

**SéminDoc 06/05/2009**  
**LIRMM – Montpellier**

# ***Visualiser un texte par un nuage arboré***

Philippe Gambette (équipes MAB/AIGco)  
en collaboration avec Jean Véronis (Aix-en-Provence)



# Plan

---

- Nuages de tags et de mots
- Nuages de tags améliorés
- Nuages arborés
- Etapes de construction
- Contrôle qualité
- Choix des paramètres
- Limites de la méthode

# Nuages de tags

- Construits depuis un ensemble de tags
- Taille de police liée à la fréquence



Ce qui est habituellement cité  
comme le premier nuage de tags,  
dans  
*Microserfs* de D. Coupland,  
HarperCollins, Toronto, 1995





















# Nuages de tags/mots améliorés

Ajouter de l'information extraite du texte :

- pâleur pour exprimer la désuétude dans Amazon
- tags partagés en rouge dans del.icio.us
- regrouper les tags cooccurrents sur la même ligne  
Hassan-Montero & Herrero-Solana, InScit'06
- optimiser l'espace vide et la proximité sémantique  
Kaser & Lemire, WWW'07
- “topigraphy”: placement 2D d'après la cooccurrence  
Fujimura, Fujimura, Matsubayashi, Yamada & Okuda, WWW'08

# Nuages de tags/mots améliorés

Ajouter de l'information extraite du texte :

- pâleur pour exprimer la désuétude dans Amazon
- tags partagés en rouge dans del.icio.us
- regrouper les tags cooccurrents sur la même ligne  
Hassan-Montero & Herrero-Solana, InScit'06
- optimiser l'espace vide et la proximité sémantique  
Kaser & Lemire, WWW'07
- “topigraphy”: placement 2D d'après la cooccurrence  
Fujimura, Fujimura, Matsubayashi, Yamada & Okuda, WWW'08

## Most Popular Tags (What's this?)

Welcome to the Amazon.com tag cloud. Tags are labels customers can use to classify a product. More frequently used tags are larger and more recently used tags will appear **darker**.

1080p action adventure american history animation anime art baby best canceled tv shows  
biography blu-ray book business canon children childrens books christian christianity  
christmas classic classic movie classic rock classical music comedy comics cookbook cooking  
defectivebydesign digital camera disney drama dvd erotica exercise family fantasy fiction stress fun  
games gift idea graphic novel harry potter hd dvd hdtv health hip hop historical fiction historical  
romance history horror humor inspirational ipod jazz kids kindle love magic manga  
meditation memoir metal movie mp3 player music mystery nonfiction paranormal  
romance pc game philosophy photography playstation 3 poetry politics progressive rock psychology  
religion rock romance rpg science science fiction self-help sex soundtrack  
spirituality suspense thriller toys travel tv series vampire vampire romance video  
games wii women world war ii xbox 360



# Nuages de tags/mots améliorés

Ajouter de l'information extraite du texte :

- pâleur pour exprimer la désuétude dans Amazon
- tags partagés en rouge dans del.icio.us
- regrouper les tags cooccurrents sur la même ligne
- optimiser l'espace vide et la proximité sémantique
- “topigraphy”: placement 2D d'après la cooccurrence

Hassan-Montero & Herrero-Solana, InScit'06

Kaser & Lemire, WWW'07

Fujimura, Fujimura, Matsubayashi, Yamada & Okuda, WWW'08





# Nuages de tags/mots améliorés

Ajouter de l'information extraite du texte :

- pâleur pour exprimer la désuétude dans Amazon
- tags partagés en rouge dans del.icio.us
- regrouper les tags cooccurents sur la même ligne

Hassan-Montero & Herrero-Solana, InScit'06

- optimiser l'espace vide et la proximité sémantique

Kaser & Lemire, WWW'07

- “topigraphy”: placement 2D d'après la cooccurrence

Fujimura, Fujimura, Matsubayashi, Yamada & Okuda, WWW'08



# Extraire l'information sémantique d'un texte

- analyse littéraire :

*approche philologique* : se concentrer sur le texte

Brody

- analyse du discours :

analyse arborée, graphe de cooc., projection géodésique

Brunet (Hyperbase), Viprey (Astartex)

- fouille de texte :

graphe sémantique

Grimmer (Wordmapper), Viegas et al. (IBM Many Eyes)

- traitement des langues naturelles :

désambiguïsation

Véronis (Hyperlex)



# Extraire l'information sémantique d'un texte

- analyse littéraire :

*approche philologique* : se concentrer sur le texte

Brody

- analyse du discours :

analyse arborée, graphe de cooc., projection géodésique

Brunet (Hyperbase), Viprey (Astartex)

- fouille de texte :

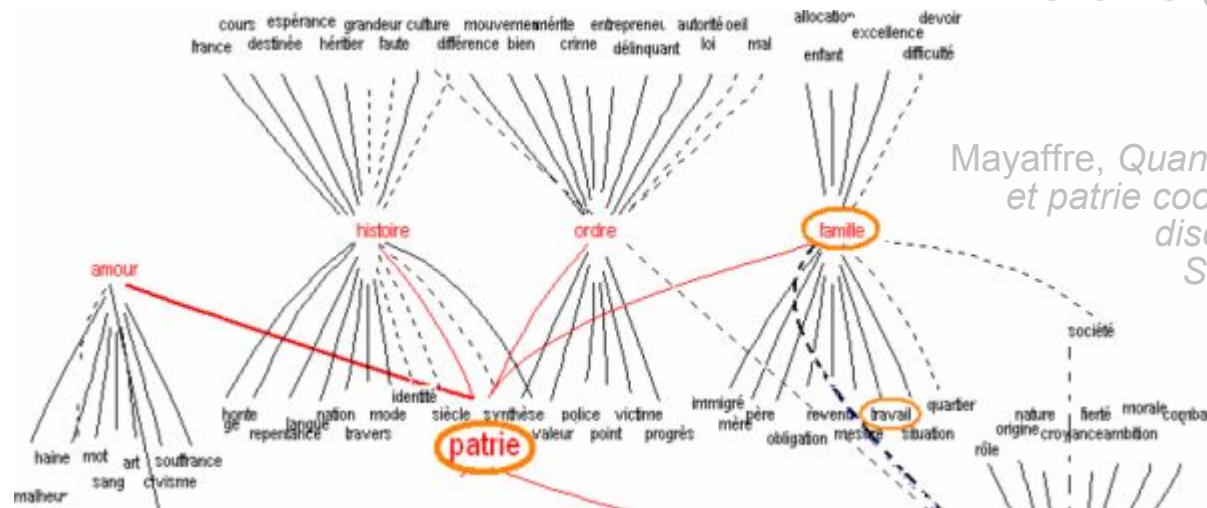
graphe sémantique

Grimmer (Wordmapper), Viegas et al. (IBM Many Eyes)

- traitement des langues naturelles :

désambiguïsation

Véronis (Hyperlex)



Mayaffre, *Quand travail, famille, et patrie cooccurrent dans le discours de Nicolas Sarkozy*, JADT'08

# Extraire l'information sémantique d'un texte

- analyse littéraire :

*approche philologique* : se concentrer sur le texte

Brody

- analyse du discours :

analyse arborée, graphe de cooc., projection géodésique

Brunet (Hyperbase), Viprey (Astartex)

- fouille de texte :

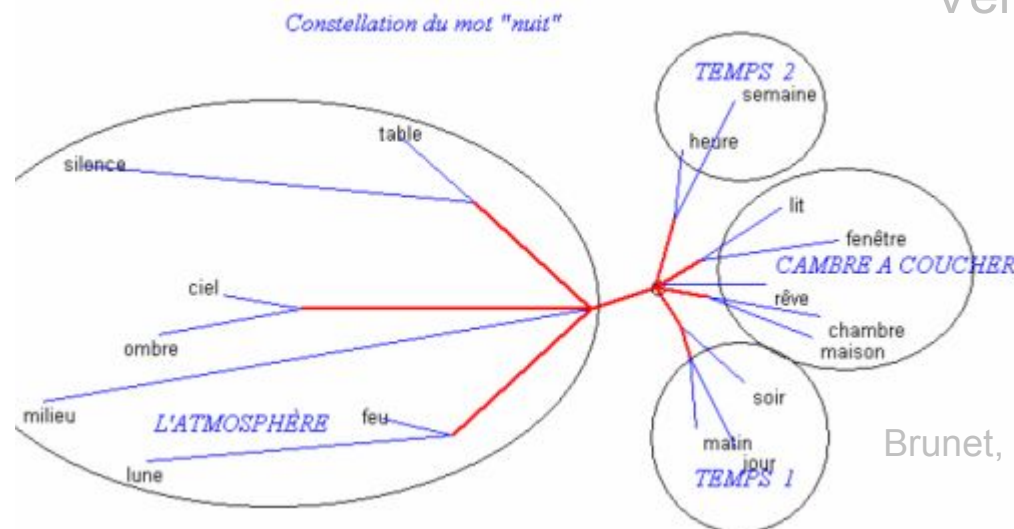
graphe sémantique

Grimmer (Wordmapper), Viegas et al. (IBM Many Eyes)

- traitement des langues naturelles :

désambiguïsation

Véronis (Hyperlex)



Brunet, *Les séquences (suite)*,  
JADT'08



# Extraire l'information sémantique d'un texte

- analyse littéraire :  
*approche philologique* : se concentrer sur le texte

Brody

- analyse du discours :  
analyse arborée, graphe de cooc., projection géodésique

