

Web Science

Master 1 IFI



Andrea G. B. Tettamanzi

Université de Nice Sophia Antipolis

Département Informatique

andrea.tettamanzi@unice.fr

CM - Séance 3

PageRank et comment Google transforme des mots en argent

Introduction

- Quelques chiffres sur Google (mise à jour : 25 mai 2014)
 - Capitalisation de marché : ~ 375 milliards de dollars
 - Chiffre d'affaires : 62 milliards de dollars
 - Bénéfices avant intérêts, impôts, etc. : 18.6 milliards de \$
 - Environ 54 000 employés
- Par comparaison
 - PIB du Luxembourg : 57 milliards de dollars
 - Si Google était un pays, il serait le 70ème par PIB sur 193
 - En 2013, Google était 15ème au monde par capitalisation parmi les entreprises cotées en bourse
- Pas mal pour un « simple » moteur de recherche...

\$590/s !!!

La clé du succès

- Le succès de Google se base sur deux algorithmes :
 - **PageRank**
 - **AdWords + AdSense**
- Le premier permet de trier les résultats des recherches :
 - valeur d'usage
 - il a imposé Google comme leader du marché
- Le deuxième génère l'impression de messages publicitaires ciblés aux intérêts du public d'une page Web :
 - valeur d'échange
 - AdWords permet d'acheter du trafic, AdSense de le vendre

