



Techniques d'accès à l'information

Présentation

- o Adrian.Chifu@irit.fr
- o Thèse
 - o Datamining et apprentissage l'adaptation des modèles de recherche
 - o Allocation ministérielle de recherche + ATER UPS
 - o Activité:
 - o Laboratoire: IRIT
 - o DCCE (Doctorant Contractuel Chargé d'Enseignement) + ATER
- o Thématiques de recherche
 - o Recherche d'information
 - o Apprentissage
 - o TAL (Traitement Automatique des Langues)
 - o Bases de données

Organisation du cours

- o Volume horaire
 - o 16h CM/TD
 - o 8h TP par groupe (discussion)
- o Modalités de contrôle des connaissances
 - o Examen/Evaluation de TP (30%)
 - o Examen final (70%)

Introduction

- o Objectifs du cours
 - o Présenter les notions et concepts de TAI
 - o Comment permettre à un individu d'accéder à l'information pertinente (répondant à ses besoins)?
 - o Acquérir des bases pour accéder plus efficacement à l'information pertinente
 - o Où et comment chercher?

Introduction (1/2)

- o Thèmes abordés
 - o Traitement documentaire
 - o Quelques définitions
 - o Représenter des documents (indexation)
 - o Requêtes (besoin d'information)
 - o Sélectionner les documents pertinents (recherche)
 - o Outils d'accès à l'information
 - o Moteur de recherche
 - o Systèmes de recommandation
 - o Outils de classification

Introduction (2/2)

- o Thèmes abordés
 - o Accéder à l'information
 - o Méthodologie de recherche
 - o Utilisation avancée des moteurs de recherche
 - o Visibilité sur le Web
 - o Bases du référencement
 - o Exploiter les médias sociaux (plan professionnel et plan personnel)



Le Traitement
Documentaire

Points abordés

- o Quelques définitions
 - o Notions: documents; collection; requête
 - o L'accès à l'information
- o Représenter les documents
 - o Indexation manuelle
 - o Indexation automatique
 - o Stockage et représentation
- o Sélectionner les documents pertinents

Documents, collections, requêtes

- o Document (unité documentaire)
 - o Unité minimale répondant aux besoins de l'utilisateur
 - o Article, page Web, cours, ligne, ...
- o Collection (fond documentaire)
 - o Ensemble de documents
 - o => Collection = { document }
- o Requête
 - o Demande, besoin d'information, mots clés
 - o « prix Nokia Lumia », requêtes SQL, ...
- o ! Pas forcément du texte !

Types d'informations

- o Distinguer le fond (contenu) et la forme
- o Le contenu
 - o Texte
 - o Graphique
 - o Image
 - o Son
 - o Vidéo
 - o Référence, lien hypertexte
 - o ...

Types d'informations

- o La forme
 - o 3 types
 - o structurée, non structurée, semi structurée
 - o Structurée
 - o Données factuelles sous forme tabulaire ou agrégative
 - o Données élémentaires typées
 - o Traite des questions précises
 - o Langage strict
 - o Peu lisible pour l'humain
 - o Interprétable par un système informatique
 - o Exemples
 - o Bases de données relationnelles
 - o Code source

Types d'informations

o C/C++:

```
int somme, a, b;  
somme = a+b;
```

o Requête SQL:

```
CREATE TABLE facture (numero  
Integer, date Date, montant  
Double);
```

Types d'informations

- o La forme
 - o Non structurée
 - o Autres informations
 - o Structure
 - o Non définie
 - o Trop complexe pour être stockés dans une base de données
 - o Exemple
 - o Article de journal

Types d'informations

- o La forme
 - o Semi-structurée
 - o Pas de schéma
 - o Auto descriptives
 - o Pas de séparation entre la donnée et son type
 - o L'interprétation d'une donnée est faite indépendamment de toute autre information
 - o Exemple
 - o Document XML

Types d'informations

```
<?xml version="1.0" encoding="ISO-8859-1"?>  
<annuaire>
```

```
<personne class = "etudiant">  
<nom>Pillou</nom>  
<prenom>Jean-Francois</prenom> <telephone>555-  
123456</telephone>  
<email>jfpillou@toto.fr</email>  
</personne>
```

```
<personne>  
...  
</personne>  
</annuaire>
```

Types d'informations

- o Les métadonnées
 - o Informations complémentaires sur le document
 - o Contexte
 - o Exemples
 - o Date de rédaction, de publication
 - o Auteur(s)
 - o Support de publication
 - o ...

Accès à l'information pertinente

- o Permettre l'accès à l'information
 - o Fournir des réponses pertinentes, même si elles sont approximatives voire incomplètes, à des requêtes (besoins d'information) exprimées en langage naturel

Représenter les documents

- o Comment représenter les documents?

La CSG pour financer la dépendance ?

lefigaro.fr Mis à jour le 19/09/2012 à 08:36 | publié le 19/09/2012 à 08:01

🗨 Réactions (2)

 J'aime

0

 Tweeter

10

 +1

 Share

 Recommander

Le budget 2013 ne comportera «aucune disposition» visant spécifiquement les retraités, a promis Jean-Marc Ayrault ce matin sur RTL. Le premier ministre répondait à une question sur la suppression de l'abattement de 10% sur les revenus pour «frais professionnels», dont bénéficient aussi les retraités, et sur l'alignement du taux de CSG payé par les retraités sur celui des actifs.

