

Comment exploiter les commentaires d'internautes pour la recommandation automatique

Damien Poirier

Paris, le 11 juin 2012

Contexte et problématique

Description de la méthodologie

- Extraction des textes

- Inférence de notes sur les textes

- Recommandation par filtrage collaboratif

Évaluation

Conclusion

Contexte et problématique

Description de la méthodologie

Extraction des textes

Inférence de notes sur les textes

Recommandation par filtrage collaboratif

Évaluation

Conclusion

Contexte

Mise en œuvre d'un système de recommandation personnalisée



Deux choix : filtrage collaboratif ou basé sur le contenu.

Contexte

Quelle que soit l'approche utilisée, des données sont nécessaires :

- ▶ Données d'usages pour le filtrage collaboratif
- ▶ Données descriptives pour le filtrage basé sur le contenu

Que faire lorsqu'elles ne sont pas disponibles ?

Contexte

Le Web regorge de ressources, principalement textuelles, qui contiennent un grand nombre d'avis et d'opinions d'utilisateurs sur des sujets de toutes sortes.

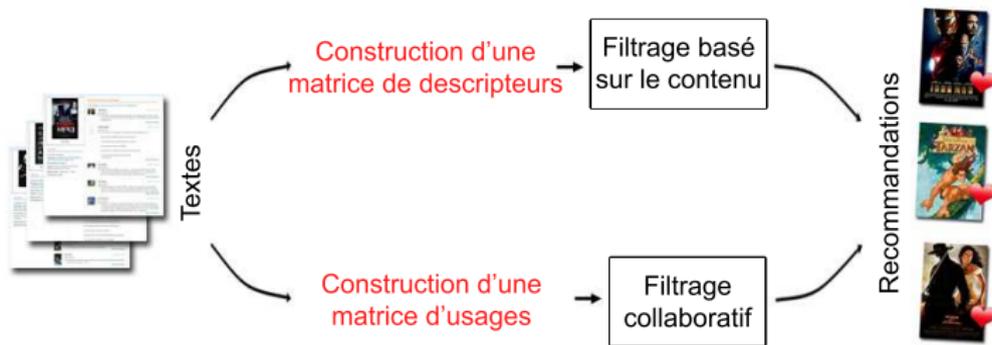
(blogs, forums, boîtes à réactions, etc.)



Problématique

Objectif principal

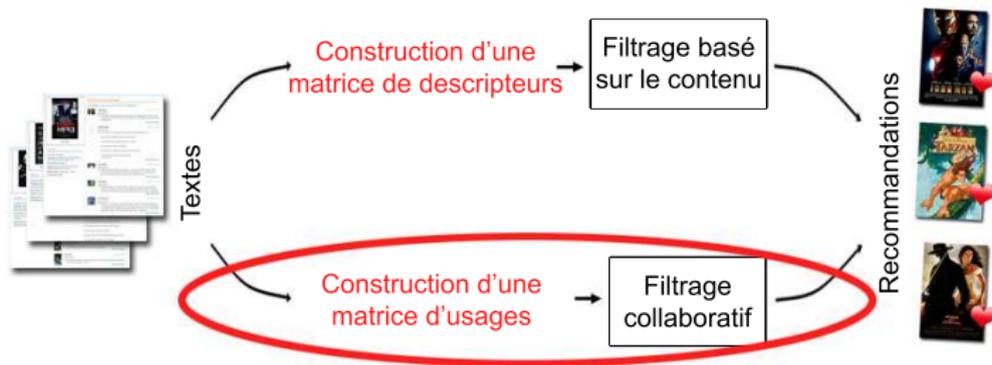
Montrer que l'on peut extraire, à partir de textes communautaires non structurés, des informations utiles aux moteurs de recommandation existants.



Problématique

Objectif principal

Montrer que l'on peut extraire, à partir de textes communautaires non structurés, des informations utiles aux moteurs de recommandation existants.



