

Partie 8

Évaluation des performances

Vincent Labatut

Laboratoire Informatique d'Avignon – LIA EA 4128

vincent.labatut@univ-avignon.fr

2019/20

M2 ILSÉN

UE Ingénierie du document et de l'information

UCE3 Indexation & Recherche d'information

Plan de la séance

- 1 Ressources et notion de pertinence
- 2 Évaluation de résultats non-ordonnés
 - Mesures basiques
 - Précision et rappel
 - F -mesure
- 3 Évaluation de résultats ordonnés
 - Courbe précision rappel
 - Mesures dérivées

Section 1

Ressources et notion de pertinence

Évaluation en recherche d'information

Ressources

- Ressources **nécessaires**
 - **Collection** de documents
 - **Besoins** informationnels
 - À convertir en requêtes
 - Ex. : Information sur l'hypothèse que boire du vin rouge réduit plus le risque d'attaque cardiaque que boire du vin blanc.
 - $\text{vin} \wedge \text{rouge} \wedge \text{blanc} \wedge \text{cœur} \wedge \text{attaque} \wedge \text{efficace}$
 - On note $\mathcal{Q} = \{q_1, \dots, q_Q\}$ l'ensemble des Q requêtes d'évaluation
 - **Pertinence** de chaque document
 - = **vérité terrain** (eng : Golden standard, ground-truth)
- **Taille** des données
 - Taille D du corpus de documents
 - Nombre Q de besoins informationnels considérés
 - Les deux affectent la fiabilité des résultats

Évaluation en recherche d'information

Notion de pertinence

- **Pertinence**
 - Évaluée par rapport aux besoins, et non pas aux requêtes
 - But = satisfaction de l'utilisateur final
 - Ex. `python` = animal ou langage de programmation ?
- **Besoins** partagés en plusieurs ensembles
 - Ensemble de **développement**
 - Utilisé seulement pour affiner les réglages du système
 - Ensemble de **test**
 - Utilisé pour obtenir la performance générale du système
 - Ne pas développer et tester un système sur les mêmes données
- **Exemples** de données de test :
 - **TREC**, **CLEF**, etc. (cf. [MRS08] Section 8.2)

Section 2

Évaluation de résultats non-ordonnés

Évaluation de résultats non-ordonnés

Types de résultats

Vrais positifs (eng : True Positives)

Nombre TP de documents pertinents **correctement renvoyés** par le moteur de recherche.

Vrais négatifs (eng : True Negatives)

Nombre TN de documents non-pertinents **correctement ignorés** par le moteur de recherche.

Faux positifs (eng : False Positives)

Nombre FP de documents non-pertinents **incorrectement renvoyés** par le moteur de recherche.

Faux négatifs (eng : False Negatives)

Nombre FN de documents pertinents **incorrectement ignorés** par le moteur de recherche.

Évaluation de résultats non-ordonnés

Taux de succès

Matrice de confusion :

	Pertinent	Non-pertinent
Renvoyé	TP	FP
Ignoré	FN	TN

Taux de succès (eng : Accuracy)

Proportion de documents **correctement classifiés** :

$$Acc = \frac{TP + TN}{TP + TN + FP + FN} = \frac{TP + TN}{D}$$

- Utilisé en **classification** (générale)
- **Pas adapté** à RI en raison du nombre élevé de négatifs
- Ex. : **rejet systématique** de tous les documents

Évaluation de résultats non-ordonnés

Précision et rappel

Matrice de confusion :

	Pertinent	Non-pertinent
Renvoyé	TP	FP
Ignoré	FN	TN

Précision (eng : Precision)

Proportion de documents **pertinents** parmi les documents **renvoyés** par le moteur de recherche (erreur de type I) :

$$Pre = \frac{TP}{TP+FP}$$

Rappel (eng : Recall)

Proportion de documents **renvoyés** par le moteur de recherche parmi les documents **pertinents** du corpus (erreur de type II) :

$$Rec = \frac{TP}{TP+FN}$$

- Ces trois mesures sont définies sur $[0;1]$

Évaluation de résultats non-ordonnés

Exemple pour Précision et Rappel

On considère un corpus de $D = 100$ documents, dont on sait que 20 sont pertinents pour une requête donnée.

