

Extraction de sentiments et d'opinions basée sur des règles

Sigrid Maurel*, Paolo Curtoni*, Luca Dini*

* CELI France, SAS
12-14, rue Claude Genin
38000 Grenoble
{maurel, curtoni, dini}@celi-france.com
<http://www.celi-france.com>

Résumé. Nous présentons ici trois méthodes différentes pour effectuer une classification automatique de textes d'opinion. La première méthode est symbolique, la seconde statistique et la dernière, hybride, est une combinaison des deux premières. Nous montrons comment la combinaison des méthodes symbolique et statistique permet de tirer parti des avantages des deux méthodes, à savoir la robustesse de l'apprentissage automatique statistique et la possibilité de configuration manuelle offerte par la méthode symbolique, permettant une utilisation dans des applications réelles. Les textes classés par ces méthodes viennent de sources informationnelles non structurées de type forum sur Internet.

1 Introduction

1.1 Motivation

Cet article s'intéresse à la classification de textes d'opinion en langue française. Dans ce cas précis, la classification a pour objectif l'analyse de sentiments exprimés dans différents types de textes comme par exemple dans des forums de discussion sur Internet où les internautes échangent des avis et s'entraident. Les textes issus de forums sur Internet constituent des sources d'informations spontanées et récentes, incontournables pour acquérir, au jour le jour, des connaissances sur les consommateurs, pour anticiper leurs besoins et leurs attentes afin de tenter d'améliorer la relation client/fournisseur. En analysant ces textes d'opinion le fournisseur d'un produit ou d'un service peut mieux réagir aux desiderata de ses clients, le client peut de son côté s'inspirer des sentiments et opinions d'autres clients sur le produit auquel il s'intéresse et profiter ainsi d'une aide à la décision (acquérir ou non le produit, choisir plutôt le produit A ou le produit B, etc.).

Comme le montrent de nombreux travaux de socio- et psycho-linguistique (Sproull et Kiesler, 1991), la communication médiée par ordinateur favorise l'expression des émotions, sentiments et opinions souvent contrôlés ou réprimés dans des cadres de communication plus traditionnels visant à étudier le point de vue des consommateurs (interviews face à face, enquêtes fermées, enquêtes ouvertes, etc.). De là, naît l'intérêt des analystes pour ces sources d'informations.

Extraction de sentiments et d'opinions

Les corpus utilisés pour le développement des systèmes de classification sont composés de textes (ou *threads*, fils de discussion) provenant de forums sur Internet qui parlent entre autres de tourisme, de jeux vidéo et d'imprimantes. Un texte (ou message) dans un forum contient un jugement argumenté de l'auteur du message, positif, négatif ou parfois mitigé, sur un sujet donné. Mais il contient aussi des parties exemptes de sentiments, comme c'est le cas par exemple dans la description du jeu vidéo sur lequel porte la critique. L'objectif de l'analyse est donc d'identifier avec précision les parties pertinentes pour la classification automatique du texte dans son entier.

Une des difficultés de la classification en *positif* et *négatif* réside dans la nécessité d'une bonne analyse syntaxique du texte, analyse qui peut se révéler particulièrement difficile dans des cas de coordination entre plusieurs parties d'une phrase, d'anaphore ou de coréférence (la reprise d'un argument présent plus loin dans le document). Une autre difficulté du langage naturel pour l'analyse automatique de sentiments réside dans les contextes intentionnels, pour lesquels l'expression d'opinion n'est pas un vrai sentiment. C'est le cas dans une phrase comme :

« Je croyais que la France était un beau pays. »

(Dini et Mazzini, 2002) ont montré le lien qui existe entre les structures syntaxiques et sémantiques d'une phrase et l'expression de l'opinion qu'elle véhicule. Ainsi l'analyse de la phrase par *paquets de mots* donne des résultats peu satisfaisants alors qu'une analyse syntaxique du texte peut aider à trouver les expressions qui contiennent des opinions. Les deux phrases suivantes contiennent les mêmes *paquets de mots* sans pour autant exprimer les mêmes sentiments. En effet, la première phrase contient un sentiment positif alors que la deuxième est négative :

« Je l'ai apprécié pas seulement à cause de ... »

« Je l'ai pas apprécié seulement à cause de ... »

Dans le cadre de sa participation à la campagne d'évaluation DEFT'07 (c.f. section 6 pour plus de détails), CELI France a mis au point trois méthodes pour classer les textes des différents corpus. La première est une méthode symbolique qui inclut un système d'extraction d'information adapté aux corpus. Elle est basée sur des règles d'un analyseur syntaxico-sémantique. Cet analyseur contient un lexique de mots qui véhiculent des sentiments sur lesquels réagissent les règles de la grammaire. La deuxième est une méthode statistique basée sur des techniques d'apprentissage automatique. Enfin, la dernière, SYBILLE, est une méthode hybride qui combine les techniques des deux précédentes pour aboutir à des résultats très précis.

L'analyse des textes se fait au niveau de la phrase, les sentiments d'un document sont extraits phrase par phrase, et c'est seulement ensuite qu'une valeur globale est attribuée au message entier. Ceci permet d'extraire une information contextuelle qui est donc très précise.

Les sections suivantes présentent brièvement l'état de l'art et les corpus utilisés pour s'attarder ensuite sur les trois méthodes développées et en fournir une première évaluation. Une section sera dédiée à l'interface graphique SYBILLE pour présenter les possibilités que celle-ci offre aux utilisateurs, avant de conclure cet article. Nous donnons en annexe des exemples et des extraits de la grammaire utilisée, ainsi qu'un schéma général du processus du traitement.

1.2 État de l'art

Aujourd'hui l'analyse de sentiments se concentre sur l'attribution d'une polarité à des expressions subjectives (les mots et les phrases qui expriment des opinions, des émotions, des sentiments, etc.) afin de décider de l'orientation d'un document (Turney, 2002; Wilson et al., 2004) ou de la valeur positive/négative/neutre d'une opinion dans un document (Hatzivassiloglou et McKeown, 1997; Yu et Hatzivassiloglou, 2003; Kim et Hovy, 2004).

Des travaux allant au-delà ont mis l'accent sur la force d'une opinion exprimée où chaque proposition dans une phrase peut avoir un fond neutre, faible, moyen ou élevé (Wilson et al., 2004). Des catégories grammaticales ont été utilisées pour l'analyse de sentiments dans (Bethard et al., 2004) où des syntagmes adjectivaux comme *trop riche* ont été utilisés afin d'extraire des opinions véhiculant des sentiments. (Bethard et al., 2004) utilisent une évaluation basée sur la somme des scores des adjectifs et des adverbes classés manuellement, tandis que (Chklovski, 2006) utilise des méthodes fondées sur un modèle pour représenter des expressions adverbiales de degré telles que *parfois*, *beaucoup*, *assez* ou *très fort*.

L'approche que nous avons adoptée pour la classification de textes d'opinion est caractérisée par une utilisation mixte d'une technologie symbolique fondée sur des règles et d'une technologie statistique reposant sur l'apprentissage automatique, approche dans laquelle la méthode symbolique a un poids plus important (Dini, 2002; Dini et Mazzini, 2002; Maurel et al., 2007, 2008). La technologie symbolique fait d'abord une analyse du texte phrase par phrase et en extrait ensuite les relations qui véhiculent des sentiments, tandis que la technologie statistique traite les textes en une seule phase et attribue un sentiment global au texte entier à la fin du traitement.

Il convient de remarquer que, contrairement à d'autres approches actuelles, la technologie de l'analyse de sentiments développée à CELI France (SYBILLE) ne se limite pas à une analyse lexicale (c'est-à-dire identification et pondération de mots positifs et négatifs), mais s'étend à une analyse syntaxique et sémantique. L'analyse syntaxique est effectuée par le biais d'une analyse robuste de surface telle que celles décrites par (Aït-Mokhtar et Chanod, 1997; Basili et al., 1999; Aït-Mokhtar et al., 2001), donnant ainsi un résultat très proche de celui produit par des grammaires de dépendance.

