

Analyses statistiques pour l'évaluation des systèmes de recherche d'information

S. Déjean, J. Mothe

www.math.univ-toulouse.fr/~sdejean

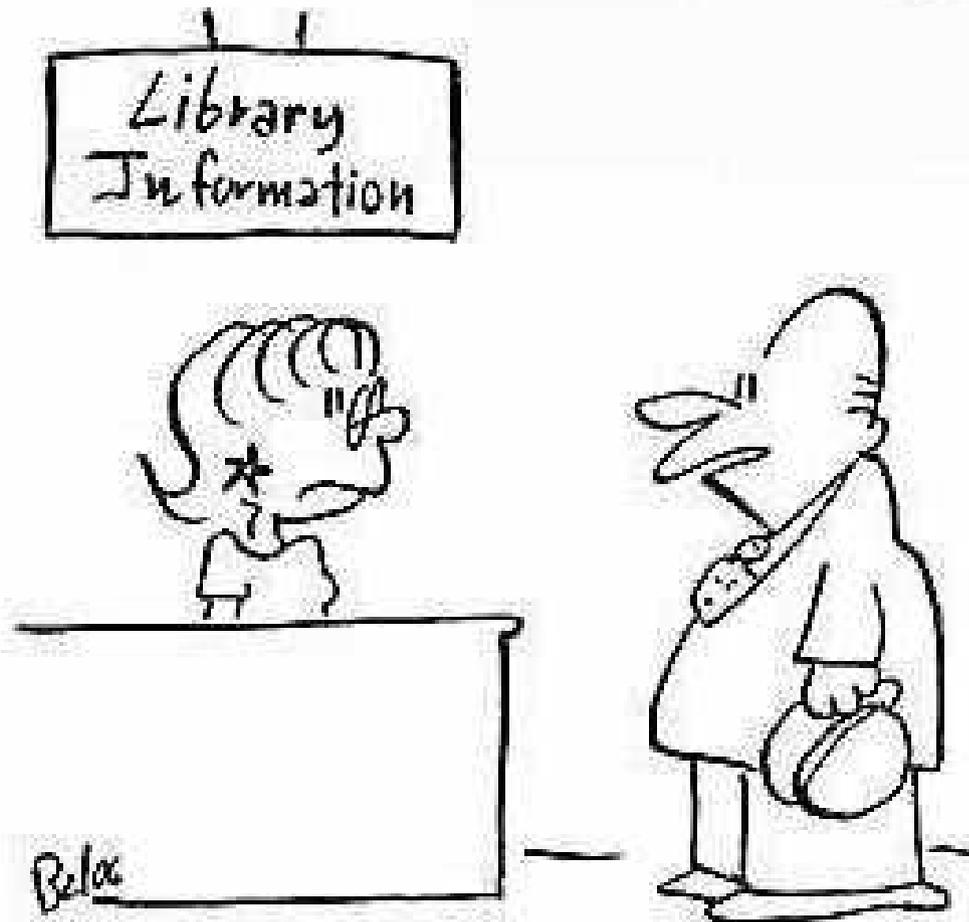


Recherche d'information



<http://boston.lti.cs.cmu.edu/classes/11-744/>

Recherche d'information

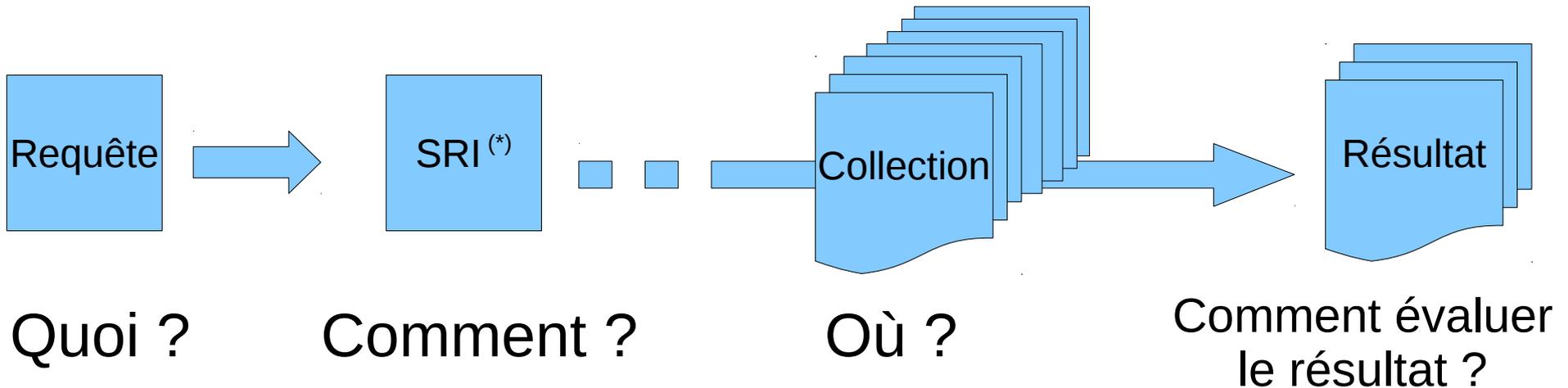


"I CAN'T FIND THE BOOKS ON
INFORMATION RETRIEVAL."

<https://oncealibrarian.wordpress.com/2012/03/12/62/>

Recherche d'information

Recherche d'Information :
domaine de la recherche qui s'intéresse à la structure, à l'analyse, à l'organisation, au stockage, à la recherche et à la découverte de l'information (G. Salton).



(*) SRI : Système de Recherche d'Information

Partie I

Mesures de performance

 A. Baccini, S. Déjean, L. Lafage, J. Mothe (2011). **How many performance measures to evaluate Information Retrieval Systems?** *Knowledge and Information System*, DOI 10.1007/s10115-011-0391-7

TREC : Text REtrieval Campaign

- Tâches : adhoc (évaluation de performances des systèmes), *filtering track* (filtrage d'informations), *cross-language track* (RI multilingue), *web track*, *question answering*...
- Requête^(*) : titre, description, partie narrative
- Les « bons » documents sont connus a priori
- *Trec_eval* : outil de calcul des performances des différents systèmes (« pertinence » des documents retournés)

<title> **Topic:** Corporate Pension Plans/Funds

(*) *Exemple de requête*

<desc> **Description:**

Document will report on problems associated with pension plans/funds such as fraud, skimming, tapping or raiding.

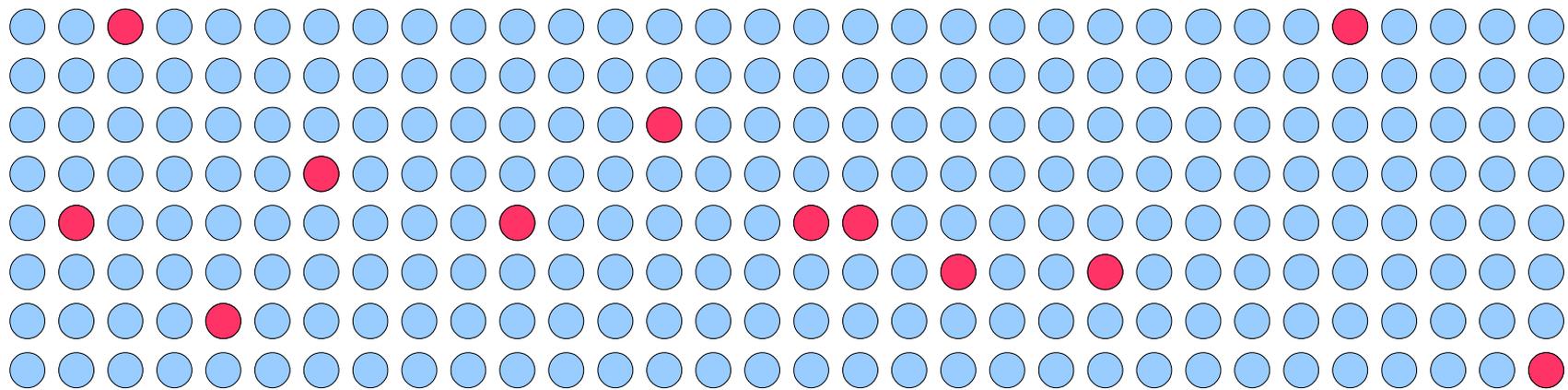
<narr> **Narrative:**

A relevant document will report on problems associated with pension plans/funds and also U.S. Government regulatory controls on pension plans. Examples of problems that are considered relevant are fraud, skimming, tapping or raiding.

Mesures de performance : rappel et précision

Rappel : nombre de documents pertinents retournés / nombre total de documents pertinents
Précision : nombre de documents pertinents retournés / nombre de documents retournés

Collection



● Document pertinent pour une requête donnée

Résultat du SRI

• 5 premiers résultats : ● ● ● ● ●

Rappel = 3/12, Précision = 3/5

• 10 premiers résultats : ● ● ● ● ● ● ● ● ● ●

Rappel = 5/12, Précision = 5/10

• ...

