

# Une introduction au Text Mining et à la sémantique

Seth Grimes

--

Président, Alta Plana

*Alta Plana*

## NEHRU PROMISES BORDER DEFENSE

Warns Peiping India Would  
Resist Further Advances  
Into Her Territory

By PAUL GRIBES  
Special to The New York Times.

NEW DELHI, India, Oct. 8.—Prime Minister Jawaharlal Nehru warned today that any further aggression by Communist China against India "will certainly be fully resisted."

"I hope it will not be necessary," he said at a crowded news conference. But he warned against "any kind of advance" by the Chinese from positions they now hold inside India's northern borders.

The Prime Minister asserted, "We do not intend to start military operations against any of these places at this stage, when we are dealing with this matter on a political level."

He said the Indian Government did not function "in an excited way with a club in hand" but instead sought to act "with determination, not with anger."

Mr. Nehru recalled that he had sent a letter to Premier Chou En-lai of Communist China Sept. 26, appealing to the Chinese to withdraw their forces so that the border disputes could be "amicably and peacefully settled."

The Prime Minister said he welcomed the "friendly message" sent by Mr. Chou in response to India's greetings on the tenth anniversary of the Communist regime Oct. 1. In this message, made public here last night, Mr. Chou stressed the "age-old friendship" between China and India and described the present difficulties as "merely an episode."

"The tone is friendly and, therefore, an improvement," Mr.