Contexte et problématique

Description de la méthodologie

Extraction des textes

Inférence de notes sur les textes

Recommandation par filtrage collaboratif

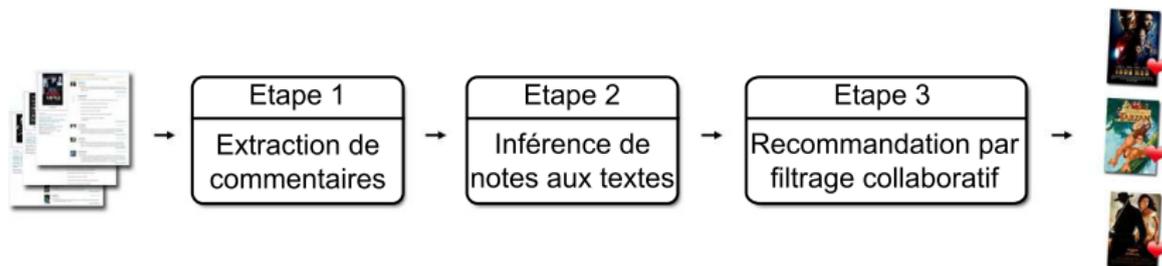
Évaluation

Conclusion

Méthodologie

Principale contribution :

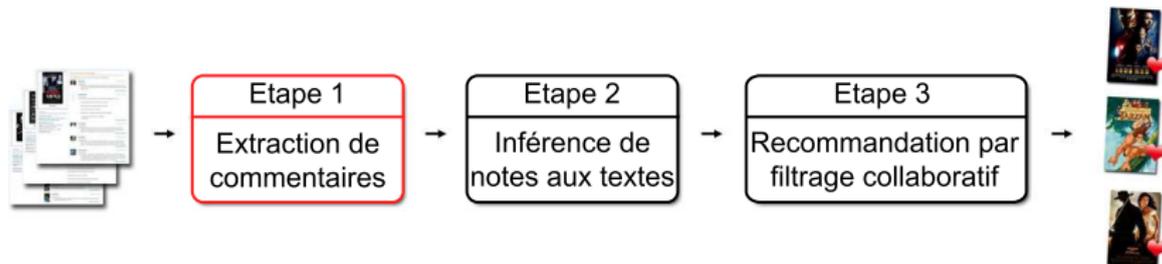
Mise en place d'une méthodologie complète, **automatisée** et en grande partie **reproductible** (étapes 2 et 3)



1^{ère} étape : Extraction des textes

Objectif de la tâche :

Acquérir des textes exprimant des opinions sur les produits de notre catalogue



Provenance des commentaires

The screenshot shows the Flixster website interface. At the top, there is a navigation bar with the Flixster logo, a search bar, and links for Sign Up and Login. Below the navigation bar are tabs for Home, Movies, Profile, Friends, Meet People, Fun & Games, and Watch Now. The main content area features a movie poster for 'Dude, Where's My Car?' on the left. To the right of the poster is a section titled 'Dude, Where's My Car? Reviews and Ratings'. This section contains a list of user reviews, each with a profile picture, a star rating, a date, and a 'Share This Review' link. The reviews are as follows:

- iluvspastley13** (January 26, 2011): 4 stars, "Really weird and random, but funny".
- JuRn** (January 5, 2011): 4 stars, "I got quite a few laughs from dis good comedy movie".
- TheFilmApache** (December 5, 2010): 1 star, "The worst comedy I have ever seen...".
- Quasim0** (November 28, 2010): 1 star, "This movie is pointless and a waste of time".
- incz** (November 25, 2010): 4 stars, "silly but still funny! :D".
- Superman1984** (November 8, 2010): 4 stars, "I Liked the part at near the end. Daddy I wanna ride that ride, me too son."

Below the movie poster, there is a 'Summary' section with the following text:

Summary

Dude, Where's My Car? Summary

Starring: Ashton Kutcher, Seann William Scott, Kristy Swanson, Jennifer Garner, Marla Sokoloff

Directed by: Danny Leiner

Genres: Comedy

Release Date: December 15, 2000

DVD Release Date: June 26, 2001

Provenance des commentaires

The screenshot shows the Flixster website interface. At the top, there is a navigation bar with 'Home', 'Movies', 'Profile', 'Friends', 'Meet People', 'Fun & Games', and 'Watch Now'. A search bar is located on the right. The main content area features a movie poster for 'Dude, Where's My Car?' on the left and a list of reviews on the right. The reviews list includes user avatars, usernames, star ratings, and review text. Red arrows originate from a red oval at the bottom right and point to specific elements in the reviews: the movie title, the reviewer's name, the review text, and the review date.

