

SPELL-CHECKING QUERIES BY COMBINING LEVENSHTTEIN AND STOILOS DISTANCES

**Zied Moalla^{1, 2}, Lina F. Soualmia^{1, 3}, Élise Prieur-Gaston¹
Thierry Lecroq¹, Stéfan J. Darmoni¹**

¹ CISMéF, Rouen University Hospital & TIBS, LITIS EA 4108, University of Rouen, France

² MIRACL, Sfax University, Tunisie

³ LIM&Bio EA 3969, Sorbonne Paris Cité, France



NETTAB 2011

Clinical Bioinformatics
October 12-14, 2011, Pavia, Italy

Content

- Context
- Introduction
- Materials and methods
 - Levenshtein distance
 - Stoilos distance
- Results
- Conclusion
- Perspectives

Context

CiSMeF Catalogue et Index des Sites Médicaux de langue Française

Doc'CISMeF
Outil de recherche en médecine

Avancée

tous types et documents
 uniquement les recommandations professionnelles
 uniquement les documents concernant l'enseignement
 uniquement les documents et associations concernant les patients

Listes et index : alphabétique, thématique, types de ressources

Portail Terminologique de Santé
Consulter le MeSH et les autres terminologies de santé

Connexion S'inscrire

Ce catalogue s'adresse en priorité aux professionnels de santé. On y trouve également des informations destinées aux patients et à leurs familles. Pour toutes remarques ou questions, veuillez nous contacter. voir aussi protection des données personnelles.

Catalog & Index of Health Resources in French on the Internet

CISMef = quality controlled health gateway for French institutional health resources

Doc'CISMeF: a search tool

- to search within the catalog CISMef more than 82,000 documents
- specific of the health resources available on the Internet, such as association, patient information, community networks

3 types of users:

- Patients
- Students
- Clinicians

Introduction

- Increase in the number of users querying different search engines
- Internet became a major source of health information
- Medical vocabularies are difficult to handle by non-professionals
- *"Did you mean:" of Google or "Also try:" of Yahoo*

Introduction

- **Purpose:** Spelling correction for medical queries in French.
- **Method:** Spelling correction based on comparing the query with a dictionary.
- **Tools:** The string distance of Stoilos and the Levenshtein edit distance to correct spelling errors. We propose here to combine them.

String distance: Levenshtein

- Minimum number of edit operations (insertion, deletion, substitution) to transform one string into the other

String distance: Levenshtein

- The Normalized Levenshtein (*LevNorm*) in the range [0, 1]

$$\text{LevNorm}(c_1, c_2) = \frac{\text{Lev}(c_1, c_2)}{\text{Max}(\text{length}(c_1), \text{length}(c_2))}$$

- Example :

$\text{LevNorm}(\text{Trigonocepahlie}, \text{Trigonocephalie}) = 2/15 = 0.133$

$\text{Lev}(\text{Trigonocepahlie}, \text{Trigonocephalie}) = 2$

$\text{max}(\text{length}(\text{Trigonocepahlie}), \text{length}(\text{Trigonocephalie})) =$
 $\text{max}(15, 15) = 15$

String distance: Stoilos

- The similarity among two entities is related to their commonalities as well as to their differences. Thus, the similarity should be a function of both these features.

$$Sto(s_1, s_2) = Comm(s_1, s_2) - Diff(s_1, s_2) + Winkler(s_1, s_2)$$

String distance: Stoilos

- The function of commonality computes the longest common substrings between 2 strings

$$Comm(s_1, s_2) = \frac{2 \times \sum_i \text{length}(\text{MaxComSubString}_i)}{\text{length}(s_1) + \text{length}(s_2)}$$

- Example: $s_1 = \text{'Trigonocephalie'}$ et $s_2 = \text{'Trigonocephalie'}$

$\text{length}(\text{MaxComSubString}_1) = \text{length}(\text{Trigonocep}) = 10$

$\text{length}(\text{MaxComSubString}_2) = \text{length}(\text{lie}) = 3$

$Comm(\text{Trigonocepahlie}, \text{Trigonocephalie}) = 13/15 = 0.866$

String distance: Stoilos

- Based on the length of the unmatched strings that have resulted from the initial matching step

$$Diff(s_1, s_2) = \frac{uLen_{s_1} \times uLen_{s_2}}{p + (1-p) \times (uLen_{s_1} + uLen_{s_2} - uLen_{s_1} \times uLen_{s_2})}$$

