

CLASSIFIEUR DE MÉTHODES DE RECHERCHE

**Faciliter les revues de littérature
mixtes en santé**

CLASSIFIEUR DE MÉTHODES DE RECHERCHE

Université de Montréal

Département d'Informatique et de Recherche Opérationnelle

Laboratoire de recherche appliquée en linguistique informatique (RALI)

- Alexis Langlois, étudiant MS
- Jian-Yun Nie, PhD



Université McGill

Département de médecine de famille

- Quan Nha Hong, OT, MSc, PhD
- Pierre Pluye, MD, PhD



University College London Institute of Education

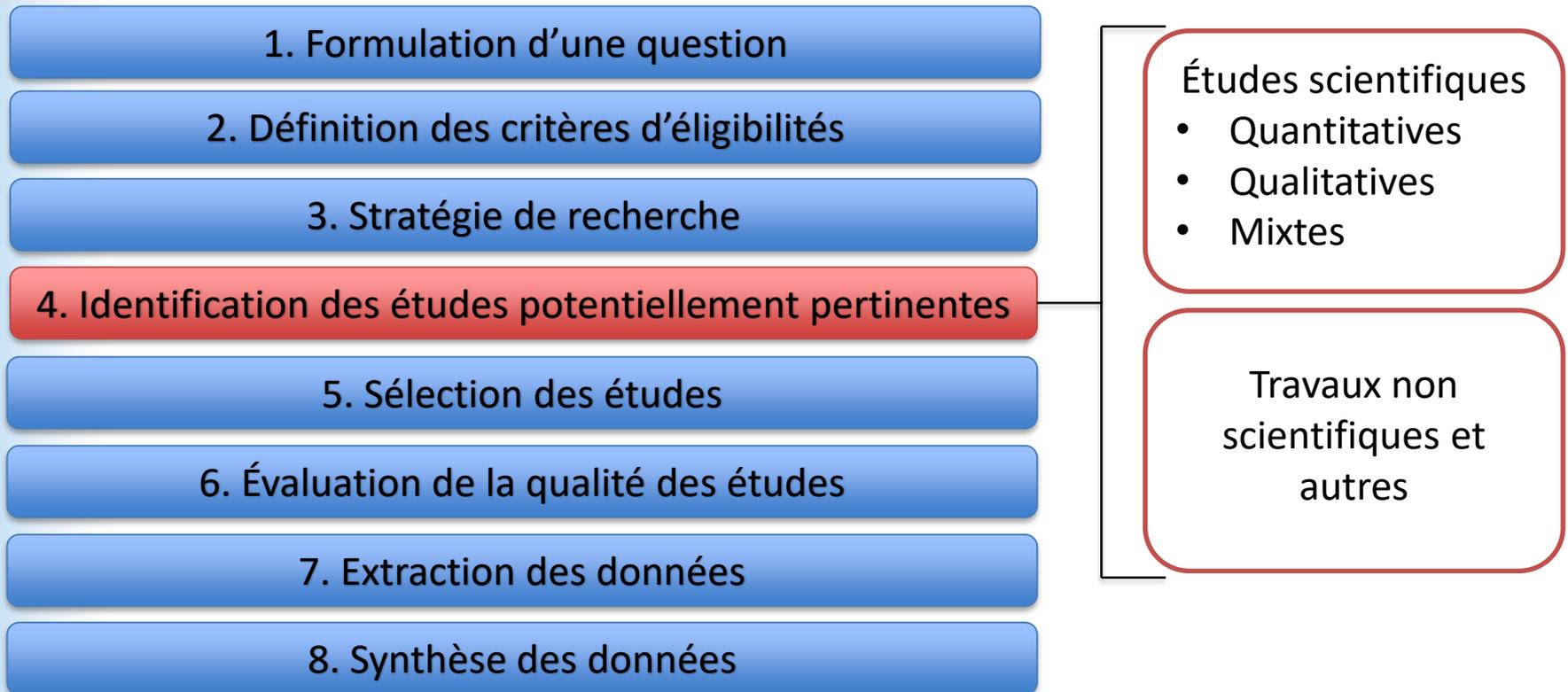
EPPI-Centre

- James Thomas, PhD



REVUES DE LITTÉRATURE MIXTES

CONTEXTE



REVUES DE LITTÉRATURE MIXTES

CONTEXTE

Méthodes quantitatives

- Essai randomisé contrôlé
- Étude non-randomisée
 - Étude de cohorte
 - Essai contrôlé
- Étude descriptive
 - Série de cas
 - Rapport de cas
 - Enquête transversal
- ...