2 Les corpus

Les données de type forums de discussion sur Internet s'articulent comme un flux d'interactions, comme par exemple : demande-réponse, argument-contre argument, commentaire-désaccord, etc. Ce flux est distribué sur une dimension temporelle qui nécessite un traitement chronologique du fil de discussion. Contrairement aux corpus utilisés par (Wilson et al., 2004), il n'est pas nécessaire ici d'identifier la personne à qui est associé un sentiment, car dans 95 % des cas, les discours analysés sont des discours à la première personne. Un exemple de flux d'interactions est donné en figure 1.

Les corpus utilisés sont assez différents les uns des autres, que ce soit par la taille des corpus eux-mêmes que par la taille de chaque *thread* (fil de discussion). Nous avons utilisé les corpus de DEFT'07 auxquels nous avons ajouté des corpus collectés sur Internet. Ces corpus nous ont permis d'augmenter la diversité des sujets et de répondre aux exigences de nos clients, selon les domaines demandés. Nous avons donc des textes de domaines très différents, entre autres

Extraction de sentiments et d'opinions

Avis sur les châteaux de la Loire en France

angie-443*5, posté le 08-10-2006 à 16:18:50:
J'ai besoin de vos conseil s.v.p. Je vais passer une ou deux journée dans la vallée de la Loire. Y-a-t-il un château en Loire avec un jardin semblable à celui de Versailles (en beauté et en superficie)? J'aime aussi l'aspect extérieurs des châteaux, plus que l'intérieur. Ce qui me plaît d'une ville est tout d'abord ses rues piétonnes, animées et pittoresques, ses charmantes places et ses promenades.

[...]

BaLadeur, posté le 13-10-2006 à 11:23:43:
Je partage l'avis d'Aston sur de nombreux points. Villandry est quelconque mais son jardin transformé en potager géant vaut le détour. Chenonceau est certainement le plus photogénique donc le plus connu et il le mérite largement. Si tu recherches la monumentalité comme à Versailles, la magnificence en plus, il faut absolument voir Chambort. Enfin s'il faut ne visiter qu'une ville ce sera Tours.

[...]

zeus77, posté le 21-10-2006 à 21:59:33:
A Amboise j'aime beaucoup le manoir du Clos-Lucé qui fut la dernière maison de Léonard de Vinci. Le parc est très agréable. Enfin un château où l'on pourrait vivre! Quel changement par rapport aux châteaux royaux. Un château que j'aime bien aussi c'est celui Du Moulin à Lassay sur Croisne entre Contres et Romorantin.

[...]

FIG. 1 – Exemple d'un flux d'interactions de messages, du domaine du tourisme. L'orthographe et la ponctuation n'ont pas été modifiées.

du domaine de la restauration rapide, du nucléaire, de l'alimentation infantile, etc.¹ Certains corpus sont structurés, d'autres contiennent beaucoup de messages en style *texto*, et le nombre de fautes d'orthographe présentes dans les messages varie aussi beaucoup.

Le comité d'organisation de DEFT'07 a pris soin de nettoyer ses corpus (Grouin et al., 2007). Ainsi, les fins de ligne ont été normalisées, les caractères encodés en ISO-Latin, et les textes ont été annotés manuellement.² Les corpus de DEFT'07 contiennent des critiques de films, de livres et de spectacles, des tests de jeux vidéo, des relectures d'articles scientifiques (de différentes conférences sur l'intelligence artificielle) et des notes de débats parlementaires (sur la loi de l'énergie).

Les textes de nos corpus portent essentiellement sur le tourisme (en France et ailleurs dans le monde), les jeux vidéo (critiques et problèmes) et les imprimantes (conseils d'achats). Ils comprennent d'un côté des aides à la solution de problèmes, mais aussi des avis sur des produits achetés et des lieux visités. Chaque *thread* contient les messages des auteurs participant aux forums sur un sujet donné.

Les fautes d'orthographe³ dans les textes des corpus posent parfois des problèmes d'analyse. Heureusement, les règles syntaxiques de la grammaire (voir la prochaine section sur la méthode symbolique) sont dans la plupart des cas assez tolérantes pour permettre l'accord entre un nom et un adjectif, ou un nom et un verbe même si le *e* ou le *s* manque. Mais malheureusement, il y a aussi des messages tellement mal écrits dans les corpus (par exemple en style *texto*, c'est-à-dire avec beaucoup d'abréviations) que l'analyse peut échouer.⁴

3 Méthode symbolique

Comme nous l'avons dit plus haut, la méthode symbolique se base sur une analyse syntaxique du texte faite par un analyseur fonctionnel et relationnel (c.f. les travaux sur l'analyse syntaxique et sémantique de (Basili et al., 1999; Aït-Mokhtar et al., 2001; Dini, 2002; Dini et Mazzini, 2002; Dini et Segond, 2007)). Cet analyseur traite, phrase par phrase, un texte donné en entrée et en extrait, pour chaque phrase, les relations syntaxiques présentes. Il s'agit de relations syntaxiques fonctionnelles de base, telles que le modifieur d'un nom, d'un verbe, sujet et objet d'une phrase, ainsi que de relations plus complexes telles que la coréférence entre deux syntagmes au sein d'une même phrase.

L'utilisateur a la possibilité d'élaborer une grammaire à sa guise et d'ajouter de nouvelles règles afin d'extraire les relations auxquelles il s'intéresse. Pour ce faire, il peut modifier les règles d'extraction de relations (par exemple ajouter des règles pour de nouvelles relations), augmenter/diminuer les traits sur les mots dans le lexique qui agissent sur les règles, enlever certaines parties du traitement, etc.

1. Ces derniers ne seront pas abordés plus profondément dans cet article, mais les ressources sont disponibles dans notre système SYBILLE.

2. En ce qui concerne nos propres corpus, ils sont encodés en UTF-8 et nous n'avons effectué aucun nettoyage. Tous les corpus sont disponibles au format XML.

3. Nous avons fait le choix de garder les textes tels quels, donc de ne pas appliquer un correcteur automatique d'orthographe ou un lexique d'abréviations. Ce choix s'explique par la volonté de garder toutes les caractéristiques stylistiques présentes dans les textes. Nous considérons qu'une uniformisation des entrées à ce moment-là du processus nous ferait perdre de l'information utile.

4. D'après ce que nous avons pu observer, ces messages sont heureusement en minorité dans les corpus et ne modifient pas les résultats de façon significative.

Extraction de sentiments et d'opinions

Un algorithme permet de calculer un indice de confiance qui servira à la méthode hybride (c.f. section 5) pour déterminer le résultat final.

La polarité positive ou négative attribuée au message entier⁵ dépend du rapport entre la quantité de relations d'opinions positives et négatives. Une majorité de relations d'opinions positives détermine une polarité positive du message, tandis qu'une majorité de relations d'opinions négatives provoque une polarité négative.

Un schéma général du processus est visualisé par le tableau 2 dans l'annexe 9.1.

3.1 Grammaire

La grammaire utilisée a été initialement développée afin d'extraire les relations de sentiments exprimés dans une phrase dans le cadre d'un projet sur le tourisme en France. Elle a été ensuite modifiée et améliorée en vue de la participation à DEFT'07 (c.f. section 6, (Maurel et al., 2007)). Dans un deuxième temps, la grammaire a été divisée en deux parties : une première partie de base (la grammaire *générique*) s'appliquant à tous les textes qui contiennent des sentiments, et une deuxième partie pour chaque domaine différent, selon le sujet du corpus : tourisme, jeux vidéo, imprimantes, etc. Les différences se situent essentiellement dans les lexiques appliqués, chaque domaine ayant ses propres mots et expressions.