Nehru said. "How far that represents any basic change in the situation, I cannot say. I hope it does."

The Prime Minister, dressed in white homespun with his usual red rose in a buttonhole, answered questions for an hour and a quarter.

Prime Minister Nehru belittled a report that Premier Khrushchev, when visiting Los Angeles, told of the Soviet interception of a highly confidential message from President Eisenhower to Mr. Nehru. The Prime Minister said the only messages he had received from the United States President in the last two or three months were "friendly and formal letters—nothing secret."

Mr. Nehru also expressed the following views:

¶That he is opposed to disbanding the armistice control commissions in Laos, Cambodia, and Vietnam as long as it is possible to carry out the Geneva agreements of 1954 that ended the Indochinese war.

¶That President Charles de Gaulle's latest proposal to end the Algerian strife was "certainly a marked advance" from previous French suggestions, particularly since it "acknowledged the right of self-determination, which was the basic thing."

¶That he disapproved thoroughly of a recent "extraordinary" resolution of the Indian Communist party that would have India weaken considerably her stand on the border issue.

¶That the Soviet Union should be "warmly congratulated" for its rocket around the moon.

Mr. Nehru, who will be 70 years old Nov. 14, was asked whether "you feel a sense of fulfillment in your life or a sense of frustration."

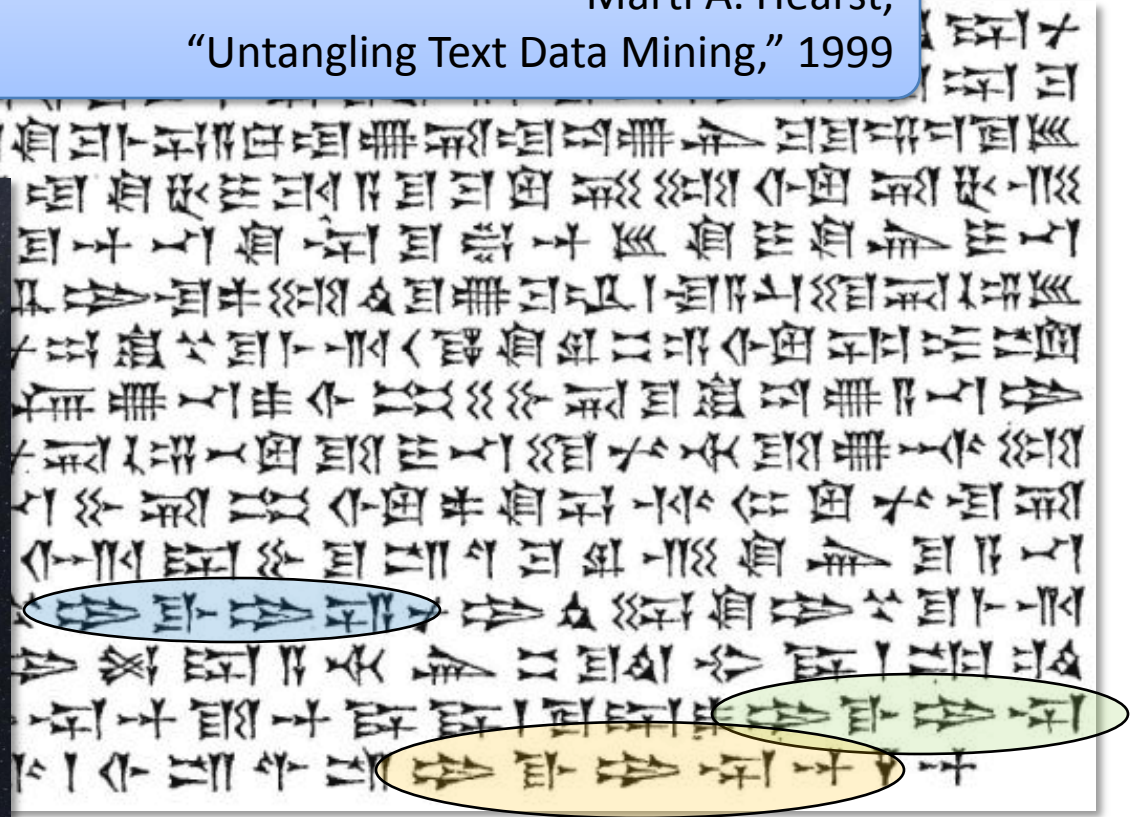
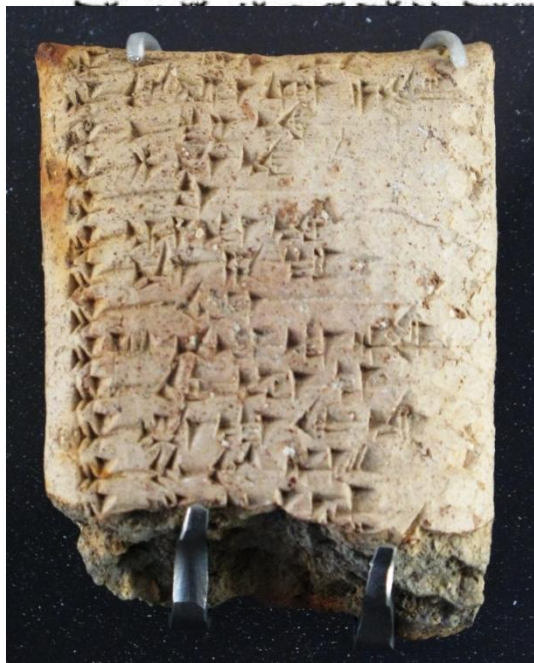
He replied:

"I have absolutely no sense of frustration in my life. I hope my face shows that. If the question is, have I achieved all I want to, my answer is no. But partial achievement comes from time to time."

# Du texte à l'information

«Le texte exprime une gamme vaste et riche d'information, mais encode cette information dans une forme qui est difficile à déchiffrer automatiquement.»

-- Marti A. Hearst,  
"Untangling Text Data Mining," 1999

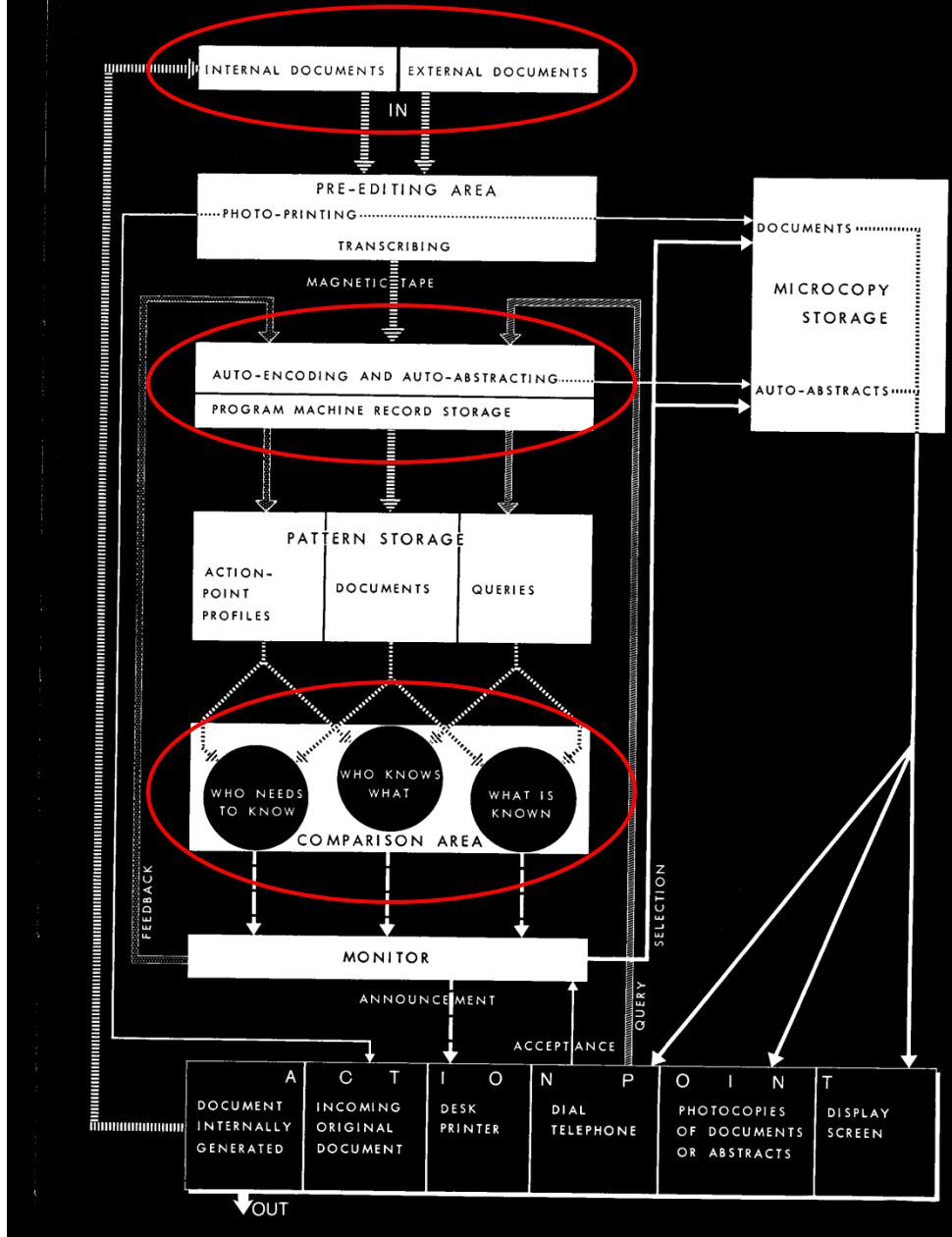


*Alta Plana*

Input et traitement des documents

Gestion des connaissances

Extraction d'information



Hans Peter Luhn, "A Business Intelligence System," IBM Journal, October 1958

Figure 1 A Business Intelligence System

# Analyse statistique du contenu

«L'information statistique obtenue par la fréquence et la distribution des mots est utilisée par la machine afin de calculer une mesure relative de leur importance.»

Significant words in descending order of frequency (common words omitted).

46	<i>nerve</i>	12	<i>body</i>	6	<i>disturbance</i>	4	<i>accumulate</i>
40	<i>chemical</i>	12	<i>effects</i>	6	<i>related</i>	4	<i>balance</i>
28	<i>system</i>	12	<i>electrical</i>	5	<i>control</i>	4	<i>block</i>
22	<i>communication</i>	12	<i>mental</i>	5	<i>diagram</i>	4	<i>disorders</i>
19	<i>adrenalin</i>	12	<i>messengers</i>	5	<i>fibers</i>	4	<i>end</i>
18	<i>cell</i>	10	<i>signals</i>	5	<i>gland</i>	4	<i>excitation</i>
18	<i>synapse</i>	10	<i>stimulation</i>	5	<i>mechanisms</i>	4	<i>health</i>
16	<i>impulses</i>	8	<i>action</i>	5	<i>mediators</i>	4	<i>human</i>
16	<i>inhibition</i>	8	<i>ganglion</i>	5	<i>organism</i>	4	<i>outgoing</i>
15	<i>brain</i>	7	<i>animal</i>	5	<i>produce</i>	4	<i>reaching</i>
15	<i>transmission</i>	7	<i>blood</i>	5	<i>regulate</i>	4	<i>recording</i>
13	<i>acetylcholine</i>	7	<i>drugs</i>	5	<i>serotonin</i>	4	<i>release</i>
13	<i>experiment</i>	7	<i>normal</i>			4	<i>supply</i>
13	<i>substances</i>					4	<i>tranquilizing</i>
Total word occurrences in the document: . . . . .							2326
Different words in document:							
Total of different words . . . . .							741
Less different common words . . . . .							170
Different non-common words . . . . .							571
Ratio of all word occurrences to different non-common words . . . . .							~4:1
Non-common words having a frequency of occurrence of 5 and over:							
Total occurrences . . . . .							478
Different words . . . . .							39

Hans Peter Luhn,

*"The Automatic Creation of Literature Abstracts,"*

IBM Journal, April 1958

# Limites de l'analyse statistique

«Cette argumentation assez simple sur la 'signification' ignore les aspects linguistiques tels que la grammaire et la syntaxe... Aucune attention n'est accordée aux rapports logiques et sémantiques établis par l'auteur.»

