



PatentsView Disambiguation Inventor Workshop

Tom Magerman (tom.magerman@econ.kuleuven.be)

ECOOM (Centre for R&D Monitoring)

INCENTIM (International Centre for Studies in Entrepreneurship and Innovation Management)

K.U.Leuven Managerial Economics, Strategy & Innovation

September 24, 2015, USPTO Washington D.C.



Overview

1. Who we are
2. What we do
3. Why we did not deliver results
4. Why we are here



Who we are

INCENTIM, research division of University of Leuven, Belgium
Bart Van Looy, Tom Magerman, Julie Callaert

In collaboration with:

Università Commerciale Luigi Bocconi, Milano, Italy
Gianluca Tarasconi

Nordic Institute for Studies in Innovation, Research and Education
Eric Iversen



Overview

1. Who we are
2. What we do
3. Why we did not deliver results
4. Why we are here



What we do

- Involved in many innovation studies, including patent statistics and indicator development, for local, federal and European government
- Pioneered in the 1990's with large scale patent and publication databases for innovation research
- Also active in methodological research (use of machine learning, data mining and text mining)



What we do

- Close collaboration with EPO Worldwide Patent Statistical Database (PATSTAT) team (auditing, quality control, new indicator development)
- Development of patentee name harmonization, NUTS-allocation and sector allocation (in collaboration with EUROSTAT)
- Development of NPR classification and linkage with WOS/SCOPUS



Overview

1. Who we are
2. What we do
- 3. Why we did not deliver results**
4. Why we are here





- Manually labeled datasets are not only important for validation purposes (calculation of precision and recall)
- But also for model development and training (quickly assess potential routes for improvements and fine tune parameter and threshold settings)
- Needs to include both positive cases (similar inventors labeled as identical) as negative cases (similar inventors labeled as non-identical) (learning from negative cases is as important as learning from positive cases)
- Needs to be exhaustive (if any case is validated, all potential related inventors need to be assessed)

- OE labeled dataset (Akinsanmi) with 98,762 labeled USPTO records corresponding to inventors of optoelectronic patents (Akinsanmi et al., 2014; Ventura et al., 2015), but without any link to the raw data, and hence not usable if you want to improve current methods;
- ALS labeled dataset with 42,376 labeled USPTO patent-inventor records corresponding to a subset of 4,801 academics in the life sciences with patents (Azoulay et al., 2007; Azoulay et al., 2012), but again without any link to the raw data, and hence not usable to improve current methods;
- IS labeled dataset with 3,845 unique Israeli inventors (Trajtenberg and Shiff, 2008), with a link to the raw data, but as this dataset is limited to Israeli inventors (and rather small), we also consider this labeled dataset as not usable to improve current methods;
- E&S labeled dataset with 96,104 labeled patent-inventor records with 14,293 unique engineers and scientists (Chunmian, Ke-wei and Ping, 2015), with a link to the raw data, hence a good starting point to improve current methods;
- EPO labeled dataset (Lissoni et al., 2010) with 1,498 PATSTAT person records of EPO patent application from scientists affiliated to French universities, and 843 PATSTAT person records of EPO and WIPO patent applications from scientists affiliated to EPFL. As these sets are small, focused on French/EPFL scientists, and linked to EPO patent applications, we also consider this labeled dataset as not usable to improve current methods.

- No explicit labeled negative cases
- Ample implicit labeled negative cases (having a different disambiguation ID is not enough, we need to know whether 'Zimmermann' is equal to 'Zimmerman')
- Non-exhaustive validation ('Lake Rickie' is linked to 'Rickie C.' and 'Rickie Charles', but not to 'Rick' and 'Rick C.')
- Not isolated cases, but fundamental problem: only 10 cases were found where the same last name and first name was linked to a different disambiguation ID (homonyms)

- How to calculate precision?
- How to train new models?
(not possible to quickly assess routes for improvements; difficult to learn where to stop merging inventors)
- Cheating?



Cheating is wrong. Cheating is wrong.
Cheating is wrong. Cheating is wrong.
Cheating is wrong. Cheating is wrong.
Cheating is wrong. Cheating is wrong.
Cheating is wrong. Cheating is wrong.
Cheating is wrong. Cheating is wrong.





