

Le(s) chinois du Shun-Pao

Pierre Magistry
pierre.magistry@inalco.fr

7 Décembre 2021
Grenoble



Le Corpus

申報

Le Shun-pao, mon «terrain» à Shanghai il y a 150 ans

- Le premier quotidien imprimé en sinogrammes
- 750M de sinogrammes
- 1872 – 1949, un période de grandes mutations
 - fin de l'Empire, début de la République de Chine
 - abandon du chinois classique au profit des langues régionales
 - définition d'une langue nationale unique
- une première une qui ne ressemble pas vraiment à la dernière

申報

號九十七百二千五萬二第

報景陳人行發

張大貳報本日今

元萬八幣國售份每
(費郵加另埠外)

號八三電 七號八電 九路郵地
號四四 號六號 號三〇口上

內政 財政 軍事 外交 教育 體育 衛生 勞務 社會 新聞

立法院今開秘密會議

討論國家總預算案

物價問題亦將提出檢討

【本報南京一日電】立法院長孫科，一日下午五時招待記者，編請該院會議機密之保持，以協助。孫氏曾聲明，招待記者在根據院會決議，代表立委同人與記者交換該院開會之意見，決無其他用意，請勿誤會。旋孫氏該院依法舉行秘密會議，諸經旁聽，乃為避免將不應公開之事項，而致發生不良影響。且消息之片斷散布，易滋誤解，為國家安全，社會安定，故不得不仿照一般行政性會議之秘密舉行。六月廿四日臨時會議，涉及軍事機密，應防止洩露予敵人得悉，而所有足以刺殺物價，危害一般人民生活及國家財政之消息，自當儘量避免登載。舉以個人之立場認為秘密會議容易引起猜測，使人心理不安，應公開為宜，而一國議會在不時舉行秘密會議實屬例外，登公開應政，取信人民，乃民主國家根本原則，但如國家運籌謀略時，為國家安全，自不得不舉行秘密會議。孫氏復謂明日之秘密會議，乃為聽取政府關於下半年預算編製之報告，當經政府報告所欲，俾立法院可得充分瞭解，故有此項措施。據謂：會議雖極秘密，但會議情形若當然可知悉，但請放心中心，俟有決定後，再行報道。至此，據以英國前財政總領事之辭職為例，指出民主國家議會討論預算案之決議。英例預算案在下院提出後即可公諸全國，然我國行憲伊始，立法院籌備慎重，乃決定俟有決議時，再對社會公開，此為避免刺殺人心，影響物價而採取之非常時期措施。同時，會中對物價問題，亦將予以檢討，亦盼記者幫忙，協同保守秘密。孫氏語末，強調民主國家中新聞界與議會關係密切，我國民主政治之發展，尤賴記者等與立法同人加強聯繫，共同努力。最後並歡迎記者多人建議，對以後採訪該院新聞，儘量予以可能方便。

【本報南京一日電】立法院第一會期第二次秘密會議，定二日晨九時舉行，討論卅七年下半年度中央政府收入支出總預算案。普通部門總預算草案及施行條例草案，業經行政院附送財政方針，施政計劃與預算案送立法院，特別收入支出預算案，亦已另案提出，政院為院長將備會計長徐鈞，財長王雲五列席，說明有關總預算案編製之各事項。同時，並報告物價發展問題，央行總辦徐鈞亦被邀請列席，以備質詢。

【本報南京一日電】國民黨籍立委，一日下午三時在中央黨部舉行談話會，就政院行將提出之總預算案交換意見，負責主言之翁文瀾、張厲生、王振五、徐鈞等均列席，將總預算案提出作詳盡說明，各立委亦提供若干意見，至下午五時許始散。

協定中雖未明文規定，但經濟作尚有根據參加馬歇爾計劃之美國，作關於改革幣制之商談，美國無意干涉我國之內政，故由各該國自行謀求其所應之幣制改革。雲氏稱：經濟合作局之成立，與投資銀行等，故該局投資，而受我國復原工作不能滿意，則合作局可停止其援助，或請國匡正錯誤幣制之局勢。雲氏就實行變遷協定後表示正式大洋洋彼岸同志之努力如何而定。