Trec_eval : 135 mesures de performance

http://www-nlpir.nist.gov/projects/trecvid/trecvid.tools/trec_eval_video/README

trec_eval is the standard tool used by the TREC community for evaluating an ad hoc retrieval run, given the results file and a standard set of judged results

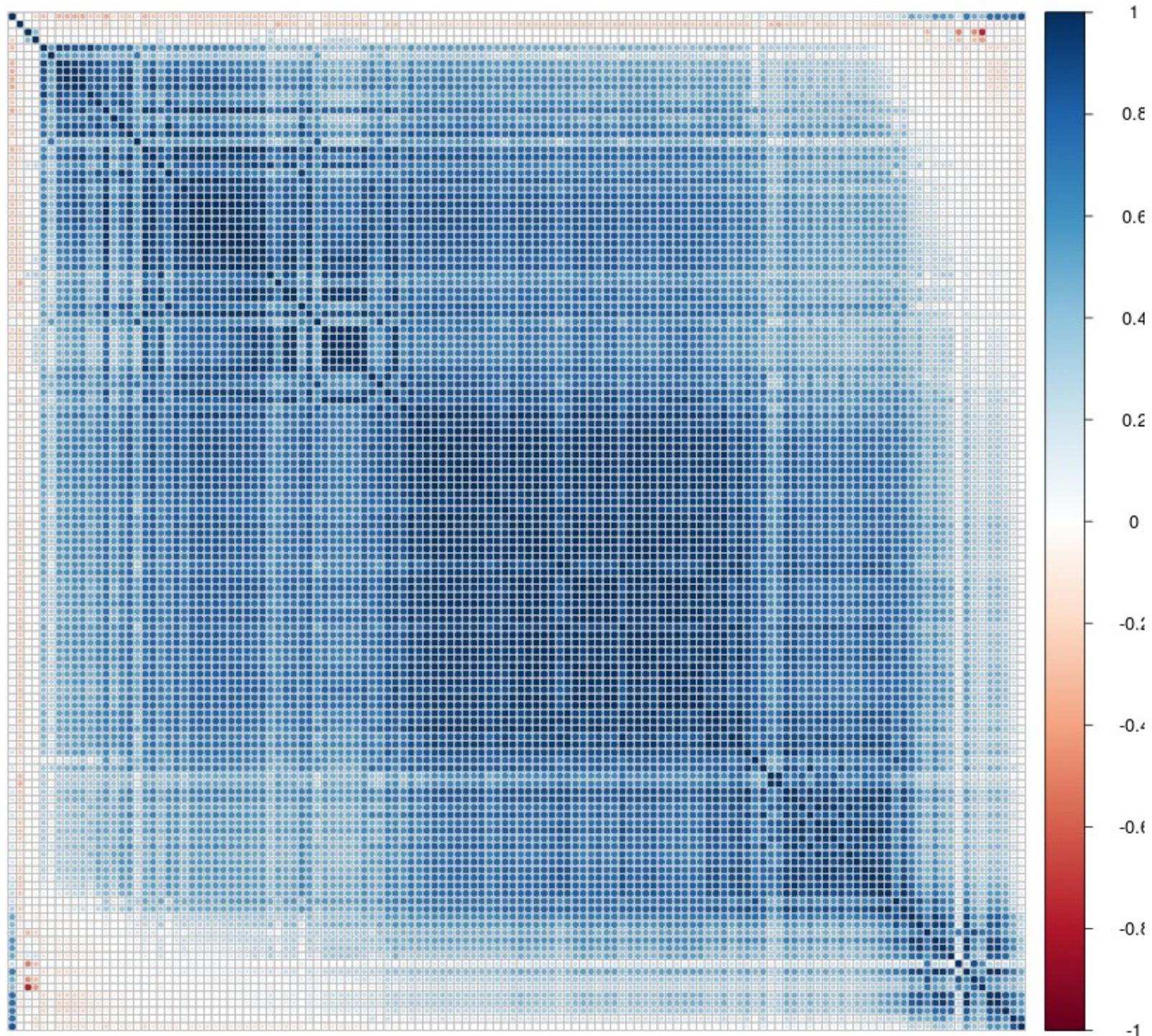
| | |
|--------------------------------|--|
| <i>map</i> | Mean Average Precision (MAP) |
| <i>gm_ap</i> | Average Precision. Geometric Mean, $q_score = \log(\text{MAX}(\text{map}, .00001))$ |
| <i>R-prec</i> | R-Precision (Precision after R (= num-rel for topic) documents retrieved) |
| <i>bpref</i> | Binary Preference, top R judged nonrel |
| <i>recip_rank</i> | Reciprical rank of top relevant document |
| <i>ircl_prn.0.00</i> | Interpolated Recall - Precision Averages at 0.00 recall |
| [...] | |
| <i>ircl_prn.1.00</i> | Interpolated Recall - Precision Averages at 1.00 recall |
| <i>P5</i> | Precision after 5 docs retrieved |
| [...] | |
| <i>P1000</i> | Precision after 1000 docs retrieved |
| <i>exact_prec</i> | Exact Precision over retrieved set |
| <i>exact_recall</i> | Exact Recall over retrieved set |
| <i>3-pt_avg</i> | Average over 3 points of recall-precision graph |
| <i>avg_doc_prec</i> | Rel doc precision averaged over all relevant docs (NOT over topics) |
| <i>exact_relative_prec</i> | Exact relative precision |
| <i>avg_relative_prec</i> | Average relative precision |
| <i>exact_unranked_avg_prec</i> | Exact Unranked Average Precision |
| <i>map_at_R</i> | Average Precision over first R docs retrieved |
| <i>int_map</i> | Interpolated Mean Average Precision |
| <i>exact_int_R_rcl_prec</i> | Exact R-based-interpolated-Precision |
| <i>int_map_at_R</i> | Average Interpolated Precision for first R docs retrieved |
| <i>bpref_topnonrel</i> | Binary Preference, top 100 judged nonrel |
| <i>bpref_top5Rnonrel</i> | Binary Preference, top 5R judged nonrel |
| <i>bpref_top25p2Rnonrel</i> | Binary Preference, top 25 + 2*R judged nonrel |
| <i>bpref_retail</i> | Binary Preference, Only retrieved judged rel and nonrel |
| <i>bpref_5</i> | Binary Preference, top 5 rel, top 5 nonrel |
| <i>bpref_10</i> | Binary Preference, top 10 rel, top 10 nonrel |
| <i>bpref_num_all</i> | Binary Preference, Number not retrieved before (all judged) |
| <i>bpref_num_ret</i> | Binary Preference, Number retrieved after |
| <i>recall5</i> | Recall after 5 docs retrieved |
| [...] | |

- Un extrait des 135 mesures fournies par *trec_eval*
- La plupart sont des combinaisons de Rappel et Précision

Redondance des mesures

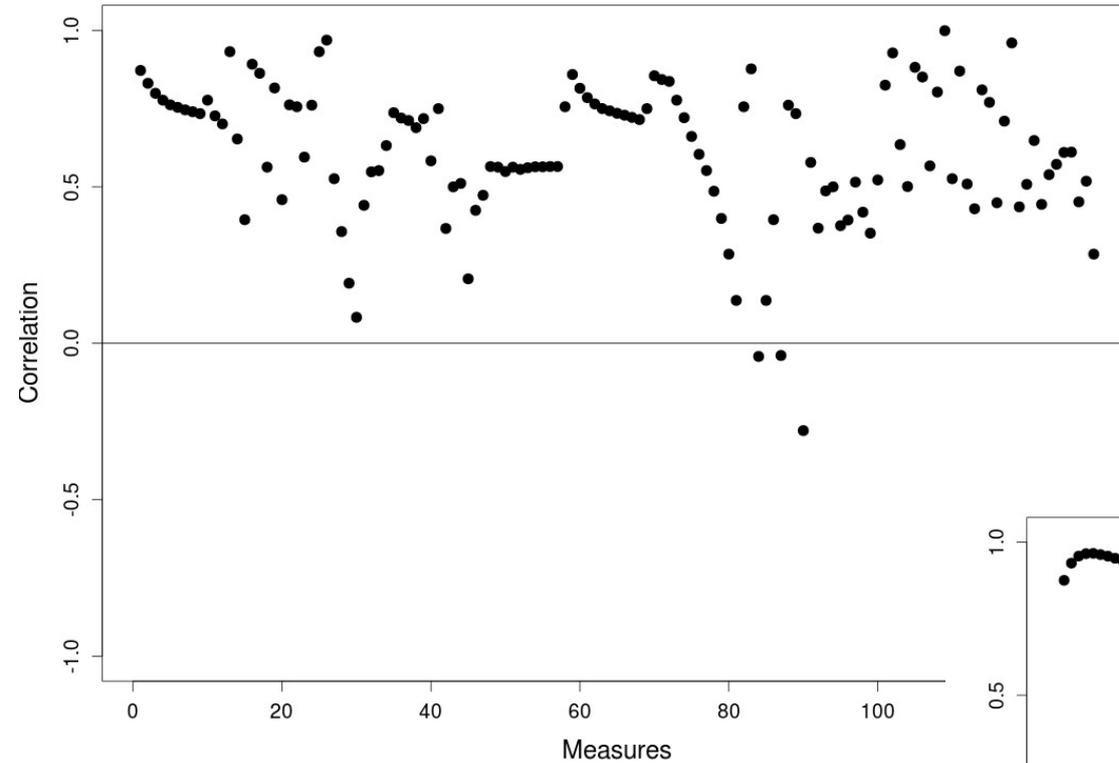
- L'interprétation globale des 135 mesures est délicate et pas forcément pertinente
- Certaines mesures sont, par nature, redondantes (corrélées) : P5, P10, P15...
- La caractérisation d'un SRI doit pouvoir se faire à partir d'un nombre restreint de mesures
- Étude réalisée sur des données issues de 7 campagnes TREC : 23 518 runs (requête × système) pour lesquels on dispose des 135 mesures