$s_1 =$ 'Trigonocep~~ah~~lie' and $s_2 =$ 'Trigonocep~~hal~~ie' and $p = 0.6$

$uLen_{s_1} = 2/15$ and $uLen_{s_2} = 2/15$

So $Diff(s_1, s_2) = 10/787 = 0.0254$

String distance: Stoilos

- The Winkler correction:

$$\text{Winkler}(s_1, s_2) = L \times p' \times (1 - \text{Comm}(s_1, s_2))$$

$s_1 = \text{'Trigonocepahlie'}$ and $s_2 = \text{'Trigonocephalie'}$

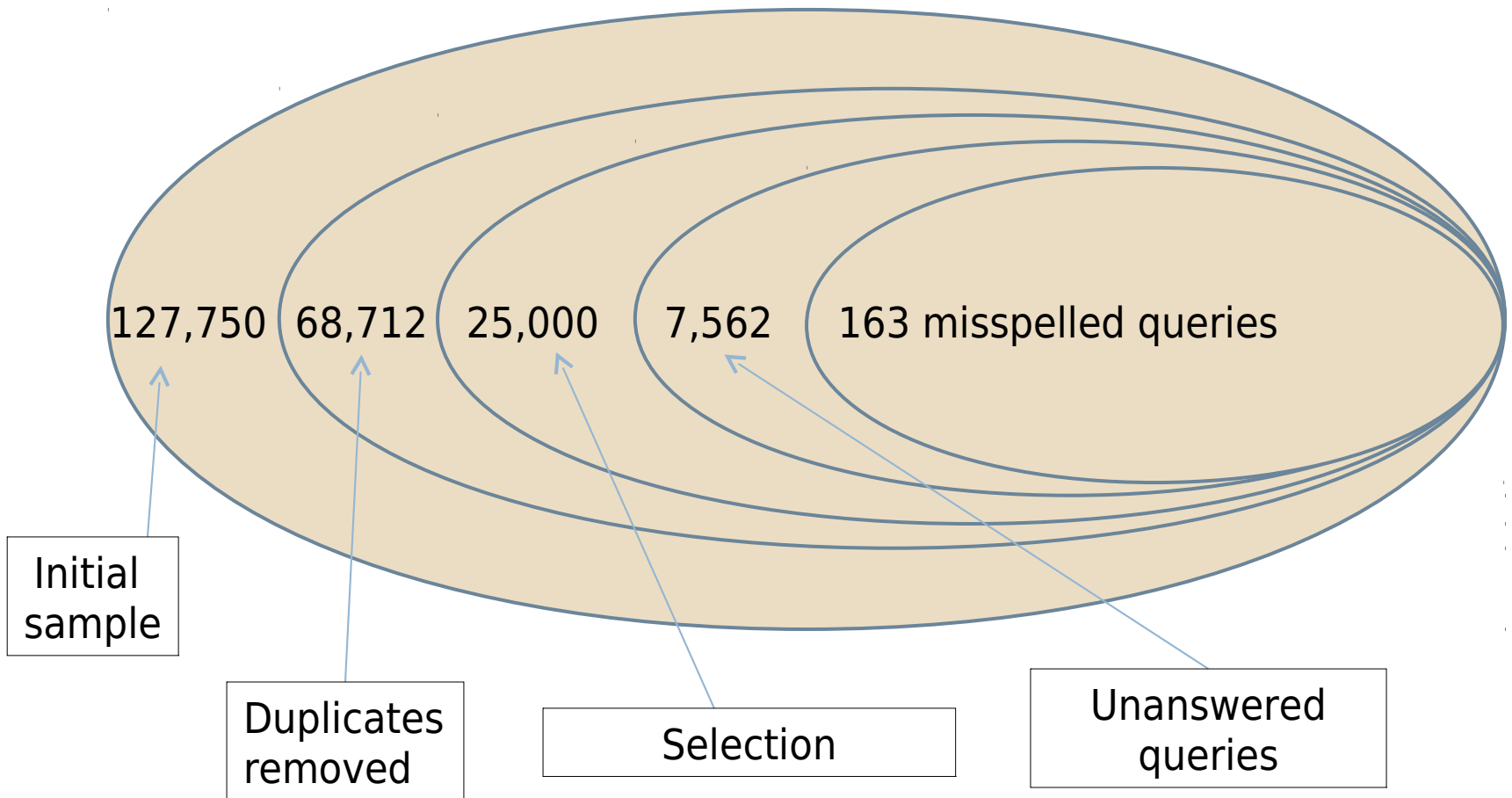
$L = 4$ and $p' = 0.1$

So $\text{Winkler}(s_1, s_2) = 4/75 = 0.053$

- Altogether

$$\begin{aligned} \text{Sto}(\text{Trigonocepahlie}, \text{Trigonocephalie}) &= 13/15 - 10/787 + 4/75 \\ &= 0.894 \end{aligned}$$

Materials: Queries



Choice of thresholds

Levenshtein and Stoilos string distances require a choice of thresholds to obtain a manageable number of propositions of correction to the user.

So we have tested this number for **163 misspelled queries**.