-- Hans Peter Luhn, 1958

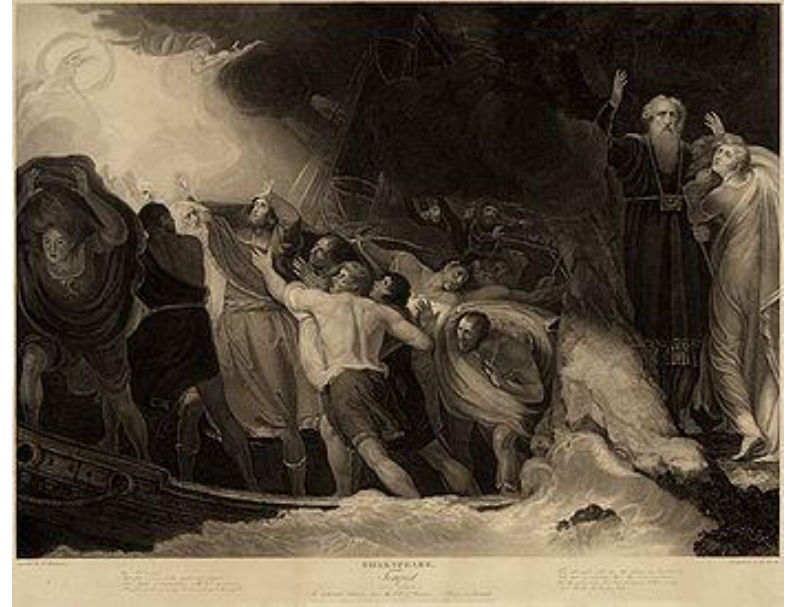
# Inconnu = nouveauté

*Miranda:*

O, merveille! Combien de belles créatures vois-je ici réunies! Que l'humanité est admirable! O splendide Nouveau Monde, Qui compte de pareils habitants !

*Prospero:*

**C'est nouveau pour toi.**



Le naufrage dans **La Tempête**, Acte I, Scène 1, dans une gravure de 1797 basée sur une peinture de George Romney.

# Liens sémantiques

New York Times,  
8 septembre, 1957

Anaphore /  
coréférence:  
"They"

## SCIENCE IN REVIEW

### Chemistry Is Employed in a Search for New Methods to Conquer Mental Illness

By ROBERT K. PLUMB

By coincidence this week-end in New York City marks the end of the annual meeting of the American Psychological Association and the beginning of the annual meeting of the American Chemical Society.

Psychologists and chemists have never had so much in common as they now have in new studies of the chemical basis for human behavior. Exciting new finds in this field were also discussed last week in Iowa City, Iowa, at the annual meeting of the American Physiological Society and at Zurich, Switzerland, at the Second International Congress for Psychiatry.

Two major recent developments have called the attention of chemists, physiologists, physicists and other scientists to mental diseases: It has been found that extremely minute quantities of chemicals can induce hallucinations and bizarre psychic disturbances in normal people, and mood-altering drugs (tranquilizers, for instance) have made long-institutionalized people amenable to therapy.

Money to finance research on the physical factors in mental illness is being made available. Progress has been achieved toward the understanding of the chemistry of the brain. New goals are in sight.

At the psychiatrists meeting in Zurich last week, four New York City physicians urged their colleagues to broaden their concept of "mental disease," and to probe more deeply into the chemistry and metabolism of the human body for answers to mental disorders and their prevention.

#### Blood May Tell

Dr. Felix Marti-Ibanez and three brothers, Dr. Mortimer D. Sackler,

Dr. Raymond R. Sackler and Dr. Arthur M. Sackler cited evidence that the blood chemistry of victims of schizophrenia is different from that of normal people. Perhaps multiple biological factors are responsible for this chemical change, they suggested.

Mental disease is a "developmental process" and long duration of a disorder may result in "permanent alteration of anatomy and physiology," they said. They urged that trials of new drugs which affect the brain should be concentrated on complex studies of the mechanism of action of the drugs. The variety of substances capable of producing profound mental effects is a new armory of weapons for use in investigating biological mechanisms underlying mental disease, they said.

The sources of behavioral disturbance are many and they may come from external as well as internal forces, the four reported. This concept has already proven practical, for instances, when it enabled psychiatrists to predict that the administration of ACTH and cortisone could produce psychosis.

"It led some years ago to the development of a blood test which was 80 per cent accurate in the identification of schizophrenic patients," they said. "It permitted us on physiologic grounds to deny that the psychoneuroses and the psychoses were lesser and greater degrees of the same disease process, and, in fact, to affirm that they represented opposite and even mutually exclusive directions of physiologic disturbances," they said.

Chemicals now available should be used not only to bring relief to mentally sick but also to uncover

the biological mechanisms of the disease processes themselves. "Only then will the metabolic era mature and bring to fruition man's long hoped for salvation from the ravages of mental disease," they reported.

#### Chemistry of the Brain

At the psychologist's meeting here, a technique for tracing electrical activity in specific portions of the animal brain was described by researchers of California reported in cat and electric animals which

In this reported sequence its various the brain may be located. Furthermore, the electrical pathways so traced out can be blocked temporarily by the use of chemicals. This poses new possibilities for studying brain and side the California sized.

The try have peutic courage for men that kn ary field last we Washin

This the Na Health informati Literatu and an bers ca technic People vited to or other letters have in Informa Silver



# Analyse de la tonalité



“Kind” = genre, variété, pas une indication de sentiment.