Données extraites :
Id Utilisateur - Id Film - Commentaire - Note

Movie Information:
Dude, Where's My Car?
Summary
Dude, Where's My Car? Summary
Starring: Ashton Kutcher, Seann William Scott, Kristy Swanson, Jennifer Garner, Marla Sokoloff
Directed by: Danny Leiner
Genres: Comedy
Release Date: December 15, 2000
DVD Release Date: June 26, 2001

Reviewer	Rating	Comment	Date
iluvspastley13	4/4	Really weird and random, but funny	January 26, 2011
JuRn	4/4	I got quite a few laughs from dis good comedy movie	January 5, 2011
TheFilmApache	4/4	The worst comedy I have ever seen...	December 5, 2010
suasim0	4/4	This movie is pointless and a waste of time	November 28, 2010
incz	4/4	silly but still funny! :D	
Superman1984	4/4	I Liked the part at near the end. Daddy I wanna ride that ride. me too son.	November 8, 2010

Description du corpus

Chiffres :

- ▶ Nombre d'entrée : > 3 millions de commentaires
(sous la forme de 4-uplets Id utilisateur-Id film-Texte-Note)
- ▶ Utilisateurs concernés : \simeq 100 000
- ▶ Films concernés : \simeq 10 000

Particularités des textes :

- ▶ Nombre moyen d'unités (mots + ponctuations + symboles + chiffres, etc.) par commentaire : 15
- ▶ Emploi du **Discours électronique Médié** (DEM) (Panckhurst¹)

Exemples :

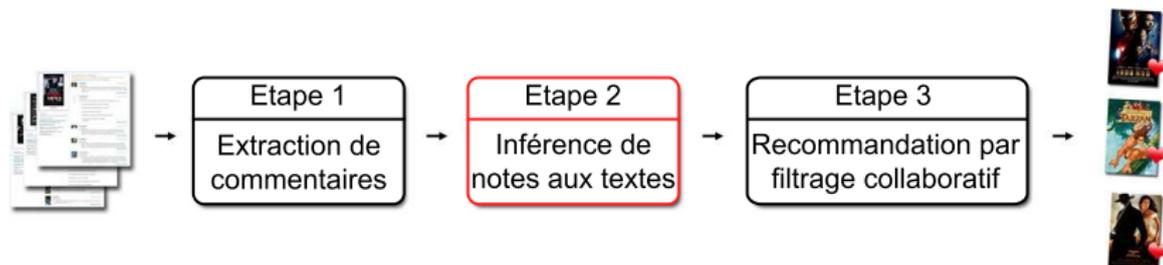
" weeeeeiiiiiiiiirrrrrddddd movie "
" i liked this cuz this shit can really happen. "
" I see it @ work everyday so this movie sucks! "
" OMG I LOVED THIS MOVIE!!!!!!!!!!!!!!!!!!!!<3 "

1. Le discours électronique médié : bilan et perspectives, 2006

2^{ème} étape : Inférence des notes

Objectif de la tâche :

Attribuer des notes aux textes afin d'obtenir une matrice d'usages



Classification d'opinion

Objectif de la classification d'opinion

Classer les textes selon l'opinion exprimée par l'auteur
(positif, négatif ou neutre)

Deux grands types de méthodes (Pang et al.²) :

- ▶ Les méthodes basées sur les lexiques
- ▶ Les méthodes basées sur l'apprentissage automatique

2. Opinion mining and sentiment analysis, *Found. Trends Inf. Retr.*, 2008

Classification d'opinion : 1^{ère} approche

Méthodes basées sur les lexiques

Objectif : Construire des lexiques contenant des mots liés à l'expression de l'opinion et classer les textes en fonction de la présence de ces mots.

