel : +33 (0)5.56.84.25.75
Fax : +33 (0)5.56.84.86.47

<http://gretha.u-bordeaux4.fr/>

Cahiers du GREThA (derniers numéros – last issues)

- 2012-17 : BERTHE Alexandre, FERRARI Sylvie, *Ecological inequalities: how to link unequal access to the environment with theories of justice?*
- 2012-18 : SALLE Isabelle, YILDIZOGLU Murat, *Efficient Sampling and Metamodeling for Computational Economic Models*
- 2012-19 : POINT Patrick, *L'évaluation des services des écosystèmes liés aux milieux aquatiques. Éléments de méthodologie.*
- 2012-20 : SALLE Isabelle, ZUMPE Martin, YILDIZOGLU Murat, SENEGAS Marc-Alexandre, *Modelling Social Learning in an Agent-Based New Keynesian Macroeconomic Model*
- 2012-21 : ABDELLAH Kamel, NICET-CHENAF Dalila, ROUGIER Eric, *FDI and macroeconomic volatility: A close-up on the source countries*
- 2012-22 : BECUWE Stéphane, BLANCHETON Bertrand, CHARLES Léo, *The decline of French trade power during the first globalization (1850-1913)*
- 2012-23 : ROUILLON Sébastien, *An Economic Mechanism to Regulate Multispecies Fisheries*
- 2012-24 : LISSONI Francesco, MONTOBBIO Fabio, *The ownership of academic patents and their impact. Evidence from five European countries*
- 2012-25 : PETIT Emmanuel, TCHERKASSOF Anna, GASSMANN Xavier, *Sincere Giving and Shame in a Dictator Game*
- 2012-26 : LISSONI Francesco, PEZZONI Michele, POTI Bianca, ROMAGNOSI Sandra, *University autonomy, IP legislation and academic patenting: Italy, 1996-2007*
- 2012-27 : OUEDRAOGO Boukary, *Population et environnement : Cas de la pression anthropique sur la forêt périurbaine de Gonsé au Burkina Faso.*
- 2012-28 : OUEGRAOGO Boukary, FERRARI Sylvie, *Incidence of forest income in reducing poverty and inequalities: Evidence from forest dependant households in managed forest' areas in Burkina Faso*
- 2012-29 : PEZZONI Michele, LISSONI Francesco, TARASCONI Gianluca, *How To Kill Inventors: Testing The Massacrator[®] Algorithm For Inventor Disambiguation*

La coordination scientifique des Cahiers du GREThA est assurée par Sylvie FERRARI et Vincent FRIGANT. La mise en page est assurée par Anne-Laure MERLETTE.