Le moteur évalué renvoie les 10 documents suivants pour cette même requête :

DocId	105	695	23	8	87	147	694	3691	68	9999
Pertinent	.	✓	✓	✓	.	.	✓	✓	.	✓

Mesures obtenues :

- $TP = 6$
- $FP = 4$
- $FN = 14$
- $TN = 76$
- $Acc = (TP + TN) / D = (6 + 76) / 100 = 82 / 100 = 0,82$
- $Pre = TP / (TP + FP) = 6 / (6 + 4) = 6 / 10 = 0,6$
- $Rec = TP / (TP + FN) = 6 / (6 + 14) = 6 / 20 = 0,3$

Évaluation de résultats non-ordonnés

Précision vs. Rappel

- Ces mesures sont **complémentaires**
- Perspective **utilisateur** : l'une est généralement prioritaire
 - Ex. : Utilisateur Web → préfère une **précision élevée**
 - Ex. : Expert légal → préfère un **rappel élevé**
- Ces mesures sont **liées** :
 - Sélectionner **tous** les documents :
 - $Pre \approx 0, Rec = 1$
 - Rejeter **tous** les document sauf un (vrai) positif
 - $Pre = 1, Rec \approx 0$
- Si **augmentation** du nombre de documents retournés :
 - Rappel ne peut pas décroître
 - Nombre de *TN* ne peut pas augmenter
 - Précision décroît généralement
 - Si augmentation du nombre de *FP*

Évaluation de résultats non-ordonnés

Définition de la F -mesure

Moyenne harmonique

La moyenne harmonique M_H est l'**inverse** de la moyenne **arithmétique** M_A des **inverses** des valeurs x_1, \dots, x_n moyennées :

$$M_H(x_1, \dots, x_n) = \frac{1}{M_A(x_1, \dots, x_n)} = \frac{1}{\frac{1}{n} \times (\frac{1}{x_1} + \dots + \frac{1}{x_n})} = \frac{n}{\frac{1}{x_1} + \dots + \frac{1}{x_n}}.$$

- Propriétés :
 - Toujours $<$ à moyennes géométrique et arithmétique
 - Plus proche de min que de moyenne arithmétique

F -mesure (eng : F -measure)

Il s'agit de la **moyenne harmonique** de la précision Pre et du rappel Rec :

$$F = \frac{2}{\frac{1}{Pre} + \frac{1}{Rec}} = 2 \frac{Pre \cdot Rec}{Pre + Rec}$$

- Pourquoi moyenne harmonique :
 - Faire ressortir une Pre ou un Rec exagérément bas
 - Ex. : $Pre = 1$ et $Rec = 0 \rightarrow$ moy. arithm. de 0,5 mais $F \approx 0$.

Évaluation de résultats non-ordonnés

Généralisation de la F -mesure

F -mesure généralisée

Version **pondérée** utilisant un **paramètre** β pour contrôler l'**importance relative** de la précision Pre et du rappel Rec :

$$F_{\beta} = \frac{1}{\frac{\alpha}{P} + \frac{1-\alpha}{R}} = \frac{(\beta^2 + 1) \cdot Pre \cdot Rec}{\beta^2 \cdot Pre + Rec},$$

où $\alpha \in [0;1]$ et $\beta^2 = (1 - \alpha)/\alpha$.

- **Comportement** :
 - $\beta = 1 \rightarrow$ équilibre entre Pre et Rec
 - $\beta < 1 \rightarrow Pre$ a plus d'importance
 - $\beta > 1 \rightarrow Rec$ a plus d'importance
- Utilisée quand on veut favoriser l'une des deux mesures (ex. cas d'utilisation p.11)

Évaluation de résultats non-ordonnés

Exemple de F -mesure

Même exemple que précédemment (p.10) : on considère un corpus de 100 documents, dont on sait que 15 sont pertinents pour une requête donnée.

Le moteur évalué renvoie les 10 documents suivants pour cette même requête :

DocId	105	695	23	8	87	147	694	3691	68	9999
Pertinent	.	✓	✓	✓	.	.	✓	✓	.	✓

Mesures obtenues :

- $TP = 6$; $FP = 4$; $FN = 9$; $TN = 81$
- $Pre = 0,6$; $Rec = 0,3$; $Acc = 0,82$
- $F = 2 \frac{Pre \cdot Rec}{Pre + Rec} = 2 \frac{0,6 \cdot 0,3}{0,6 + 0,3} = 2 \frac{0,18}{0,9} = 0,4$