Plan

- PageRank
- AdWords + AdSense
- TD

1ère Partie

PageRank

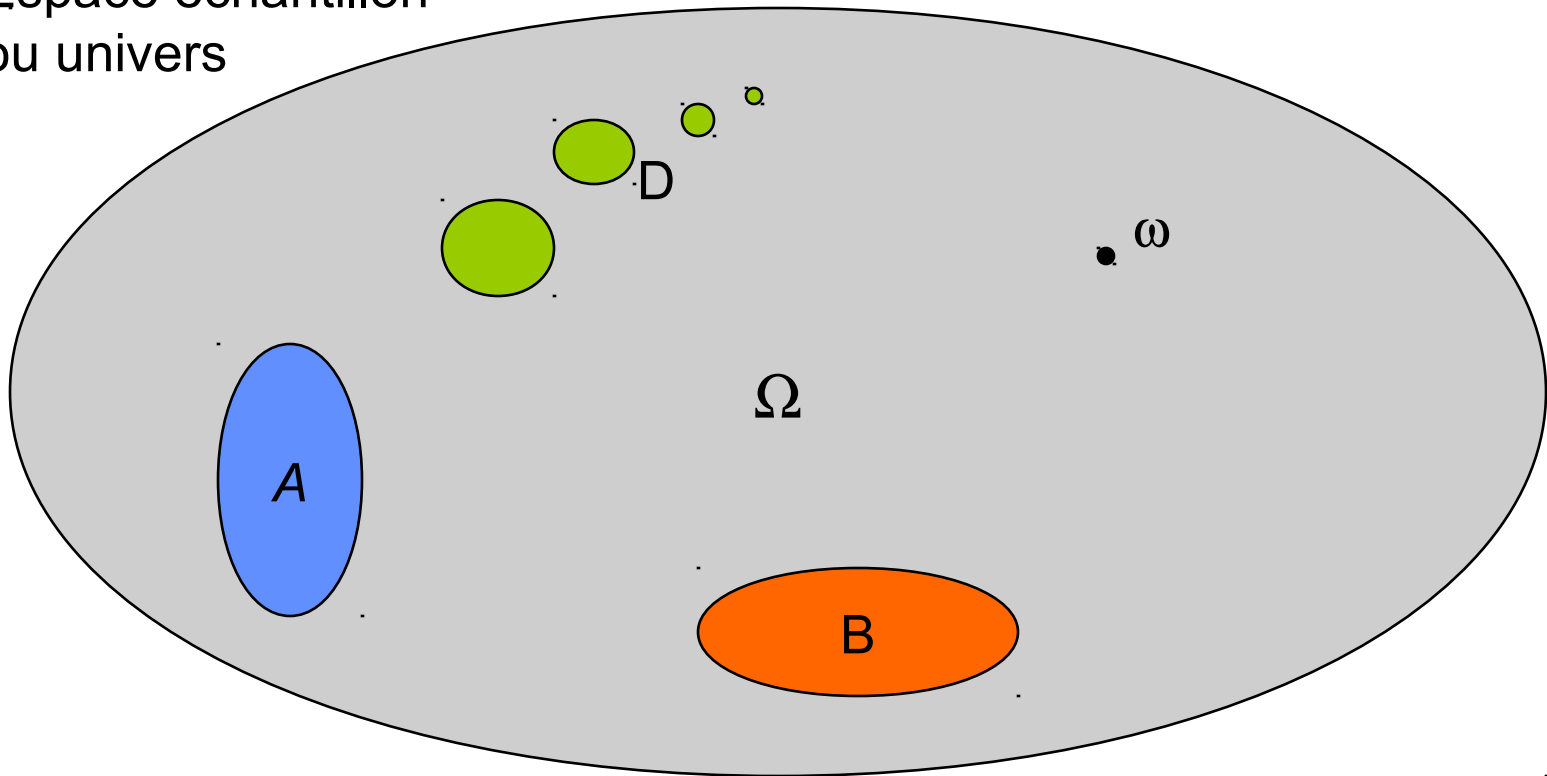


Intuition de fond

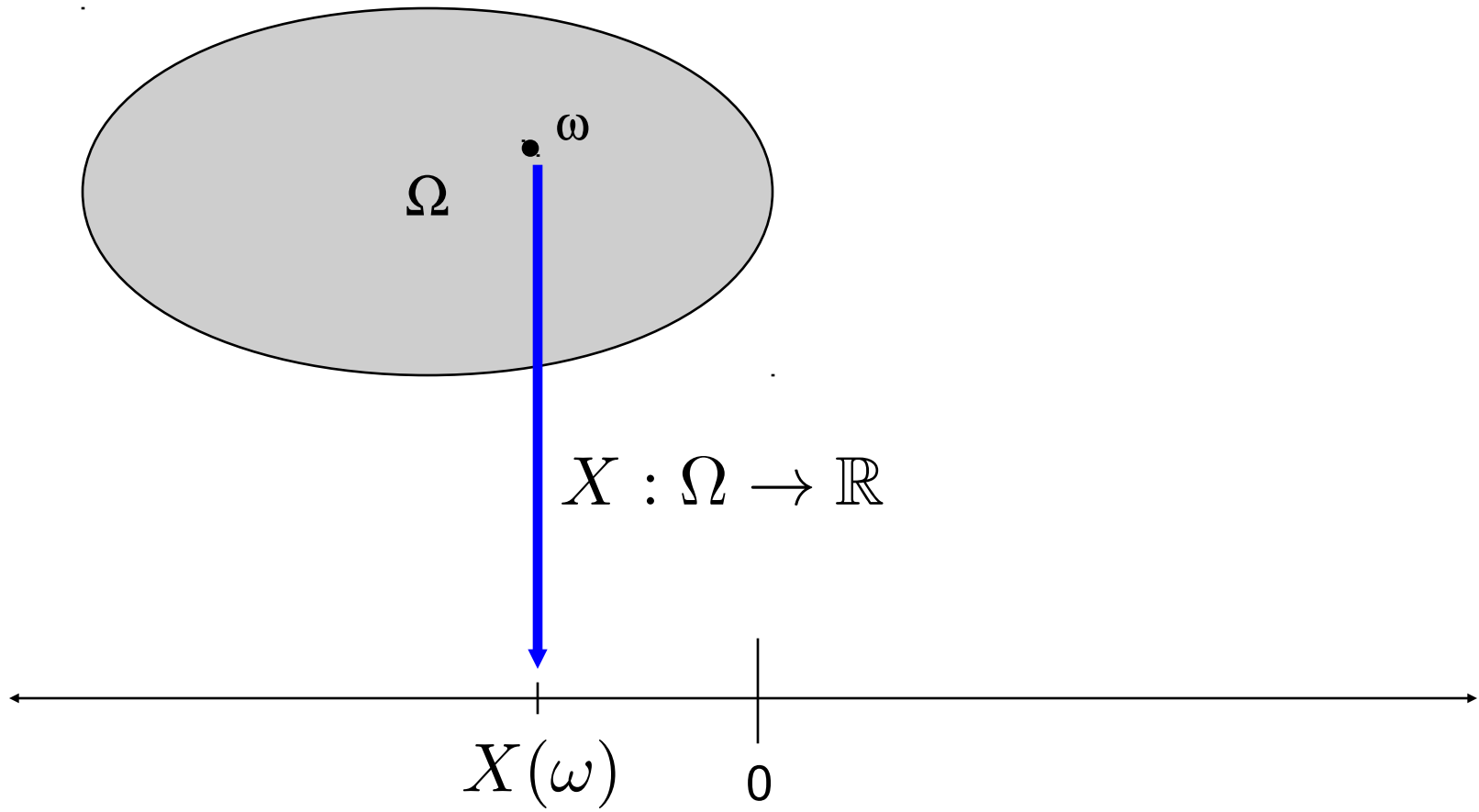
- Le Web comme graphe orienté
 - Ses nœuds sont les pages HTML
 - Ses arcs sont les hyperliens ` . . . `
- Quelles pages visiterait un **surfeur aléatoire** ?
 - Le surfeur aléatoire commencerait par une page arbitraire
 - Il sauterait d'une page à la suivante en cliquant sur un hyperlien choisi au hasard
- Idée : mesurer l'importance d'un page par la probabilité qu'elle soit visitée à l'instant t par un surfeur aléatoire !
- Cette probabilité est la fréquence de visite de la page

