

Hackathon "Les champs de Sirene"

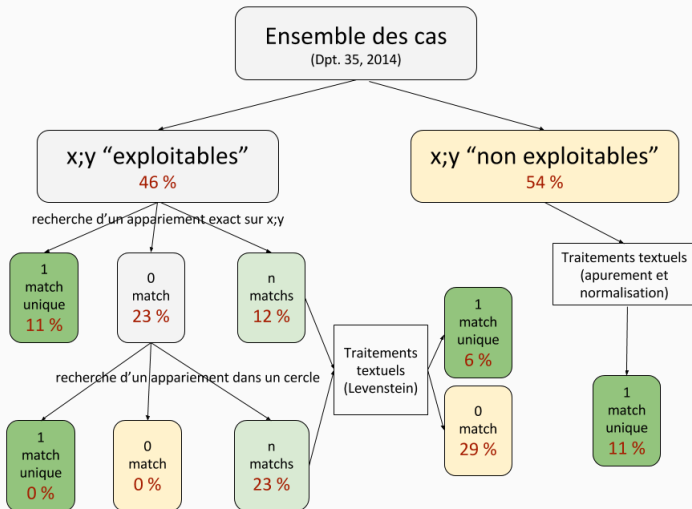
Restitution de l'équipe des géocodeurs

Cédric Couralet, Gaël de Peretti, Arlindo Dos Santos, Maëlle Fontaine, Juliette Fourcot, Joaquim Morais

Insee

- mobilisation en premier lieu des coordonnées géographiques
- prise en compte de la marge d'incertitude liée au processus de géocodage
- intérêt : limiter le nombre de cas à soumettre aux approches textuelles

Démarche



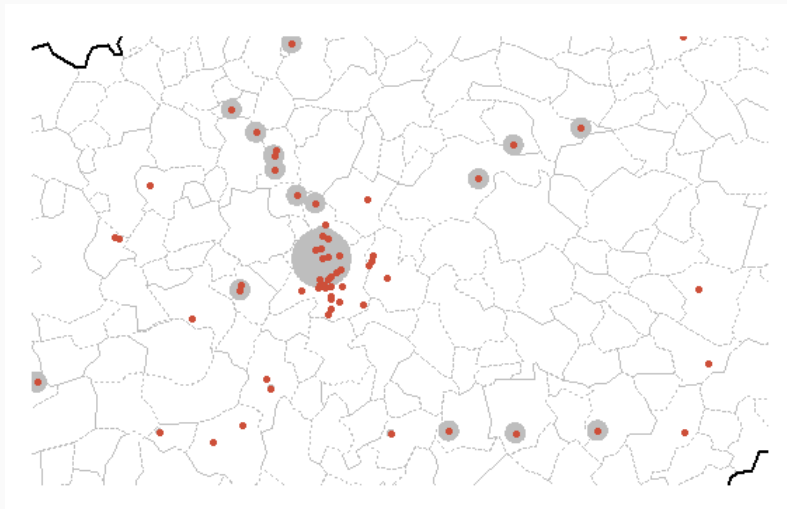


Figure 1: Exemples de cercles de rayons variables en fonction de la qualité du x,y dans Sirius

- Apurement de la raison sociale déclarée (mentions faites à la forme juridique, au nom de la commune) : permet d'augmenter les possibilités de matching sur cette variable
- Calcul de la distance de Damerau–Levenshtein entre **RS_X** et les champs Sirius pour accepter celle en dessous d'un certain seuil
⇒ Accepte les fautes de frappes, et permet de trouver certaines entreprises individuelles (le nom de famille apparait dans l'adresse)
- Si plusieurs résultats, prise en compte du n° et libellé de voie
- Si pas de résultats, ne pas prendre en compte **RS_X** mais l'activité et l'adresse

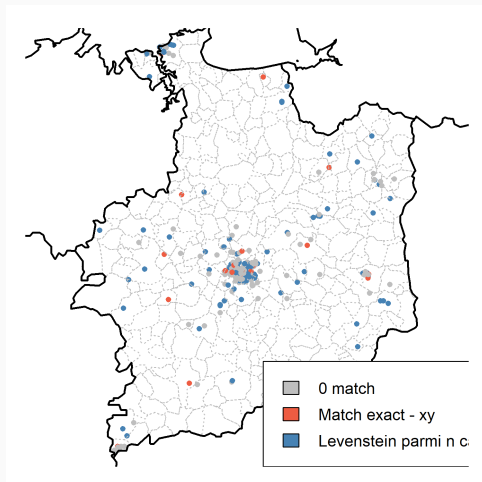


Figure 2: Exemples de cas à problème en Ille-et-Vilaine

- Une mobilisation du $x;y$ peut conduire à de nouveaux matchs et réduire fortement la combinatoire lors de l'identification textuelle ...
- ... mais cela repose sur une bonne qualité de celui-ci (ou une bonne caractérisation de la marge d'erreur).
- Prolongements possibles :
 - parmi les n matchs possibles, en éliminer en utilisant l'activité, la commune ...
 - recherche de $x;y$ alternatifs pour les $x;y$ "non-exploitable" ou ceux ne satisfaisant pas certains contrôles (sur la commune par exemple)