Ainsi les mots se rattachant à la vitesse (*lent, rapide, etc.*) ont des polarités différentes selon qu'ils qualifient une imprimante ou un voyage. De même, comme le montrent les phrases ci-dessous, l'adjectif *effrayant* est plutôt perçu comme positif dans une description romanesque alors qu'il est perçu comme négatif dans le domaine des assurances ou du tourisme :

« Dans *Ghost*, les habitants du village sont vraiment effrayants ! »

« C'est effrayant de voir comment la côte est de plus en plus bétonnée. »

En général, une relation de sentiment a deux arguments : le premier est l'expression linguistique qui véhicule le sentiment en question, le deuxième est la cause ou l'objet du sentiment (si la cause est exprimée dans la phrase). Ceci donne pour la phrase

« J'aime beaucoup Grenoble. »

la relation SENTIMENT_POSITIF (*aimer, Grenoble*). L'attribut POSITIF de la relation, c'est-à-dire la valeur de sa classe, indique qu'il s'agit d'un sentiment positif dont l'objet est *Grenoble*. Dans le cas d'une phrase comme

« Je déteste !!!! »

la relation n'aura qu'un seul argument : SENTIMENT_NEGATIF (*détester*), dans la mesure où l'objet du sentiment n'est pas exprimé dans la phrase.

3.2 Fonctionnement de la grammaire

L'objectif de la grammaire est d'extraire le plus d'informations possible dans le *thread*, en particulier les sentiments positifs et négatifs, les lieux et produits. Pour ceci, les *threads* sont analysés phrase par phrase. Chaque phrase peut contenir zéro, une ou plusieurs relations de sentiment. Il est tout à fait possible d'avoir des relations de sentiments positifs et négatifs dans une même phrase :

5. L'attribution d'un sentiment global au message entier est utilisée dans des contextes spécifiques, comme par exemple pour l'évaluation DEFT'07 (c.f. section 6). Sinon nous n'attribuons pas de sentiment global mais gardons les sentiments attribuées à chaque phrase.

« En qualité d'impression, la Epson est meilleure, en texte comme en photo, malheureusement c'est aussi la plus chère. »

⇒ SENTIMENT_POSITIF (meilleur, Epson)

⇒ SENTIMENT_NEGATIF (cher, ce)

Les parties de la grammaire qui varient selon le corpus se distinguent essentiellement par le lexique de mots qui reçoivent les traits *positif* et *négatif* correspondant aux valeurs des classes des textes (et par leur liste de termes, c.f. section 3.5). Par exemple, le lexique de la grammaire du *tourisme* contient les mots *joli* et *beau* :

« Ce monument est vraiment *beau*. »

Pourtant, dans un corpus qui porte sur le cinéma, les livres ou les jeux vidéo, ces mêmes mots n'expriment pas toujours des sentiments. Ils ont donc été supprimés du lexique de la grammaire des *jeuxvidéo* parce qu'ils produisent trop de relations éronnées :

« Cela dépendra moins de vous que de l'imbécillité contagieuse des ennemis qui attendent sagement derrière un petit muret, leur *beau* visage buriné dépassant allègrement. »

Comme on le voit dans la phrase précédente, dans ce contexte, les mots de type *joli* ou *beau* sont utilisés pour décrire une action ou un personnage, mais pas un sentiment. La difficulté réside dans le fait de pouvoir distinguer les parties subjectives des parties objectives d'un texte. La description d'une action peut contenir des phrases avec des sentiments, donc subjectives, qui se réfèrent au déroulement de l'histoire. Cependant ces phrases devront être considérées comme étant objectives pour l'évaluation.

Des exemples de règles de grammaire se trouvent dans l'annexe 9.2.

3.3 Lexique de sentiments

L'analyse du texte se base sur les mots du lexique qui ont reçu des traits spécifiques marquant le sentiment positif ou négatif. Il s'agit pour la plupart de verbes (*aimer*, *apprécier*, *détester*, ...) et d'adjectifs (*magnifique*, *superbe*, *insupportable*, ...), mais aussi de quelques noms communs (*plaisir*, ...) et d'adverbes (*malheureusement*, ...). Par exemple, quand une relation de modifieur du nom est extraite (*paysage magnifique*) et que le modifieur (*magnifique*) porte le trait *sents*, la relation de sentiment (⇒ SENTIMENT_POSITIF (*magnifique*, *paysage*)) est extraite ensuite entre le nom et son modifieur. Après cette phase d'analyse, il y a évidemment des règles plus complexes pour extraire les relations des phrases plus compliquées.

Le lexique a été défini par un linguiste au fur et à mesure de l'avancé de chaque projet. A chaque fois qu'un mot intéressant est apparu dans les textes qui n'était pas encore dans le lexique il a été ajouté à ce dernier, selon le domaine du texte. Pour chaque domaine il y a le même lexique de base et ensuite un lexique spécifique qui contient les mots du domaine en question.

L'attribut de la relation (*positif* ou *négatif*) d'un sentiment sera inversé quand une négation est présente dans la phrase, comme par exemple :

« J'aime pas du tout les randonnées en montagne ! »

⇒ SENTIMENT_NEGATIF (*aimer*, *randonnée*)

« Ce n'est pas un mauvais restaurant. »

⇒ SENTIMENT_POSITIF (*mauvais*, *restaurant*)

Extraction de sentiments et d'opinions

Quand cela est possible, les pronoms *qui* et *que* se rapportant à une entité présente ailleurs dans la même phrase, seront remplacés par cette même entité :

« Grenoble est une ville qui vaut vraiment le détour hiver comme été. »
⇒ SENTIMENT_POSITIF (valoir, ville)

Certains noms communs ainsi que des verbes de type interrogatif ont reçu un trait (*no-sents*) pour empêcher l'extraction de relations. Dans *Je cherche un bon hôtel.*, *Bon voyage!* ou *Bonne journée!* il ne s'agit pas de sentiments proprement dit exprimés par l'auteur du texte, mais plutôt de souhaits comme on peut les trouver surtout au début ou à la fin de messages. C'est pour cette raison que nous essayons d'éviter d'extraire ces relations.

Les noms de lieu et de produit ont également des traits spéciaux pour pouvoir extraire d'autres relations qui seront potentiellement intéressantes dans le futur. Voici un extrait du lexique où les mots reçoivent des traits en plus de ceux qu'ils portent déjà (la valeur 1 ajoute ce trait au mot, la valeur 0 l'enlève).

Chaque mot qui peut véhiculer un sentiment reçoit le trait *sents*, puis le trait *positif* ou *négatif* selon sa polarité. D'après la taxonomie d'(Ogorek, 2005) (c.f. la section suivante 3.4) sont ajoutées des valeurs de sentiment plus fines comme à l'aise, détendu, etc. Les mots qui ne doivent pas entrer en relation de sentiment reçoivent le trait *no-sents*. Les traits *genre* et *plateforme* servent à extraire d'autres relations intéressantes dans le domaine des *jeuxvidéo*.

Lexique:

```
agréable = {sents=1, positif=1, à l'aise=1}
sympathique = {sents=1, positif=1, détendu=1}
aimer = {sents=1, positif=1, enchanté=1}
conseiller = {sents=1, positif=1, conseil=1}
plaisir = {sents=1, positif=1, enchanté=1}
décevant = {sents=1, négatif=1, triste=1}
cher = {sents=1, négatif=1, cher=1}
regretter = {sents=1, négatif=1, triste=1}
malheureusement = {sents=1, négatif=1, triste=1}
appétit = {no-sents=1}
vacance = {no-sents=1}
chercher = {no-sents=1}
aventure = {genre=1}
PC = {plateforme=1}
```

La taille du lexique varie selon le domaine d'application. Le lexique de la grammaire de base des sentiments contient environ 250 mots (noms, verbes, adjectifs, etc.) avec des traits de sentiment (*positif* et *négatif*). À ce lexique de base, s'ajoutent environ 150 mots dans le domaine du *tourisme*, et environ 250 mots dans le domaine des *jeuxvidéo*.

3.4 Annotation manuelle de textes

La configuration de la grammaire générique a été faite sur la base d'un travail d'annotation manuelle (à l'aide du logiciel Protégé 3.2⁶ avec le plugin Knowtator⁷) de *threads* venant du domaine du tourisme. Ce corpus du *tourisme* contient une centaine de *threads* annotés (avec comme sujet différentes régions et destinations en France). Chaque *thread* est composé de messages des utilisateurs du forum ; la longueur varie entre dix et 55 messages par document. Un message peut ne contenir qu'une phrase ou plusieurs paragraphes. L'annotation de ce corpus avec Protégé et Knowtator a été faite dans la lignée des travaux de (Riloff et al., 2005, 2006; Wiebe et Mihalcea, 2006).