Section 3

Évaluation de résultats ordonnés

Évaluation de résultats ordonnés

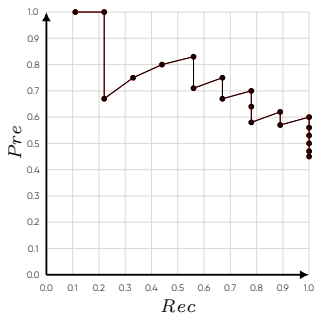
Définition de la courbe Précision-Rappel

- Notations :
 - $Pre(k)$: précision obtenue en considérant les k premiers documents
 - $Rec(k)$: pareil pour le rappel
- Courbe Précision-Rappel : méthode de construction
 - 1 Ordonner les documents par **score décroissant**
 - 2 Pour $k = 1, \dots, D$, calculer $Pre(k)$ et $Rec(k)$
 - 3 Représenter la performance obtenue pour chaque valeur de k par un point dans le repère $Pre-Rec$.
- Propriétés :
 - $Rec(k)$ est une fonction **non-décroissante** : augmenter le nombre de documents considérés ne peut pas diminuer le rappel
 - courbe en **dents-de-scie** (Pre diminuant à chaque document incorrect)

Évaluation de résultats ordonnés

Exemple de courbe Précision-Rappel

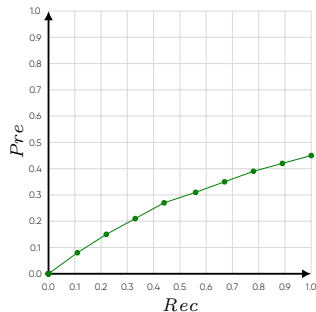
Rang	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	19	20
Pertinent	✓	✓	.	✓	✓	✓	.	✓	.	✓	.	.	✓	.	✓
k	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	19	20
TP	1	2	2	3	4	5	5	6	6	7	7	7	8	8	9	9	9	9	9	9
FP	0	0	1	1	1	1	2	2	3	3	4	5	5	6	6	7	8	9	10	11
FN	8	7	7	6	5	4	4	3	3	2	2	2	1	1	0	0	0	0	0	0
$Pre(k)$	1/1	2/2	2/3	3/4	4/5	5/6	5/7	6/8	6/9	7/10	7/11	7/12	8/13	8/14	9/15	9/16	9/17	9/18	9/19	9/20
	1,00	1,00	0,67	0,75	0,80	0,83	0,71	0,75	0,67	0,70	0,64	0,58	0,62	0,57	0,60	0,56	0,53	0,50	0,47	0,45
$Rec(k)$	1/9	2/9	2/9	3/9	4/9	5/9	5/9	6/9	6/9	7/9	7/9	7/9	8/9	8/9	9/9	9/9	9/9	9/9	9/9	9/9
	0,11	0,22	0,22	0,33	0,44	0,56	0,56	0,67	0,67	0,78	0,78	0,78	0,89	0,89	1,00	1,00	1,00	1,00	1,00	1,00



Évaluation de résultats ordonnés

Exemple de courbe Précision-Rappel

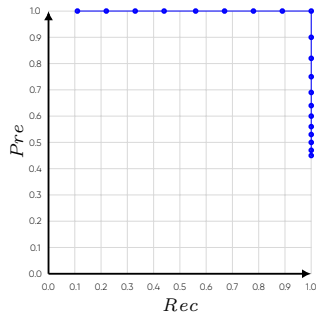
Rang	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	19	20
Pertinent	✓	✓	✓	✓	✓	✓	✓	✓	✓
k	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	19	20
TP	0	0	0	0	0	0	0	0	0	0	0	1	2	3	4	5	6	7	8	9
FP	1	2	3	4	5	6	7	8	9	10	11	11	11	11	11	11	11	11	11	11
FN	9	9	9	9	9	9	9	9	9	9	9	8	7	6	5	4	3	2	1	0
$Pre(k)$	0/1 0,00	0/2 0,00	0/3 0,00	0/4 0,00	0/5 0,00	0/6 0,00	0/7 0,00	0/8 0,00	0/9 0,00	0/10 0,00	0/11 0,00	1/12 0,08	2/13 0,15	3/14 0,21	4/15 0,27	5/16 0,31	6/17 0,35	7/18 0,39	8/19 0,42	9/20 0,45
$Rec(k)$	0/9 0,00	0/9 0,00	0/9 0,00	0/9 0,00	0/9 0,00	0/9 0,00	0/9 0,00	0/9 0,00	0/9 0,00	0/9 0,00	0/9 0,00	1/9 0,11	2/9 0,22	3/9 0,33	4/9 0,44	5/9 0,56	6/9 0,67	7/9 0,78	8/9 0,89	9/9 1,00