Événements

Espace échantillon
ou univers



Variables aléatoires



Processus Stochastiques

Une suite de variables aléatoires

$$X_1, X_2, \dots, X_t, \dots$$

Chacune dotée de sa propre distribution de probabilité.

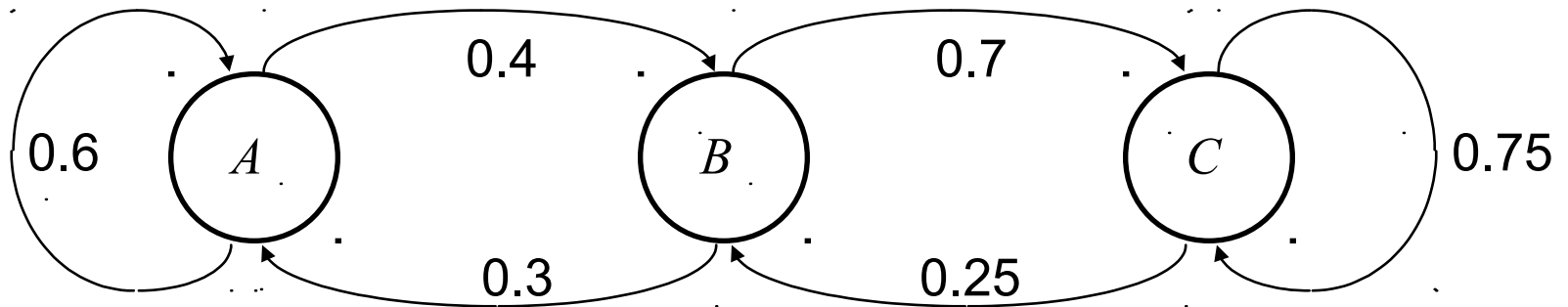
Notation: $\{X_t(\omega)\}_{t=0,1,\dots}$