【本報南京一日電】立法院長王寵惠，一日下午四時召集司法正式就職。前院長居正於十前到京，等候王院長蒞臨後，在院長室內交談約記，旋在大家堂內，舉行簡單之交換儀式，該職員均參加。行禮如儀後，

【本報南京一日電】立法院長王寵惠，一日下午四時召集司法正式就職。前院長居正於十前到京，等候王院長蒞臨後，在院長室內交談約記，旋在大家堂內，舉行簡單之交換儀式，該職員均參加。行禮如儀後，

王寵惠就職

勉僚屬共圖

申報

擬製造局新刻西學書十三種

今夫人欲習新奇之技藝必當有秘密之傳
 通迎機啟悟方能力窮其奧窔精析夫毫其
 圖以為測算之學而測算無不精微也按甘
 之人指點之益則雖有規矩迹象焉耳矣而
 焉故有知其所當然而不能知其所以然者
 則何以為考究之資孟子有言曰羿之教射
 矩學者亦必以規矩此物此志也滬上自
 行之然執事於其間者大抵皆藉西師為之
 知其所以然也間有受其書而伏誦之者則
 不可不繙譯出之以為工師之規矩儒生之
 數年之功費數萬之資現已印行之書則有
 夫西人技藝之精甲於天下洩造化而通神
 購之歐洲今既設局製造而又有繙出各書
 法之原起於積算勾絃圓徑製器所宗推其
 規約指三書萬物之質有疑有化凝全乎于
 原化學鑑原二書氣分輕養遇熱而融蒸釜
 機發軔汽機必以二書錐鑿雜施詎傷地脈
 若金石識別開煤要法二書海中船越帆蓬
 則有若御風要術航海簡法二書地球之下
 方時則有若地學淺釋一書火攻之毒至今
 則有若製火藥法一書以上十三種精深微
 皆戢上稱過者以之作為工師之規矩儒生
 有為加意於富國強兵之道取其所長集其
 將可以復絕古今矣豈不美哉豈不美哉

吳門張少卿校書花燭詞 并序

吳門名校書張少卿者系出毘陵問小家之

Objectifs

Histoire assistée par ordinateur

- Améliorer la structuration du corpus
- Faciliter la recherche d'information

Linguistique de corpus

- observer la naissance du chinois moderne dans les données de la presse quotidienne
- caractériser les évolutions (ou ruptures)

Cadre et circonstances

2018-2019, postdoc dans le cadre de l'ERC

- difficultés pour obtenir le corpus
- pas de structure, ni de points de comparaisons
- intérêt à «faire du deep»

2021, reprise du sujet

- comparer à d'autres corpus
- liberté dans les méthodes (avec un retour à Salem, 1991)
- outils et méthodes

Deux moments

2018-2019, postdoc dans le cadre de l'ERC

- Traitement de données massives et brutes
- Modèles de langues et plongements contextualisés

2021, reprise du sujet

- Gros corpus un peu plus structuré
- Ajout de (petits) corpus de référence pour comparaison
 - chinois archaïque / médiéval / pré-moderne
 - corpus équilibré de chinois standard (TW)
- Utilisation de l'Analyse Factorielle des Correspondances (AFC)