Évaluation de résultats ordonnés

Exemple de courbe Précision-Rappel

Rang	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	19	20
Pertinent	✓	✓	✓	✓	✓	✓	✓	✓	✓
k	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	19	20
TP	1	2	3	4	5	6	7	8	9	9	9	9	9	9	9	9	9	9	9	9
FP	0	0	0	0	0	0	0	0	0	1	2	3	4	5	6	7	8	9	10	11
FN	8	7	6	5	4	3	2	1	0	0	0	0	0	0	0	0	0	0	0	0
$Pre(k)$	1/1	2/2	3/3	4/4	5/5	6/6	7/7	8/8	9/9	9/10	9/11	9/12	9/13	9/14	9/15	9/16	9/17	9/18	9/19	9/20
	1,00	1,00	1,00	1,00	1,00	1,00	1,00	1,00	1,00	0,90	0,82	0,75	0,69	0,64	0,60	0,56	0,53	0,50	0,47	0,45
$Rec(k)$	1/9	2/9	3/9	4/9	5/9	6/9	7/9	8/9	9/9	9/9	9/9	9/9	9/9	9/9	9/9	9/9	9/9	9/9	9/9	9/9
	0,11	0,22	0,33	0,44	0,56	0,67	0,78	0,89	1,00	1,00	1,00	1,00	1,00	1,00	1,00	1,00	1,00	1,00	1,00	1,00



Évaluation de résultats ordonnés

Principe de l'interpolation

Précision interpolée

La précision interpolée $\tilde{Pre}(k)$ est la précision **maximale** obtenue en renvoyant k **documents ou plus**. Formellement :

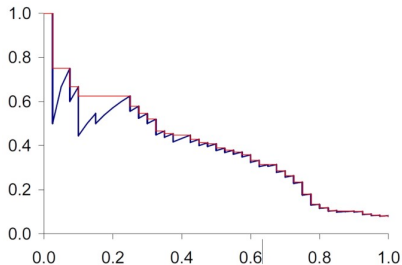
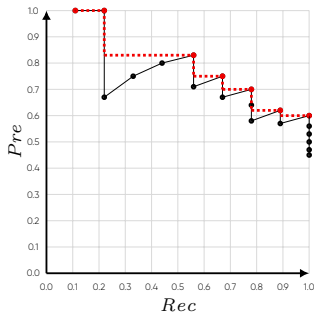
$$\tilde{Pre}(k) = \max_{k \leq k' \leq D} Pre(k').$$

- Courbe Précision-Rappel interpolée
 - Construction : on utilise \tilde{Pre} à la place de Pre
 - Intérêt :
 - En jouant sur k , on peut parfois augmenter Rec sans diminuer Pre
 - La courbe Précision-Rappel interpolée permet de détecter ces situations

Évaluation de résultats ordonnés

Exemple de courbe Précision-Rappel interpolée

Rang	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	19	20
Pertinent	✓	✓	.	✓	✓	✓	.	✓	.	✓	.	.	✓	.	✓
k	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	19	20
TP	1	2	2	3	4	5	5	6	6	7	7	7	8	8	9	9	9	9	9	9
FP	0	0	1	1	1	1	2	2	3	3	4	5	5	6	6	7	8	9	10	11
FN	8	7	7	6	5	4	4	3	3	2	2	2	1	1	0	0	0	0	0	0
Pre(k)	1/1	2/2	2/3	3/4	4/5	5/6	5/7	6/8	6/9	7/10	7/11	7/12	8/13	8/14	9/15	9/16	9/17	9/18	9/19	9/20
Rec(k)	1,00	1,00	0,67	0,75	0,80	0,83	0,71	0,75	0,67	0,70	0,64	0,58	0,62	0,57	0,60	0,56	0,53	0,50	0,47	0,45
	0,11	0,22	0,22	0,33	0,44	0,56	0,56	0,67	0,67	0,78	0,78	0,78	0,89	0,89	1,00	1,00	1,00	1,00	1,00	1,00