L'annotation inclut les informations de cause/objet, d'intensité et de l'émetteur du sentiment. Dans

« J'aime énormément Grenoble. »

aimer véhicule le sentiment, *Grenoble* est l'objet du sentiment et *je* est l'émetteur du sentiment. L'adverbe *énormément* exprime l'intensité, le sentiment ici est plus intense que dans la phrase

« J'aime bien Grenoble. »

Cette phase d'annotation sert ensuite aussi à la méthode statistique pour l'élaboration d'un modèle à l'aide de l'entraînement du système sur les textes (c.f. le tableau schématique dans l'annexe 9.1). Les phrases annotées comme positives seront séparées des phrases négatives et un modèle statistique est créé ainsi. L'annotation correcte et précise des phrases est donc très importante pour les deux méthodes de traitement.

L'annotation pour le *tourisme* ne contient pas seulement les deux valeurs *positif* et *négatif* pour classer les sentiments, mais est détaillée beaucoup plus finement (c.f. par exemple les travaux de (Mathieu, 2000, 2006)). Le schéma d'annotation choisi est même plus fin et on voit donc que la classification des sentiments que l'on propose permet un grand nombre de modalités et va au-delà de la simple opposition positif-négatif.

En effet, nous avons repris la taxonomie d'(Ogorek, 2005) qui propose 33 sentiments différents (17 positifs et 16 négatifs) auxquels nous avons ajouté les pseudo-sentiments comme *bon-marché*, *conseil*, *cher* et *avertissement*, car dans le domaine du *tourisme* il y a beaucoup de messages concernant les prix des prestations dont les auteurs des messages sont contents (ou pas).

Les sentiments de la taxonomie d'Ogorek sont classés en groupes⁸ comme AMOUR-DÉSIR (*amour*, *envie*, *tendresse*, *désir*), JOIE (*enchanté*, *excité*, *heureux*, *joyeux*), TRISTESSE- DÉTRESSE (*découragé*, *bouleversé*, *démoralisé*, *triste*), COLÈRE-DÉGOÛT-MÉPRIS (*colère*, *mépris*, *désapprobation*), etc.

Le tableau 1 ci-après montre un extrait de la taxonomie utilisée lors de l'annotation manuelle des *threads* du *tourisme*. Pour chaque occurrence de sentiment, qui peut être un verbe, un adjectif ou même une expression, est indiqué le sentiment correspondant en question, son groupe d'après Ogorek et sa polarité.

6. <http://protege.stanford.edu/>

7. <http://bionlp.sourceforge.net/Knowtator/index.shtml>

8. Sauf les pseudo-sentiments concernant les prix et conseils introduits par notre équipe comme *gratuit*, etc.

Extraction de sentiments et d'opinions

pol.	groupe	sentiment	termes
ATTRACTION	JOIE	<i>enchanté</i>	accueillant, (à) admirer, j'affectionne, aimable, j'aime (aussi/beaucoup/bien/énormément/mieux/particulièrement), (très) amical, (des plus) apprécié/ j'apprécie/à apprécier, à visiter, à/aller voir, ..., chouette, vaut le coup, à découvrir, délicieux, mérite/ vaut le détour, ..., (un de mes/endroit) préféré/je préfère
		<i>excité</i>	mon best off, un bijou, célèbre, coup de cœur/foudre, à couper le souffle, éblouissant, (super) excellent, extra, extraordinaire, exceptionnel, fabuleux, fameux, fantastique, ... , incontournable, incroyable, indescriptible, indispensable, inévitablement, inoubliable, ... , magnifique, à ne pas manquer, un must, pas oublier, le pied, pas rater, ...
	-	<i>bon-marché</i>	abordable, bon marché, pas (trop/très) cher, gratuit, le meilleur rapport qualité/prix, moins onéreux, (beaucoup plus) raisonnable, sans se ruiner

RÉPULSION	COLÈRE, DÉGOÛT, MÉPRIS	<i>mépris</i>	affreux, n'a rien d'attirant, (mes) aversions, bétonné, chiant, (trop) commercial, (le plus) déplu/déplaît, désagréable, détester, effarant, épouvantable, ..., j'ai horreur de ça, impossible, insupportable, trop/très/ultra/hyper touristique
		<i>dés-appro- bation</i>	pas agréable, moins apprécié, bof, c'est bondé, bourré de touristes, vaut pas le coup, (très) envahi, étouffant, éviter, moins fan, (très) fréquenté, pas incontournable, pas inévitable, pas (si) intéressant, moins joli, (vraiment) moche, beaucoup/trop de monde, oublier, pénible, plaît pas
	TRISTESSE, DÉTRESSE	<i>triste</i>	pas chaleureux, pas cool, déçu/déception/décevant, (quel) dommage, pas fameux, pas impressionnant, pas indispensable, sans plus, je regrette, pas sensationnel, pas super, pas sympa, triste

TAB. 1 – Extrait de la taxonomie utilisée pour l'annotation manuelle. La polarité ATTRACTION contient les 17 sentiments positifs et la polarité RÉPULSION les 16 sentiments négatifs d'(Ogorek, 2005). Ils sont groupés en cinq groupes respectifs auxquels s'ajoutent les pseudo-sentiments bon-marché, conseil, cher et avertissement dont nous avons besoin pour l'annotation de textes du domaine du tourisme où il y a beaucoup de messages concernant les prix des prestations et dont les auteurs des messages sont satisfaits (ou pas). Absents de la taxonomie d'Ogorek, ils ne sont pas associés à un groupe de sentiment spécifique.

3.5 Listes de termes

Une liste de termes a été élaborée pour chaque corpus. Chaque liste contient les noms qui sont propres au domaine du corpus. En voici quelques exemples : la liste du corpus *tourisme* contient les mots *ville*, *auberge*, *lieu*, ... ; celle du corpus *jeuxvidéo* englobe *jeu*, *graphisme*, *soft*, Les listes ne regroupent pas les synonymes, nous avons préféré définir des listes détaillées où chaque mot est noté. Grâce à ces listes, des relations erronées, c'est-à-dire dont le deuxième argument n'est pas dans la liste parce qu'il n'appartient pas au domaine, peuvent être refusées.

Par exemple dans le corpus *jeuxvidéo*, cette mesure s'applique à la plupart des relations extraites de la partie résumé du jeu, etc. Considérons la phrase suivante

« Le héros a passé une *magnifique* journée. »

Le mot *journée* n'étant pas dans la liste, la relation `SENTIMENT_POSITIF` (*magnifique*, *journée*) est refusée. Ceci est correct dans la mesure où cette phrase ne contient pas un sentiment exprimé par l'auteur du message, mais fait partie du résumé de l'histoire. En revanche, dans le corpus *tourisme*, le même adjectif *magnifique* exprime bien un sentiment comme c'est le cas dans la phrase :

« Cette *ville* est vraiment *magnifique*. »

⇒ `SENTIMENT_POSITIF` (*magnifique*, *ville*)

Chaque liste ne contient que des noms communs. Toutes les relations qui contiennent des noms propres sont gardées telles quelles :

« Bref, inutile de dépenser le moindre euro pour ce *Yetisports* qui n'en vaut vraiment pas la chandelle. »

⇒ `SENTIMENT_NEGATIF` (*valoir*, *Yetisports*)

Les relations extraites à partir de mots qui ne sont pas des noms sont gardées elles aussi :

« Je n'aime pas *aller* au cinéma. »

⇒ `SENTIMENT_NEGATIF` (*aimer*, *aller*)

Ces listes ont été élaborées automatiquement à partir des textes des corpus d'entraînement avec une méthode basée sur le *bootstrapping*⁹. Elles contiennent tous les noms qui ont été extraits en deuxième argument d'une relation et qui ont été jugés utiles. La phrase qui contient les arguments de cette relation doit se trouver dans un message dont la valeur de classe (positif ou négatif) est la même que l'attribut de la relation extraite, c'est-à-dire, un nom dans une relation de sentiment positif doit se trouver dans un message avec la valeur positif de la classe dans le corpus d'entraînement.