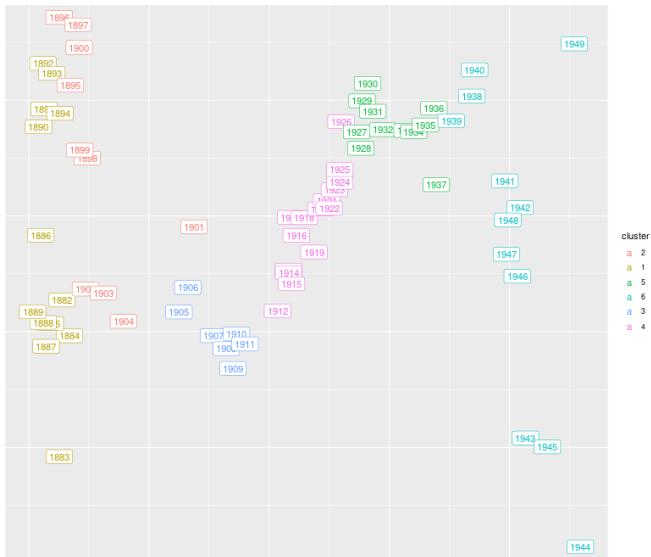
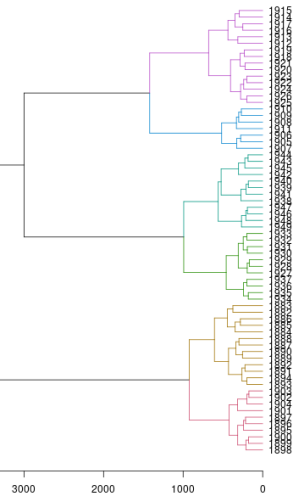
Perplexité de modèles de langue

Objectif

Détecter des périodes, ou des ruptures dans le temps.

Procédure

- diviser le corpus par année
- entraînement de 77 modèles
 - niveau du sinogramme
 - sans ponctuation
- construction d'une matrice de distances entre les années
- MDS et clustering hiérarchique pour visualiser



Résultats

En comparant les sorties de différents algorithmes de clustering.

- 1894 – 1904
- 1905 – 1911
- 1912 – 1922
- 1927 – 1938
- 1938 – 1949

Plongements lexicaux contextualisés

Objectifs

Suivre des compétitions entre formes

Procédure

- sur chacune des 5 périodes précédemment définies
- entraînement d'un Bi-LSTM, à la *flair* (Akbik et al., 2018)
- choix d'une liste d'items en concurrence à échantillonner
- pour chaque occurrence, on obtient un vecteur associé à la forme, dans un contexte et à une date précise.
- projeter et clusteriser ces vecteurs

Character-based representation of words in context

from (Akbik et al., 2018)

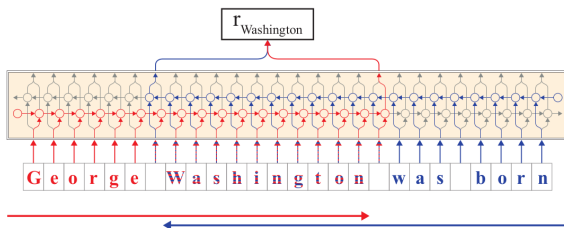


Figure 2: Extraction of a contextual string embedding for a word ("Washington") in a sentential context. From the forward language model (shown in red), we extract the output hidden state after the last character in the word. This hidden state thus contains information propagated from the beginning of the sentence up to this point. From the backward language model (shown in blue), we extract the output hidden state before the first character in the word. It thus contains information propagated from the end of the sentence to this point. Both output hidden states are concatenated to form the final embedding.

list of 字

很 狠 甚 什 極

time period

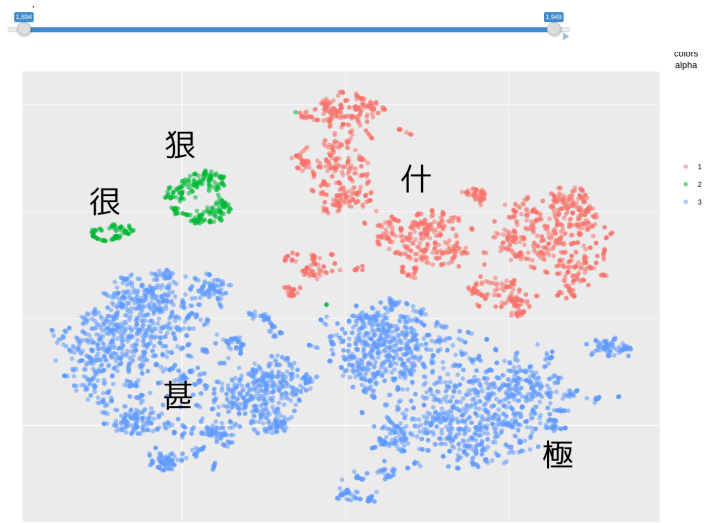
1894-1904

color by

字 cluster

number of clusters

3



Résultats (Magistry, 2019)