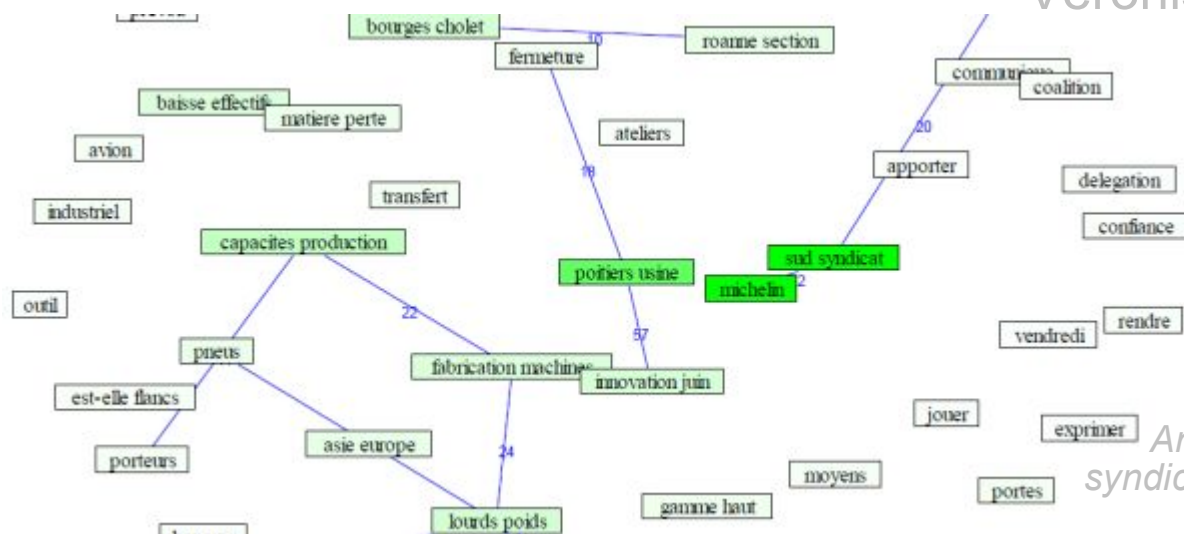
Brunet (Hyperbase), Viprey (Astartex)

- fouille de texte :  
graphe sémantique

Grimmer (Wordmapper), Viegas et al. (IBM Many Eyes)

- traitement des langues naturelles :  
désambiguïsation

Véronis (Hyperlex)



Peyrat-Guillard,  
*Analyse du discours  
syndical sur l'entreprise*,  
JADT'08





# Extraire l'information sémantique d'un texte

- analyse littéraire :

*approche philologique* : se concentrer sur le texte

Brody

- analyse du discours :

analyse arborée, graphe de cooc., projection géodésique

Brunet (Hyperbase), Viprey (Astartex)

- fouille de texte :

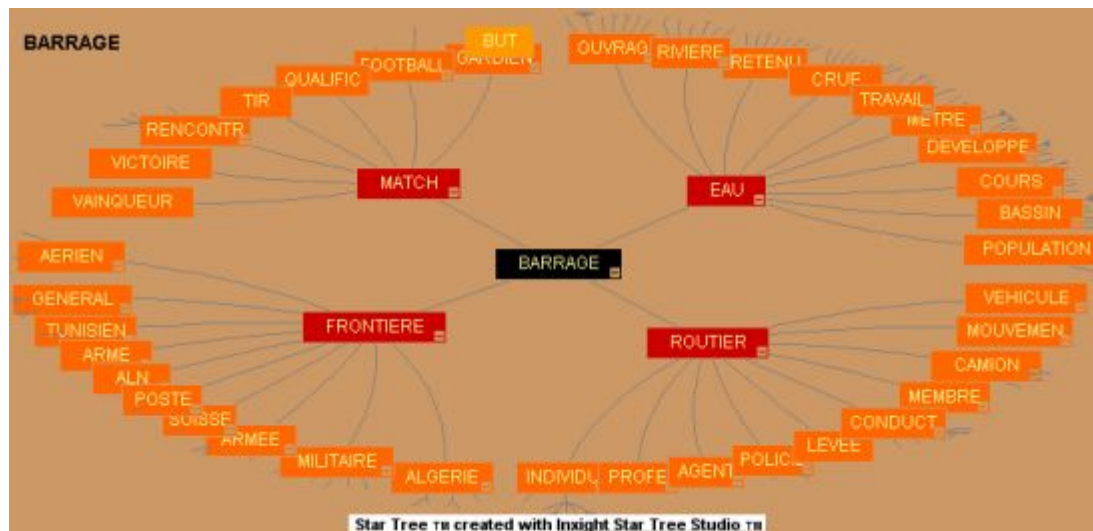
graphe sémantique

Grimmer (Wordmapper), Viegas et al. (IBM Many Eyes)

- traitement des langues naturelles :

désambiguïsation

Véronis (Hyperlex)



Désambiguïsation du mot  
"barrage".

Véronis, *HyperLex: Lexical Cartography for Information Retrieval*, 2004





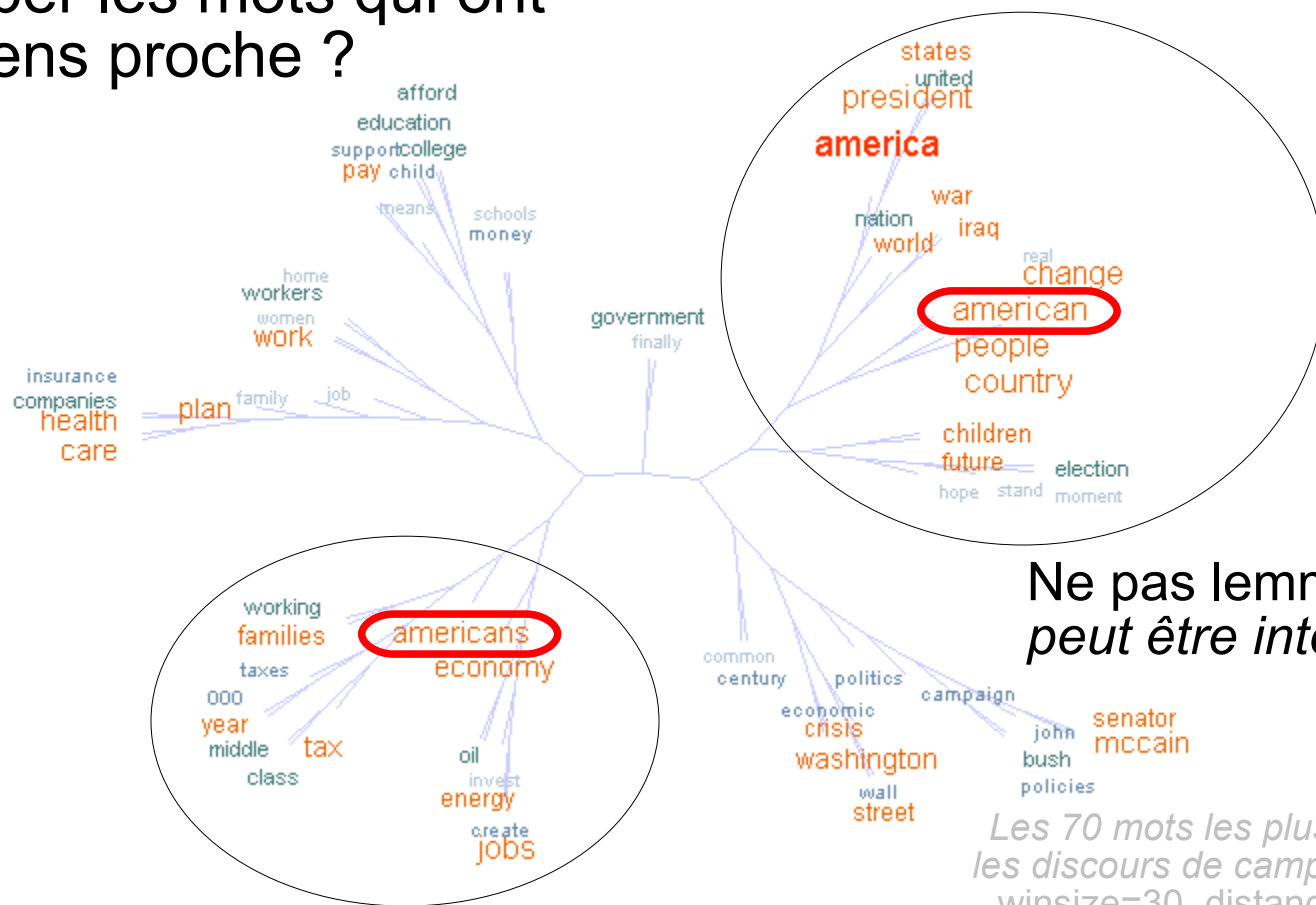




# Construction – extraction des mots

Extraire les mots avec leur fréquence :

- lemmatisation?
- grouper les mots qui ont un sens proche ?



# Construction – matrice de dissimilarité

---

De nombreuses formules de distance sémantique utilisent la cooccurrence.

# Construction – matrice de dissimilarité

De nombreuses formules de distance sémantique utilisent la cooccurrence.