À l'intérieur de chaque liste, les termes sont ordonnés en fonction de la fréquence avec laquelle ils ont satisfait les conditions d'extraction de relations. L'utilisation de ces listes permet d'augmenter le F-score¹⁰ des résultats d'environ 5-10 %, selon le corpus. Ces listes évitent surtout l'apparition de résultats faux positifs qui ne sont pas de vrais résultats souhaités, car ils se rattachent, la plupart du temps, à la partie objective du texte.

9. Un *bootstrap* est un petit programme d'amorçage qui permet d'en lancer un plus gros.

10. La qualité et la fiabilité du résultat sont calculées à partir du quotient entre la précision et le rappel, la précision étant le nombre des résultats corrects obtenus par rapport à tous les résultats obtenus, et le rappel étant le nombre des résultats corrects obtenus par rapport à tous les résultats corrects.

4 Méthode statistique

Pour la méthode statistique, nous utilisons une technique d'apprentissage automatique qui se base sur les travaux de (Pang et al., 2002; Pang et Lee, 2004, 2005)¹¹. Nous l'avons adaptée aux corpus de langue française. Nous l'avons testée d'une part sur les *threads* du corpus sur le *tourisme*, et d'autre part sur les textes des corpus de DEFT'07. (Pang et Lee, 2004) proposent deux axes de classification possibles, soit dans l'opposition subjectif-objectif, soit dans la distinction des opinions subjectives dans l'opposition positif-négatif.

(Pang et Lee, 2004) améliorent la classification de l'axe positif-négatif en supprimant d'abord du texte toutes les phrases objectives et en faisant la classification seulement sur la partie subjective. Cette extraction correspond dans leurs expérimentations à 60 % du texte original. Nous n'avons pas retenu cette façon de faire car nous disposons pas d'un corpus d'entraînement ayant des parties subjectives et objectives bien distinctes. Nous avons choisi de faire la séparation de texte subjectif-objectif à l'aide de la méthode symbolique (c.f. section 3), qui permet d'obtenir finalement des résultats plus nuancés. Les extraits peuvent être vus comme de bons résumés du texte au niveau des sentiments qu'ils expriment.

La méthode statistique se base sur des n -gram de caractères. Pour les projets sur la langue française (le tourisme, les jeux vidéo et DEFT'07) nous avons choisi $n = 12$. Comme pour la méthode symbolique, un indice de confiance est attribué aux textes. Il permet de comparer le résultat avec celui de la méthode symbolique pour en conclure le résultat final avec la méthode hybride. Pour l'entraînement des textes, les techniques de *support vector machines* (SVM) et de *naive bayes* (NB) ont été utilisées. Les résultats sont légèrement meilleurs avec NB, mais ceci reste négligeable.

Des expérimentations ont été faites avec un des corpus de DEFT'07 qui contient des critiques de livres et films, en prenant seulement la/les première(s) et/ou la/les dernière(s) phrase(s) du message. Nous sommes partis de l'hypothèse que le jugement de l'auteur dans une critique de livre ou de film se trouve la plupart du temps en début ou en fin du message, la place du milieu étant vraisemblablement occupée par le résumé du livre ou du film. Les résultats de classification positif ou négatif avec cette technique sont meilleurs qu'en prenant le message en entier. Pourtant, cette technique n'a finalement pas été retenue, car elle ne sera pas facilement reproductible sur des messages provenant d'autres domaines que la critique de film et de livre, où il n'y a pas forcément un résumé au milieu du message.

L'entraînement du module statistique est donc réalisé uniquement sur les phrases de chaque *thread* qui ont été sélectionnées par la méthode symbolique, qui contiennent donc des sentiments, et selon les valeurs de leur classe (positive ou négative) attribuées à chaque corpus par l'annotation manuelle. Les résultats sont ensuite confrontés aux résultats de la méthode symbolique pour donner un résultat final pour chaque message.

5 SYBILLE, la méthode hybride

La méthode hybride est une combinaison des deux méthodes précédentes (c.f. sections 3 et 4). Elle prend en entrée les sorties des deux autres méthodes et calcule d'après les indices de confiance de chaque résultat, une moyenne qui sera traduite en positif ou négatif.

11. <http://www.cs.cornell.edu/home/lllee/papers.html>

La classification définitive est calculée avec les deux méthodes symbolique et statistique. Les résultats respectifs sont confrontés pour obtenir une classification finale. La pondération exacte varie selon plusieurs facteurs, notamment la précision de l'annotation manuelle et la taille du corpus.

La méthode statistique permet de faire une première fouille dans les textes pour obtenir les messages positifs et négatifs. Ensuite, l'utilisateur qui a configuré la grammaire peut modifier et améliorer celle-ci pour obtenir de meilleurs résultats. Le travail prend alors la forme d'un cycle où les résultats s'améliorent constamment.

L'analyse du *thread* se fait au niveau des phrases et permet d'améliorer le résultat en ajoutant ou supprimant par exemple des mots à la liste de termes. Ceci a l'avantage de montrer exactement quelles phrases du document expriment un sentiment.

C'est une approche qui permet de garder la robustesse de l'apprentissage automatique de la méthode statistique et d'orienter en même temps la base de l'entraînement sur une configuration manuelle de la méthode symbolique. Ceci permet de corriger de façon significative les erreurs de l'apprentissage automatique et d'intégrer les spécificités du cahier des charges, c'est-à-dire les particularités de chaque corpus (à l'aide de lexiques et de listes de termes différents selon le domaine d'application).

La méthode hybride a été évaluée, notamment au moment du concours DEFT'07 (c.f. la section 6), avec la mesure du F-score. Le F-score utilisé dans nos expériences est calculé de la manière suivante :

$$F_{score}(\beta) = \frac{(\beta^2 + 1) \times Précision \times Rappel}{\beta^2 \times Précision + Rappel}$$

avec $\beta = 1$.

Elle a été utilisée pour trois des quatre corpus DEFT'07 et a donné les meilleurs résultats pour les corpus *jeuxvidéo* avec un F-score de 0,71, contre 0,54 (méthode symbolique) et 0,70 (méthode statistique) et *relectures* avec un F-score de 0,54, contre 0,48 (méthode symbolique) et 0,51 (méthode statistique).¹² Pour le corpus *débats politiques* seule la méthode statistique a été utilisée.

6 Première évaluation

DEFT'07¹³ (le DÉfi Fouille de Texte) est une campagne d'évaluation dont le thème était en 2007 la classification de textes d'opinion, présents dans différents types de textes. Plusieurs groupes de recherche (laboratoires universitaires ou entreprises privées) ont pu tester leurs systèmes de classification sur les mêmes textes. Dans la phase initiale, chaque groupe inscrit a reçu les deux tiers de chacun des quatre corpus différents qui avaient comme sujet des critiques de films et de livres, des tests de jeux vidéo, des relectures d'articles scientifiques et des notes de débats parlementaires. Pour les trois premiers corpus, une note à trois valeurs (positif,

12. Pour le corpus *aVoiRLire* le meilleur résultat a été obtenu par la méthode statistique avec un F-score de 0,52, contre 0,51 (méthode hybride) et 0,42 (méthode symbolique). Ce corpus n'est probablement pas assez uniforme (il parle de livres, films actuels au cinéma, disques, films plus anciens enregistrés, ...) pour pouvoir faire une liste de termes plus performante.

13. <http://deft07.limsi.fr/>

moyen ou négatif) a été attribuée à chaque texte par le comité des organisateurs, une note à deux valeurs seulement (positif ou négatif) pour le dernier corpus.

Après un certain temps pendant lequel chaque groupe a mis au point son ou ses systèmes de classification, un troisième tiers de chaque corpus a été envoyé pour faire les tests dont les résultats ont dû être soumis quelques jours plus tard. Dix équipes ont participé à l'édition 2007, CELI France est arrivée à la troisième place ¹⁴, les deux premières équipes sont issues de laboratoires de recherche universitaires.

La grammaire de l'analyseur a été paramétrée pour répondre aux besoins des différents corpus DEFT'07, du point de vue lexical mais aussi pour résister aux fautes d'orthographe répétitives. Le point le plus important à modifier a été la classification du message entier qui peut contenir plusieurs sentiments avec une seule valeur globale, et en particulier l'introduction de la notion de sentiment moyen. Dans notre approche standard, au niveau des phrases, les sentiments sont positifs ou négatifs. Il n'est pas nécessaire d'utiliser des sentiments moyens dans le domaine du tourisme, dans la mesure où la taxonomie utilisée (c.f. section 3.4) permet de nuancer suffisamment.