Matrice de corrélation



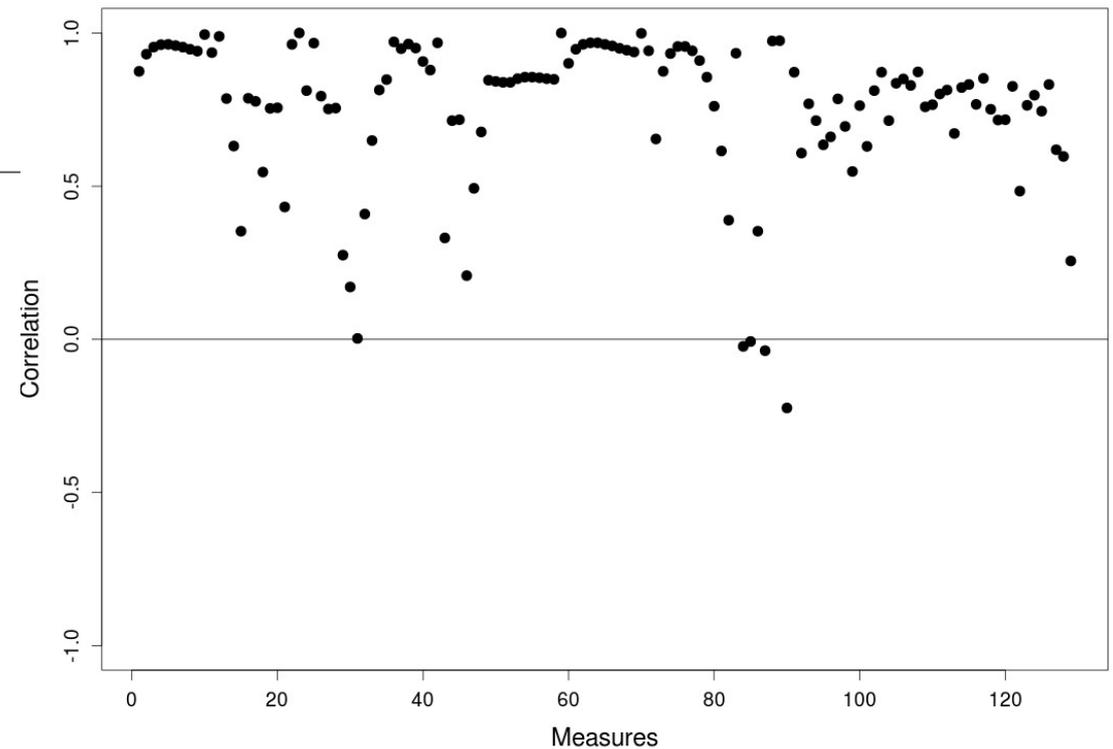
Corrélations (suite)

P5



Représentation des corrélations entre 2 indicateurs particuliers (P5 et MAP) et l'ensemble des autres indicateurs.

map



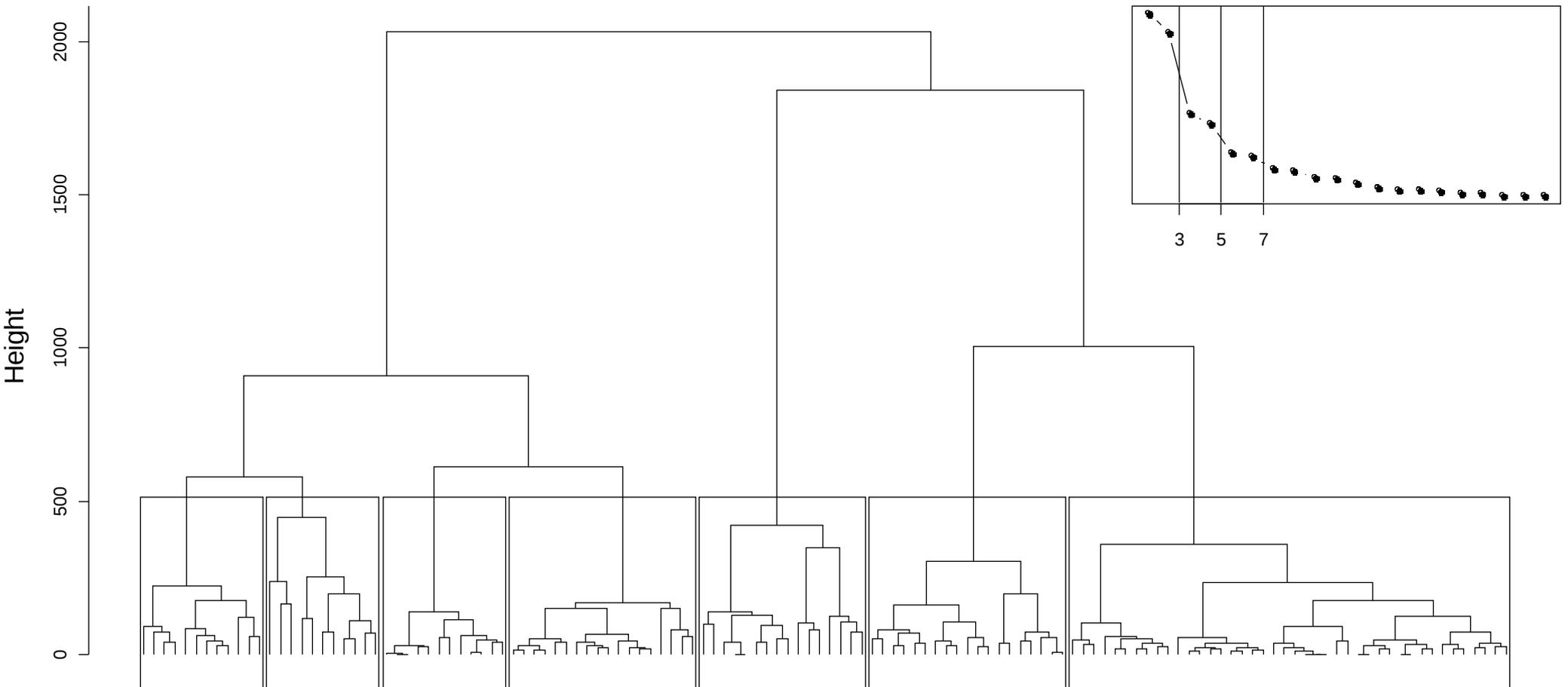
Analyse multivariée

- **Objectifs :**
 - Identifier des groupes de mesures homogènes
 - Définir un représentant par groupe
- **Méthodes :**
 - Classification hiérarchique + k-means
 - Heuristique pour le choix de représentant
- **Validation :**
 - Comparaison de classements élaborés soit avec l'ensemble complet des mesures, soit avec un ensemble restreint

Analyse multivariée

- **Objectifs :**
 - Identifier des groupes de mesures homogènes
 - Définir un représentant par groupe
- **Méthodes :**
 - Classification hiérarchique + k-means
 - Heuristique pour le choix de représentant
- **Validation :**
 - Comparaison de classements élaborés soit avec l'ensemble complet des mesures, soit avec un ensemble restreint

Classification des mesures de performance



Classification hiérarchique (distance euclidienne + critère de Ward) + k-means
→ 7 groupes homogènes (dont un contenant des mesures non informatives :
nombre de documents retournés par exemple)

Interprétation des groupes

Cluster 1 (23 measures)

relative unranked avg prec30 - relative unranked avg prec20 - relative prec30 - map at R - relative unranked avg prec15 - relative prec20 - P10 - relative prec15 - int 0.20R.prec - relative unranked avg prec10 - X0.20R.prec - ircl prn.0.10 - P20 - P15 - relative prec10 - bpref 10 - P10 - relative unranked avg prec5 - relative prec5 - P5 - bpref 5 - recip rank - ircl prn.0.00

Précision avec un faible nombre de documents retournés

Cluster 2 (16 measures)

P100 - P200 - unranked avg prec500 - unranked avg prec1000 - bpref num ret - P500 - bpref num all - P1000 - unranked avg prec1000 - exact relative unranked avg prec - bpref num possible - int 1.0 - .1.0 0.0 0.0 - exact relative unranked avg prec - bpref num possible

Précision avec un nombre élevé de documents retournés

Cluster 3 (12 measures)

bpref top10Rnonrel - bpref retnonrel - relative unranked avg prec500 - avg relative prec - recall500 - relative prec1000 - bpref top10Rnonrel - relative unranked avg prec1000 - exact relative prec - recall2000 - relative prec1000 - exact relative prec

Rappel avec un nombre élevé de documents retournés

Cluster 4 (45 measures)

X1.20R.prec - ircl prn.0.30 - X1.40R.prec - int map - X1.00R.prec - R.prec - int 1.20R.prec - exact int R rcl prec - int 1.00R.prec infAP - avg doc prec - map - X11.pt avg - X1.60R.prec - int 0.80R.prec - int 1.40R.prec - X0.80R.prec - old bpref top10pRnonrel - ircl prn.0.40 - X1.80R.prec - int 1.60R.prec - X6.pt avg bpref - X2.00R.prec - bpref top25p2Rnonrel - old bpref - bpref top10pRnonrel - int 1.80R.prec - int 0.60R.prec - int 2.00R.prec - bpref top25pRnonrel - X0.60R.prec - bpref top50pRnonrel - bpref top5Rnonrel - ircl prn.0.20 - ircl prn.0.50 - int 0.40R.prec - X0.40R.prec - int map at R - ircl prn.0.60 - unranked avg prec30 - ircl prn.0.70 - ircl prn.0.80 - unranked avg prec200 - unranked avg prec100

Mesures « moyennes »

Cluster 5 (18 measures)

bpref topnonrel - fallout recall 42 - fallout recall 28 - fallout recall 56 - rcl at 142 nonrel - fallout recall 71 - fallout recall 85 - relative unranked avg prec100 - fallout recall 99 - fallout recall 113 - relative prec100 - fallout recall 127 - relative unranked avg prec200 - fallout recall 142 - recall100 - relative prec200 - recall200 - bpref retall