Chaînes de Markov

Un processus stochastique $\{X_t(\omega)\}_{t=0,1,\dots}$

est une chaîne de Markov si et seulement si, pour tout t ,

$$\Pr[X_t = x \mid X_0, X_1, \dots, X_{t-1}] = \Pr[X_t = x \mid X_{t-1}]$$



Matrice de transition

$$\mathbf{T} = \begin{bmatrix} \Pr(X_t = x_1 \mid X_{t-1} = x_1) & \dots & \Pr(X_t = x_n \mid X_{t-1} = x_1) \\ \Pr(X_t = x_1 \mid X_{t-1} = x_2) & \dots & \Pr(X_t = x_n \mid X_{t-1} = x_2) \\ \vdots & & \vdots \\ \Pr(X_t = x_1 \mid X_{t-1} = x_n) & \dots & \Pr(X_t = x_n \mid X_{t-1} = x_n) \end{bmatrix}$$

\mathbf{T} est une matrice stochastique :

$$\forall i, \quad \sum_{j=1}^n \Pr(X_t = x_j \mid X_{t-1} = x_i) = 1$$

Définition « idéalisée » de PageRank

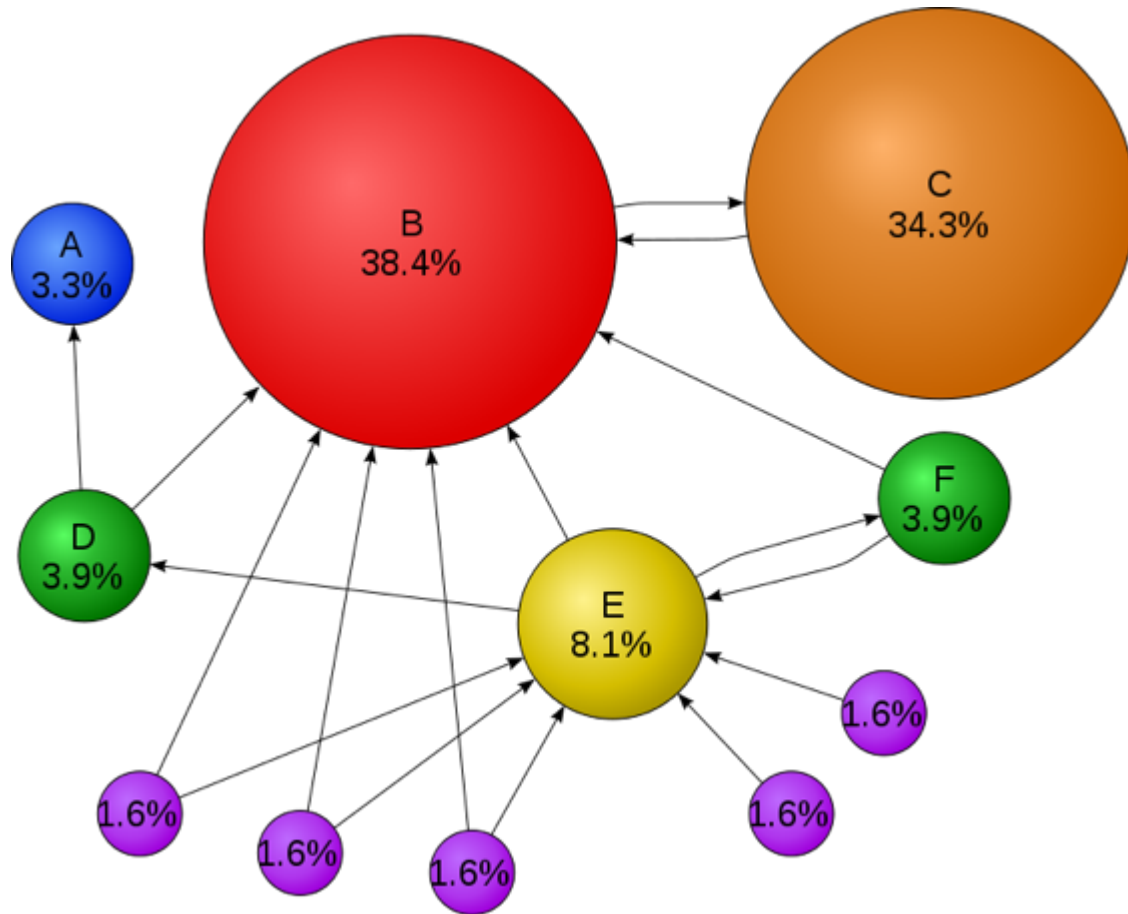
$q_i = N^o$ liens sortant de la page i

$$\mathbf{H} = (h_{ij})$$

$$h_{ij} = \begin{cases} 1/q_i & \text{s'il existe un lien de } i \text{ à } j; \\ 0 & \text{sinon.} \end{cases}$$

$$\pi_j = \sum_i \pi_i h_{ij} \iff \pi = \pi \mathbf{H}$$

Exemple



Hypothèse de fond

Une page Web est importante
si elle est référencée par
d'autres pages importantes

Analyse de la définition

- Il y a trois facteurs qui déterminent le PageRank d'une page :
 - Le nombre de liens qui pointent vers elle ;
 - La propension des pages d'origine de ces liens à diriger le surfeur vers elle, c'est-à-dire, le nombre total de leurs liens sortants ;
 - Le PageRank des pages d'origine des liens
- Le modèle idéalisé a deux problèmes :
 - Les pages sans liens sortants (*dangling*), qui captureraient le surfeur
 - Le surfeur peut aussi rester piégé dans un *bucket* (= sceau), une composante joignable et fortement connectée, sans arcs sortants vers le reste du graphe

Modèle réel : matrice Google

- Les lignes de la matrice \mathbf{H} avec tous les éléments égaux à zéro, qui correspondent aux pages sans liens sortants, sont remplacées par une distribution uniforme ou arbitraire.
- Soit \mathbf{S} la matrice ainsi modifiée.
- Pour résoudre le problème des *buckets*, Brin et Page proposent de remplacer la matrice \mathbf{S} par la matrice Google :

$$\mathbf{G} = \delta \mathbf{S} + (1 - \delta) \mathbf{E}$$

damping factor δ \leftarrow *matrice de téléportation* \mathbf{E}

$$\mathbf{E} = \begin{bmatrix} 1/n & 1/n & \cdots & 1/n \\ \vdots & \vdots & & \vdots \\ 1/n & 1/n & \cdots & 1/n \end{bmatrix}$$

Interprétation de la matrice Google

- La définition de la matrice Google peut être expliquée comme suit
 - Avec probabilité δ , le surfeur aléatoire suit le lien suivant
 - Avec probabilité $1 - \delta$, le surfeur se lasse de suivre des liens et dirige le navigateur vers un nouveau URL, qui pourrait n'avoir rien à voir avec la page courante.
 - Dans ce cas, le surfeur est « téléporté » à cette nouvelle page
- Les inventeurs de PageRank proposent un facteur d'amortissement $\delta = 0.85$:
 - En moyenne, après avoir suivi 5 liens, le surfeur choisit une nouvelle page au hasard.
- Le vecteur PageRank est donc π tel que

$$\pi = \pi \mathbf{G}$$

Existence et unicité du vecteur PageRank

- Le vecteur π est un vecteur propre de \mathbf{G} de valeur propre 1.
- La matrice \mathbf{S} est stochastique, ainsi que la matrice \mathbf{E} .
- La matrice \mathbf{G} est donc stochastique.
- Si \mathbf{G} est stochastique, l'équation $\pi = \pi\mathbf{G}$ a au moins une solution.
- D'après le théorème de Perron-Frobenius, si \mathbf{A} est une matrice carrée non-négative irréductible, alors il existe un vecteur \mathbf{x} tel que $\mathbf{x}\mathbf{A} = r\mathbf{x}$, où r est le rayon spectral de \mathbf{A} .
- La matrice \mathbf{S} est vraisemblablement réductible mais, grâce à la matrice de téléportation, \mathbf{G} ne l'est sûrement pas.
- En plus, puisque \mathbf{G} est stochastique, son rayon spectral est 1.
- Par conséquent, un vecteur PageRank > 0 existe et est unique.