Les sentiments moyens pour DEFT'07 n'ont pas été extraits à l'aide de mots dans le lexique avec un trait moyen, mais d'après des structures de phrase. Par exemple à une phrase qui contient un sentiment positif et un sentiment négatif coordonnés par *mais* est attribué un sentiment moyen à la place :

« Ce jeu est *amusant* au début **mais** *ennuyant* la deuxième semaine. »

Quelques mots clés (surtout des adverbes comme *malgré*, *pourtant*, ...) sont utilisés pour aider à classer un texte qui contient des phrases avec des sentiments positifs et négatifs (c.f. les travaux de (Sándor, 2005)). Le texte entier est alors classé comme moyen.

7 L'interface graphique SYBILLE

Pour conclure cet article quelques figures présentent notre interface graphique SYBILLE, une interface qui aide l'analyste et le client à naviguer parmi les messages analysés pour en prendre connaissance. L'interface que nous utilisons est une dérivation du navigateur *Longwell* du MIT ¹⁵.

Les deux figures (2 et 3) suivantes montrent l'interface graphique du système SYBILLE, dans le domaine des *imprimantes*. Sur la figure 2 en haut à droite il y a un champ dans lequel l'utilisateur peut faire une recherche (1) de messages qui contiennent un mot de son choix ; sinon, il peut ne choisir que les messages positifs ou négatifs (2). Une autre façon de faire une recherche serait de se limiter aux messages qui n'évoquent qu'une marque précise (3), ou en dessous un domaine d'application plus spécialisé (4), ou encore un mot précis d'un domaine. On offre aussi l'option de sélectionner un forum donné parmi tous ceux qui ont été analysés. Les options de recherche peuvent être combinées à volonté pour limiter le nombre de réponses souhaitées.

La figure 3 montre en détail la relation de sentiment qui est indiquée avec ses arguments (5), la phrase qui contient le sentiment et un lien (*external link* (6)) vers le *thread* entier qui permet de visualiser le contexte.

14. <http://deft07.limsi.fr/actes.php#palma>

15. <http://simile.mit.edu/wiki/Longwell>

The screenshot shows the SYBILLE interface with the following elements:

- Header:** CEU France logo and navigation tabs: Order, Commands, List View, Data Mining, Map View (Unavailable for Sybille), Graph View, Timeline View.
- Filter Criteria:**
 - Domaine: "VITESSE" (remove) [add more]
 - Mots Domaine: "lent" (remove) [add more]
- Search Bar:** "Type here to search" with a magnifying glass icon (1).
- Filters:**
 - attitude:** negative (271), positive (130) (2).
 - Secteur:** HP (42), Epson (36), Canon (30), Brother (13) (3).
 - Domaine:** "HARDWARE" (207), "IMPRESSION" (170), "QUALITE" (86), "IMAGE" (81), "GRAPHISME" (77), "SYSTEME" (28) (4).
 - Mots Domaine:** "imprimante" (139)
- Message Details (materiel/20070509):**
 - attitude:** negative
 - Domaine:** IMPRESSION, VITESSE, HARDWARE
 - materiel:**
 - Secteur:** Brother
 - forum:** forum.hardware
 - text:** NEGATIF ~ lent ~ Brother ~ | Je sais que les Brother bas de gamme sont lentes, mais c'est une question de prix.
 - Expediteur:** linuxafficion
 - Sujet:** Multifonction rapport qualité/prix : la nouvelle canon MP500 ? - Imprimantes - Hardware - Périphériques - FORUM HardWare.fr
 - [external link]**
 - Show Referers

FIG. 2 – L'interface graphique SYBILLE, ici pour le domaine des imprimantes. Plusieurs moyens différents permettent de naviguer dans les résultats des messages analysés.

The screenshot shows a detailed view of a message (materiel/20060618) with the following details:

- attitude:** positive (5)
- Domaine:** IMAGE, IMPRESSION, VITESSE
- materiel:**
- Secteur:** Epson
- forum:** forum.hardware
- text:** POSITIF ~ superbe ~ impression ~ | J'en suis très content, elle est très rapide pour le texte et les impressions photos sont superbes et je la consommation d'encre très raisonnable comparé à mon ancienne Epson.
- Expediteur:** Steph657
- Sujet:** Canon Fixma 6600D - Qualite PARFAITE !!! - Imprimantes - Hardware - Périphériques - FORUM HardWare.fr
- [external link]** (6)
- Show Referers

FIG. 3 – Exemple détaillé d'une relation de sentiment positif, toujours dans le domaine des imprimantes.

Extraction de sentiments et d'opinions

Le domaine des prochaines figures sont les *jeux vidéo*. Le paramétrage choisi ici sont les deux marques *Ubisoft* et *Electronic Arts*, c'est-à-dire on aura sélectionné tous les messages qui parlent soit d'*Ubisoft* soit d'*Electronic Arts*, sur toute la période analysée.

Les figures (4 et 5) montrent la chronologie (*TimelineView*) des messages. Ceci permet de visualiser combien de messages ont été publié par unité de temps. Par exemple si un nouveau produit est lancé sur le marché on peut voir sur la *TimelineView* si les internautes le commentent directement et s'il suscite plutôt des réactions positives ou négatives. Chaque point de couleur représente un message, bleu pour les messages positifs et rouge pour les messages négatifs. Un clic sur un des points permet de voir le texte du message.

Sur la figure 4 on peut voir la période du mois d'août 2007. Dans les cases en-bas on peut saisir des mots-clés contenus dans les messages pour les faire ressortir, le message marqué en jaune contient le mot *super*. Dans cet échantillon il n'y a qu'un seul message contenant ce mot-clé. L'interface permet de saisir plusieurs mots-clés qui s'afficheront avec des couleurs différentes.



FIG. 4 – La *TimelineView* montre la chronologie des messages. Chaque point représente un message, ici pour la période du mois d'août 2007. Les cases *Highlight* permettent de mettre en évidence des mots-clés dans les messages.

La figure 5 montre en détail le message marqué en jaune. Il contient le mot *super* dans la phrase suivante :

« Histoire d'illustrer un peu plus les *superbes* effets de fumées dégagées... »

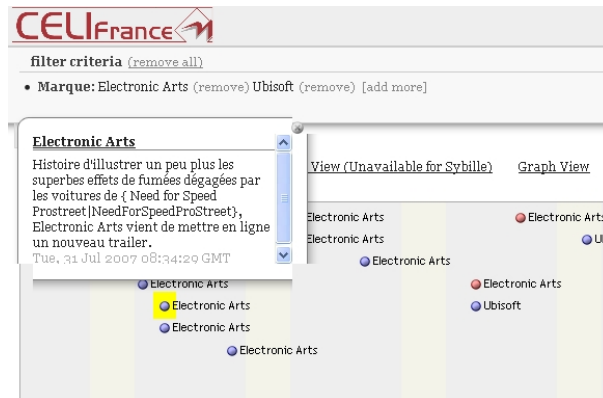


FIG. 5 – La TimelineView avec en détail le message qui contient le mot *super*.

Les trois figures suivantes (toujours avec les paramètres des deux marques *Ubisoft* et *Electronic Arts*) montrent graphiquement les statistiques les plus importantes. Ces graphiques sont intéressants pour l'utilisateur qui peut se faire rapidement une idée du contenu du forum analysé en ce qui concerne la fréquence d'interaction, les produits détectés et l'équilibre entre messages positifs et négatifs.

Les sujets discutés sont montrés dans la figure 6. On peut voir que le sujet GRAPHISME n'a pas provoqué de messages négatifs dans cet échantillon, contrairement au sujet JOUABILITÉ qui reçoit plus de messages négatifs que positifs, le seul dans ce cas. D'après nos observations il y a très souvent (beaucoup) plus de messages positifs que négatifs, donc l'inverse peut avoir une signification intéressante.