Évaluation de résultats ordonnés

Précision interpolée moyenne

Notation : $Rec_j(k)$, $Pre_j(k)$ et $\tilde{Pre}_j(k)$ dénotent respectivement le rappel, la précision et la précision interpolée obtenus sur les k premiers documents renvoyés pour la requête q_j

Points d'interpolation

Pour une requête q_j , c'est l'**ensemble** $\{k_{ij}\}_{i \in \{0, \dots, 10\}}$ des 11 **valeurs** de k telles que :

$$Rec_j(k_{ij}) = i/10, \text{ où } i \in \{0, \dots, 10\}.$$

Précision interpolée moyenne

Moyenne arithmétique des **précisions interpolées** obtenues pour l'ensemble \mathcal{Q} des **requêtes** testées, en considérant leur i ème point d'interpolation :

$$IAP_i = \frac{1}{Q} \sum_{j=1}^Q \tilde{Pre}_j(k_{ij}).$$

Évaluation de résultats ordonnés

Courbe de précision interpolée moyenne

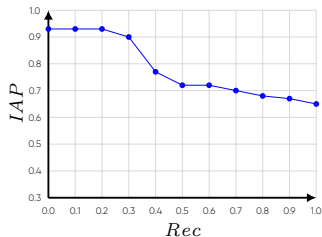
Exemple : on utilise 5 requêtes pour l'évaluation

Rang	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	19	20
Pert. q_1	✓	✓		✓	✓	✓		✓		✓			✓		✓					
Rec ₁ (k)	0,11	0,22	0,22	0,33	0,44	0,56	0,56	0,67	0,67	0,78	0,78	0,78	0,89	0,89	1,00	1,00	1,00	1,00	1,00	1,00
Pre ₁ (k)	1,00	1,00	0,83	0,83	0,83	0,83	0,75	0,75	0,70	0,70	0,64	0,62	0,62	0,60	0,60	0,56	0,53	0,50	0,47	0,45
Pert. q_2	✓	✓	✓					✓	✓	✓			✓			✓				
Rec ₂ (k)	0,10	0,20	0,30	0,30	0,30	0,30	0,40	0,50	0,60	0,70	0,70	0,80	0,90	0,90	0,90	1,00	1,00	1,00	1,00	1,00
Pre ₂ (k)	1,00	1,00	1,00	0,75	0,70	0,70	0,70	0,70	0,70	0,70	0,69	0,69	0,69	0,64	0,63	0,63	0,59	0,56	0,53	0,50
Pert. q_3		✓		✓		✓		✓	✓	✓		✓		✓						
Rec ₃ (k)	0,00	0,10	0,10	0,20	0,20	0,30	0,40	0,50	0,60	0,60	0,70	0,70	0,80	0,90	1,00	1,00	1,00	1,00	1,00	1,00
Pre ₃ (k)	0,67	0,67	0,67	0,67	0,67	0,67	0,67	0,67	0,67	0,67	0,67	0,67	0,67	0,67	0,67	0,63	0,59	0,56	0,53	0,50
Pert. q_4	✓	✓	✓					✓	✓	✓			✓							
Rec ₄ (k)	0,13	0,25	0,25	0,25	0,25	0,38	0,50	0,63	0,63	0,63	0,63	0,63	0,75	0,88	1,00	1,00	1,00	1,00	1,00	1,00
Pre ₄ (k)	1,00	1,00	0,67	0,63	0,63	0,63	0,63	0,63	0,62	0,62	0,62	0,62	0,62	0,62	0,57	0,53	0,50	0,47	0,44	0,42
Pert. q_5		✓			✓	✓	✓	✓	✓	✓										
Rec ₅ (k)	0,17	0,33	0,33	0,33	0,50	0,67	0,83	1,00	1,00	1,00	1,00	1,00	1,00	1,00	1,00	1,00	1,00	1,00	1,00	1,00
Pre ₅ (k)	1,00	1,00	0,75	0,75	0,75	0,75	0,75	0,75	0,67	0,60	0,55	0,50	0,46	0,43	0,40	0,38	0,35	0,33	0,32	0,30