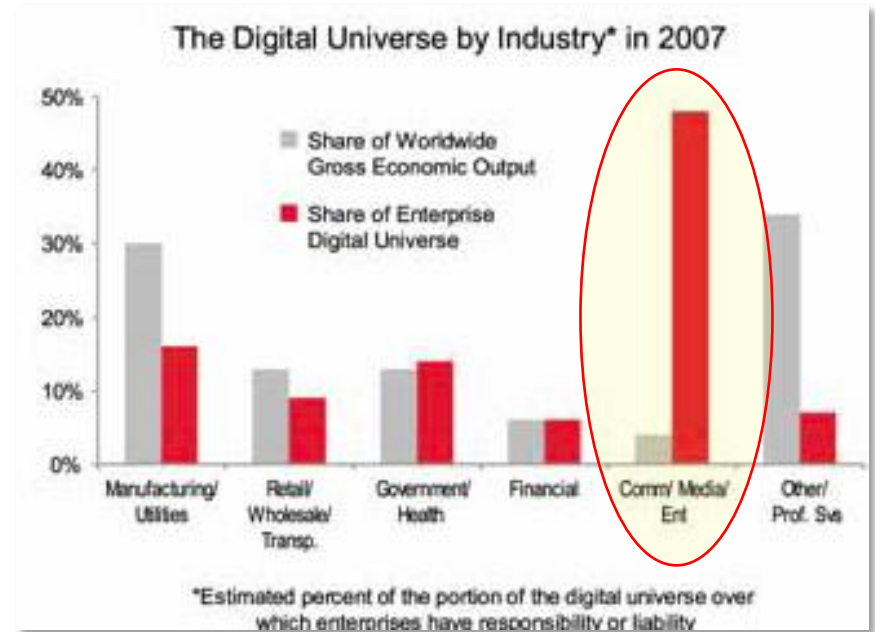
Référence externe

Répétitions non filtrées

# L'univers du contenu numérique

“La télédiffusion, les médias et les industries du loisir recueillent environ 4% des revenus globaux, mais produisent déjà, dirigent, ou supervisent 50% de l'univers numérique.”

**Environ 70% de l'univers numérique est créé par les individus.**



“The Diverse and Exploding Digital Universe,”  
(IDC, 2008)

# Le défi de « l'information non-structurée »

- Les sites Web, articles de journaux et magazines, images, vidéos.
- Les blogs, forums, et médias sociaux.
- Le courrier électronique, les notes et transcriptions de centres d'appel; les conversations enregistrées.
- Les sondages, formulaires de satisfaction, demandes d'indemnité et de garantie.
- Les documents, rapports, écrits scientifiques.
- Et tout autre type de document imaginable.

**Est-ce que la recherche à elle seule est suffisante?**

# Résultats de recherche

Quelles sont la qualité, la valeur et la crédibilité des résultats de recherche?

The screenshot shows a Google search interface with the query "Is the Hilton New York a good hotel?". The results page includes a sponsored link for Hilton Hotels New York, a search result from Hilton.com, a sponsored link for New York Best Hotels, and a search result from TripAdvisor. Red circles highlight the words "sophisticated" and "conveniently" in the Hilton.com snippet, and "TripAdvisor" and "good deal" in the TripAdvisor snippet. Arrows point from these circles to callout boxes.

Google  Search [Advanced Search](#) [Preferences](#)

Web Results 1 - 10 of about 1,590,000 for [Is the Hilton New York a good hotel?](#). (0.41 seconds)

**Hilton Hotels New York** Sponsored Link  
[www.Hilton.com](#) Book direct on **Hilton.com** For Our Best Rates, Guaranteed!

**New York Hotels - Hilton New York Hotel NYC**  
Hilton New York City Hotel - Rockefeller Center Hotel - NYC Hotels ... The Hilton New York Hotel is a sophisticated Manhattan hotel conveniently located in ...  
[www1.hilton.com/en\\_US/hil/hotel/NYCNHHH-Hilton-New-York-New-York/index.do-77k](#) - Cached - Similar pages

[Accommodations](#) [Video Tour](#)  
[Directions](#) [Photos](#)  
[Reservations](#) [Dining](#)  
[Services & Amenities](#) [Employment](#)

[More results from hilton.com »](#)

**Hotels by Hilton - Hotel Reservations, Deals, and Room Rates**  
Find and book hotel rooms online at the official site of Hilton brand hotels. ... Nevada, New Hampshire, New Jersey, New Mexico, New York, North Carolina ...  
[www1.hilton.com/](#) - 87k - Cached - Similar pages

**Hilton New York (New York City, NY) - Hotel Reviews - TripAdvisor**  
Hilton New York: Save up to 50% off Hotels. Everyday Smart Deal! ... I stayed in this hotel over a weekend and got a good deal on Priceline for \$110 a night ...  
[www.tripadvisor.com/Hotel\\_Review-g60763-d611947-Reviews-Hilton\\_New\\_York-New\\_York\\_City\\_New\\_York.html](#) - 116k - Cached - Similar pages

L'opinion de l'hôtel

L'opinion de l'invité... à propos de Priceline

À qui profite la recherche?