Méthodes qualitatives

- Entretien
- Biographie
- Ethnographie
- Étude de cas qualitative
- Étude descriptive ou interprétative
- Phénoménographie
- ...

Travaux non scientifiques

- Éditorial
- Lettre
- Présentation
- Commentaire
- Débat
- Travail méthodologique
- ...

REVUES DE LITTÉRATURE MIXTES

FILTRE MIXTE (BOOLEEN)

- Combinaison de mots clés pour MEDLINE
(focus group* OR ethno* ...)

OR (random*.mp. OR control*.mp. ...)

NOT (letter OR comment ...)

- Temps du triage réduit d'environ 50%

- Collection de 4 500 documents

- 1 400 scientifiques
- 3 100 non scientifiques

- Résultats

- Taux de succès ≈ 57%
- **Rappel ≈ 89%**
- Précision ≈ 60%
- Spécificité ≈ 55%

Inconvénients

Rigidité (forme fléchie du langage naturel?)

Taux de succès ≈ lancer d'une pièce

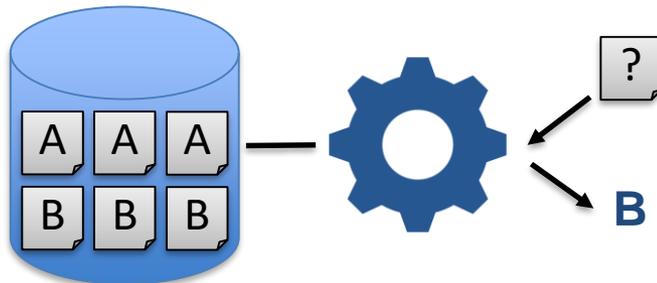
Dépendance aux bases de données bibliographiques

Maintenance et ajustements complexes

APPRENTISSAGE MACHINE

CLASSIFICATION AUTOMATIQUE

- Apprentissage **statistique** guidé par un **algorithme**
- Apprentissage supervisé / **Classification binaire**
 - Scientifique ou non scientifique?
 - Nécessite une quantité raisonnable d'**entrées annotés**
 - **Généralisation** des connaissances à-priori
- Prédiction → Catégorie ou probabilité
- Définition explicite des **caractéristiques** d'une entrée



APPRENTISSAGE MACHINE

ALGORITHMES

On cherche une fonction mathématique $x \rightarrow y$ qui généralise bien l'ensemble d'entraînement (entrées déjà annotées)

Plusieurs algorithmes utilisés

- Machine à vecteurs de support
- Réseau de neurones
- Classifieur de Bayes
- Arbre de décision
- Méthode des k proches voisins
- Histogramme
- Perceptron
- ...

APPRENTISSAGE MACHINE

CARACTÉRISTIQUES

1. Quelles sont les caractéristiques et les particularités d'un résumé de texte scientifique et non scientifique?
2. Comment les exploiter par un algorithme?

Indexation

TF-IDF



{mot 1, mot 2, mot 3, empirique}



{0.11, 0.25, 0.12, 1}

Meta-caractéristiques

Concepts

{0.11, 0.25, 0.12, 1}



{0.11, 0.12, 0.18, 0.8, 0.3, 0.5, 1}

Sélection

Info Gain

{0.11, 0.25, ..., 1}

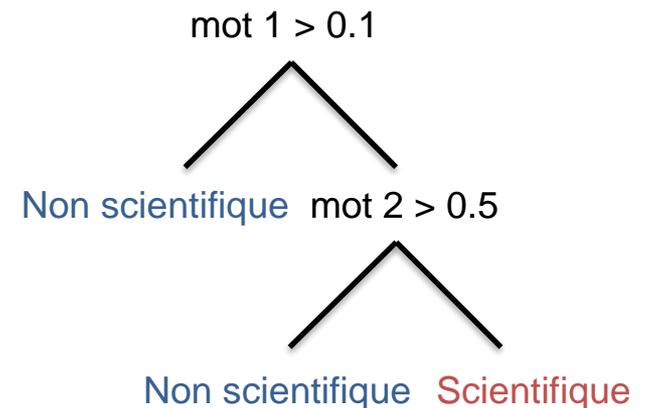
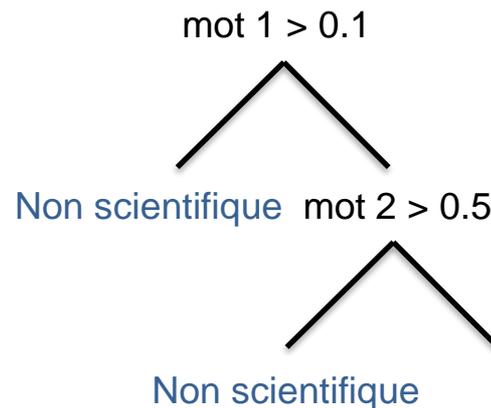
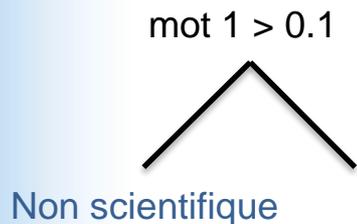