Pt. interp.	k_{11}	k_2	k_3	k_4	k_5	k_6	k_7	k_8	k_9	k_{10}	k_{11}
Rec(k_{ij})	0,00	0,10	0,20	0,30	0,40	0,50	0,60	0,70	0,80	0,90	1,00
Pre ₁ (k_{i1})	1,00	1,00	1,00	0,83	0,83	0,83	0,83	0,75	0,70	0,62	0,60
Pre ₂ (k_{i2})	1,00	1,00	1,00	1,00	0,70	0,70	0,70	0,70	0,69	0,69	0,63
Pre ₃ (k_{i3})	0,67	0,67	0,67	0,67	0,67	0,67	0,67	0,67	0,67	0,67	0,67
Pre ₄ (k_{i4})	1,00	1,00	1,00	1,00	0,63	0,63	0,63	0,62	0,62	0,62	0,62
Pre ₅ (k_{i5})	1,00	1,00	1,00	1,00	1,00	0,75	0,75	0,75	0,75	0,75	0,75
IAP _i	0,93	0,93	0,93	0,90	0,77	0,72	0,72	0,70	0,68	0,67	0,65

Construction de la courbe : calcul des valeurs IAP_i pour les 11 points d'interpolation

(i.e. Précisions associées aux valeurs de rappel $Rec = 0 ; 0,1 ; 0,2 ; \dots ; 1$)



Évaluation de résultats ordonnés

Définition de la moyenne des précisions moyennes

La **moyenne des précisions moyennes** (MAP) repose sur **deux niveaux** de moyennage

Points de moyennage

Pour une requête q_j , c'est l'**ensemble** $\{k_{ij}\}_{1 \leq i \leq m_j}$ des m_j valeurs de k telles que le k_{ij} ème document soit **pertinent**.

Précision moyenne (eng : Average precision)

Pour une requête q_j , c'est la **moyenne** arithmétique des précisions obtenues pour l'ensemble des points de moyennage :

$$AP_j = \frac{1}{m_j} \sum_{i=1}^{m_j} Pre_j(k_{ij}).$$

Évaluation de résultats ordonnés

Propriétés de la moyenne des précisions moyennes

Moyenne des précisions moyennes

La MAP (eng : *Mean Average Precision*) est la **moyenne** arithmétique des **précisions moyennes** AP_j obtenues pour l'ensemble Q des **requêtes** testées :

$$MAP = \frac{1}{Q} \sum_{j=1}^Q AP_j.$$

- **Intérêt** : permet de se ramener à une seule valeur
- **Approximation** de l'aire (moyenne) sous la précision interpolée
- **Note** : de nombreuses autres mesures existent, cf. bibliographie

Section 4

Conclusion

Concepts abordés dans cette partie

- Courbe Précision-Rappel
- Courbe PR Interpolée
- Points d'interpolation
- Précision interpolée moyenne
- Précision
- Rappel
- F -mesure
- MAP

Lectures recommandées

- [MRS08] *Introduction to Information Retrieval*, chapitres 6 & 8.
- [BCC10] *Information Retrieval : Implementing and Evaluating Search Engines*, chapitre 12.
- [BR11] *Modern Information Retrieval : The Concepts and Technology behind Search*, chapitre 4.
- [AG13] *Recherche d'information - Applications, modèles et algorithmes*, chapitres 2 & 3.
- [CMS15] *Search Engines : Information Retrieval in Practice*, chapitres 8.

Références bibliographiques I

- [AG13] M.-R. Amini et É. Gaussier. *Recherche d'information – Applications, modèles et algorithmes*. Paris, FR : Eyrolles, 2013. url : <https://www.eyrolles.com/Informatique/Livre/recherche-d-information-9782212673760/>.
- [BR11] R. Baeza-Yates et B. Ribeiro-Neto. *Modern Information Retrieval : The Concepts and Technology behind Search*. 2nd Edition. Boston, USA : Addison Wesley Longman, 2011. url : <http://people.ischool.berkeley.edu/~hearst/irbook/>.
- [BCC10] S. Büttcher, C. L. A. Clarke et G. V. Cormack. *Information Retrieval : Implementing and Evaluating Search Engines*. Cambridge, USA : MIT Press, 2010. url : <http://www.ir.uwaterloo.ca/book/>.
- [CMS15] W. B. Croft, D. Metzler et T. Strohman. *Search Engines : Information Retrieval in Practice*. Pearson, 2015. url : <http://www.search-engines-book.com/>.

Références bibliographiques II

- [MRS08] C. D. Manning, P. Raghavan et H. Schütze. *Introduction to Information Retrieval*. New York, USA : Cambridge University Press, 2008. url : <http://www-nlp.stanford.edu/IR-book/>.