Texte

fenêtre  
glissante S

→ Pas de  
glissement  $s$

largeur  $w$

matrices de cooccurrence

$O_{11}$ ,  $O_{12}$ ,  $O_{21}$ ,  $O_{22}$

	$v \in S$	$v \notin S$
$u \in S$	$O_{11}$	$O_{12}$
$u \notin S$	$O_{21}$	$O_{22}$



matrice de dissimilarité  
sémantique

*chi squared, mutual information, liddel, dice, jaccard, gmean, hyperlex, minimum sensitivity, odds ratio, zscore, log likelihood, poisson-stirling...*

# Construction – matrice de dissimilarité

Transformations à appliquer pour obtenir une dissimilarité :

- transformer la similarité en dissimilarité
- normalisation linéaire pour les matrices positives pour avoir des distances dans l'intervalle  $[0,1]$
- normalisation affine pour les matrices avec des nombres positifs et négatifs, pour avoir des distances dans  $[\alpha,1]$  (par exemple  $\alpha=0.1$ )

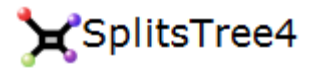


# Construction – construction de l'arbre

Plusieurs méthodes possibles :

- Neighbor-Joining

Saitou & Nei, 1987



- Variantes d'Addtree

Barthelemy & Luong, 1987

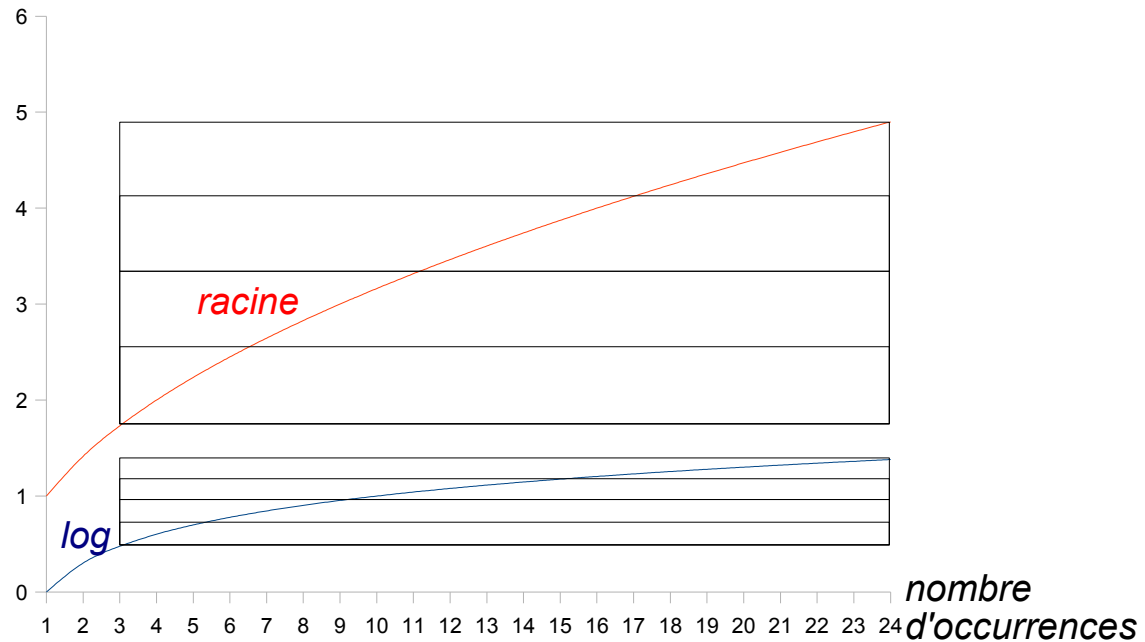
- Heuristique de quadruplets

Cilibrasi & Vitanyi, 2007

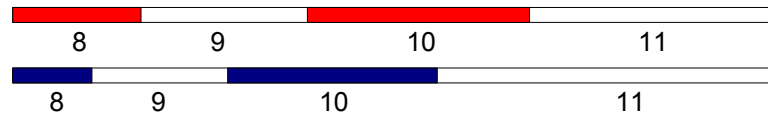
# Construction – décoration de l'arbre

Choix des tailles de polices :

- **fréquence** (appliquer un log !)



taille de police :



# Construction – décoration de l'arbre

Choix des tailles de polices :

- **fréquence** (appliquer un log !)

ACTEURS ADORE ADORÉ AIME AIMÉ ALLÉE ALLER AMOUR  
AMOUREUX ANNÉE ANS ARRIVE BEAU BREF CINÉMA CITÉ  
COEUR COMPLÈTEMENT COMPREND COMPRIS COUP CROIT DÉBUT  
DEVIENT DRÔLE ENFANTS ENVIE ÉVIDEMMENT FAMILLE FEMME  
FILLE **FILM** FILMÉ FILMS FILS FIN FOIS FONT  
FORMIDABLE GÉNIAL GENS GRAND GRANDE GUERRE  
HISTOIRE HOMME JEAN JEUNE JOLI JOUE JOUÉ LONG  
MAGNIFIQUE MAISON MARCHÉ MARI MARRANT MÉCHANT  
MEILLEURS MÈRE MIGNON MOMENT MONDE MORT PART  
PASSE PAUVRE PÈRE PETIT PETITE PETITS PEUR PLEIN  
PREMIER RÉSUMÉ RIGOLO RIRE RÔLE SAIS SAIT SEMAINE SUBLIME  
SUPER SYMPA TELLEMENT TEMPS TOMBE TROUVE TUE TUER  
TYPE VIE VIEUX VOIR VOIT VONT VRAI VRAIE VRAIMENT  
VU

*racine*

ACTEURS ADORE ADORÉ AIME AIMÉ ALLÉE ALLER  
AMOUR AMOUREUX ANNÉE ANS ARRIVE BEAU BREF  
CINÉMA CITÉ COEUR COMPLÈTEMENT COMPREND COMPRIS  
**COUP** CROIT DÉBUT DEVIENT DRÔLE ENFANTS ENVIE  
ÉVIDEMMENT FAMILLE FEMME FILLE **FILM** FILMÉ  
FILMS FILS FIN FOIS FONT FORMIDABLE GÉNIAL  
GENS **GRAND** GRANDE GUERRE HISTOIRE HOMME  
JEAN **JEUNE** JOLI JOUE JOUÉ LONG MAGNIFIQUE  
MAISON MARCHÉ MARI MARRANT MÉCHANT MEILLEURS MÈRE  
MIGNON MOMENT MONDE MORT PART PASSE PAUVRE  
PÈRE **PETIT** PETITE PETITS PEUR PLEIN PREMIER RÉSUMÉ  
RIGOLO RIRE RÔLE SAIS SAIT SEMAINE SUBLIME  
SUPER SYMPA TELLEMENT TEMPS TOMBE TROUVE TUE  
TUER TYPE VIE VIEUX **VOIR** VOIT VONT VRAI VRAIE  
VRAIMENT VU

*log*

Merci, la loi de Zipf !

# Construction – décoration de l'arbre

Choix des tailles de polices :

- **fréquence** (appliquer un log !)

ou

- **classement des fréquences** (distribution exponentielle)