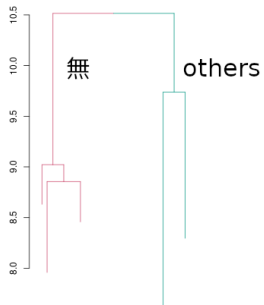
adverbes intensifieurs de verbes d'états

極、很、狠、甚、什

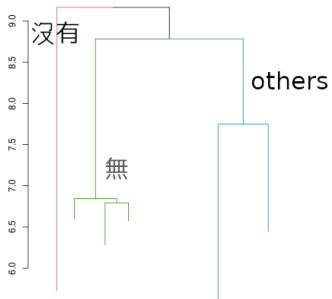
- 甚 et 什 variantes graphiques ? (juste dans 什麼)
- 很 et 狠 variantes graphiques ?
- 什 surtout dans des emprunts au manchu.
- figements de 極 au fil des époques.

négation des existentiels

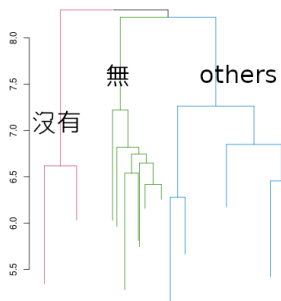
- Coblin (2000) décrit 未曾, 不曾 et 沒有 comme équivalents, avec “an apparent freedom of choice”
- En fait 沒有 n'arrive que plus tard et semble distinct, sauf dans sur la fin où c'est 無 (classique) qui s'isole.



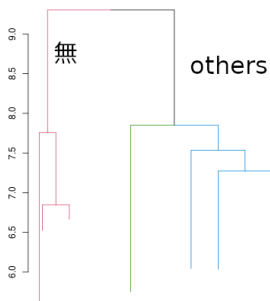
1894 ~ 1904



1905 ~ 1911



1912 ~ 1922



1939 ~ 1949

nouveaux corpus, nouvelle hypothèse

corpus historique Academia Sinica

- 上古 48 livres (-8 – +2)
- 中古 36 livres (2 – 6)
- 近世 13 livres (9 – 16)

Sinica Treebank

350k tokens, années 90

- news
- oral
- manuels scolaires
- magazines

Dans le Shun-pao

- colinguisme du début à la fin (selon les rubriques)
- particularité de la première page ?

Analyse des correspondances

(Salem, 1991) Les séries textuelles chronologiques

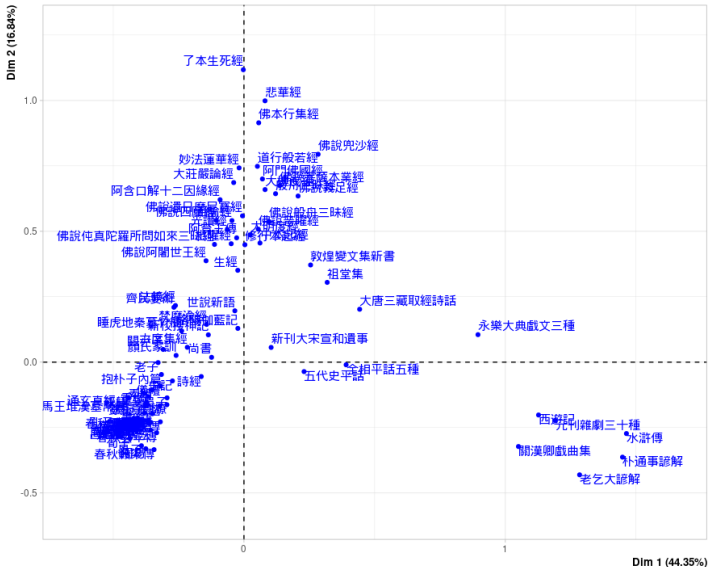
prétraitements

- sélection de mots grammaticaux dans les corpus de références
- compte d'occurrences par document
- du Shun-pao, pour chaque année je sample les premières et les avant-dernières pages (séparément, on obtient donc deux points par an).