## Comment pouvons-nous faire mieux?

«Nous avons en place plusieurs instruments -- des technologies Web 2.0...»

“The Diverse and Exploding Digital Universe,”  
(IDC, 2008)

# Le Web 2.0

«Le Web 2.0 est la révolution d'affaires dans l'industrie de l'informatique provoquée par l'évolution de l'Internet comme une plate-forme.»

-- Tim O'Reilly, 2004

Le Web 2.0 marque un «mouvement des sites Internet personnels aux blogs et l'agrégation de blogs, de la publication à la participation, ... un processus en cours et interactif ... vers les liens basés sur le balisage.»

-- Terry Flew, "New Media: An Introduction," 2008

**Le Web 2.0 est interactif, personnalisé, collaboratif, dynamique.**

## Mais comment pouvons-nous faire mieux?

«Nous avons en place plusieurs instruments -- des technologies Web 2.0... logiciels qui fouillent les données non structurées et le Web Sémantique -- pour apprivoiser l'univers numérique. **Fait de manière adéquate, nous pouvons transformer la croissance d'information en croissance économique.**»

“The Diverse and Exploding Digital Universe,”  
(IDC, 2008)

# Text mining: de l'information à l'intelligence

Le text mining permet une recherche plus intelligente, basée sur les objectifs de l'utilisateur, par exemple la réponse aux questions –

The image displays several overlapping search engine results pages (SERPs) for different queries, demonstrating text mining and information intelligence. The queries and their corresponding results are:

- Live Search: population peru**
  - Web 1-10 of 6,550,000 results - [Advanced](#)
  - See also: [Images](#), [Video](#), [News](#), [Maps](#), [More](#) ▼
  - » [Peru](#) Population, total: 29,041,593
  - 2008 estimate · United States C
  - Is this useful? [Yes](#) | [No](#)
  - [Peru Population](#)
  - Related searches: [Peru Travel](#), [Peru Food](#)
- Google: map massachusetts**
  - Web [Maps](#) [Images](#)
  - Map of Massachusetts and surrounding areas (Connecticut, New York, New Jersey).
- Google: illinois unemployment rate**
  - Web Results 1 - 10 of about 698,000 for [illinois unemployment rate](#). (0.21 seconds)
  - Unemployment rate, Illinois**
  - 9.4% of the labor force - Not seasonally adjusted - Mar 2009
  - Source: U.S. Bureau of Labor Statistics
  - www.google.com/publicdata
  - unemployment rate**
  - Local Area **Unemployment** Statistics: LAUS. **Illinois** and Chicago Metropolitan Area **Unemployment** Rates. Current Monthly **Unemployment** Rates. Year-to-Date Data ...
  - [lmi.ides.state.il.us/laus/lausmenu.htm](#) - 20k - [Cached](#) - [Similar pages](#)
  - 20\*30 = 600**
  - Yahoo! [Shortcut](#) - [About](#)
  - [Renault 20/30](#) - [Wikipedia, the free encyclopedia](#)
  - The Renault 20 and Renault 30 are two executive cars produced by the French ... Renault 20 had two single rectangular headlights whereas the Renault 30 had ...
  - [en.wikipedia.org/wiki/Renault\\_20/30](#) - [Cached](#)
  - [Active 20-30 USA & Canada](#)
  - ... for young adults between the ages of 20 and 39. Provides young adults with an ...
  - 30 20**
  - Save on 30 20 and More! Buy, Bid, or Make an Offer now.
  - [www.eBayMotors.com](#)
  - 20 30**
  - Find Bargain Prices On 20 30.
  - [www.BizRate.com](#)

# Du Web 2.0 au Web 3.0

## **Pour trouver encore plus facilement:**

«Le Web sémantique est un web de données, d'une certaine manière comme une base de données globale.»

-- Tim Berners-Lee, 1998

---

**Web 3.0 = Web 2.0 + Web sémantique + outils sémantiques.**

---

## **Thèmes fréquents du Web 3.0:**

- Contenus enrichis sémantiquement.
- Données reliées (« linked data »)
- Sensibilité au contexte.
- Sensibilité à la localisation géographique.