Overview

1. Who we are
2. What we do
3. Why we did not deliver results
- 4. Why we are here**



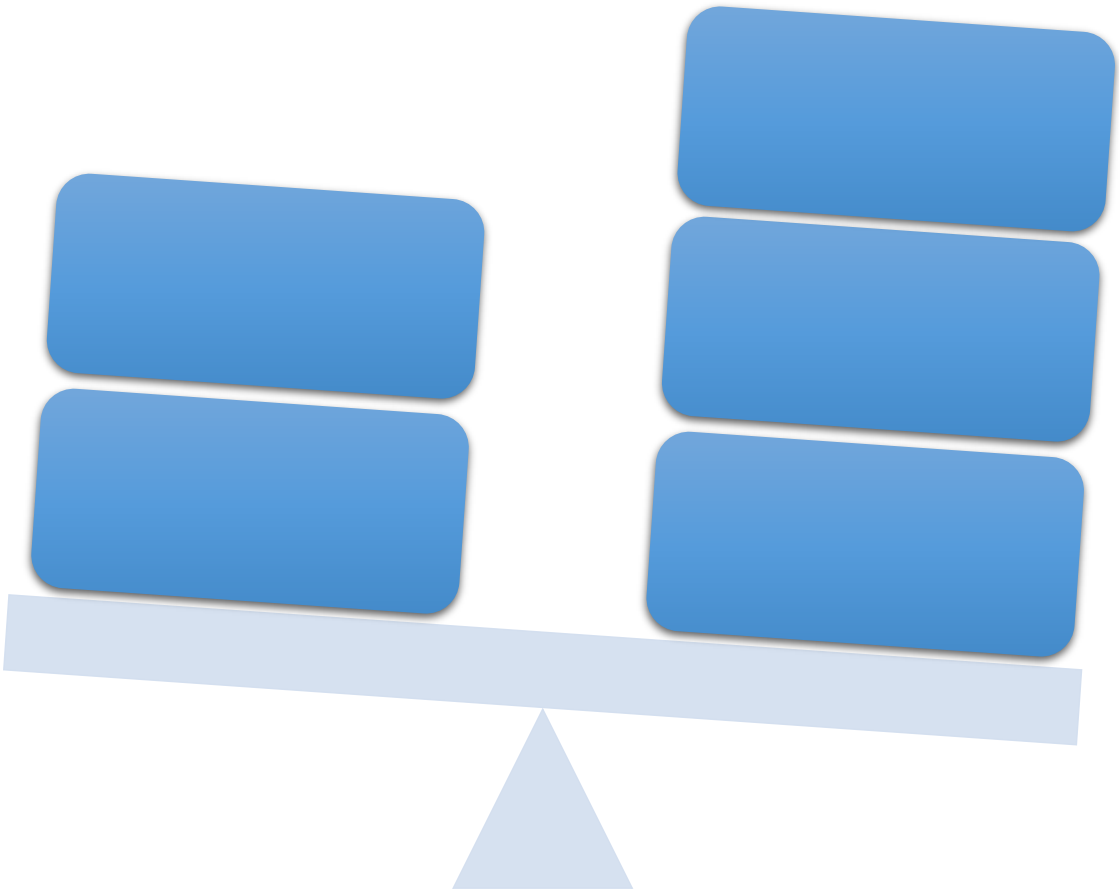
==

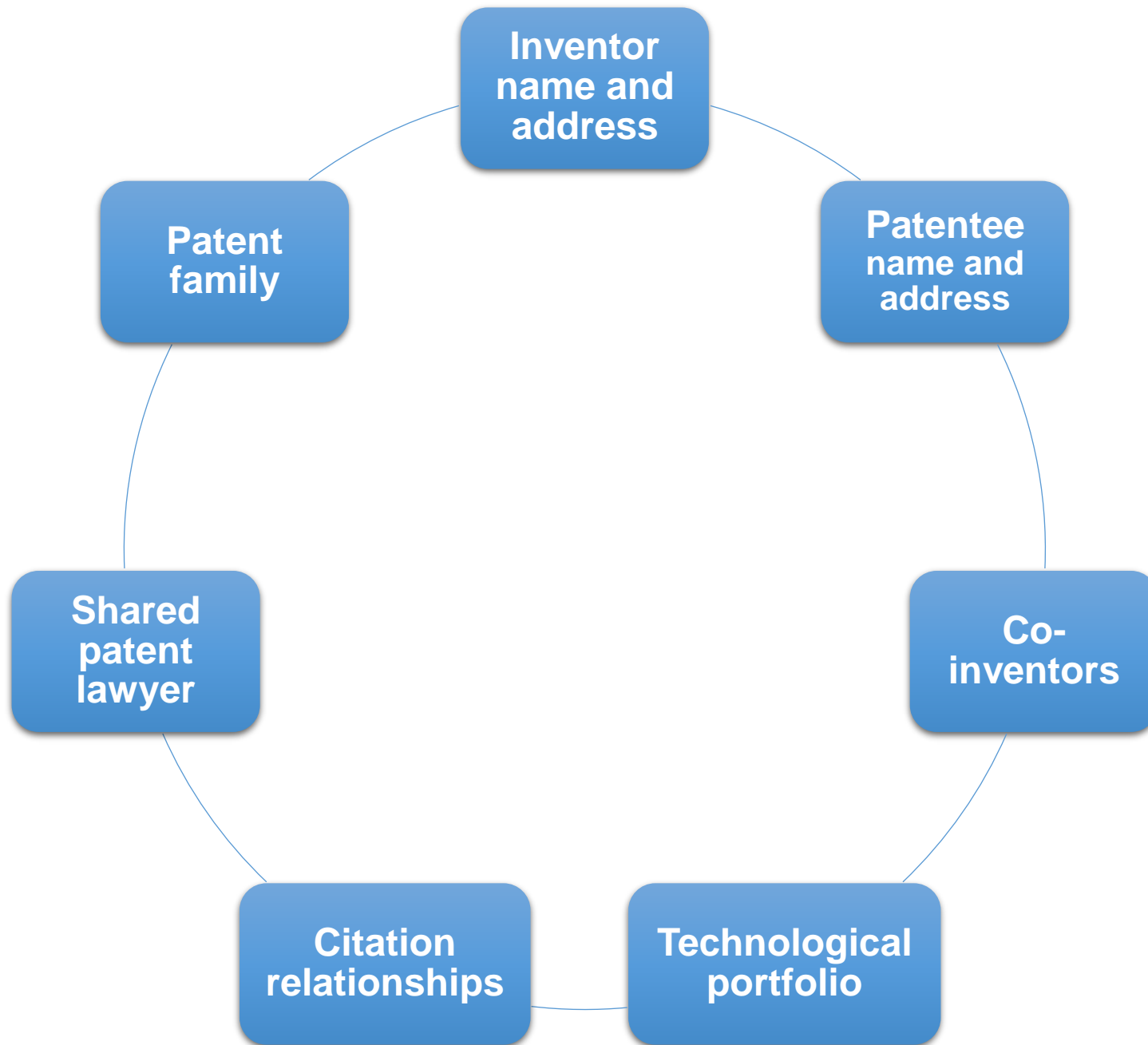




Model

Feature selection





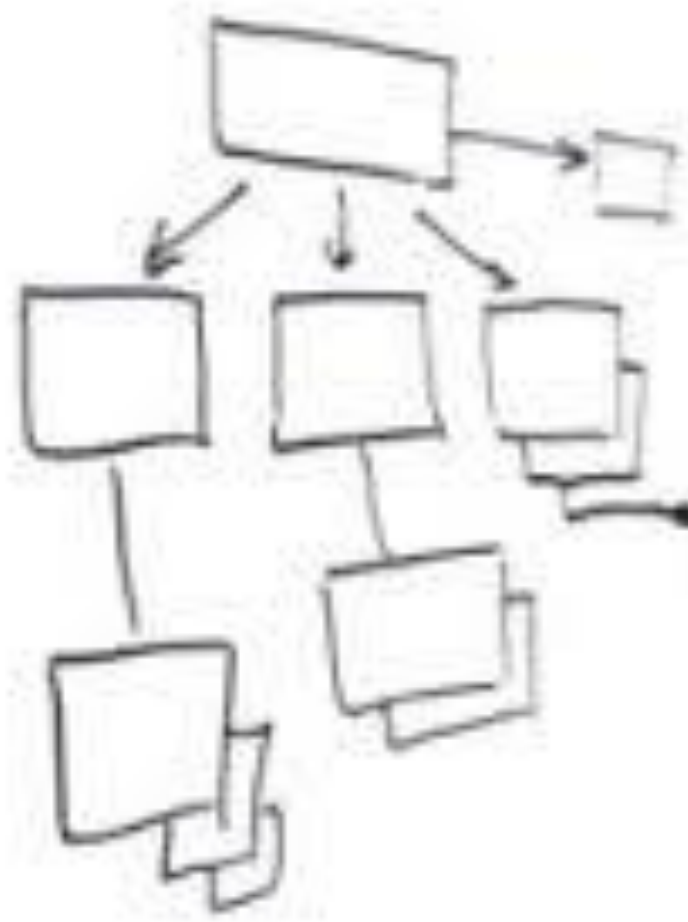


Share our insights

- We are convinced no complex models are needed; no complex (non-linear) relationships between features
- Compared to existing methods we think more fine grained control is needed for the classification process

Some examples of current practices that can be improved:

- Calculate similarities for all features, derive weighted sum, and classify if weighted sum exceeds threshold (difficult to completely compensate low outcomes on one feature with strong outcomes on other features)
- Present all features in a vector space and cluster based on cosine similarity measure (all features similarities get equal weight)

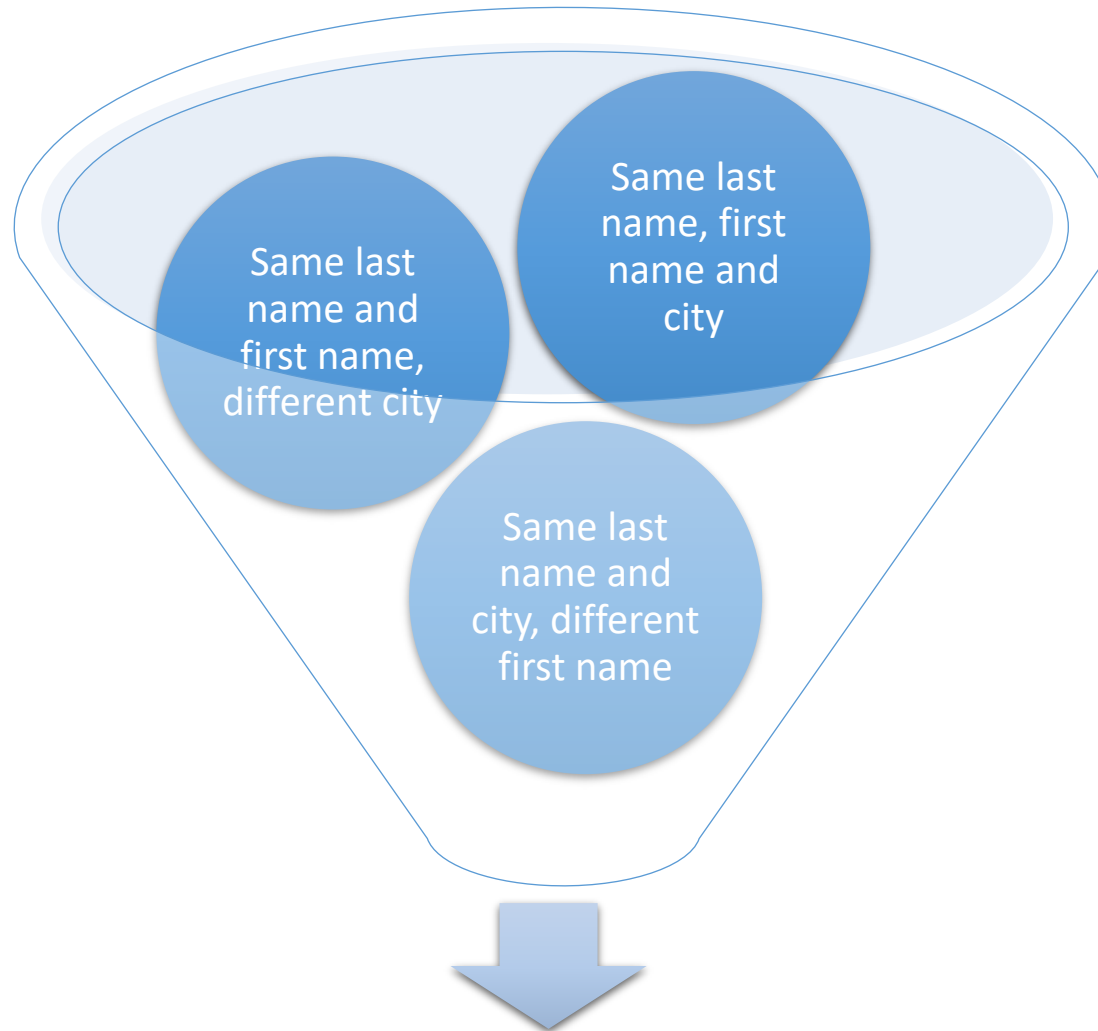




Share our insights

Non-inventor and patentee based information is crucial to assess cases (solve homonymy problem):

- Two inventors with same name, same address and different patentee: different person or same person that switched jobs?
- Two inventors with same name and different address: different person or same person that moved?



Cluster by last name, calculate all pairwise distances for the given similarities within each cluster, retain combinations with at least one positive similarity
=> 3.7 million combination records (inventors with the same last name and at least one positive similarity for the given similarities)



Lack of additional relevant similarities (co-inventor, patentee name, technology subclass, patent lawyer, US patent citation)

- In the subset of inventors with same last name and first name, but different patentee, 22% of the cases labeled as identical have no additional similarity information.
- How can we ever classify these cases?