La figure 7 visualise la fréquence des produits qui ont fait l'objet de discussion. Les marques *EA Sports* et *Maxis* sont évoquées beaucoup plus souvent dans le forum que les autres marques. Il s'agit dans cette figure des marques que le système a pu identifier, il y a bien sûr toujours des messages qui contiennent des noms de marques qui n'ont pas été identifiées à cause de fautes de frappe ou abréviations inconnues.

Et finalement, la figure 8 présente la répartition du nombre de messages par mois, en rouge les messages positifs et en bleu les négatifs. Ce graphe permet de constater immédiatement qu'il y a presque tout le temps (beaucoup) plus de messages positifs que de messages négatifs, et de voir les périodes de pics dues probablement à la sortie de nouveaux jeux.

Extraction de sentiments et d'opinions

Main Topics

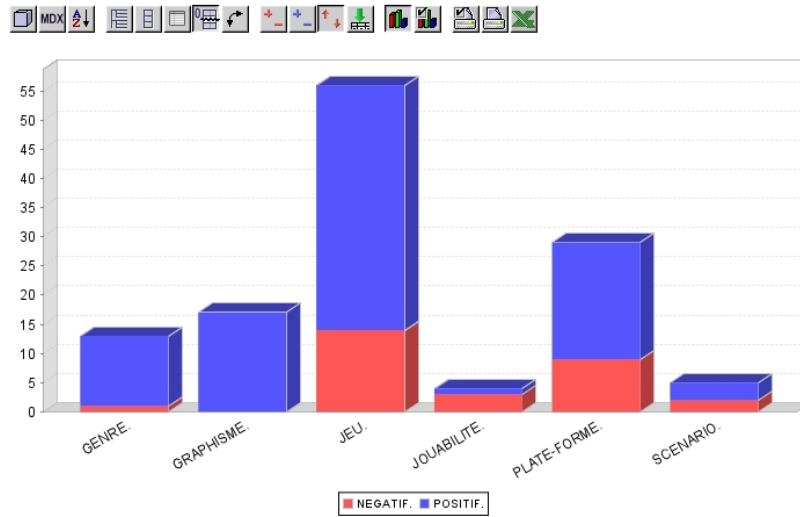


FIG. 6 – Ce diagramme présente les sujets discutés dans le domaine des jeuxvidéo.

Buzz Brand Market Share

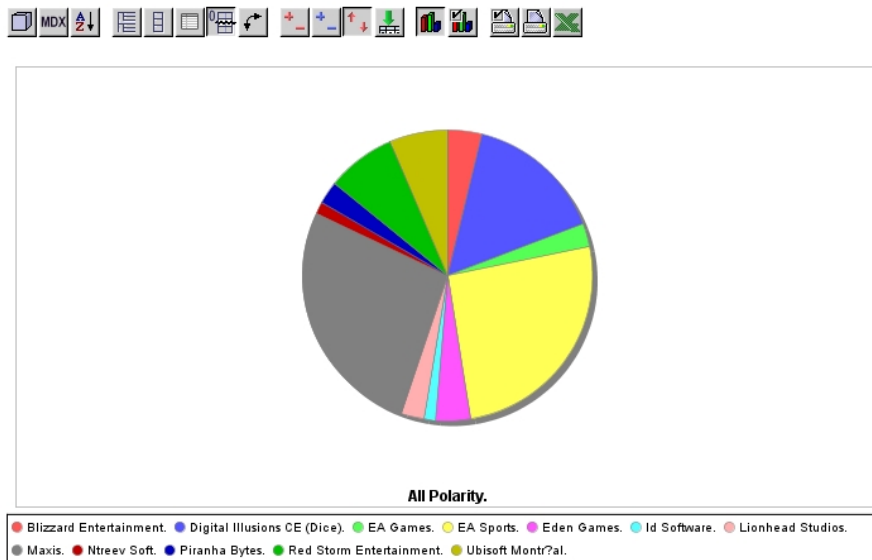


FIG. 7 – Le diagramme ci-dessus visualise la fréquence des produits évoqués dans les messages.

Monthly opinion trend

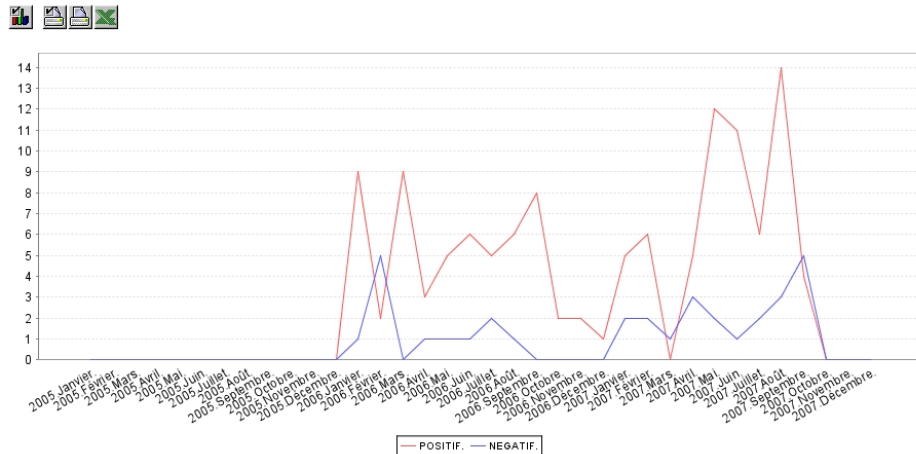


FIG. 8 – Ce graphe présente la répartition du nombre de messages par mois sur l'ensemble de la période analysée, en rouge les messages positifs (plus nombreux) et en bleu les négatifs.

8 Conclusion

Dans cet article nous avons présenté comment l'utilisation de grammaires syntaxiques, un outil du traitement automatique du langage naturel, peut améliorer la qualité d'un système d'extraction de sentiments. Nous avons décrit une méthode symbolique avec une grammaire adaptée au domaine des textes, et une méthode d'apprentissage automatique. L'évaluation de notre système de classification SYBILLE a montré que la combinaison des méthodes symbolique et statistique donne des résultats plus précis que chacune des méthodes employée séparément.

L'intérêt de la méthode hybride repose sur la prise en compte des contextes d'application de ses résultats. Il est bien connu que la méthode purement symbolique a souvent pour le client un coût d'entrée plutôt élevé. Cette considération est liée au temps de configuration, de repérage ou de création de lexiques spécifiques, de taxonomies etc.

L'utilisation d'une méthode hybride permet, au contraire, de minimiser les coûts de configuration, en réduisant une partie du travail à l'annotation de textes, une tâche qui dans la plupart des cas peut être réalisée par le client lui-même. Les algorithmes d'apprentissage automatique sont alors en mesure de donner des premiers jugements au niveau du texte entier.

Ce qui est le plus important, c'est qu'avec ce type de système statistique on peut ajouter, selon la méthode exposée dans cet article, une couche *symbolique* au fur et à mesure, de plus en plus importante dès que les exigences d'une application deviennent plus précises. On peut par exemple superposer une couche d'identification de jugement, qui permet d'avoir une visibilité sur les jugements sans devoir lire le texte dans son entier. On peut identifier certains patrons sémantiques qui sont d'importance capitale pour une application donnée et qui doivent

avoir la priorité sur les résultats statistiques (par exemple le souci de sécurité exprimé par les internautes sur un certain modèle de voiture).

Les exemples pourraient être multipliés. Ce qui apparaît avant tout intéressant, c'est que la démarche hybride est importante non seulement pour des raisons scientifiques de performance (le meilleur résultat entre les technologies que nous avons adoptées) mais, aussi et surtout pour des raisons de développement et d'acceptation par le marché.