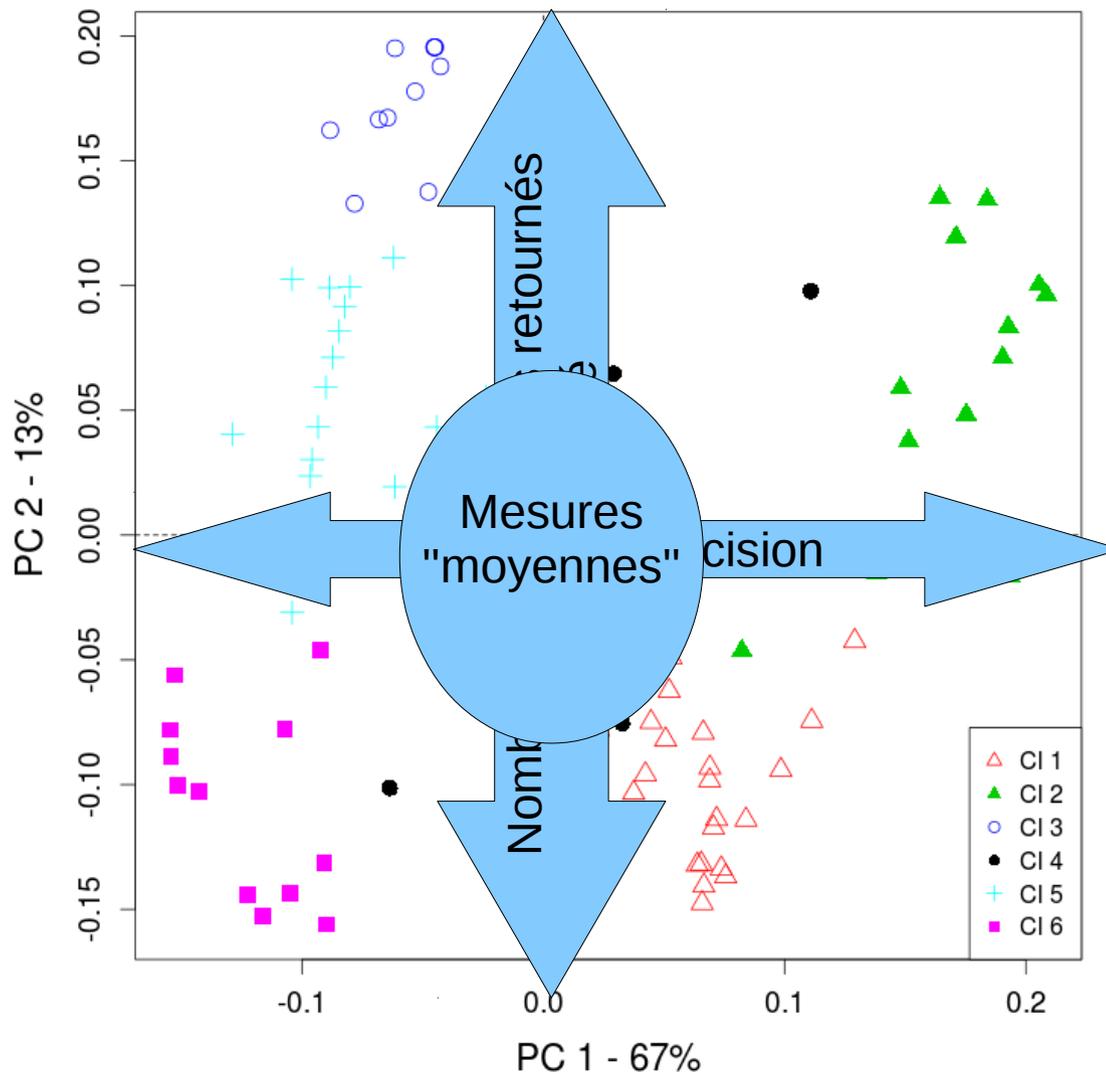
Rappel (???)

Cluster 6 (13 measures)

fallout recall 14 - unranked avg prec10 - unranked avg prec5 - ircl prn.0.60 - recall100 - recall10 - unranked avg prec10 - recall130 - ircl prn.1.00 - recall120 - recall115 - unranked avg prec5 - recall110 - recall15

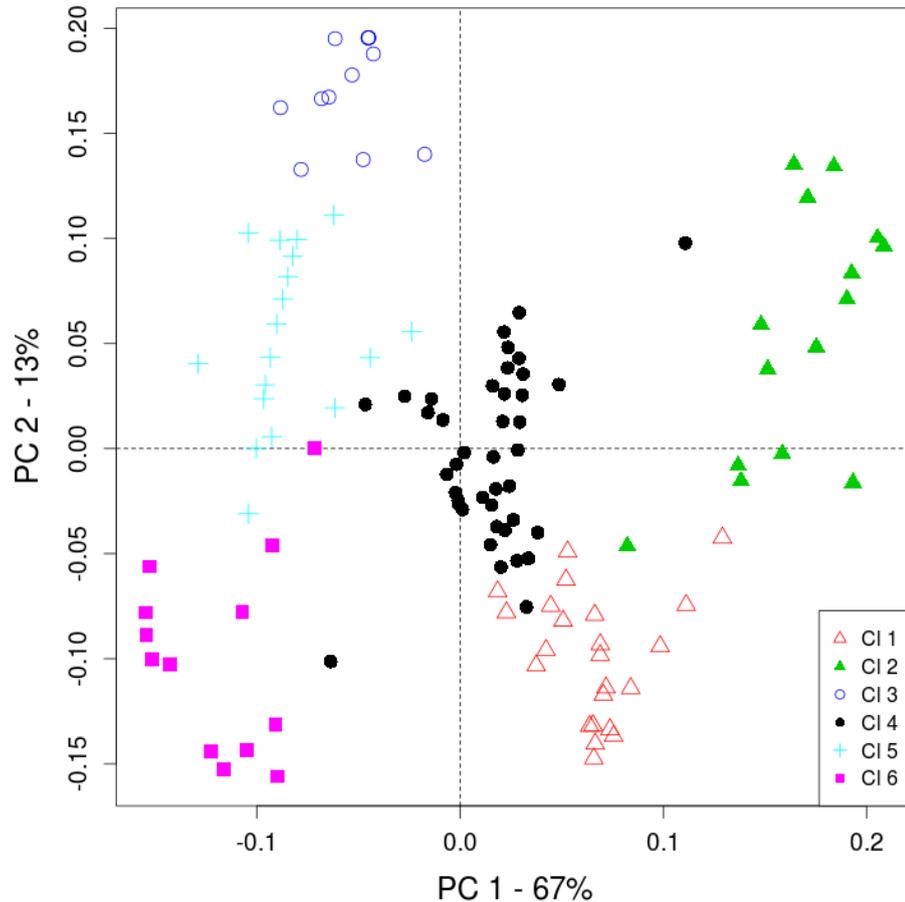
Rappel avec un faible nombre de documents retournés

Analyse en Composantes Principales



△ CI 1 Précision avec un faible nombre de documents
▲ CI 2 Précision avec un nombre élevé de documents
○ CI 3 Rappel avec un nombre élevé de documents
● CI 4 Mesures « moyennes »
+ CI 5 Rappel (???)
■ CI 6 Rappel avec un faible nombre de documents

Choix d'un représentant



- Les groupes apparaissent homogènes et relativement bien distincts les uns des autres
- Proposition : le représentant d'un groupe doit occuper une position « centrale » dans son groupe
- Heuristique :
 - 1) Calculer le barycentre de chaque groupe
 - 2) Calculer la distance de chaque mesure du groupe au barycentre
 - 3) Prendre comme représentant la mesure la plus proche du barycentre
- Liste proposée :
P30 - P100 - Exact recall -
MAP - bpref retail - recall130

Analyse en Composantes Principales –
Représentation des mesures

Vers un classement des systèmes

$$\begin{pmatrix} X_{11} & \dots & X_{1p} \\ \cdot & & \cdot \\ \cdot & & \cdot \\ \cdot & & \cdot \\ X_{n1} & \dots & X_{np} \end{pmatrix}$$

X_{ij} : score obtenu par le système i pour l'indice de performance j .

$i = 1$ à n , nombre de systèmes

$j = 1$ à 135 (version complète)

$j = 1$ à 6 (version réduite)

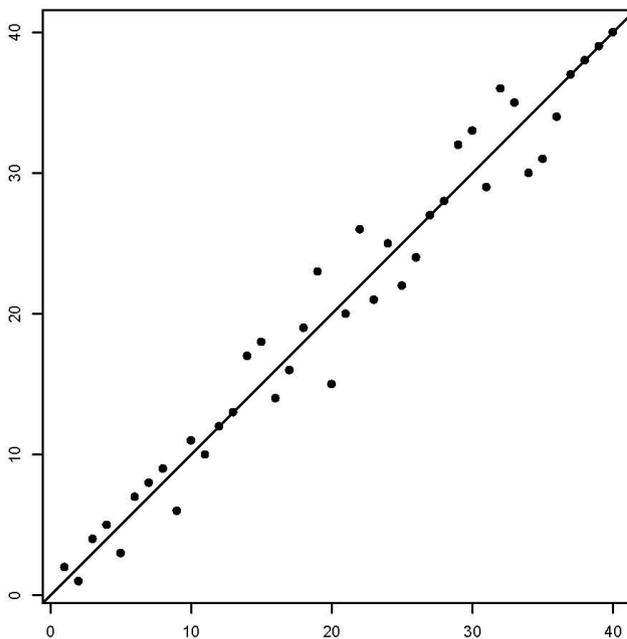
$$\begin{array}{c} \overline{X}_1 \\ \cdot \\ \cdot \\ \cdot \\ \cdot \\ \overline{X}_n \\ \underbrace{\hspace{1.5cm}} \\ \text{score moyen} \end{array}$$

$$\begin{array}{c} R_1 \\ \cdot \\ \cdot \\ \cdot \\ \cdot \\ R_n \\ \underbrace{\hspace{1.5cm}} \\ \text{rang basé sur} \\ \text{le score moyen} \end{array}$$