	Method							
	Levenshtein			Stoilos			Levenshtein & Stoilos	
Thresholds	<0.2	<0.1	<0.05	>0.7	>0.8	>0.9	Lev < 0.2 Stoilos > 0.8	Lev < 0.2 Stoilos > 0.7
Nb of answers	224 1.37	76 0.46	8 0.04	1454 8.92	489 3	140 0.85	179 1.09	213 1.30

Evaluation

$$\text{Recall} = \frac{\text{Queries correctly corrected}}{\text{Queries}}$$

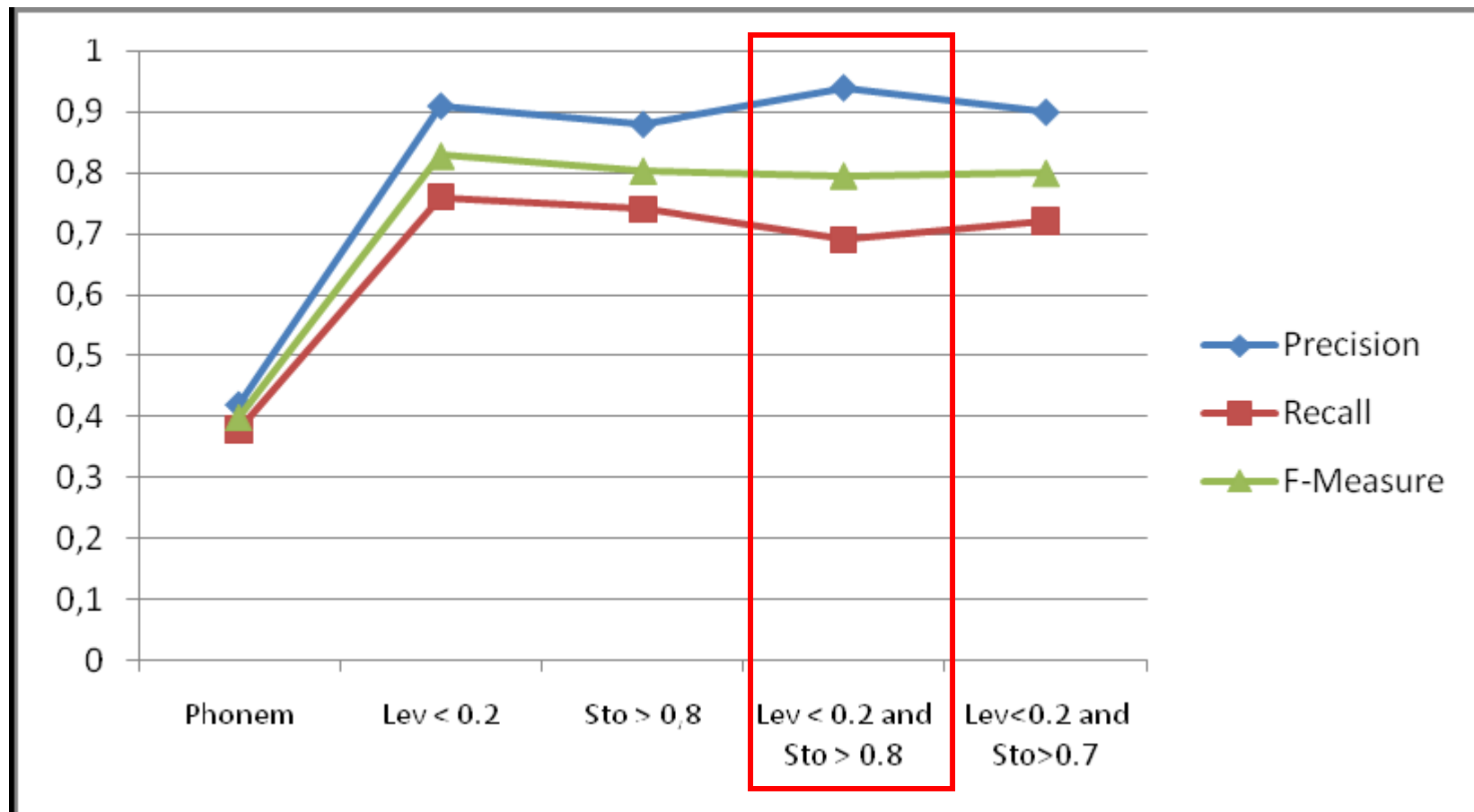
$$\text{Precision} = \frac{\text{Queries correctly corrected}}{\text{Queries corrected}}$$

$$\text{F-Measure} = \frac{2 \times \text{Precision} \times \text{Recall}}{\text{Precision} + \text{Recall}}$$

Results

Method	Recall	Precision	F-Measure
Phonetic transcription	0.38	0.42	0.399
Levenshtein < 0.2	0.76	0.91	0.8283
Stoilos > 0.8	0.74	0.88	0.8039
Levenshtein < 0.2 & Stoilos > 0.8	0.69	0.94	0.7958

Evaluation



Conclusion

- A method to automatically correct misspelled queries submitted to health search tool
- The combination of the 2 distances gives a recall of 69% and a precision of 94%
- This combination has increased the precision, but decreased the recall
- The functionality is implemented in CISMeF

Perspectives

- Misspelled queries categorized according to their number of words
- The configuration of a keyboard, by studying the distances between keys