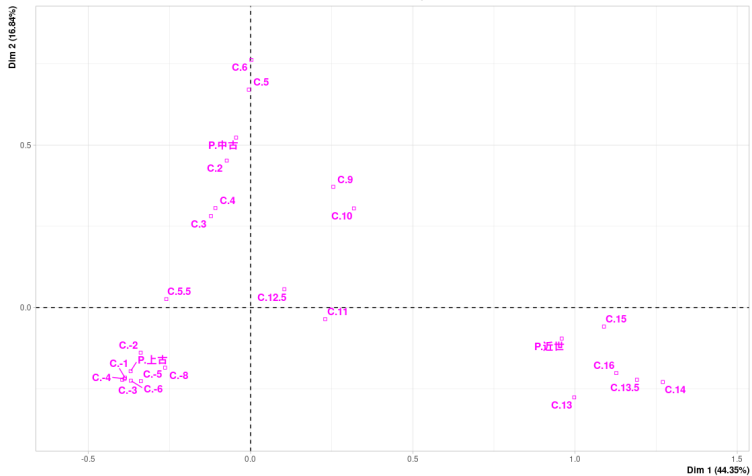
Calculs, avec FactoMineR (Lê, Josse & Husson, 2008)

- AFC
- clustering hiérarchique

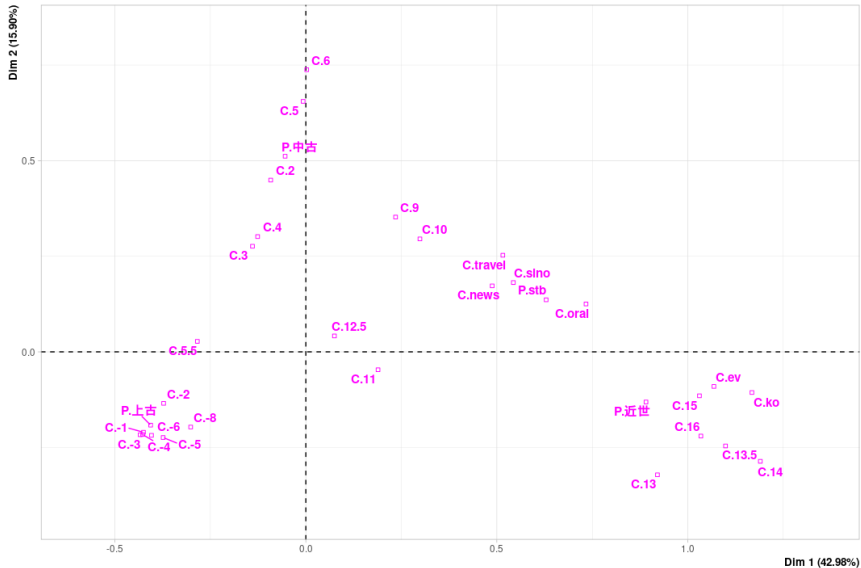
CA factor map



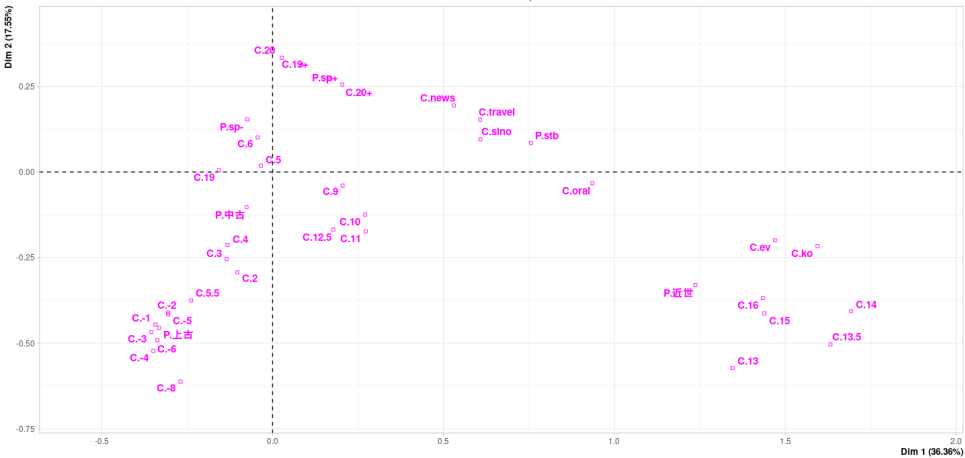