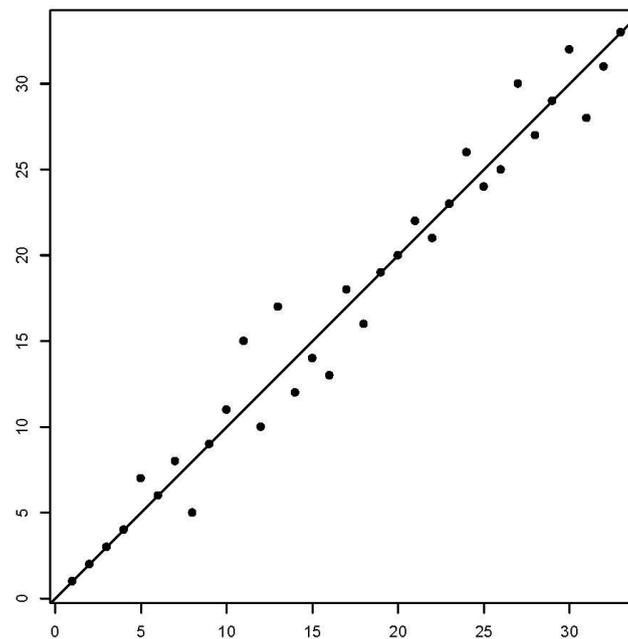
Classement de systèmes

Illustration : pour 3 requêtes particulières (#192 – TREC3, #246 – TREC4, #372 – TREC7), chaque point a pour coordonnées son rang selon le classement obtenu à partir de l'ensemble complet (horizontal) ou de l'ensemble réduit de mesures (vertical).

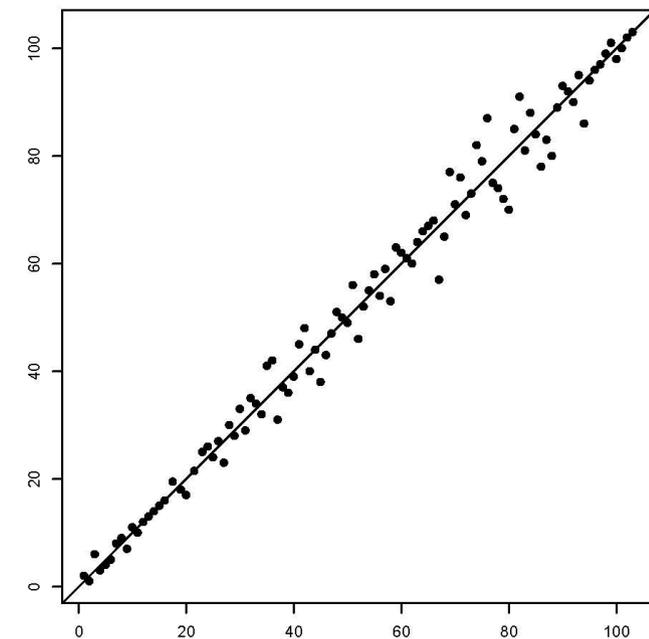
Topic 192



Topic 246



Topic 372



→ Le classement des systèmes reste globalement le même.

Conclusion Partie 1

Typologie des mesures de performances : 6 grandes catégories de mesures

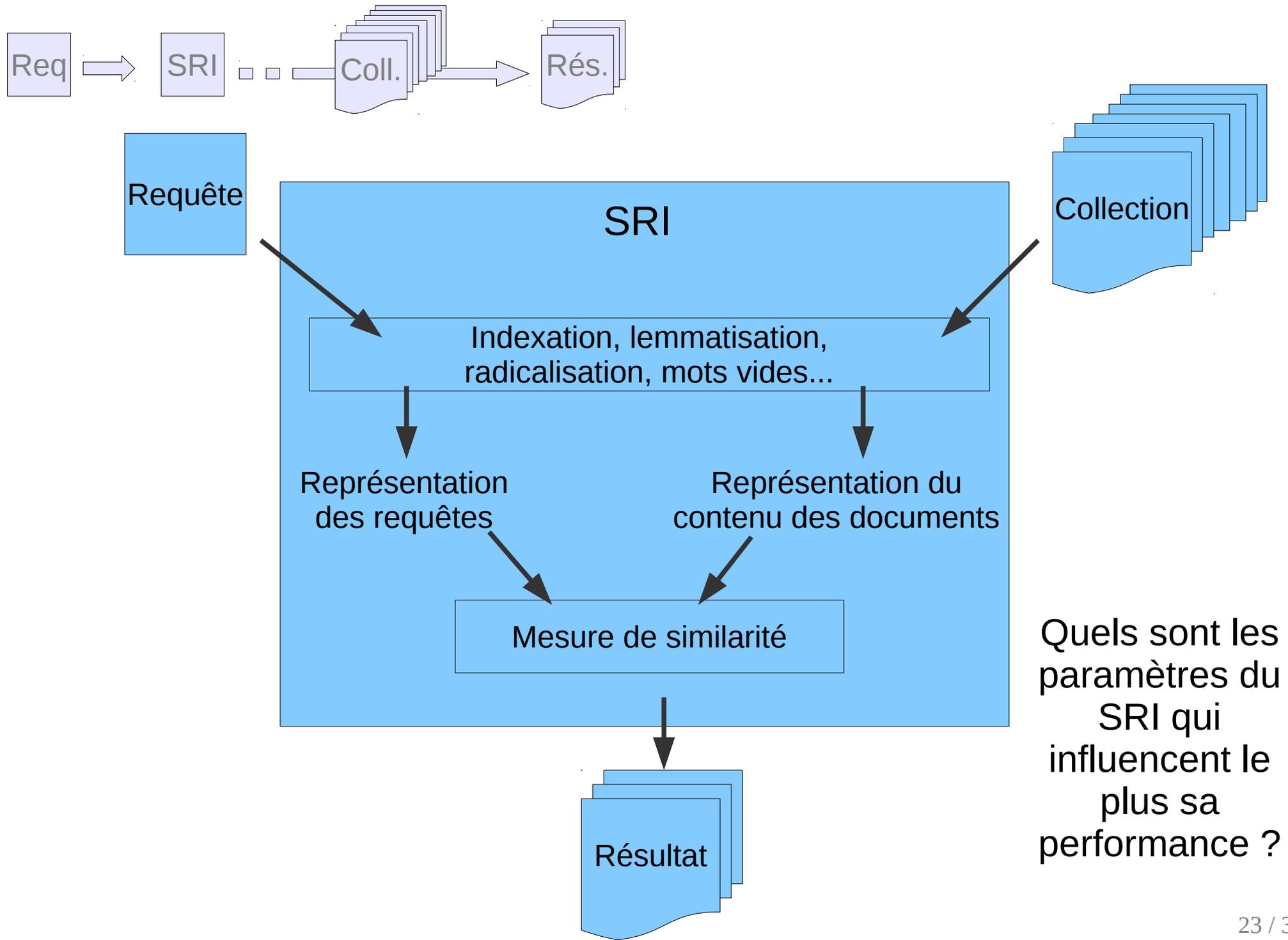
- Caractériser un SRI par l'évaluation de 6 mesures (une par groupe)
- Faciliter la comparaison de systèmes
- Cibler les mesures pertinentes (restreintes à 1 ou 2 groupes) à évaluer en vue d'une optimisation particulière d'un SRI
- Positionnement d'une nouvelle mesure vis-à-vis de cette typologie

Partie II

Paramètres d'un SRI

 J. Compaoré, S. Déjean, A.M. Gueye, J. Mothe, J. Randriamparany. **Mining Information Retrieval Results: Significant IR parameters**. Dans / In : *Advances in Information Mining and Management*, Barcelone, 23/10/2011-28/10/2011, IARIA

Paramètres d'un SRI



Quels sont les paramètres du SRI qui influencent le plus sa performance ?

La plateforme Terrier

Terrier is a highly flexible, efficient, and effective open source search engine, readily deployable on large-scale collections of documents.



<http://terrier.org/>

Terrier implements **state-of-the-art indexing and retrieval functionalities**, and provides an ideal platform for the rapid development and **evaluation of large-scale retrieval applications**.

Terrier is open source, and is a comprehensive, flexible and transparent platform for research and experimentation in text retrieval. Research can easily be carried out on standard TREC and CLEF test collections.