Lack of validation data to identify false positives

- Deriving rules to merge inventors is not difficult, but where to stop?
- The number of found potential matches is far far far higher than the number of labeled cases. Current available labeled validation data does not allow to identify false positives.



Tryout

- Try out on set with same last name and different first name (allows quick identification of false positives on mismatched first names) having a co-inventor (strong indication of similarity): 110,260 cases
- Classification rules: Co-inventor, and at least one of following criteria: shared assignee; at least 2 shared technology subclasses; shared patent lawyer, at least 2 shared US patent citations: 36,970 matches
- Precision (based on 500 cases checked): about 90% (42 cases questionable, all on Ryde Niels-Peter and Ryde Tuula, both active on nanoparticles, but with same INPADOC family ID – can be corrected by adding criterion on first name).
- Recall (within subset): 70% (4,356 cases labeled identical, of which 1,344 missed by rule)

Difficult to improve: 478 of these cases have different inventor city, non of these linked to the same patent family, first names can match (mostly), and mostly title similarity. But dropping additional restrictions yields many false positives.



Tryout

- Obtaining high precision is not a problem: we can achieve almost 100% precision with some small changes to our basic classification rules
- Obtaining high recall is a challenge: we arrive at about 70% with basic classification rules, but difficult to improve, **not** because our model is too simple, **but because 20 to 30% of sampled inventors records labeled as identical in validation set have no additional similarity besides name:** no co-inventor, no shared patentee name, no shared patent lawyer, no significant shared technology or shared US or non-US patent citation (relevant = more than one shared technology subclasses or more than one shared US or non-US patent citation)

It is almost impossible to classify these cases correctly, even with (very) complex models (no information to know whether these are different persons or identical persons who moved and/or switched job)

Conclusions:

1. Homonymy problem (find out whether similar inventor names are identical persons or different persons that moved or switched jobs) is far more challenging than synonymy problem (identify related inventor names due to spelling variation and errors: “Van Der Bilt” / “Vanderbilt”; “Zimmermann” / “Zimmerman”);
2. Need to focus on homonymy problem, synonymy problem will be solved automatically (by using same criteria to solve homonymy problem);
3. (1) + (2) => Additional similarity information derived from patent portfolio is more important than direct name and address similarity (you cannot solve the homonymy problem with name and address similarity alone);
4. Additional derived similarity info is limited: co-inventor, technological portfolio, shared patentee, shared patent lawyer, shared patent citations;
5. No complex (non-linear) relationships between derived similarity;
6. (4) + (5) => No complex methods are needed, this is a simple linear classification problem;
7. Fine-grained control over classification rule is needed (some features have more importance than other features);
8. (6) + (7) => Use classic decision tree;

9. **20% to 30% of sampled inventor records labeled identical in validation set have no relevant additional relationship, which makes it almost impossible to classify these cases, regardless the complexity of methods used;**
10. **(3) + (9) => It is impossible to obtain high precision (95% and beyond) and high recall (95% and beyond) at the same time;**
11. **Our basic classification rules yield 99% precision and 70% recall, and this will be difficult to improve (complex methods or network based features/similarities seem to be of little value because of (9));**
12. **Impossible to increase recall without jeopardizing precision (you can start ignoring homonymy problem and classify all somewhat related inventors as identical, but this will inevitably result in false positives, hence lowering precision);**
13. **Current labeled validation data makes it impossible to quickly and correctly assess precision whenever method is adapted to increase recall (when recall is increased, many additional inventor records will be merged and labeled as identical person, however, there is no way to know whether these cases are correct, as these additional matches are not labeled in the dataset => false positives might remain undetected);**
14. **It is important that the result of any method is checked for precision manually and independently, not using any of the labeled validation sets, otherwise high precision rated might be claimed that are however biased and overestimated because of the limitations of the labeled validation sets (difficult to derive true precision because validation info is not exhaustive and do not contain explicit labeled negative cases).**

Directions to go:

- 1. Better feature selection is more relevant than more complex models;**
- 2. Lack of qualitative labeled validation data hampers assessment and development of methods (currently not possible to quickly and correctly derive precision);**
- 3. Most of current labeled validation data was not developed for developing name disambiguation. Collaborative approach needed to get better validation data, or publish OE labeled dataset (Akinsanmi);**
- 4. Street address information would benefit the classification process;**
- 5. Additional derived similarity data is crucial to solve homonymy problem;**
- 6. In that respect, content similarity (title, abstract, claims) is definitely a good candidate to start with;**
- 7. Including time dimension to inventor cluster results might help too, but will not solely solve the homonymy problem.**