PageRank et Théorie de Markov

- Le modèle de marche aléatoire sur le graphe du Web, modifié avec la téléportation, induit naturellement une chaîne de Markov avec un nombre fini (même si énorme) n d'états (= pages)
- **G** est la matrice de transition de cette chaîne de Markov
- Puisque **G** est irréductible, la chaîne est ergodique et a une distribution stationnaire unique, qui correspond au vecteur PageRank π .

Calcul du vecteur PageRank

- La **méthode de la puissance** est une méthode numérique qui permet de déterminer la valeur propre de module maximal d'une matrice à coefficients réels.
- On prend un vecteur \mathbf{x} au hasard et on calcule la suite récurrente:

$$\mathbf{x}^{(0)} = \mathbf{x}, \quad \mathbf{x}^{(t+1)} = \mathbf{x}^{(t)} \mathbf{A} / \|\mathbf{A}\|$$

- Cette suite converge vers la valeur propre de module maximal de la matrice \mathbf{A}
- Pour calculer π , on part du vecteur $\mathbf{u} = (1/n, \dots, 1/n)$ et on s'arrête lorsque

$$\|\pi^{(t+1)} - \pi^{(t)}\| < \epsilon$$

2ème Partie

AdWords et AdSense

**... ou de Comment Google
transforme des mots en argent**

Enjeu

- Mars 2000 : Éclatement de la « Bulle Internet »
 - Beaucoup de *start-ups* proposant une valeur d'usage, mais pas de valeur d'échange ne résistèrent pas
 - Meilleure idée que de simplement vendre de la publicité
 - Accumulation de capital linguistique grâce à ses services
 - Exploiter ce capital
- Un algorithme qui organise automatiquement la spéculation autour des mots a permis de créer le premier marché linguistique mondial
- Trademarks : il était déjà possible d'acheter certains mots
- Google a élargi et libéralisé ce marché

1

[X-HelvetiC Tours | helvetic-tours.ch](http://helvetic-tours.ch)www.helvetic-tours.ch/xHelveticToursDes **vacances** à la mer aussi peu chères, ça fait vraiment du bien!

2

[Vacances tout compris | clubmed.ch](http://clubmed.ch)www.clubmed.ch/-15% sur vos **vacances** d'hiver 12/13 ou jusqu'à 480 CHF offerts now !

3

[Voyages Jusqu'à -70% - Offres Imbattables: Vos Voyages](http://groupon.ch/Voyages)www.groupon.ch/Voyages

Jusqu'à -70% avec Groupon. Ici !

[Vacances 2012. L'été à la mer, en France, Corse, Var, Bretagne ...](http://www.vacances.com/)www.vacances.com/Nombreuses annonces, de professionnels et de particuliers, pour les **vacances** d'hiver. Chalets, studios, appartements au ski. Week-ends, voyages, séjours en ...

↳ Location - Espagne - Location Aquitaine - Provence et Côte d'Azur

[Le calendrier scolaire - Ministère de l'Éducation nationale](http://www.education.gouv.fr)www.education.gouv.fr > ... > Le ministère > Repères, histoire et patrimoine**Vacances**, Zone A, Zone B, Zone C ... **Vacances** de la Toussaint ... Le départ en **vacances** a lieu après la classe, la reprise des cours le matin des jours indiqués.