# Exemples

Web [Images](#) [Vidéo](#) [Maps](#) [Actualités](#) [Groupes](#) [Gmail](#) [plus](#) [Connexion](#)


**Google**   [Recherche avancée](#)  
[Préférences](#)

Rechercher dans :  Web  Pages francophones  Pages : France

Web

Recherches associées : [cecilia sarkozy](#) [nicolas sarkozy juif](#)

Résultats dans l'Actualité pour **Nicolas Sarkozy**

 [La venue de Nicolas Sarkozy devant le Congrès](#) - Publié il y a 1  
La venue de **Nicolas Sarkozy**, le 22 juin, devant le Parlement ré  
constituera une première dans l'histoire de la France ...  
[Le Point](#) - [166 autres articles >](#)  
[Nicolas Sarkozy et Angela Merkel soutiennent José Manuel Ba](#)  
L'Express - [292 autres articles >](#)  
[Barack Obama et Nicolas Sarkozy affichent leur bonne entente](#)  
Les Échos - [748 autres articles >](#)

[Nicolas Sarkozy - Wikipédia](#)  
Nicolas Paul Stéphane Sarközy de Nagy-Bocsa, dit **Nicolas Sarkozy** [N  
le 28 janvier 1955 à Paris (17<sup>e</sup> arrondissement), est un homme d'État ...  
[fr.wikipedia.org/wiki/Nicolas\\_Sarkozy](#) - [En cache](#) - [Pages similaires](#)

[Elysee.fr | Présidence de la République](#)  
Les services de la Présidence - Les institutions - L'Élysée et les résident  
PRESIDENT, **Nicolas Sarkozy** ou son porte-parole vous répondent. ...  
[www.elysee.fr](#) - [En cache](#) - [Pages similaires](#)

[Nicolas Sarkozy | Ensemble tout devient possible](#)  
La campagne officielle du président de l'UMP. Contient notamment une "I  
en ligne, des débats, des informations et inscriptions.  
[www.sarkozy.fr](#) - [En cache](#) - [Pages similaires](#)

[Sarkozy Blog - Tout sur Nicolas Sarkozy](#)  
29 avr 2009 ... Donne des informations sur le personnage en reproduisant  
déclarations et des articles que lui consacre la presse.  
[sarkozyblog.free.fr](#) - [En cache](#) - [Pages similaires](#)

[Nicolas Sarkozy | Dossier d'actualité - Yahoo! Actualités](#)  
Le dossier **Nicolas Sarkozy** de Yahoo! Actualités. Les dernières dépêch  
analyses et revue de Web sur **Nicolas Sarkozy**.  
[fr.news.yahoo.com/fr/nicolas-sarkozy.html](#) - [En cache](#) - [Pages similaire](#)

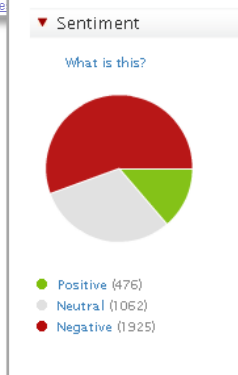
**Newsift** BETA  
FROM THE FINANCIAL TIMES GROUP

New Search | Saved Searches | Search History | Help | Feedback: Take our survey

Search Term: Nicolas Sarkozy

Close Suggestions

<b>Enter term</b> OR Select linked term to refine search	<b>Business Topic</b> Trade Policy Labor Force And Unions Natural And Environmental Risks Litigation And Settlement Environmental Policy <input type="button" value="More options"/>	<b>Organization</b> European Union Air France (AKH) European Parliament European Commission <input type="button" value="More options"/>	<b>Place</b> France Paris Europe Germany London Spain Iran <input type="button" value="More options"/>	<b>Person</b> Nicolas Sarkozy Barack Obama Angela Merkel Gordon Brown Carla Bruni Silvio Berlusconi Prince Charles <input type="button" value="More options"/>	<b>Theme</b> 65th Anniversary Gaulle Airport European Elections France Jet France Plane Automatic Message Day Landings <input type="button" value="More options"/>
--	--	--	--	--	--



**1-10 of 3463 Articles: Frequent terms above**

Sort by: **Relevance** | Date

Timeframe:

**Sarkozy tries to rescue internet law after court decision**  
EUobserver.com, June 12, 2009

Sarkozy tries to rescue internet law after court decision... The government's aim is to get the new agency envisaged in the bill, the Haute Autorité pour la Diffusion des œuvres et la Protection des Droits sur Internet (HADOPI), or High Authority on Diffusion of Works of Art and the Protection of the Rights on the Internet...

SPONSOR ADVERTISEMENT

**Subscribe now for unlimited access to FT.com**

*We live in FINANCIAL TIMES®*

# Le Web 3.0

## **Le text mining permet le Web 3.0 et le Web sémantique.**

- La catégorisation et la classification automatique du contenu.
- L'augmentation de texte: la création de métadonnées; le balisage du contenu.
- L'extraction d'informations vers les bases de données.
- L'analyse exploratoire et la visualisation.

## **Concepts techniques:**

- Les Microformats
- RDF, SPARQL
- OWL

# Text Mining: perspective des utilisateurs

J'ai récemment publié un rapport, "Text Analytics 2009: User Perspectives on Solutions and Providers" («Text Analytics 2009: les perspectives des utilisateurs sur les solutions et les fournisseurs»).