ou

- **saillance** par rapport à un corpus de référence

ou

- **charge émotionnelle** calculée par l'entropie

Eda, Uchiyama, Uchiyama, Yoshikawa, WWW 2009



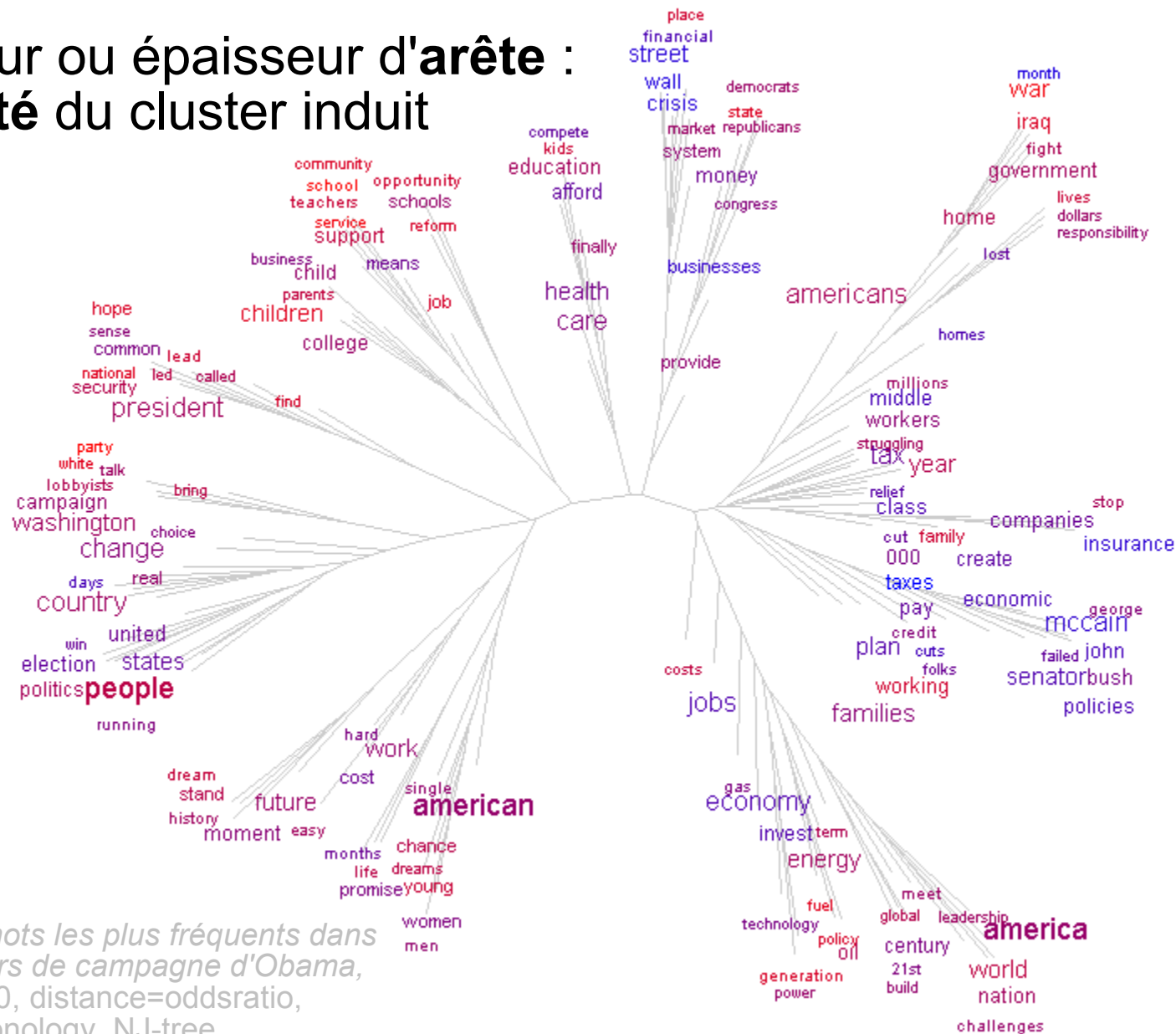






# Construction – décoration de l'arbre

couleur ou épaisseur d'arête :  
qualité du cluster induit



Les 150 mots les plus fréquents dans  
les discours de campagne d'Obama,  
winsize=30, distance=oddsratio,  
color=chronology, NJ-tree.

# Démo

**Treecloud**

Treecloud  
 Français  
 English

Distance  
 liddell  
 gmean  
 jaccard  
 dice  
 ms  
 zscore  
 hyperlex  
 chisquared  
 poissonstirling  
 loglikelihood  
 oddsratio  
 ngd  
 mi

Fenêtre glissante  
Taille:  
30  
Déplacement:  
1

Couleurs  
 yahoo  
 berry  
 chronology  
 dispersion  
 chronodisp

Longueurs d'arêtes  
 unitaires  
 réelles

Ce programme est une interface graphique pour Treecloud, et permet de construire le nuage arboré d'un texte. Téléchargez Treecloud, ainsi que son code source Python, sur <http://www.treecloud.org>.

Emplacement de Python (télécharger version 2.X sur [www.python.org](http://www.python.org))  
C:\Python26\python.exe

Emplacement de SplitsTree (télécharger version 4.X sur [www.splitstree.org](http://www.splitstree.org))  
C:\Treecloud\SplitsTree.lnk


*Le chemin de fichier ne doit pas contenir d'espace. Sinon, créez un raccourci vers SplitsTree dont le nom de fichier ne contient pas d'espace.*

Texte à visualiser : Ouvrir un fichier texte

Antidictionnaire :  
C:\Sites\GambetteLirrm\TreecloudDistribution\StoplistFrench.txt Perso

Nombre de mots du nuage arboré :  ou nombre minimal d'occurrences pour apparaître dans le nuage arboré :

Ligne de commande :  
`"C:\Python26\python.exe" C:\Sites\GambetteLirrm\TreecloudDistribution\Treecloud.py  
stoplist=C:\Sites\GambetteLirrm\TreecloudDistribution\StoplistFrench.txt  
splitstreepath=C:\Treecloud\SplitsTree.lnk unit=1 minnb=4 distance=jaccard window=30 step=1`

 Calcule le nuage arboré avec TreeCloud !

Interface graphique pour TreeCloud

# Contrôle qualité

---

Peut-on mesurer objectivement la qualité des nuages arborés ?



# Contrôle qualité

---

Peut-on mesurer objectivement la qualité des nuages arborés ?

Quelle est la meilleure méthode pour construire un nuage arboré à partir de mes données ?

# Contrôle qualité

Peut-on mesurer objectivement la qualité des nuages arborés ?

Quelle est la meilleure méthode pour construire un nuage arboré à partir de mes données ?

Variations du nuage arboré en cas de petits changements ?

➔ **bootstrap** pour évaluer :

- **stabilité du résultat**
- **robustesse de la méthode**

# Contrôle qualité

Peut-on mesurer objectivement la qualité des nuages arborés ?

Quelle est la meilleure méthode pour construire un nuage arboré à partir de mes données ?

Variations du nuage arboré en cas de petits changements ?

➡ **bootstrap** pour évaluer :

- **stabilité du résultat**
- **robustesse de la méthode**

Y a-t-il une méthode plus directe ?

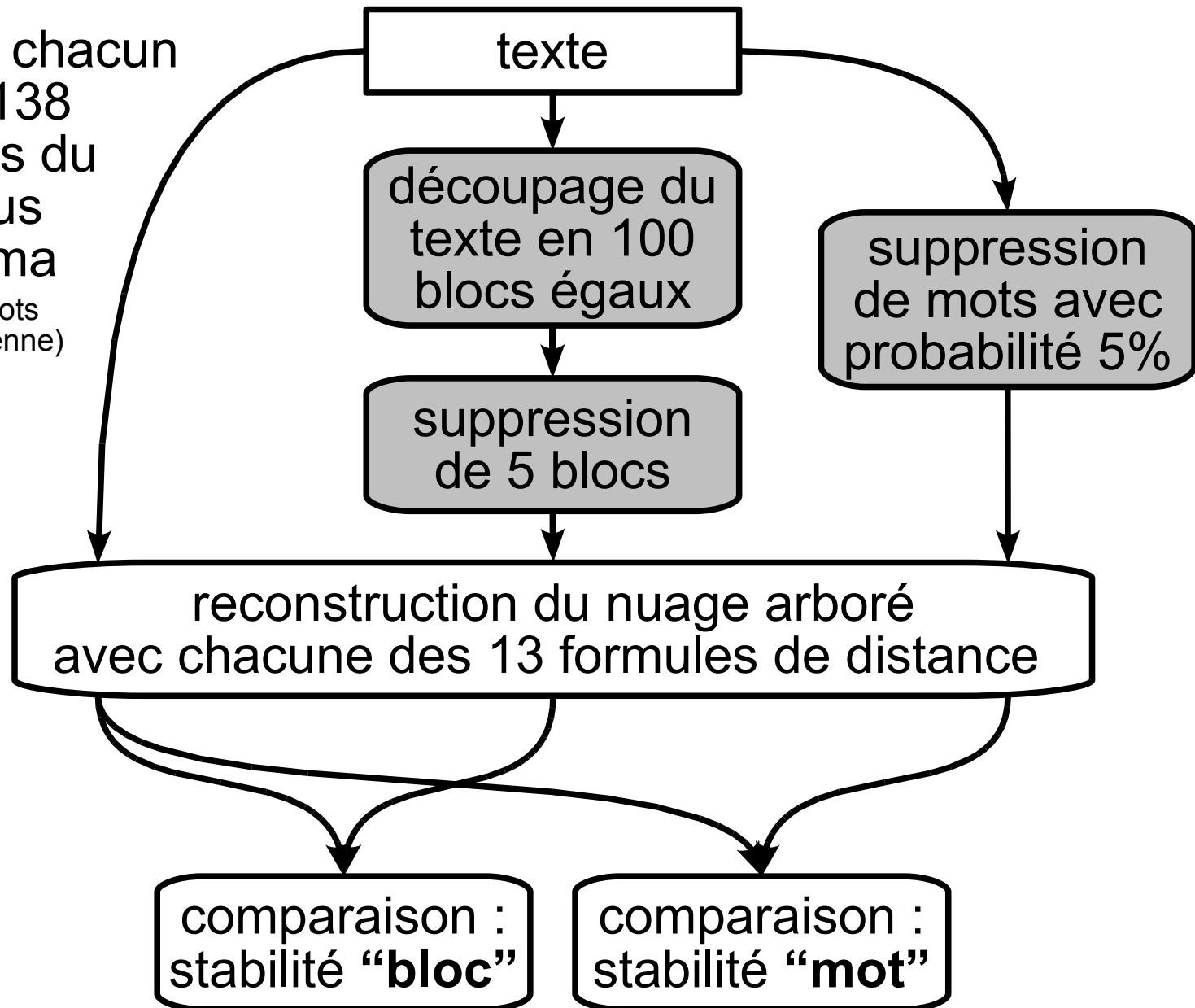
➡ **l'arboricité** montre à quel point la matrice de distance correspond à un arbre

➡ **implique stabilité ?**

Guénoche & Garreta, 2001  
Guénoche & Darlu, 2009

# Contrôle qualité – bootstrap

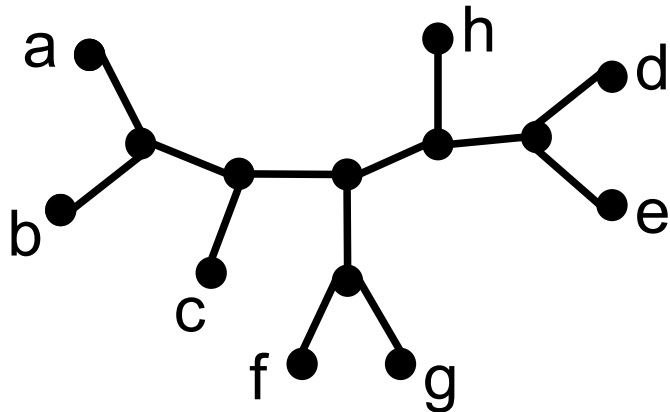
Pour chacun  
des 138  
textes du  
corpus  
Obama  
(3000 mots  
en moyenne)