Références

- Aït-Mokhtar, S. et J.-P. Chanod (1997). Subject and object dependency extraction using finite-state transducers. In P. Vossen, G. Adriaens, N. Calzolari, A. Sanfilippo, et Y. Wilks (Eds.), *Automatic information extraction and building of lexical semantic resources for NLP applications*, pp. 71–77. Association for Computational Linguistics.
- Aït-Mokhtar, S., J.-P. Chanod, et C. Roux (2001). A multi-input dependency parser. In *Actes d' IWPT'01*.
- Basili, R., M. T. Pazienza, et F. M. Zanzotto (1999). Lexicalizing a shallow parser. In *Actes de TALN'99*.
- Bethard, S., H. Yu, A. Thornton, V. Hatzivassiloglou, et D. Jurafsky (2004). Automatic extraction of opinion propositions and their holders. In *Actes d' AAAI'04*.
- Chklovski, T. (2006). Deriving quantitative overviews of free text assessments on the web. In *Actes d' IUI'06*, pp. 155–162.
- Dini, L. (2002). Compréhension multilingue et extraction de l'information. In F. Segond (Ed.), *Multilinguisme et traitement de l'information (Traité des sciences et techniques de l'information)*. Editions Hermes Science.
- Dini, L. et G. Mazzini (2002). Opinion classification through information extraction. In A. Zanasi, C. A. Brebbia, N. F. F. Ebecken, et P. Melli (Eds.), *Data Mining III*, pp. 299–310. WIT Press.
- Dini, L. et F. Segond (2007). La linguistique informatique au service des sentiments. In *Revue de l'électricité et de l'électronique*, pp. 66–77. Editions SEE.
- Grouin, C., J.-B. Berthelin, S. El Ayari, T. Heitz, M. Hurault-Plantet, M. Jardino, Z. Khalis, et M. Lastes (2007). Présentation de DEFT'07 (DÉfi Fouille de Textes). In *Actes de DEFT'07*, pp. 1–8.
- Hatzivassiloglou, V. et K. R. McKeown (1997). Predicting the semantic orientation of adjectives. In *Actes d' ACL'97*, pp. 174–181.
- Kim, S.-M. et E. Hovy (2004). Determining the sentiment of opinions. In *Actes de COLING'04*, pp. 1267–1373.
- Mathieu, Y. Y. (2000). *Les verbes de sentiment. De l'analyse linguistique au traitement automatique*. CNRS Editions.
- Mathieu, Y. Y. (2006). A computational semantic lexicon of french verbs of emotion. In J. G. Shanahan, Y. Qu, et J. Wiebe (Eds.), *Computing attitude and affect in text: Theorie and applications*, pp. 109–124. Springer.

- Maurel, S., P. Curtoni, et L. Dini (2007). Classification d'opinions par méthodes symbolique, statistique et hybride. In *Actes de DEFT'07*, pp. 111–117.
- Maurel, S., P. Curtoni, et L. Dini (2008). L'analyse des sentiments dans les forums. In *Actes de FODOP'08*, pp. 9–22.
- Ogorek, J. R. (2005). Normative picture categorization: Defining affective space in response to pictorial stimuli. In *Actes de REU'05*.
- Pang, B. et L. Lee (2004). A sentimental education: Sentiment analysis using subjectivity summarization based on minimum cuts. In *Actes d' ACL'04*, pp. 271–278.
- Pang, B. et L. Lee (2005). Seeing stars: Exploiting class relationships for sentiment categorization with respect to rating scales. In *Actes d' ACL'05*, pp. 115–124.
- Pang, B., L. Lee, et S. Vaithyanathan (2002). Thumbs up? Sentiment classification using machine learning techniques. In *Actes d' EMNLP'02*, pp. 79–86.
- Riloff, E., S. Patwardhan, et J. Wiebe (2006). Feature subsumption for opinion analysis. In *Actes d' EMNLP'06*, pp. 440–448.
- Riloff, E., J. Wiebe, et W. Phillips (2005). Exploiting subjectivity classification to improve information extraction. In *Actes d' AAAI'05*.
- Sándor, A. (2005). A framework for detecting contextual concepts in texts. In *Actes du Electra Workshop*.
- Sproull, L. et S. Kiesler (1991). *Connections: New ways of working in the networked organization*. Cambridge: MIT Press.
- Turney, P. D. (2002). Thumbs up or thumbs down? Semantic orientation applied to unsupervised classification of reviews. In *Actes d' ACL'02*.
- Wiebe, J. et R. Mihalcea (2006). Word sense and subjectivity. In *Actes d' ACL'06*, pp. 1065–1072.
- Wilson, T., J. Wiebe, et R. Hwa (2004). Just how mad are you? Finding strong and weak opinion clauses. In *Actes d' AAAI'04*.
- Yu, H. et V. Hatzivassiloglou (2003). Towards answering opinion questions: Separating facts from opinions and identifying the polarity of opinion sentences. In *Actes d' EMNLP'03*, pp. 129–136.

9 Annexes

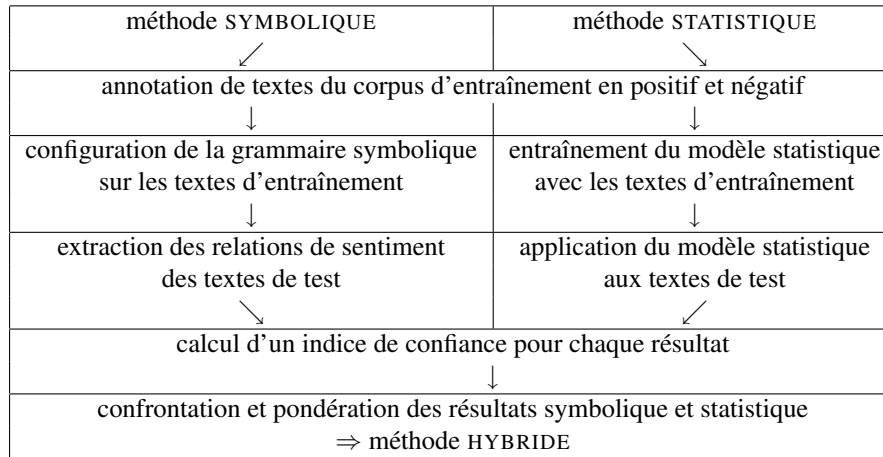
9.1 Schéma général du processus

Le tableau 2 ci-après décrit le schéma général du processus de traitement, un corpus de textes d'entraînement différents d'un corpus de textes de test est indispensable.

9.2 Extrait de la grammaire

Les règles de grammaire se basent sur les relations extraites par la grammaire de base comme le sujet et l'objet d'une phrase, ou le modifieur d'un nom, d'un verbe ou d'un adjectif,

Extraction de sentiments et d'opinions



TAB. 2 – *Le schéma général du processus.*

ou encore la coordination. Pour changer l'attribut d'une relation on efface l'attribut actuel et ajoute le nouvel attribut souhaité.

La règle suivante extrait une relation de sentiment quand un adjectif marqué avec le trait `sents` est l'objet d'un verbe dont le sujet ne porte pas le trait `no-sents`.

```
%% Le parc est très agréable.
if ( OBJET[spred:1](*2[cond:0, fut:0, subj:0], *1[adj:1,
  sents:1, positif:1]) & SUJET(*2, *3[no-sents:0]) )
--> SENTIMENT[positif=1](*1, *3).
```

La polarité d'une relation de sentiment déjà extraite est inversée par le modifieur négatif *pas* (porteur du trait `psneg`).

```
%% Un hôtel pas agréable.
if ( ^SENTIMENT[positif:1](*1[adj:1, positif:1], *2)
  & MADJ(*1, *3[psneg:1]) & ~MADV(*3, *4[psneg:1]) )
--> SENTIMENT[negatif=1, positif=0](*1, *2).
```

Une relation de coordination provoque une nouvelle relation de sentiment avec le deuxième argument de la première relation. Un test vérifie que cette relation n'est pas déjà extraite par une autre règle.

```
%% La ville est vraiment moche et affreuse.
if ( COORD(*10, *1[sents:1, negatif:1], *2)
  & SENTIMENT(*2, *3) & ~SENTIMENT(*1, *3) )
--> SENTIMENT[negatif=1](*1, *3).
```

Summary

We present three different methods to perform an automatic classification of texts which include opinions. The first method is symbolic, the second statistic and the last, hybrid, is a combination of the first two. We will show how the combination makes it possible to exploit the advantages of both methods, namely robustness of statistical machine learning and the possibility of a manual configuration given by the symbolic method allowing the use of real-life applications. The texts classified by these methods come from non structured information sources such as internet forums.