Exemples :

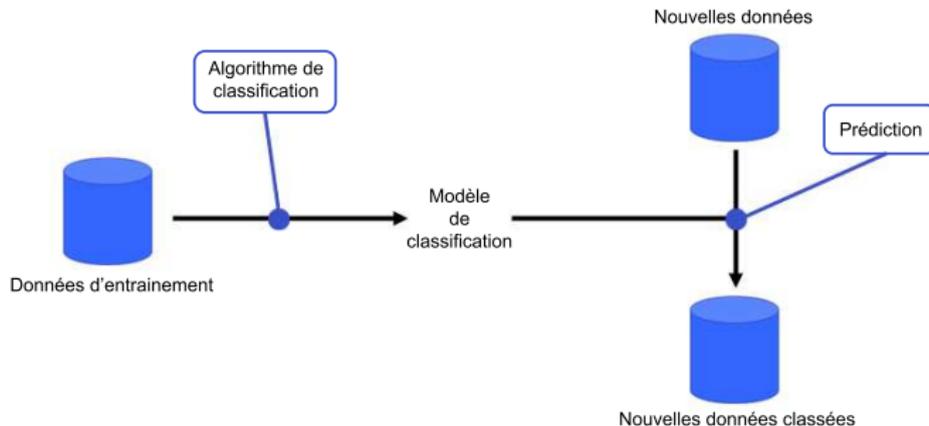
Mots positifs	awesome, great, best, love, ...
Mots négatifs	dislike, awful, hate, bad, ...

⇒ Trop de contraintes manuelles pour l'automatisation des traitements

Classification d'opinion : 2^{ème} approche

Méthodes basées sur l'apprentissage automatique

Objectif de la classification supervisée : apprendre des modèles à partir d'exemples et classer les nouveaux textes à l'aide de ces modèles.



Méthodes basées sur l'apprentissage automatique

Plusieurs choix à plusieurs niveaux :

Pré-traitements	Représentation
Minusculation	Binaire
Lemmatisation	Fréquentielle
Stemmatisation	Fréquentielle normalisée
Correction	TF-Idf
...	...
Classifieur	Nombre de classes
SVM	2
Naïve Bayes	3
Réseaux de neurones	5
Règles de décision	Échelle continue
...	...

Méthodes basées sur l'apprentissage automatique

Conditions expérimentales :

- ▶ Construction d'un corpus d'apprentissage :
 - ▶ 175 000 nouveaux commentaires déjà classés
 - ▶ Classes équilibrées
 - ▶ Aucune intersection au niveau des films et des utilisateurs
 - ▶ Pas plus de 10 commentaires rédigés par le même auteur
- ▶ Mesure utilisée :

$$F_{score} = 2 * \frac{Précision * Rappel}{Précision + Rappel}$$

- ▶ Meilleur résultat obtenu :

(Pré-traitements minimaux, représentation fréquentielle, SVM,

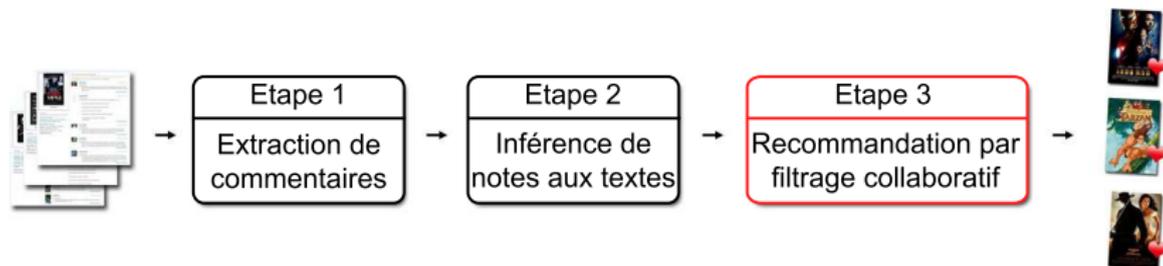
classification binaire \implies

$$F_{score} = 0,74$$

3^{ème} étape : Recommandations

Objectif de la tâche :

Établir des recommandations à partir de la matrice d'usages



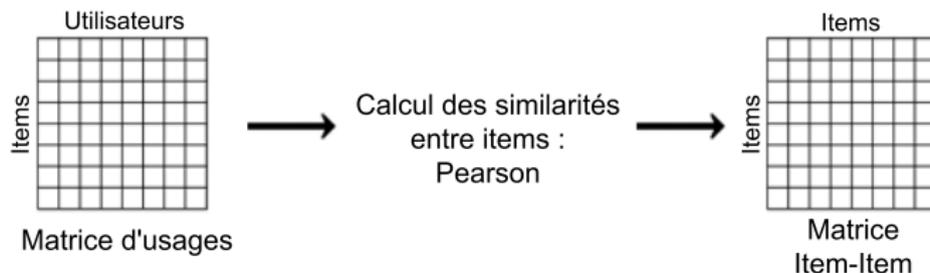
Description du moteur

- ▶ Moteur développé chez Orange Labs (Meyer³)
- ▶ Permet le filtrage collaboratif et basé sur le contenu
- ▶ Permet la recommandation contextuelle et personnalisée
- ▶ Performances au niveau de l'état de l'art