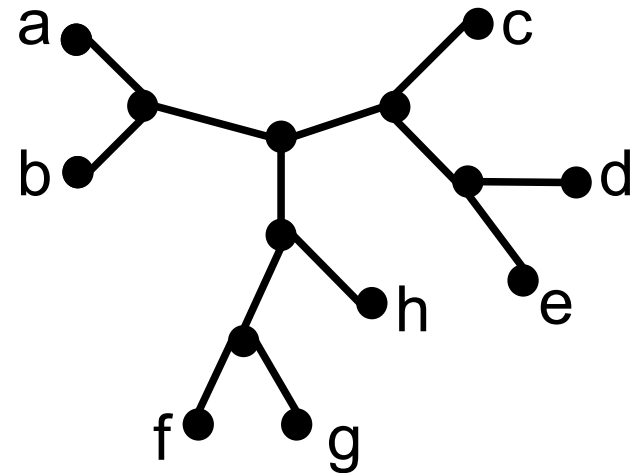


# Contrôle qualité – similarité d'arbres

arbre obtenu sur le texte



arbre obtenu sur le texte altéré

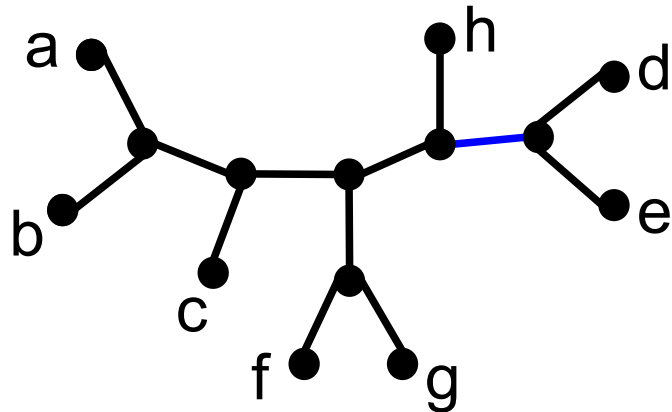


**Distance de Robinson-Foulds :**  
Nombre de splits différents

**Similarité :**  
Pourcentage de splits non triviaux identiques

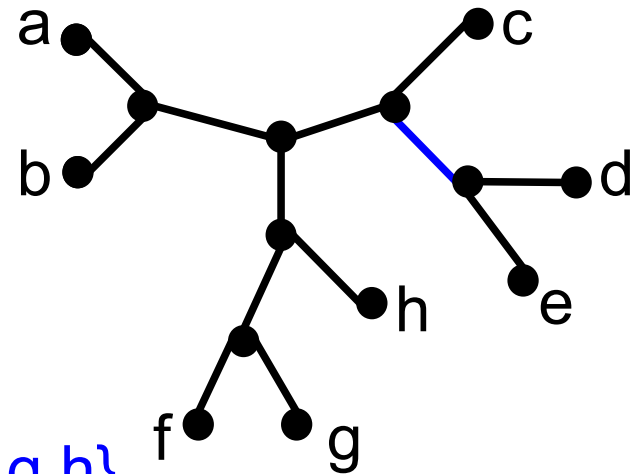
# Contrôle qualité – similarité d'arbres

arbre obtenu sur le texte



split  $\{d,e\}$  séparé de  $\{a,b,c,f,g,h\}$

arbre obtenu sur le texte altéré



**Distance de Robinson-Foulds :**

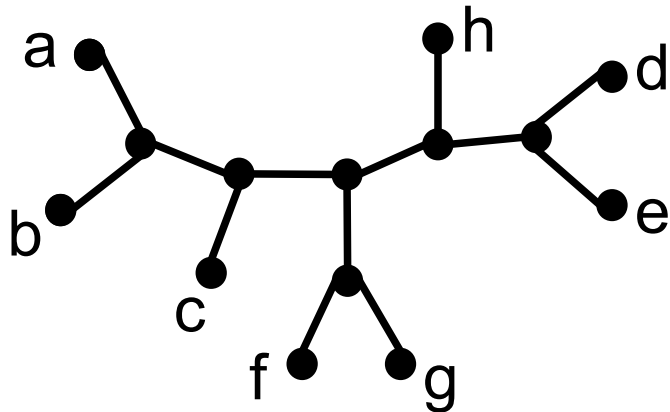
Nombre de **splits** différents

**Similarité :**

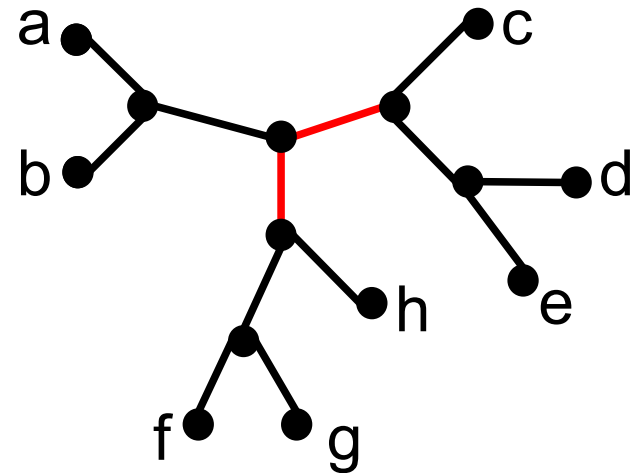
Pourcentage de splits non triviaux identiques

# Contrôle qualité – similarité d'arbres

arbre obtenu sur le texte



arbre obtenu sur le texte altéré



**Distance de Robinson-Foulds :**

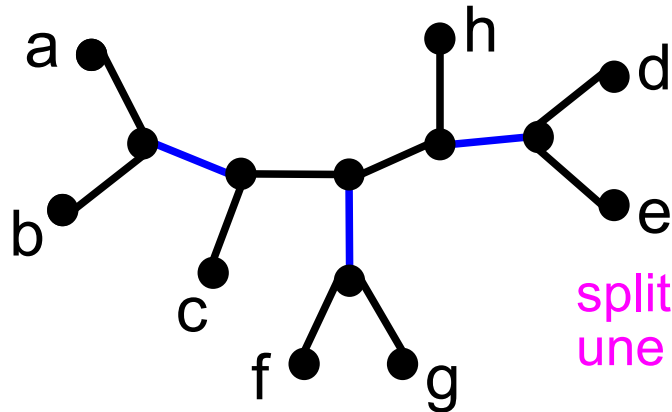
Nombre de **splits différents** : 2

**Similarité :**

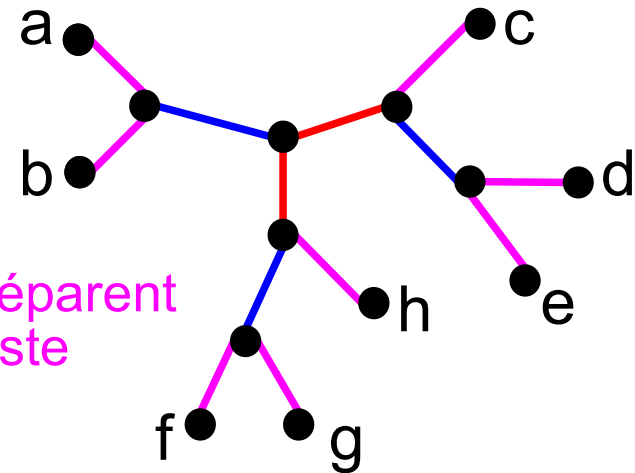
Pourcentage de splits non triviaux identiques

# Contrôle qualité – similarité d'arbres

arbre obtenu sur le texte



arbre obtenu sur le texte altéré



splits *triviaux* : séparent  
une feuille du reste

**Distance de Robinson-Foulds :**

Nombre de **splits différents** : 2

**Similarité :**

Pourcentage de **splits non triviaux identiques** : 60%

# Contrôle qualité – arboricité

**Arboricité “discrète”** d'une matrice symétrique  $M$  :

$$\text{Arb}_d(M) = \frac{1}{C_n^4} |\{\{i,j,k,l\} \text{ tels que } S_{\max} - S_{\text{med}} < S_{\text{med}} - S_{\min}\}|$$

où  $S_{\min}$ ,  $S_{\text{med}}$ ,  $S_{\max}$  sont les trois sommes  $M_{i,j} + M_{k,l}$ ,  $M_{i,k} + M_{j,l}$  et  $M_{i,l} + M_{j,k}$  rangées dans l'ordre croissant.