↳ Vacances Scolaires 2011-2012 - Le calendrier scolaire ... - MENE0914826A

[Vacances Look Voyages : séjour pas cher en famille, club tout ...](http://www.look-voyages.fr/)www.look-voyages.fr/Votre séjour en club de **vacances** Lookea, en hôtel ou en formule circuit au meilleur prix avec Look Voyages. Départ dernière minute, pas cher ou en promo ![Villages et Clubs de vacances en tout inclus Thomas Cook](http://tt.thomascook.fr/village-club-vacances/)tt.thomascook.fr/village-club-vacances/

Ambiance. Une ambiance animée en journée comme en soirée. CHAQUE SEMAINE, NOS ANIMATEURS CONCOCTENT LE PROGRAMME DU VILLAGE : ...

[Tech Travel](http://www.techtravel.ch/)www.techtravel.ch/

les meilleures offres de voyages Culturels, Escapades, Plages

[Site Officiel Air Transat](http://www.airtransat.ch/)www.airtransat.ch/

Les plus bas prix pour les vols vers le Canada. Réservez en ligne!

[Voyager Moins Cher](http://www.voyagermoinscher.com/)www.voyagermoinscher.com/

Voyage et billet d'avion Promotions et dernière minute.

[Vacances en France](http://www.declifrance.com/Vacance_France)www.declifrance.com/Vacance_FranceTrouvez vos **vacances** en France parmi plus de 5000 offres ![Vacances Go Voyages](http://www.govoyages.com/Vacances)www.govoyages.com/VacancesChoisissez vos **vacances** au prix le plus bas chez Go Voyages![Vacances à l'étranger](http://www.sejour-express.com/)www.sejour-express.com/Trouvez les moins chères sur le comparateur de **vacances**[Vacances en Club](http://www.club-vacances-express.fr/)www.club-vacances-express.fr/

Trouvez les moins chères sur le comparateur de location en Club

4

5

6

7

8

9

10

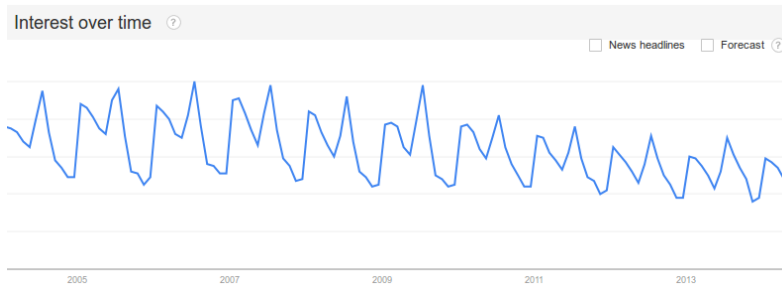
AdWords

- Mécanisme d'enchère sur des mots pour placer des annonces
- Tous les mots(-clefs) peuvent donner lieu à des enchères
- L'algorithme classe automatiquement les annonces selon un calcul en quatre étapes :
 - Enchère sur un mot (E) : l'annonceur fixe un prix max qu'il est prêt à payer en cas de clic
 - Calcul du score qualité Q de l'annonce (pertinence) : **secret !**
 - Calcul du rang de l'annonce : $R = E Q$, et classement i
 - Calcul du prix à payer en cas de clic :