En revanche, Jean-Marc Ayrault a laissé entendre que ces mesures, **préconisées par la Cour des comptes**, pourraient être prises dans le cadre d'une réforme de la dépendance. «Le jour où nous ferons cette réforme - et nous le ferons - il faudra voir comment on la finance», a déclaré le chef du gouvernement, en appelant à la «solidarité intergénérationnelle».

Le processus d'indexation

- o Construire une représentation numérique des documents
 - o Extraire des caractéristiques sur le contenu sémantique
 - o Phase primordiale pour la qualité de l'accès à l'information
 - o Analyser le contenu textuel
 - o Sélection de termes significatifs/représentatifs du contenu (contenu sémantique)
 - o => Langage d'indexation

Le processus d'indexation

- o Langage d'indexation
 - o Libre
 - o Construit au moment de l'indexation
 - o Termes directement extraits des textes
 - o Peut inclure toute la variété du langage naturel
 - o Contrôlé
 - o Construit généralement avant la phase d'indexation
 - o Termes prédéfinis et limités
- o Mixte

Indexation manuelle

- o Réalisée par des documentalistes
 - o Experts
 - o Caractériser au mieux les idées contenues dans les documents
 - o Effort cognitif et intellectuel important
- o Indexation
 - o Performante
 - o Mais subjective
 - o Dépend des connaissances des documentalistes
 - o Très couteuse

Indexation automatique

- o Indexation

- o Applicable à de grandes collections
- o Rapide
- o Limite les représentations des documents aux « entrées » (termes) utiles

Indexation automatique

- o Vue globale de l'indexation automatique
 - o 3 étapes
 - o Extraction des mots simples
 - o Normalisation des mots extraits
 - o Troncature
 - o Lemmatisation
 - o Radicalisation
 - o Pondération des mots normalisés

Indexation automatique

- o Segmenter le texte
 - o Extraire les mots
 - o Mots simples
 - o Groupe de mots ou de symboles (idéogrammes)
 - o Eliminer les signes de ponctuation

Le budget 2013 ne comportera « aucune disposition » visant spécifiquement les retraités a promis Jean-Marc Ayrault
 - o => Texte = { mot }

Le	budget	2013	ne	comportera	aucune	disposition	visant	spécifiquement
----	--------	------	----	------------	--------	-------------	--------	----------------

les	retraités	a	promis	Jean-Marc	Ayrault
-----	-----------	---	--------	-----------	---------

Indexation automatique

- o Eliminer les mots-vides
 - o Anti dictionnaire
 - o Mots n'ayant qu'un rôle syntaxique (pas de sens)
 - o Mots ordinaires ou athématique
 - o Exemples
 - o « de », « à » en français
 - o En anglais, liste de plus de 275 mots

budget	2013	comportera	aucune	disposition	visant	spécifiquement
--------	------	------------	--------	-------------	--------	----------------

retraités	promis	Jean-Marc	Ayrault
-----------	--------	-----------	---------

Indexation automatique

- o Normaliser les mots extraits
 - o Désaccentuer

Le	budget	2013	ne	comportera	aucune	disposition	visant	specifiquement
----	--------	------	----	------------	--------	-------------	--------	----------------

les	retraites	a	promis	Jean-Marc	Ayrault
-----	-----------	---	--------	-----------	---------

- o Minusculiser

le	budget	2013	ne	comportera	aucune	disposition	visant	spécifiquement
----	--------	------	----	------------	--------	-------------	--------	----------------

les	retraités	a	promis	jean-marc	ayrault
-----	-----------	---	--------	-----------	---------

Indexation automatique

o Lemmatisation

- o Verbe: ce verbe à l'infinif
- o Autre mot: le mot au masculin singulier
 - o Petit, petite, petits, petites => petit
 - o Avoir, ai, a, avons, eu, eussions => avoir

budget	2013	comporter	disposition	viser	specifique
--------	------	-----------	-------------	-------	------------

retraite	avoir	promettre	jean-marc	ayrault
----------	-------	-----------	-----------	---------

Indexation automatique

- o Troncature
 - o Couper le mot au-delà d'un certain nombre de caractères
 - o Troncature à 6

le	budget	2013	ne	compor	aucune	dispos	visant	specif
----	--------	------	----	--------	--------	--------	--------	--------

les	retrai	a	promis	jean-m	ayraul
-----	--------	---	--------	--------	--------

- o Troncature à 4

le	budg	2013	ne	comp	aucu	disp	visa	spec
----	------	------	----	------	------	------	------	------

les	retr	a	prom	jean	ayra
-----	------	---	------	------	------

Indexation automatique

- o Radicalisation
 - o Forme morphologique d'un mot

Préfixe	Radical	Désinence
----------------	----------------	------------------

- o Radicalisation

Le	budget	2013	ne	comport	aucun	disposit	vis	spécifiqu
----	--------	------	----	---------	-------	----------	-----	-----------

les	retraité	a	promi	Jean-Marc	Ayrault
-----	----------	---	-------	-----------	---------

Indexation automatique

- o Conséquences liées à la normalisation