Guénoche & Garreta, Jobim 2001

**Arboricité “continue”** :

$$\text{Arb}_d(M) = \frac{1}{C_n^4} \sum_{i < j < k < l} \frac{S_{\text{med}} - S_{\min}}{S_{\max} - S_{\text{med}}}$$

Guénoche & Darlu, Alphy 2009

# Contrôle qualité – arboricité

Arboricité “discrète” d'une matrice symétrique  $M$  :

$$\text{Arb}_d(M) = \frac{1}{C_n^4} \left| \left\{ \{i,j,k,l\} \text{ tels que } S_{\max} - S_{\text{med}} < S_{\text{med}} - S_{\min} \right\} \right|$$

*condition liée à la condition des quatre points*

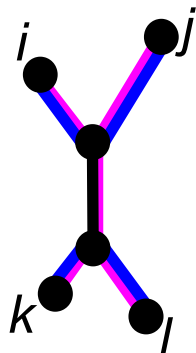
où  $S_{\min}$ ,  $S_{\text{med}}$ ,  $S_{\max}$  sont les trois sommes  $M_{i,j} + M_{k,l}$ ,  $M_{i,k} + M_{j,l}$  et  $M_{i,l} + M_{j,k}$  rangées dans l'ordre croissant.

Guénoche & Garreta, Jobim 2001

Arboricité “continue” :

$$\text{Arb}_d(M) = \frac{1}{C_n^4} \sum_{i < j < k < l} \frac{S_{\text{med}} - S_{\min}}{S_{\max} - S_{\text{med}}}$$

Guénoche & Darlu, Alphy 2009



**Condition des quatre points :**

$S_{\min} \leq \min(S_{\text{med}}, S_{\max})$  pour tout quadruplet ssi  $M$  correspond à une distance d'arbre



# Contrôle qualité – arboricité

Arboricité “discrète” d'une matrice symétrique  $M$  :

$$\text{Arb}_d(M) = \frac{1}{C_n^4} |\{\{i,j,k,l\} \text{ tels que } S_{\max} - S_{\text{med}} < S_{\text{med}} - S_{\min}\}|$$

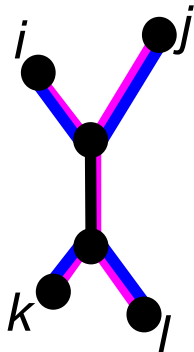
où  $S_{\min}$ ,  $S_{\text{med}}$ ,  $S_{\max}$  sont les trois sommes  $M_{i,j} + M_{k,l}$ ,  $M_{i,k} + M_{j,l}$  et  $M_{i,l} + M_{j,k}$  rangées dans l'ordre croissant.

Guénoche & Garreta, Jobim 2001

Arboricité “continue” :

$$\text{Arb}_d(M) = \frac{1}{C_n^4} \sum_{i < j < k < l} \frac{S_{\text{med}} - S_{\min}}{S_{\max} - S_{\text{med}}}$$

Guénoche & Darlu, Alphy 2009

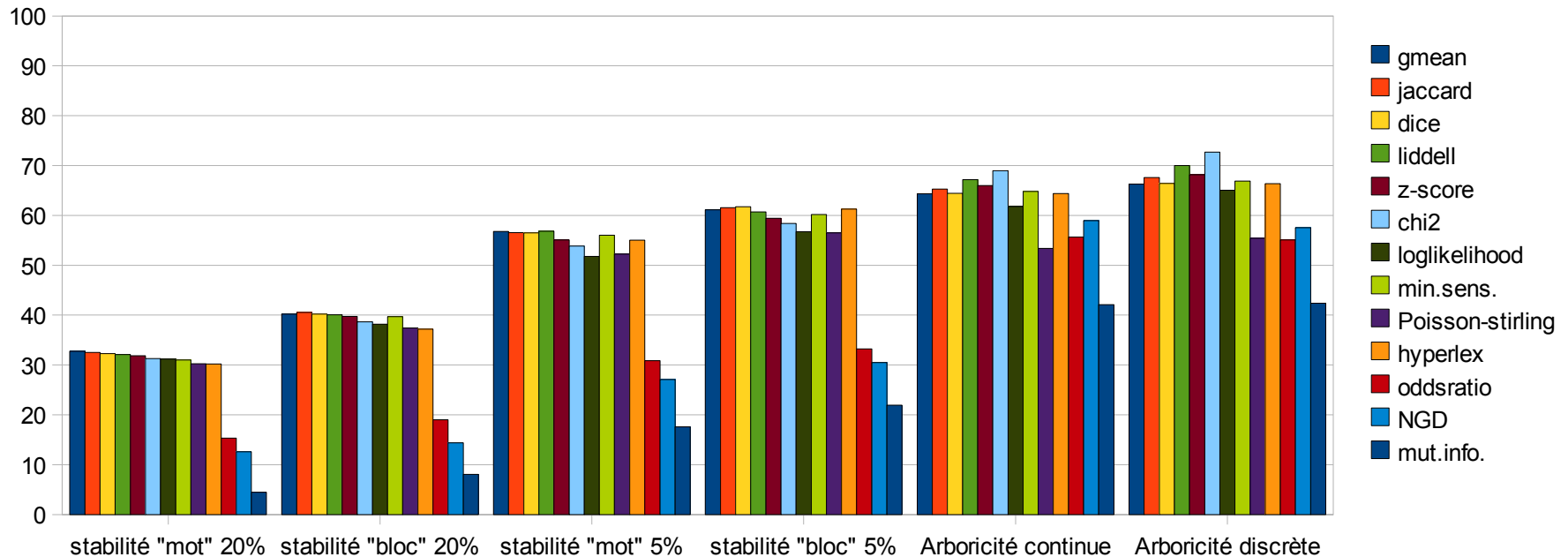


$$S_{\min} \leq S_{\text{med}} = S_{\max}$$

Une distance est **proche d'une distance d'arbre** si  $S_{\text{med}}$  est plus proche de  $S_{\max}$  que de  $S_{\min}$ .

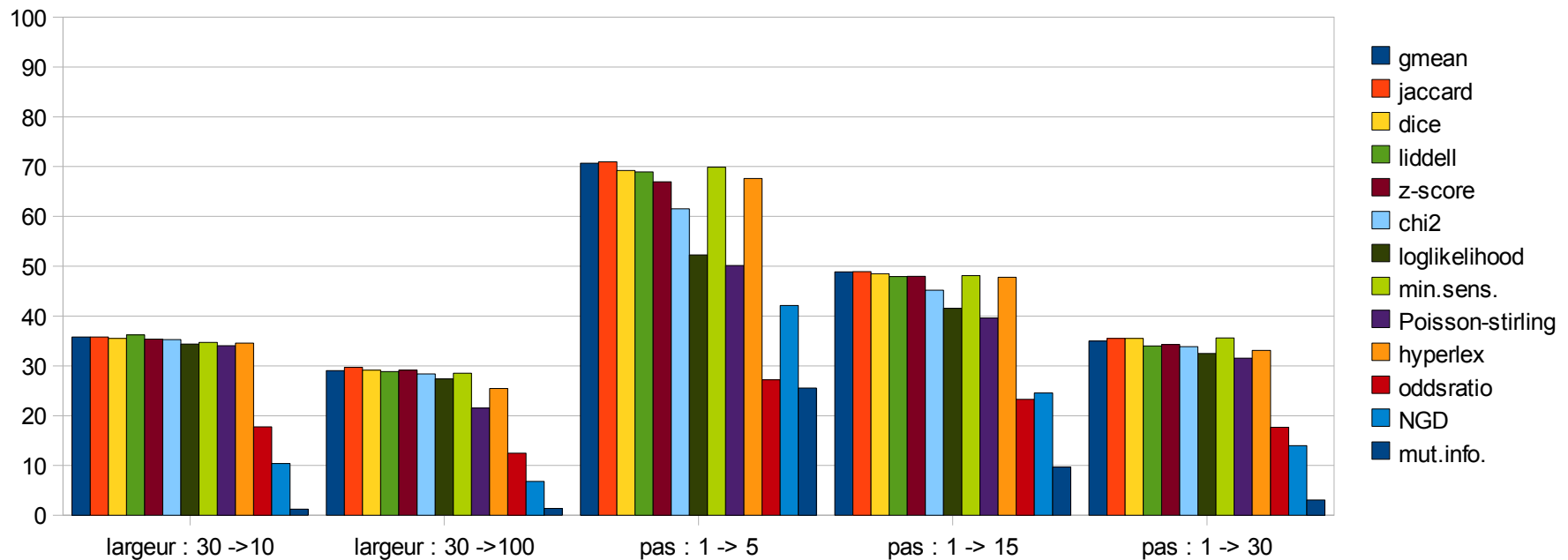
# Contrôle qualité – résultats

Qualité de bootstrap et arboricité :