Quelques paramètres

| Parameters | Meaning | Values |
|------------|--|--|
| Top | Numéro de la requête (<i>topic</i>) | 351, ..., 400 |
| Field | Champ de la requête | T, T+D, T+D+N |
| Bloc | Taille du bloc d'indexation | 0, 1, 5, 10 |
| Idf | Fréquence inverse de document | FALSE, TRUE |
| Model | Modèle de recherche (<i>retrieval model</i>) | BB2c1, BM25b0.5, DFRBM25c1.0, IFB2c1.0, InexpB2c1.0, InexpC2c1.0, InL2c1.0, PL2c1.0, TFIDF |
| Ref | Reformulation de requête | None, Bo1bfree, Bo2bfree, KLbfree |
| DocNb | Nombre de documents pour la reformulation | 0, 3, 5, 10, 50, 100, 200 |

Jeu de données

- 98 650 lignes : 1 ligne = 1 *run*, une requête (*topic*) traitée en utilisant une configuration de 6 paramètres
- 8 colonnes : 1 requête + 6 paramètres + 1 mesure de performance (map)

| # | Top | Field | Bloc | Idf | Model | Ref | DocNb | map |
|-------|-----|-------|------|-------|---------|----------|-------|--------|
| 1 | 351 | T | 1 | false | BB2c1.0 | Bo1bfree | 3 | 0.6134 |
| 2 | 352 | T | 1 | false | BB2c1.0 | Bo1bfree | 3 | 0.3412 |
| 3 | 353 | T | 1 | false | BB2c1.0 | Bo1bfree | 3 | 0.3479 |
| 4 | 354 | T | 1 | false | BB2c1.0 | Bo1bfree | 3 | 0.0662 |
| 5 | 355 | T | 1 | false | BB2c1.0 | Bo1bfree | 3 | 0.2794 |
| 6 | 356 | T | 1 | false | BB2c1.0 | Bo1bfree | 3 | 0.0460 |
| ... | | | | | | | | |
| 98645 | 445 | T | 0 | true | TFIDF | NONE | 1 | 0.1514 |
| 98646 | 446 | T | 0 | true | TFIDF | NONE | 1 | 0.2234 |
| 98647 | 447 | T | 0 | true | TFIDF | NONE | 1 | 0.1121 |
| 98648 | 448 | T | 0 | true | TFIDF | NONE | 1 | 0.0114 |
| 98649 | 449 | T | 0 | true | TFIDF | NONE | 1 | 0.0714 |
| 98650 | 450 | T | 0 | true | TFIDF | NONE | 1 | 0.3226 |

Requêtes plus ou moins faciles

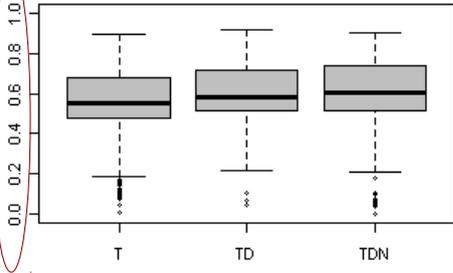
- La plupart du temps, les résultats sont moyennés sur un ensemble de requêtes alors que les résultats d'un système sont dépendants de la requête.
- Étude de l'impact des paramètres en fonction de la difficulté de la requête :
 - Requêtes « faciles »: requêtes pour lesquelles, l'AP moyenne sur les systèmes est la plus élevée ; dans notre cas, **13** requêtes avec une AP moyenne > 0.45)
 - Requêtes « difficiles »: requêtes pour lesquelles, l'AP moyenne sur les systèmes est la plus faible ; dans notre cas, **19** requêtes avec une AP moyenne < 0.045)

Analyses univariées

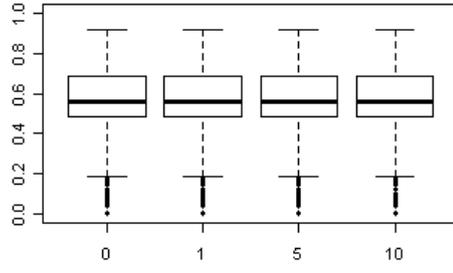
 Effect significatif (ANOVA 1 facteur)

1

Field ***

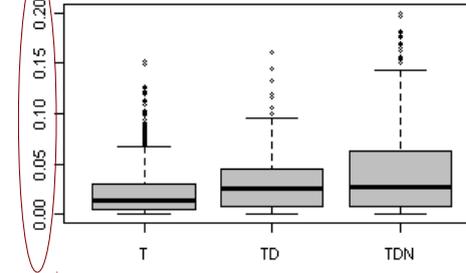


Bloc **

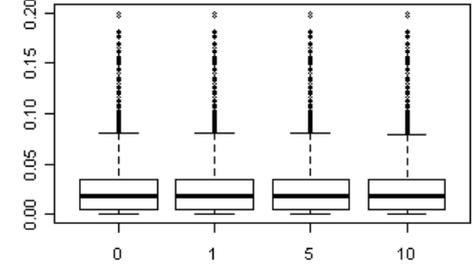


0.2

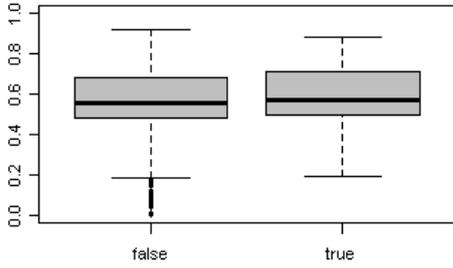
Field ***



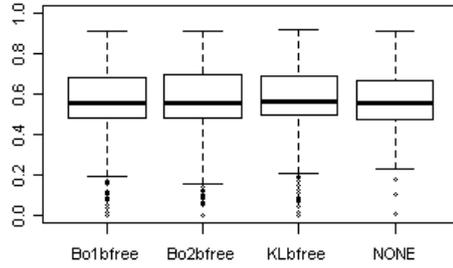
Bloc **



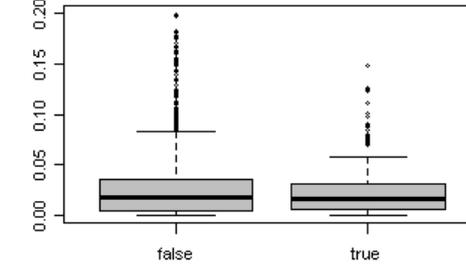
Idf ***



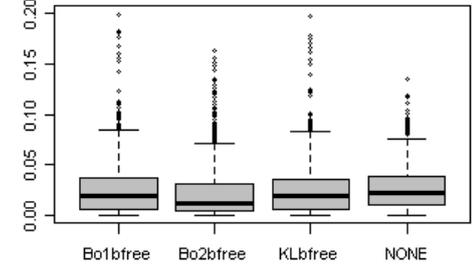
Ref **



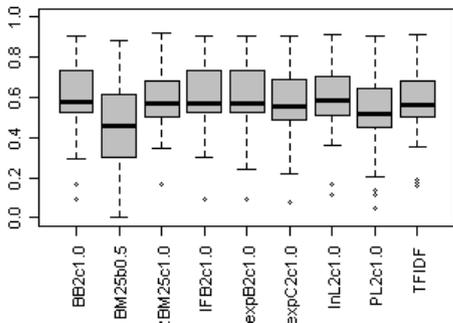
Idf ***



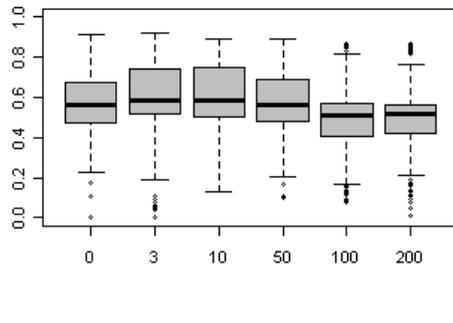
Ref ***



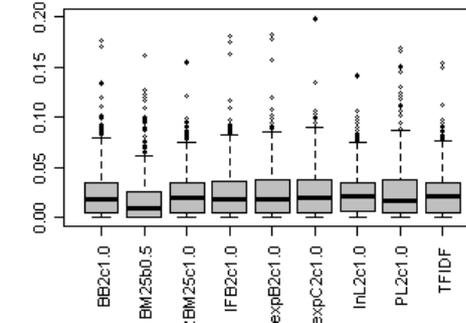
Model ***



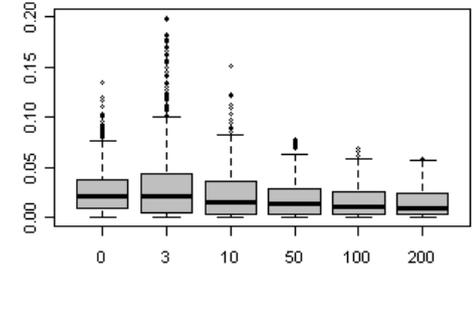
DocNb ***



Model ***



DocNb ***



Requêtes faciles

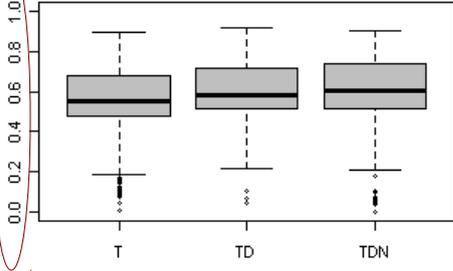
Requêtes difficiles

Analyses univariées

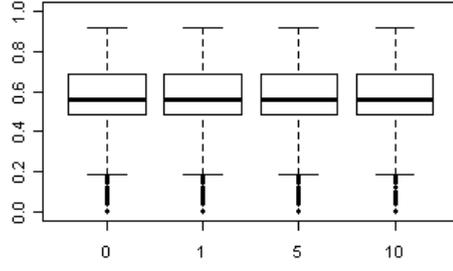
■ Effect significatif (ANOVA 1 facteur)