J'ai estimé un marché global de \$350 millions en 2008, une croissance de 40% par rapport à 2007.

J'ai présenté les résultats d'un sondage dans lequel j'ai posé les questions...

# Principales applications

Quelles sont pour vous les applications primaires où le texte joue un rôle?



# Types d'informations textuelles analysées

## Quelles informations textuelles analysez-vous ou projetez-vous d'analyser?

Les utilisateurs **actuels** ont répondu:

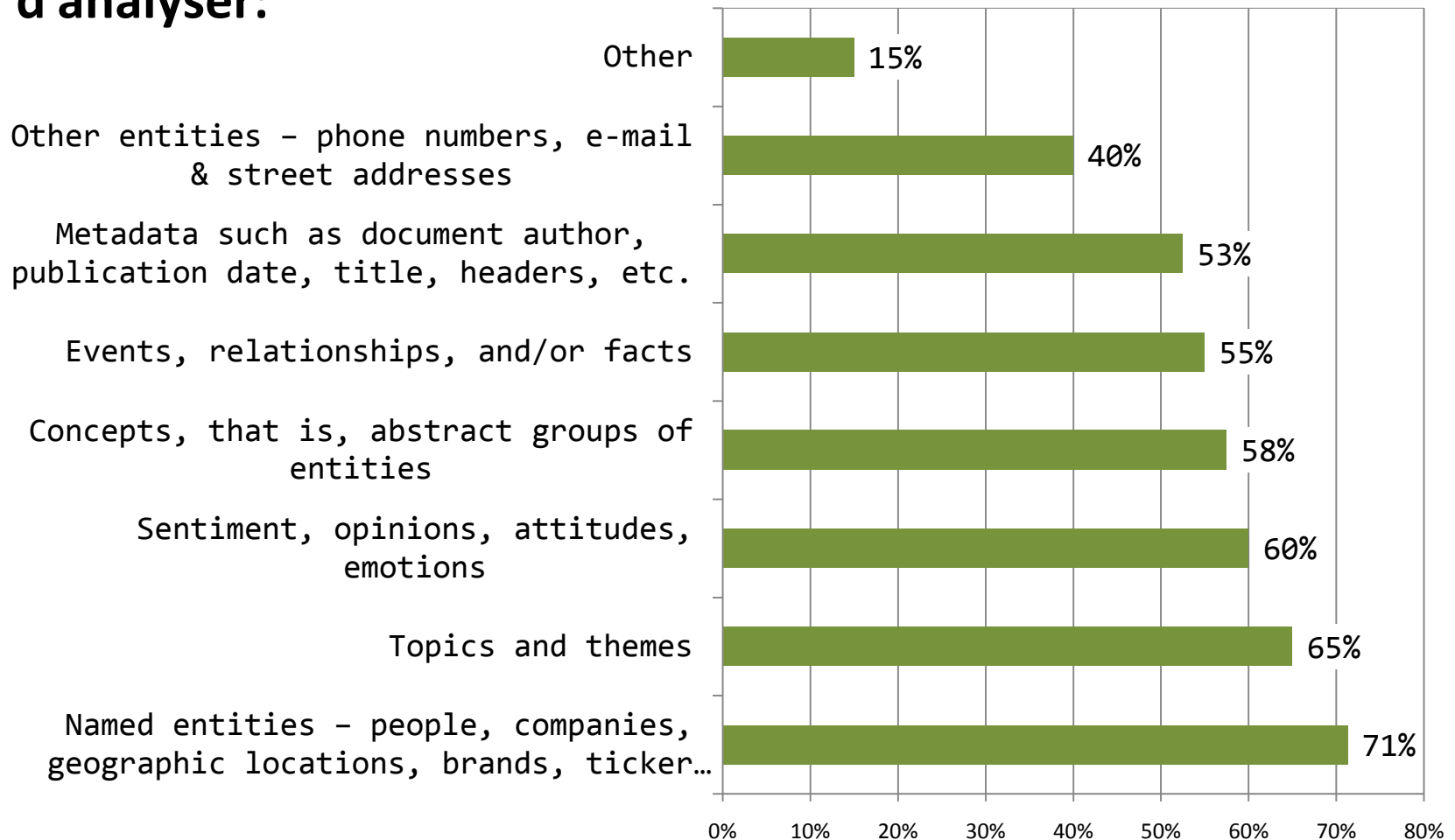
---

blogs and other social media (twitter, social-network sites, etc.)	62%
news articles	55%
on-line forums	41%
e-mail and correspondence	38%
customer/market surveys	35%

---

# Types d'informations extraites

**Avez-vous (ou comptez-vous avoir) besoin d'extraire ou d'analyser:**



# Satisfaction générale

**Veillez évaluer votre expérience générale – votre satisfaction – avec le text mining.**