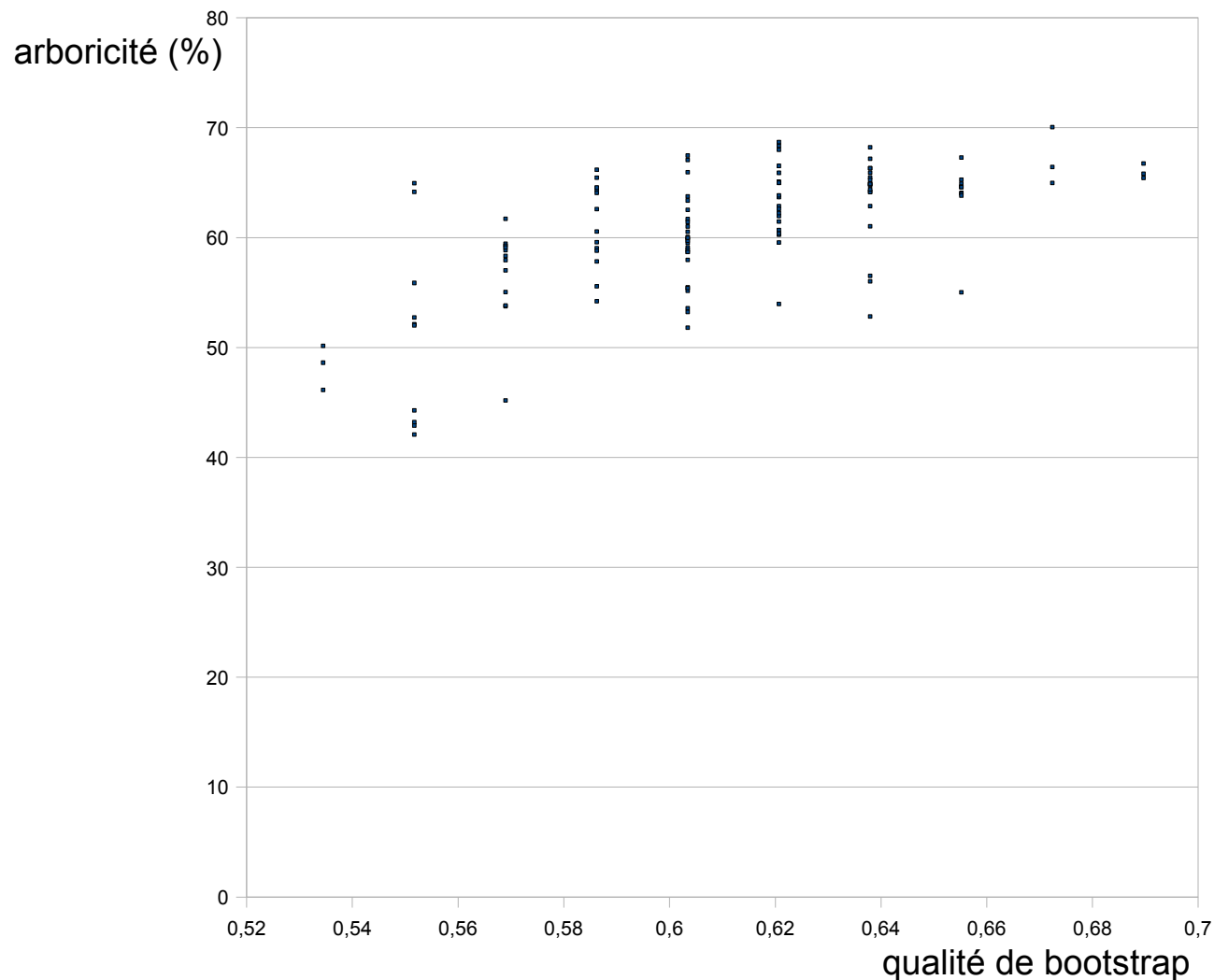
# Contrôle qualité – résultats

Robustesse au changement de paramètres de la fenêtre glissante :



# Contrôle qualité – arboricité

Relations entre “qualité de bootstrap” et arboricité :

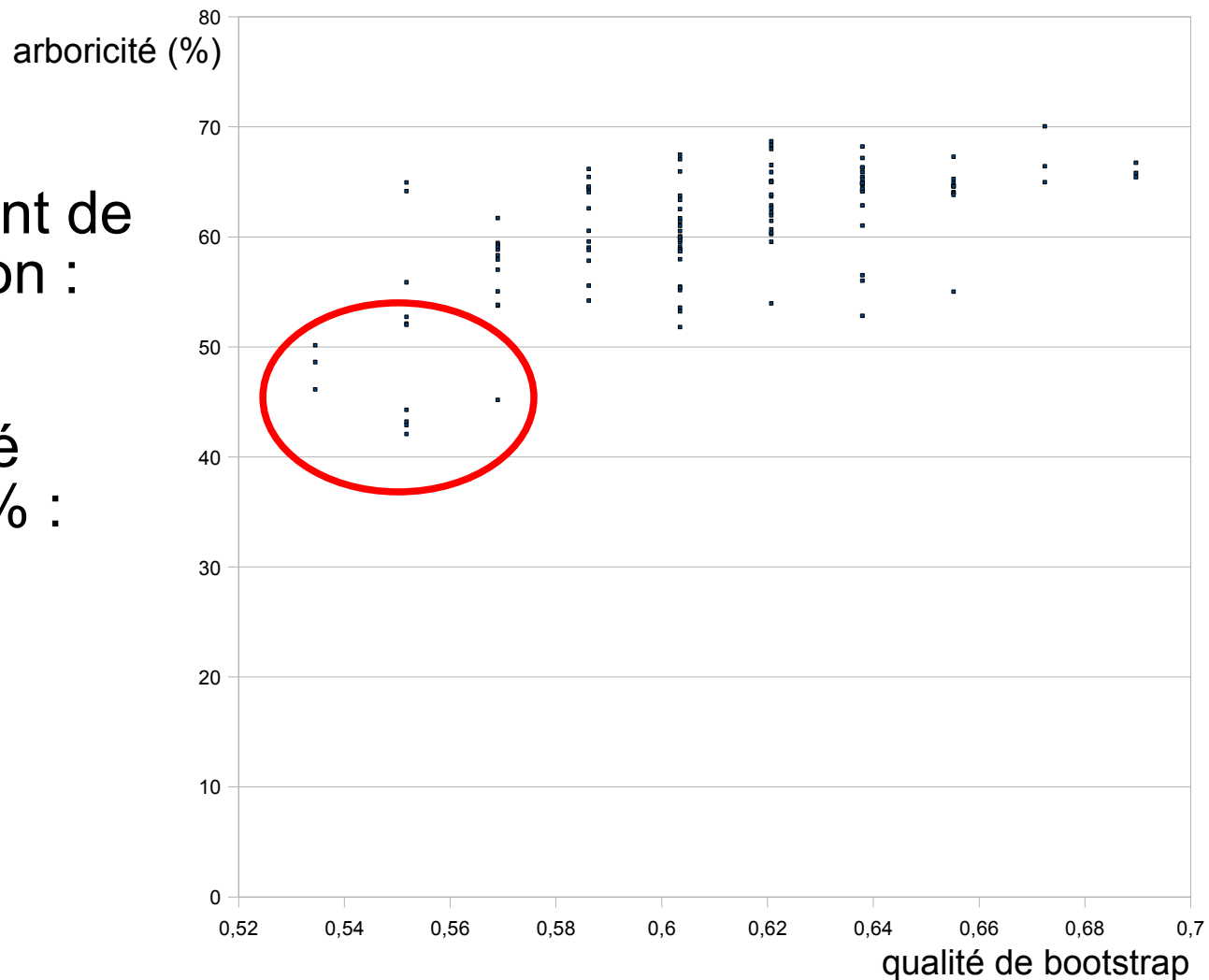


# Contrôle qualité – arboricité

Relations entre “qualité de bootstrap” et arboricité :

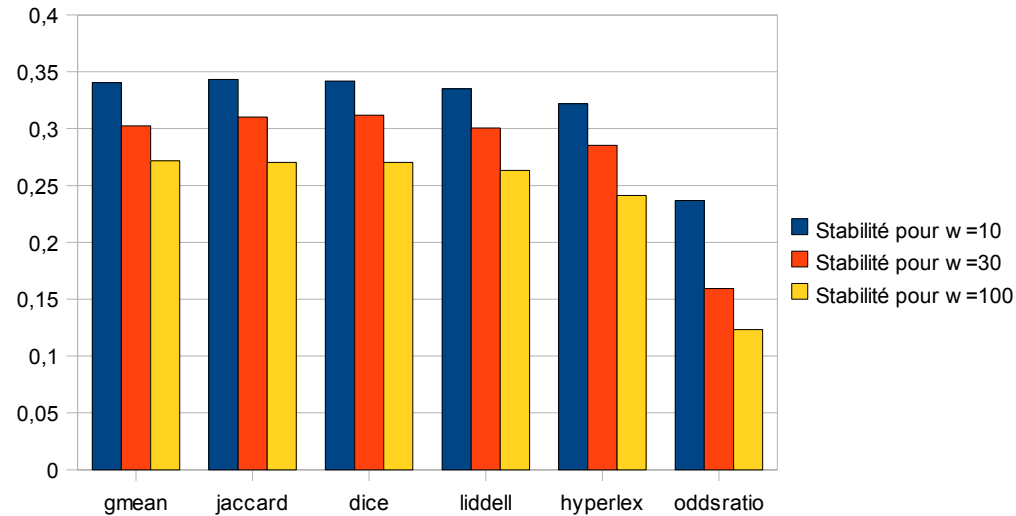
Coefficient de  
corrélacion :  
0.64

Arboricité  
sous 50% :  
**danger !**



# Choix des paramètres – fenêtre glissante

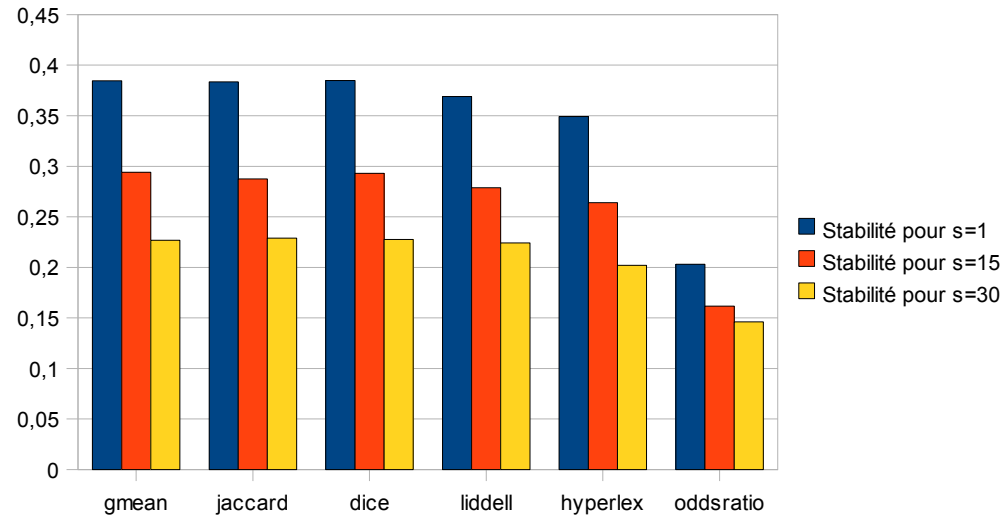
Choix de la largeur de la fenêtre glissante :