⇒ Fonctionnement en 2 étapes

3. Apport des données thématiques dans les systèmes de recommandation : hybridation et démarrage à froid, *In Extraction et gestion des connaissances (EGC'2011), Brest, France, 2011*

Construction de la matrice item-item



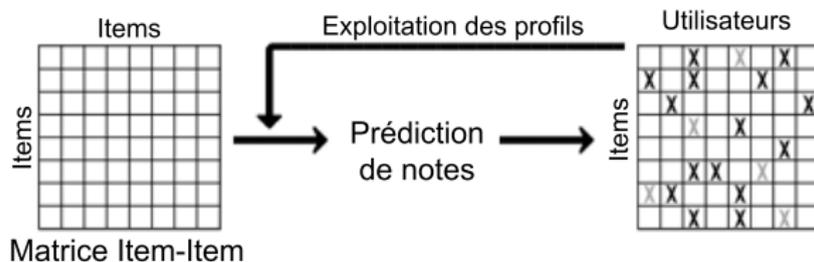
$$Pearson(i, j) = \frac{\sum_{\{u \in S_i \cap S_j\}} (r_{iu} - \bar{r}_i) \times (r_{ju} - \bar{r}_j)}{\sqrt{\sum_{\{u \in S_i \cup S_j\}} (r_{iu} - \bar{r}_i)^2 \sum_{\{u \in S_i \cup S_j\}} (r_{ju} - \bar{r}_j)^2}}$$

S_i (resp. S_j) : l'ensemble des utilisateurs ayant noté l'item i (resp. j)

r_{iu} (resp. r_{ju}) : la note donnée par l'utilisateur u sur l'item i (resp. j)

\bar{r}_i (resp. \bar{r}_j) la moyenne des notes obtenues par i (resp. j)

Prédiction des notes



$$p_{ui} = \bar{r}_i + \frac{\sum_{\{j \in S_u\}} \text{Pearson}(i, j) \times (r_{uj} - \bar{r}_j)}{\sum_{\{j \in S_u\}} \text{Pearson}(i, j)}$$

S_u : l'ensemble des notes connues données par l'utilisateur u

Contexte et problématique

Description de la méthodologie

Extraction des textes

Inférence de notes sur les textes

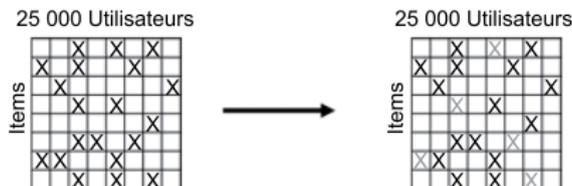
Recommandation par filtrage collaboratif

Évaluation

Conclusion

Conditions expérimentales

Données de test : extraites du challenge Netflix (Bennett et al.)⁴



X : 90% des notes dédiées aux profils utilisateurs

X : 10% des notes dédiées à la validation des notes prédites

Données d'entrée :

- ▶ Données d'usages obtenue à partir des commentaires **Flixster**
 - ▶ Notes issues de la classification d'opinion
 - ▶ Vraies notes
- ▶ Descripteurs **IMDb** pour le filtrage basé sur le contenu
- ▶ Données d'usages **Netflix** (profils) pour le filtrage collaboratif