CA factor map

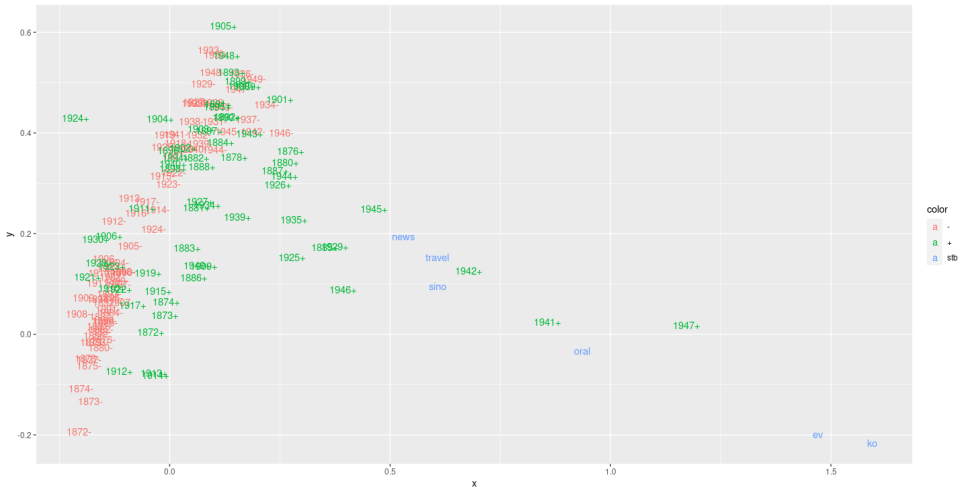


CA factor map

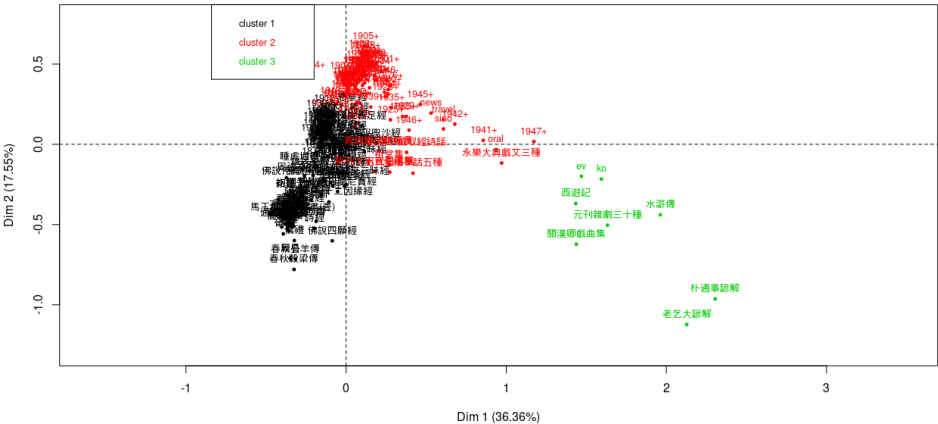


CA factor map

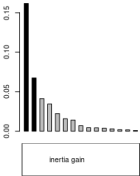
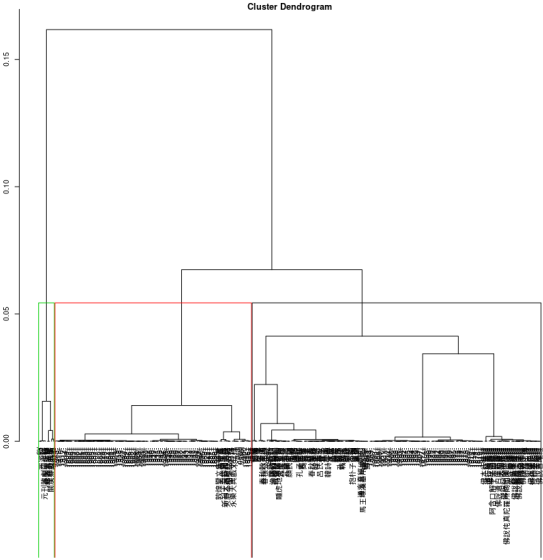