1

Field ***

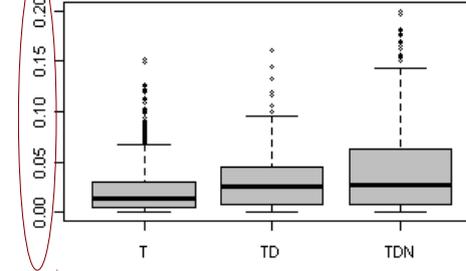


Bloc **

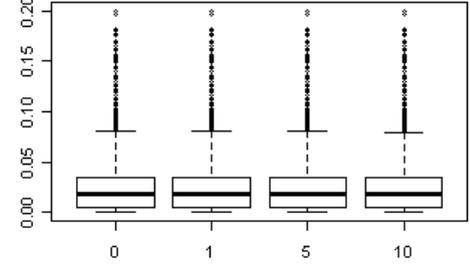


0.2

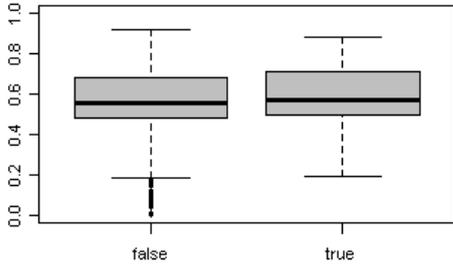
Field ***



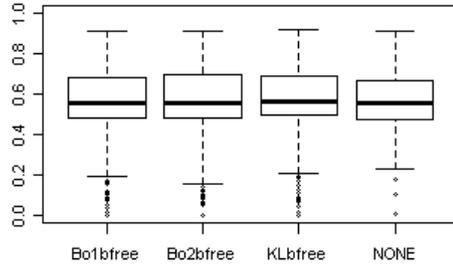
Bloc **



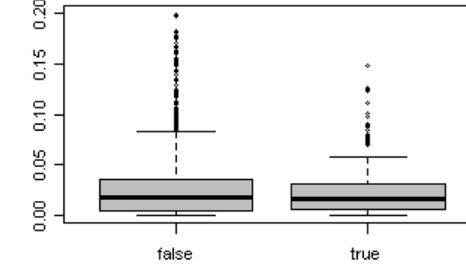
Idf ***



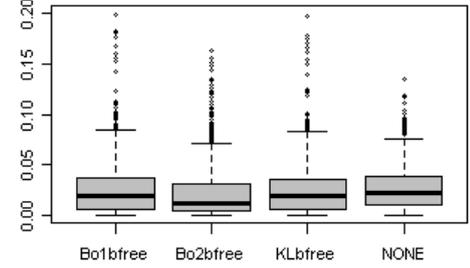
Ref **



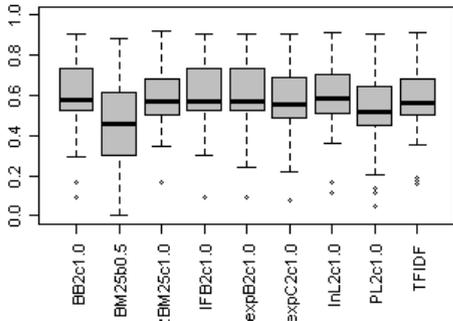
Idf ***



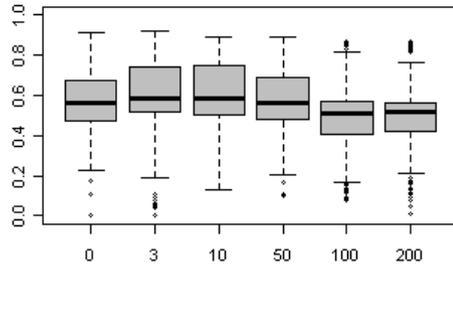
Ref ***



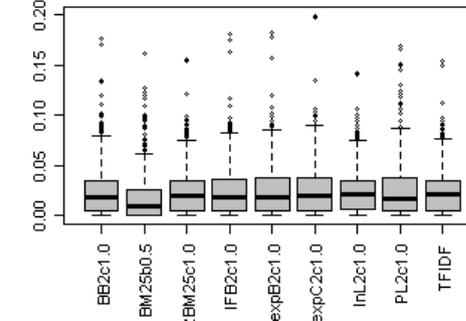
Model ***



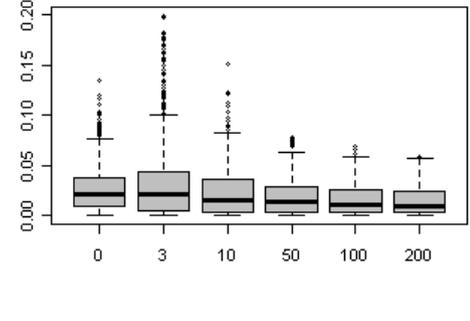
DocNb ***



Model ***



DocNb ***



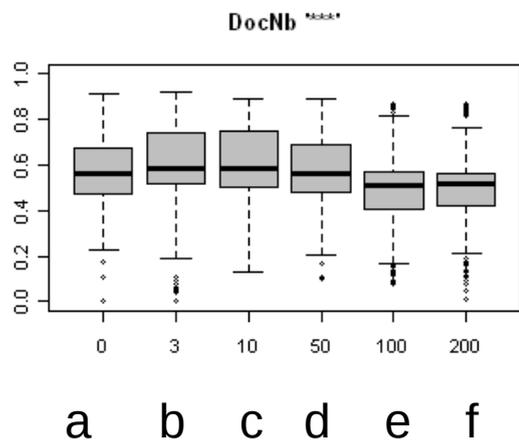
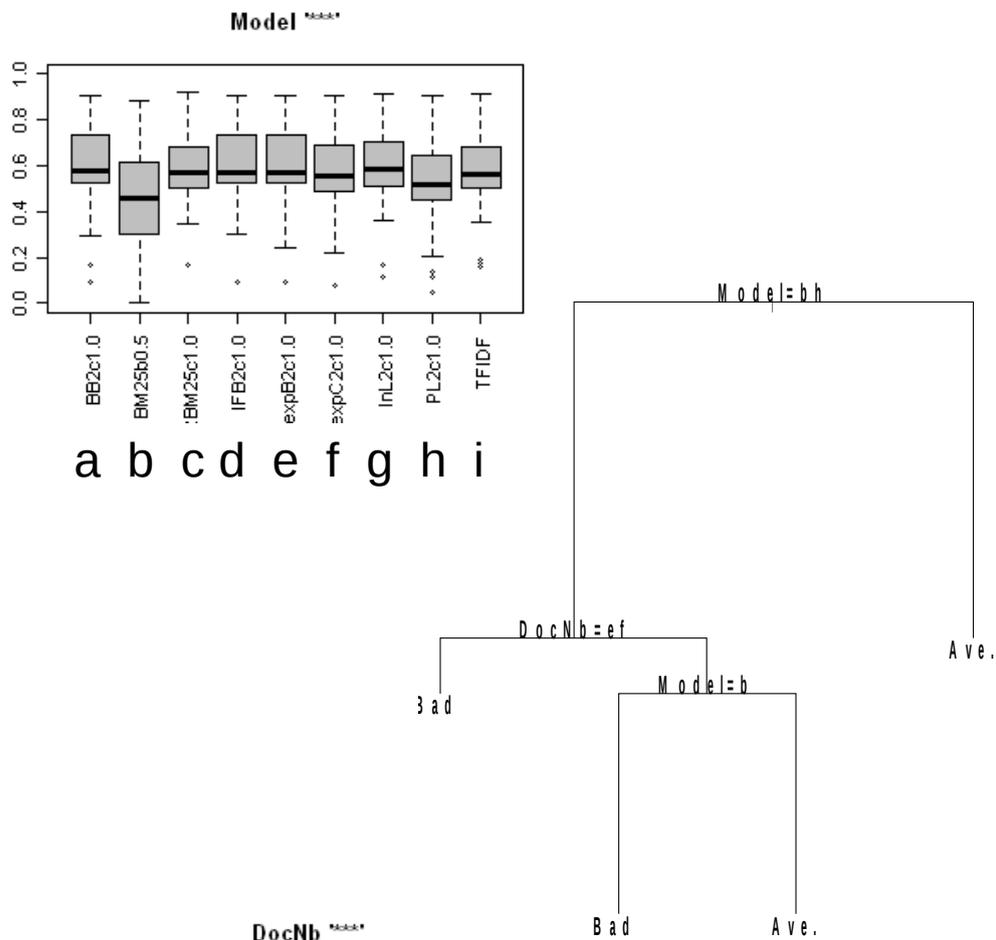
Requêtes faciles

Requêtes difficiles

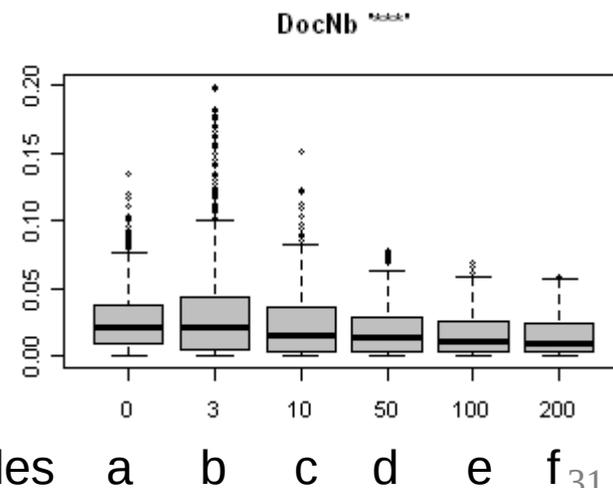
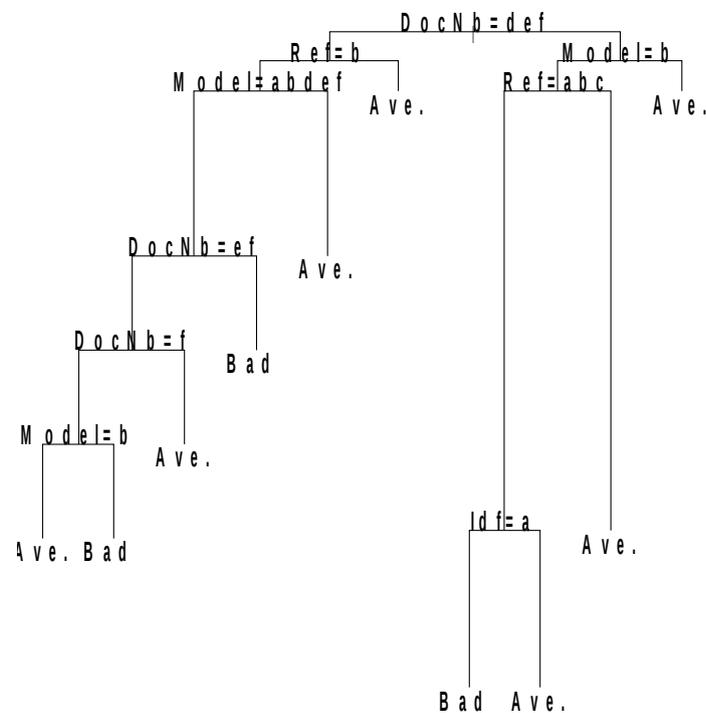
Analyse multivariée