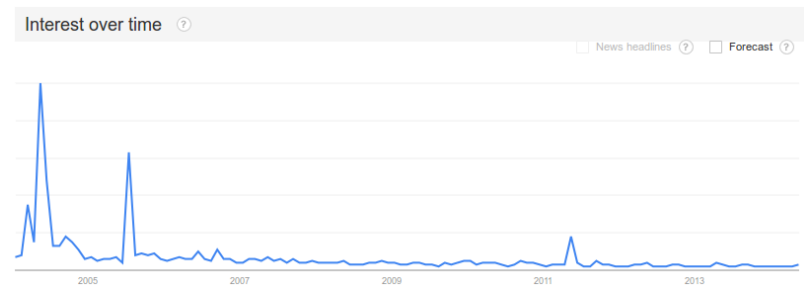
$$P_i = E_{i+1} \frac{Q_i}{Q_{i+1}}$$

Google Trends

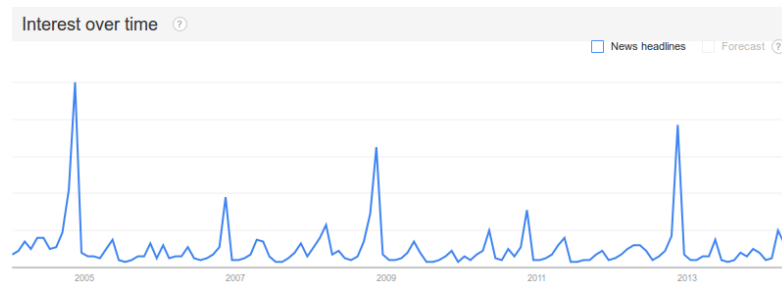
Vacances



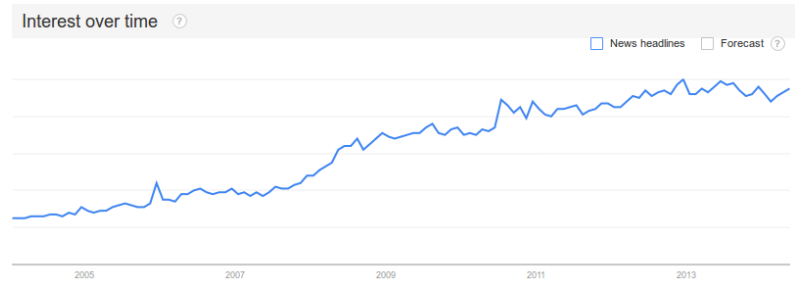
al-Qaeda



Elections



Porn



Achat et vente de trafic

Ad Words

Annonceurs

Les annonceurs font des enchères sur les mots pour en acheter le trafic



Google



Ad Sense

Propriétaires/créateurs de pages Web

Les sites Web vendent leur trafic à Google pour afficher les annonces

Avantages pour les utilisateurs

- Des services « gratuits » (recherche, info, email, maps, etc.)
- Une publicité utile, pertinente, non-invasive
- Une expérience agréable des contenus en ligne

Deux sources de bénéfice

Google search results for "bicycle shops in mountain view". The search bar shows "bicycle shops in mountain view" and "Search". Below the search bar, there are navigation links for "Everything", "Maps", "More", and "Show search tools". The main content area is titled "Local business results for bicycle shops near Mountain View, CA". It features a map on the left and a list of three businesses on the right:

- Performance Bicycle Shop** - www.performancebike.com
2124 West El Camino Real, Mountain View - (650) 964-1796
9 reviews, directions, and more »
- El Camino Bicycle Shop Jack'Ssee Off Ramp The** - offrampbikes.com
2320 West El Camino Real, Mountain View - (650) 968-2974
"Their customer service has always been great."
★★★★☆ 21 reviews, directions, and more »
- REI - Mountain View** - www.rei.com
2450 Chabot Road, Mountain View

Below the list, there are "Sponsored links" for "Local Business Listings", "Mountain Bikes Store", and "New Bikes Up To 60% Off".

Publicité sur les sites de Google, tels que

- google.com
- gmail.com
- orkut.com
- Youtube.com



A screenshot of the TechCrunch website. The header includes "Tech", "Gadgets", "Mobile", "Enterprise", and "More". The main content area features a large advertisement for "THE 88" with the text "NOW 50% SOLD! Urban living begins in the \$300,000s" and "GO TO THE885J.COM San Jose, CA". Below the ad, there are social media icons and a "Subscribe:" button.

Publicité sur les sites des clients adSense

Merci de votre attention