4. The netflix prize, *San Jose, California, ACM, 2007*

Méthode d'évaluation

Mesure utilisée : Erreur quadratique moyenne

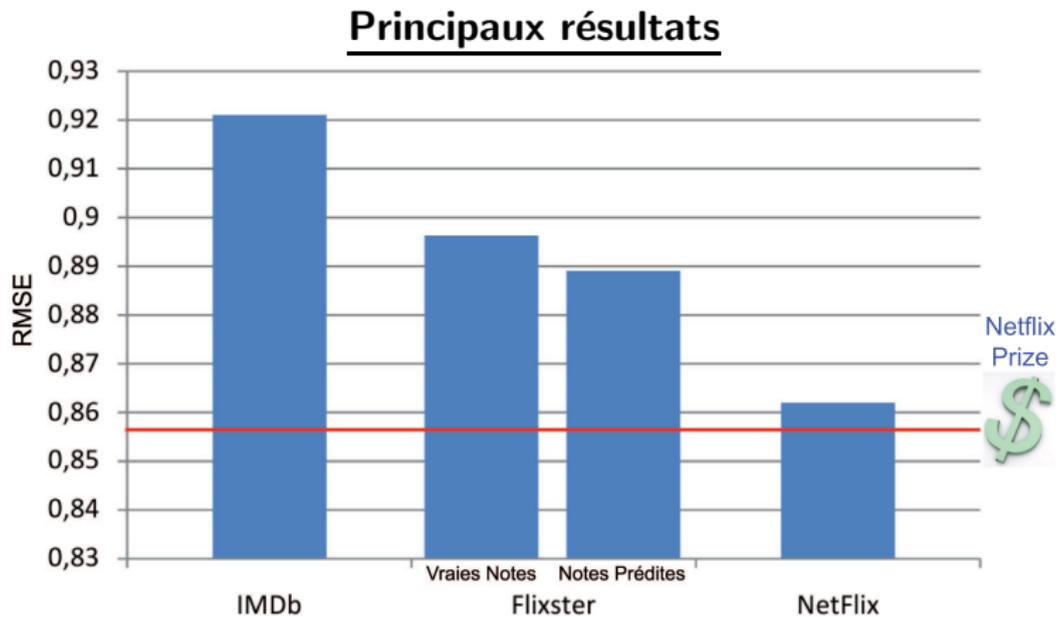
$$RMSE = \sqrt{\frac{\sum_{u,i} (p_{ui} - n_{ui})^2}{n}}$$

n_{ui} = note réelle donnée par l'utilisateur u sur l'item i

p_{ui} = note prédite par le moteur de recommandation

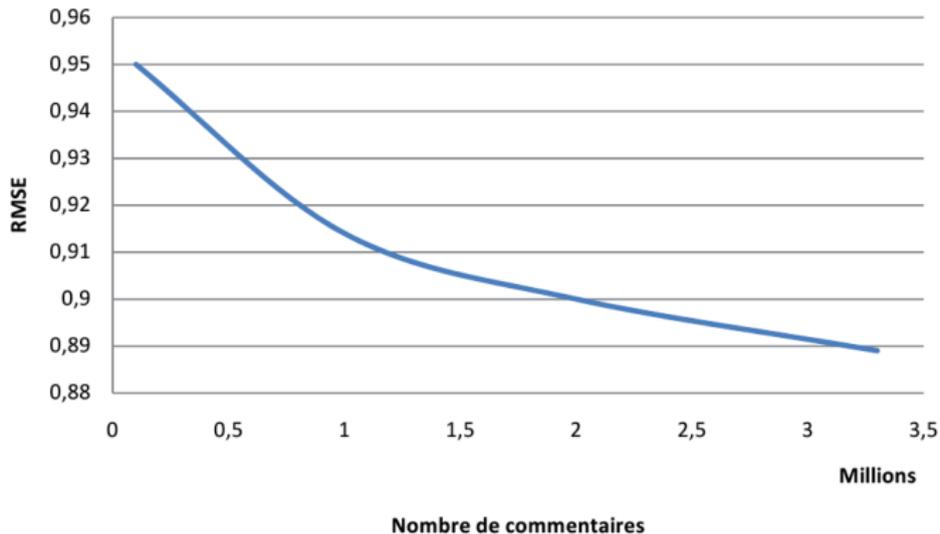
n = nombre total de notes prédites

Résultats



Résultats

Une amélioration possible et simple :
Augmenter le nombre de commentaires présents dans la base



Contexte et problématique

Description de la méthodologie

- Extraction des textes

- Inférence de notes sur les textes

- Recommandation par filtrage collaboratif

Évaluation

Conclusion

Objectif atteint !

Démonstration faite que les textes non-structurés issus de blogs peuvent pallier le manque de données en recommandation

- ▶ Le choix de transformer les textes en données d'usages (classification d'opinion) s'est avéré pertinent :
⇒ Meilleures performances qu'avec les descripteurs IMDb

Merci pour votre attention