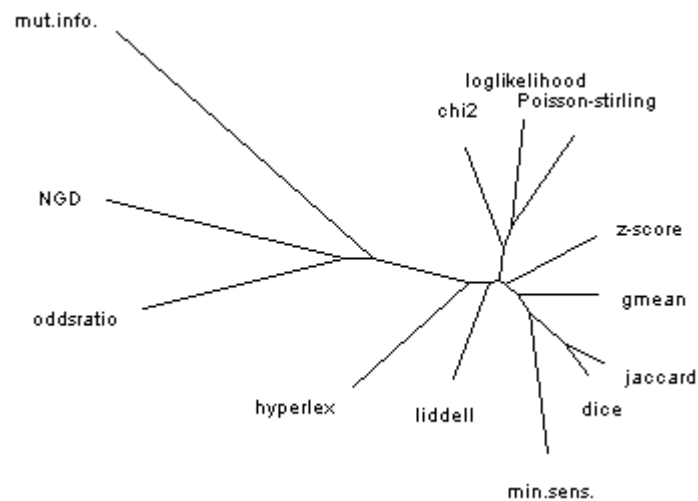
# Choix des paramètres – fenêtre glissante

Choix du pas de glissement de la fenêtre glissante :



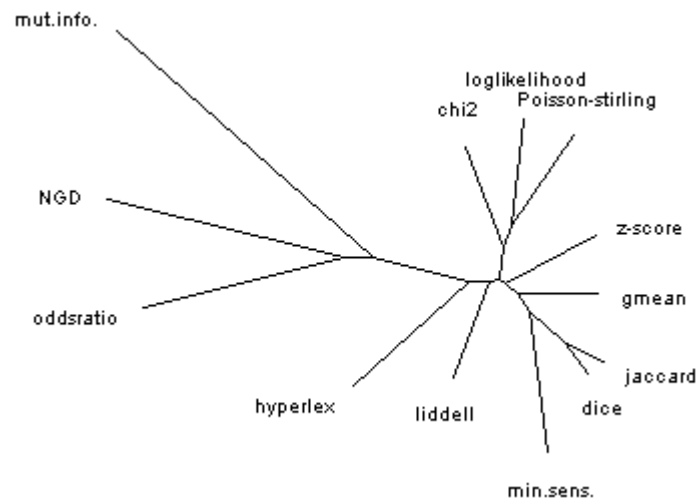
# Choix des paramètres – distance

L'arbre des distances entre arbres de distances :



# Choix des paramètres – distance

L'arbre des distances moyennes, sur les textes du corpus Obama, entre les arbres de mots, en fonction du choix de la formule de distances entre mots :



# Limites de la visualisation

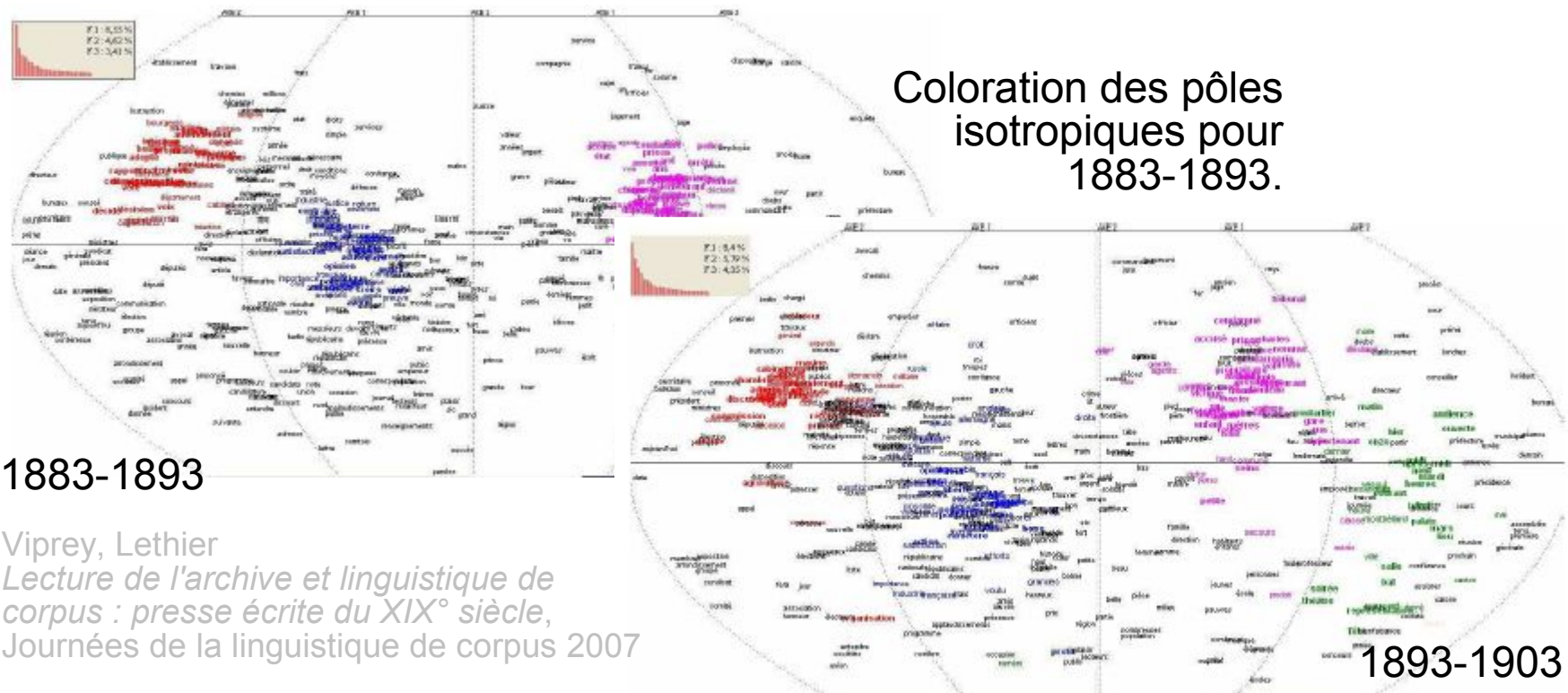
---

- problème de **place** :  
espace en  $O(n^2)$  pour représenter  $O(n)$  mots
- problème de **robustesse** :  
5% de modifications sur le texte induit 40% de modifications sur l'arbre

# Limites de la visualisation

- problème de **robustesse** :  
5% de modifications sur le texte induit 40% de modifications sur l'arbre

Robustesse de l'AFC sur gros corpus :



Viprey, Lethier  
*Lecture de l'archive et linguistique de corpus : presse écrite du XIX<sup>e</sup> siècle,*  
Journées de la linguistique de corpus 2007

# Perspectives

- Créer une interface web  
*SplitsTree : problème de droits. Scriptree ?*
- Tester d'autres méthodes de construction d'arbre
- Evaluer l'utilité des nuages arborés en analyse de textes
- Construire le nuage arboré quotidien des personnalités citées sur les blogs, avec

<http://www.treecloud.fr>

<http://www.scriptree.org>

**WIKIO** Labs

Alex Ferguson Antoine Kombouaré Arthur Levinson Audrey Tautou  
Ben Bermanke Bernard Kouchner Bixente Lizarazu Camélia Jordana  
Christina Aguilera Claude Evin Claude Puel Cristiano Ronaldo David  
Sikharoulidzé Denis Sassou Nguesso **Didier**  
**Deschamps** Didier Drogba **Eric Gerets** Eric Schmidt  
**Eric Woerth** Fabio Santos Franck Ribéry Guy Bedos Guy Stephan  
Jack Lang Jessica Biel Jiri Dlabaja Joe Biden Justin Timberlake Kate Moss Kenny  
Van Hummel Lakhdar Boumediene Lamine Ouahab Landry Chauvin Mikheil  
Saakachvili Nicolas Cage Nicolas de Tavernost Olivier Blanc Patrick Kron  
Patxi Lopez **Paul Le Guen** Robert Louis-Dreyfus Ron Artest Samir  
Nasri **Sébastien Bazin** Sébastien Chavanel Steven De Jongh  
Teodoro Obiang Virginie Guillaume Yves Colleu Yves Jégo

<http://labs.wikio.net>



**Merci pour votre attention !**

---