{0.11, 0.12, 0.8, 1}

CLASSIFIEUR DE MÉTHODES DE RECHERCHE

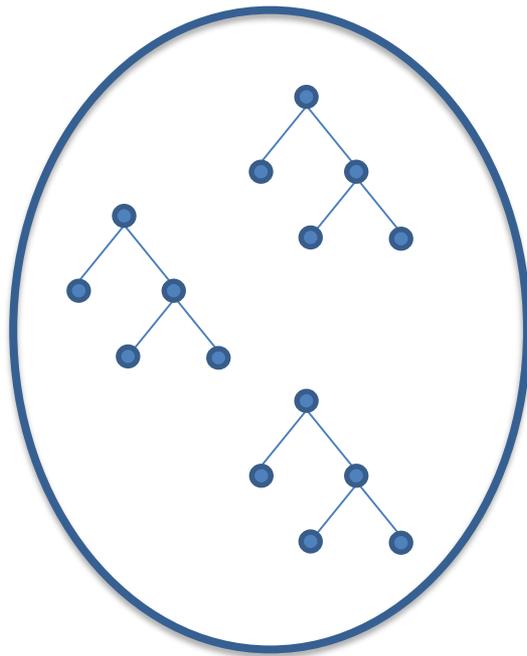
ARBRES DE DÉCISION

- Suite de règles définissant la classe d'un document:



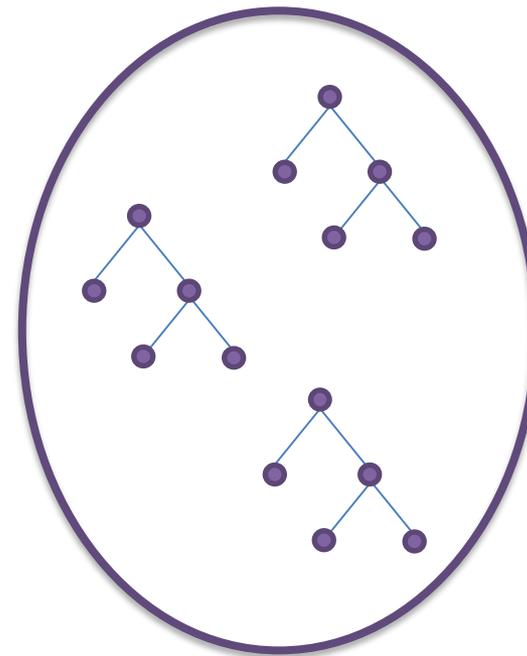
CLASSIFIEUR DE MÉTHODES DE RECHERCHE

MODÈLE FINAL



Caractéristiques
formées par les **termes**

+



Caractéristiques
formées par les **concepts**

CLASSIFIEUR DE MÉTHODES DE RECHERCHE

RÉSULTATS

Catégorie	Précision	Sensibilité	Spécificité	Taux de succès
Titres et résumés de texte				
Scientifique	81.4 %	85.2 %	90.4 %	88.7 % (+32.1)
Non scientifique	92.5 %	90.4 %	85.2 %	
Textes complets				
Scientifique	86.3 %	85.4 %	93.3 %	90.7 % (+34.1)
Non scientifique	92.8 %	93.3 %	85.4 %	

CLASSIFIEUR DE MÉTHODES DE RECHERCHE

LA SUITE

- **Accessibilité**
 - Outil exécutable et exploitable par les libraires et les chercheurs
 - Interface conviviale sur le web
- **Recherche accentuée**
 - Apprentissage profond
 - Collection de données étendue
 - Accès centralisé et automatisé aux textes intégraux
 - ...

Merci!