- Arbre de décision (***Classification And Regression Trees, CART***)
- Construction d'un arbre de décision par découpages successifs de variables explicatives selon les valeurs d'une variable à expliquer (quantitative pour la régression ou qualitative pour la classification).
- Nous avons opté pour le cadre de la classification. Les valeurs de MAP ont été converties en classes :
 - *Bad*: MAP inférieure au premier quartile
 - *Average*: MAP comprise entre le premier et le troisième quartile
 - *Good*: MAP supérieure au troisième quartile

Résultats de CART



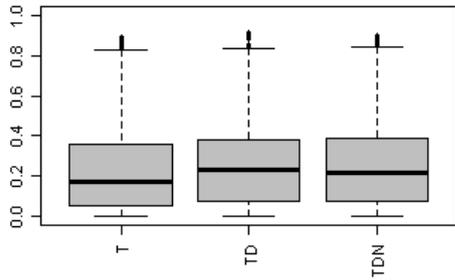
Requêtes faciles



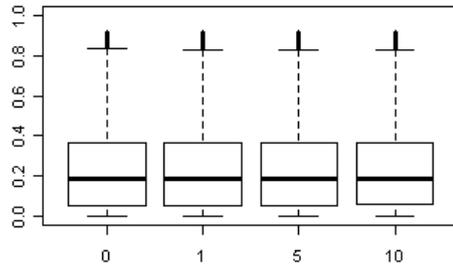
Requêtes difficiles

Sur l'ensemble des requêtes

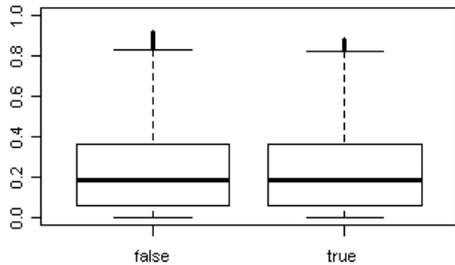
Field ^{***}



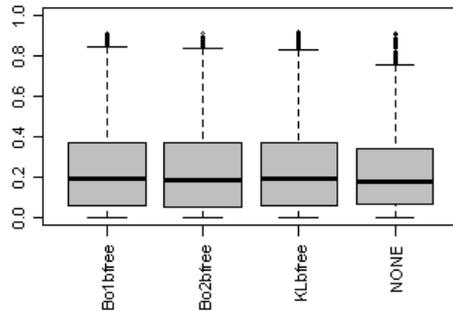
Bloc ^{**}



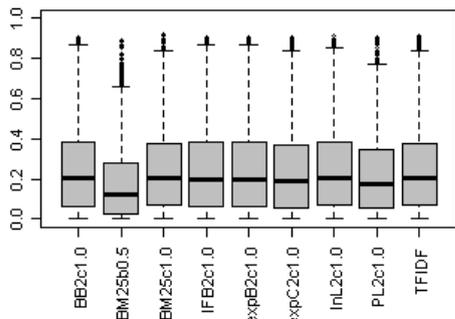
Idf ^{**}



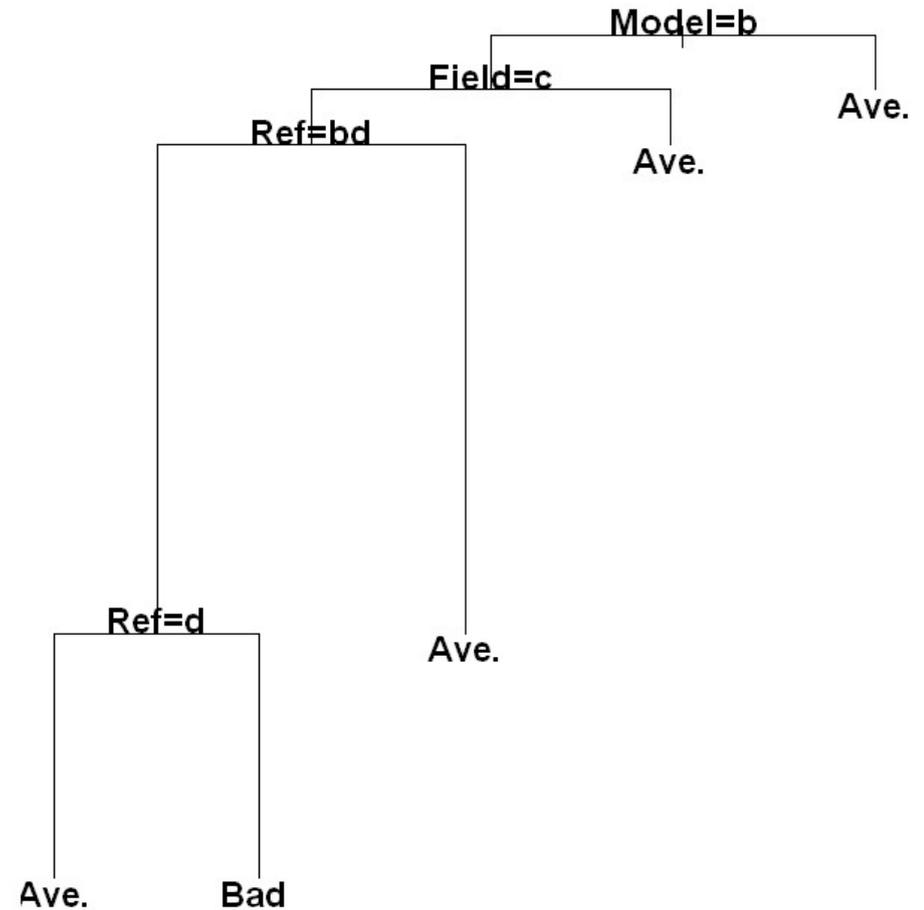
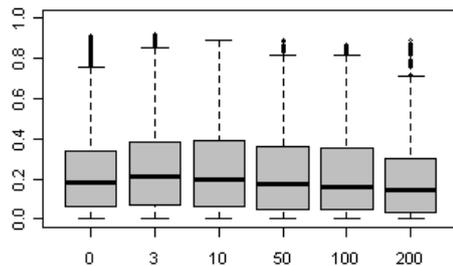
Ref ^{***}



Model ^{***}



DocNb ^{***}



Conclusion Partie II

- Certains paramètres produisent des changements significatifs dans les performances d'un SRI.
- Ces changements dépendent de la difficulté des requêtes.
- Ces travaux fournissent des pistes concernant le réglage des paramètres susceptibles d'améliorer la performance d'un SRI.
- Résultats à étendre à d'autres mesures de performance.
- Besoin d'une procédure plus systématique et *non auto-référente* pour définir les groupes de requêtes.

Suites

Besoin d'une procédure plus systématique et non auto-référente pour définir les groupes de requête

A. CHIFU, L. LAPORTE, J. MOTHE
La prédiction efficace de la difficulté
des requêtes : une tâche impossible ?
CORIA 2015 • Paris • 18-20 Mars

Prédicteurs pré-recherche

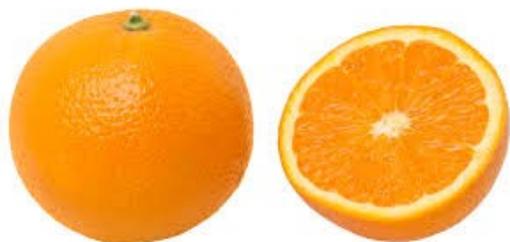
- L'ambiguïté des termes. Nombre de sens de WordNet (WNS) : prédicteur linguistique de pré-recherche, mesure de l'ambiguïté
- La discrimination des termes. Fréquence Inverse (idf) : prédicteur statistique de pré-recherche, mesure la rareté/popularité d'un terme

Prédicteurs post-recherche (?!)

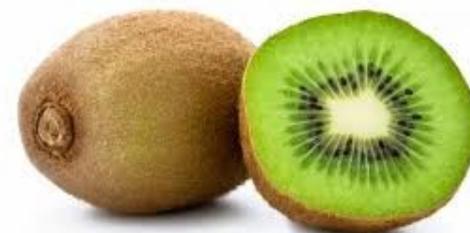
- L'homogénéité des listes de documents. Ecart-type (STD) : prédicteur post-recherche statistique, mesure le degré de variation de la liste des scores
- La divergence des listes. Retour sur la requête (QF) : un prédicteur post-recherche, calcule le chevauchement entre deux listes de documents retrouvés

Ambiguïté des termes

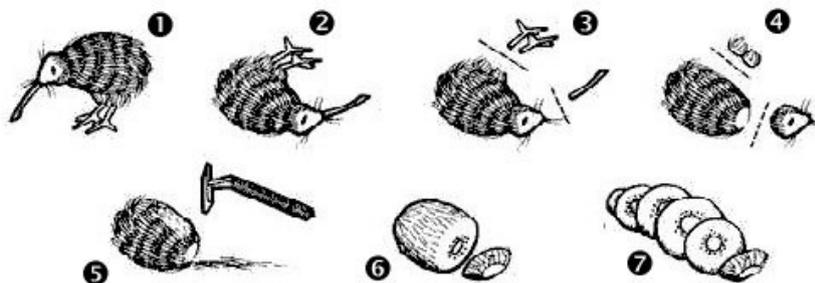
ORANGE



KIWI



How to prepare a kiwi



TERRIER


Terrier
<http://terrier.org/>



Analyses statistiques pour l'évaluation des systèmes de recherche d'information

S. Déjean, J. Mothe

www.math.univ-toulouse.fr/~sdejean