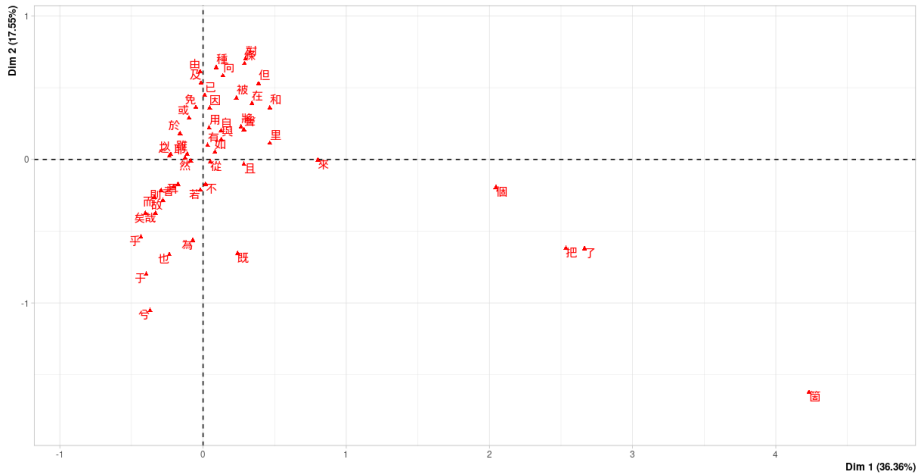
Plan factoriel



Arbre hiérarchique



CA factor map



	Classe 1 ▾	Classe 2 ▾	Classe 3 ▾
P=上古	7.54	-7.11	-1.37
哉	6.8	-6.55	-0.912
乎	6.71	-6.2	-1.63
矣	6.54	-6.09	-1.45
P=中古	6.32	-5.96	-1.09
而	5.87	-5.33	-1.68
也	5.82	-6.23	0.926
C=19	5.44	-5.13	-0.891
者	5.41	-5.34	-0.376
則	5.35	-4.95	-1.3
故	5.32	-4.91	-1.28

	Classe 1	Classe 2	Classe 3
了	-2.86	-0.373	8.97
個	-3.7	0.67	8.46
把	-2.28	-0.743	8.38
來	-1.65	-0.804	6.79
箇	-1.43	-0.946	6.59
P=近世	-4.44	1.07	5.24
里	-1.69	0.196	4.17
且	0.513	-1.98	4.02
聲	-1.03	-0.249	3.55
既	1.78	-2.97	3.21
將	-1.15	0.0549	3.05

Synthèse

Comparaison impossible

Les trois stratégies

- répondent à des questions différentes
- traitent des corpus différents
- ont des définition d'« observable » différentes

Perplexité et clustering

Objectif Détecter des sous-périodes (des ruptures)

Corpus Un sous-corpus par année

Observable chaque modèle de langue

Plongements contextualisés

Objectif Suivre une compétition entre items dans le temps

Corpus Un sous-corpus d'une période

Observable chaque **occurrence** d'un terme ciblé

Analyse factorielle des correspondances

Objectif Situer des documents les uns par rapport aux autres

Corpus ensemble de documents, enrichi de points de référence

Observable chaque document

avantages/inconvénients

Plongements contextualisés

- + niveau occurrence très pratique pour suivre une évolution
- + robuste face aux variantes
- + découverte de n-grammes figés
- besoin d'un gros corpus
- magie noire

Analyse factorielle des correspondances

- + marche sur des petits documents (avec metadonnées)
- + moins boîte-noire
- + correspondance entre documents et lexèmes
- sensible aux variantes, aveugle aux n-grammes

Merci !

Questions ? Commentaires ? Conseils ?
pierre.magistry@inalco.fr