

Le programme d'analyse par iTrameur

Chargement du corpus

The screenshot shows the iTrameur software interface. At the top, there is a navigation bar with the following items: **Chargement**, Trame, Cadre, SR/Patron, Section, Coocs, Bi-Texte, Dépendance, Sélection, Export, Aide. Below this, a light blue header reads "Création d'une nouvelle base / Importation d'une base".

Inside the dialog box, there is a text area with the following instructions:

Deux possibilités pour charger des données dans iTrameur :

1. Charger un fichier (nouvelle base) au format TXT brut, encodé en UTF-8, en ayant préalablement partitionné son contenu (*cf* onglet Aide).
2. Importer une base annotée déjà constituée (*cf* onglet Aide pour le format de cette base).

Une fois la base chargée, les données textuelles sont représentées sous la forme d'une *Carte des sections* (sections définies via le délimiteur de contexte choisi) qui apparaît au bas de cette page.

Below the text area, there are five numbered options for loading data:

| | |
|----------------------------------|--|
| 1. NOUVELLE BASE | Choisir le fichier GROSFIKIER.TXT |
| 1. IMPORTER UNE BASE | Choisir le fichier AUCUN FIKIER SÉL. |
| 2. DÉLIMITEUR DE CONTEXTE | § (si cette zone est vide, contexte=ligne) |
| 3. DÉLIMITEUR(S) | ;;"~ &#@='-.?!%*\$()[\{\}_.:+<>§\ |
| 4. BI-TEXTE | <input type="checkbox"/> (chargement d'un bitexte aligné <i>cf</i> Aide) |
| 5. DÉPENDANCE | <input type="checkbox"/> (chargement d'une base avec annotations en dépendance <i>cf</i> Aide) |

figure(1)

Pour charger le corpus dans iTrameur, il nous faudrait de choisir un fichier en cliquant le premier bouton "choisir le fichier". Comme on avait préparé ce corpus en amont nommé "Grosfichier" en combinant tous les textes après la segmentation en cas de textes chinois et aussi en ajoutant les balises. Donc, il nous suffit de choisir ce fichier en cliquant dessus.

Cadre

Tout d'abord, nous vérifions notre corpus structuré avec les paires de balises ouvrante et fermantes pour chacune partie. Le bouton "cadre" permet de donner une représentation graphique de l'organisation des parties.

The screenshot shows a software interface with a menu bar at the top containing: Chargement, Trame, **Cadre**, SR/Patron, Section, Coocs, Bi-Texte, Dépendance, Sélection, Export, Aide. Below the menu bar is a section titled "Opérations sur le Cadre" with several buttons: Cadre, Parties, PCLC, Ventilation*, Spécifs-partie*, Spécifs totales*, Mots Spécifs+, and TGF+BT+VN. Below this section is a list of 11 dump files, each enclosed in a red rectangular box:

- dump=1.txt <68965:71743>
- dump=10.txt <71744:73762>
- dump=11.txt <73763:74249>
- dump=12.txt <74250:75562>
- dump=13.txt <75563:76251>
- dump=14.txt <76252:76845>
- dump=15.txt <76846:79066>
- dump=16.txt <79067:81167>
- dump=17.txt <81168:82833>
- dump=18.txt <82834:83223>
- dump=19.txt <83224:83859>

figure(2)

Dictionnaire

Cette fonction nous représente le contenu du dictionnaire du corpus, et aussi les fonctions permettant de montrer les fréquences, la ventilation, la concordance, la carte et la sélection qui sont tous dans la même ligne du tableau.

Dictionnaire

Recherche :

| Item | Fq | Concordance | Ventilation | Carte | Sélection |
|------|------|---|---|--|---|
| , | 4300 |  |  |  |  |
| 的 | 3390 |  |  |  |  |
| 。 | 2439 |  |  |  |  |
| 、 | 1916 |  |  |  |  |
| 环境 | 1389 |  |  |  |  |
| 和 | 1000 |  |  |  |  |
| 污染 | 973 |  |  |  |  |
| 生态 | 743 |  |  |  |  |
| 是 | 577 |  |  |  |  |
| 在 | 543 |  |  |  |  |

Affichage de 1 à 10 des 11,748 items Préc. ... Suiv.

figure(3)

On voit que le mot “的” apparaît dans le corpus 3390 fois qui est le plus nombreux parmi les 11,784 mots. Ensuite, le groupe de mot “环境” qui signifie “l’environnement ” est le seconde le plus fréquent dans ce corpus. Alors, notre thématique est “la protection environnementale”. Donc, en l’occurrence, nous saisissons notre groupe de mots clés en mandarin “环保” qui veut dire “la protection environnementale”, voilà ce que nous obtenons:

Dictionnaire

Copy CSV Excel PDF Print

Recherche : 环保

| Item | Fq | Concordance | Ventilation | Carte | Sélection |
|------|-----|-------------|-------------|-------|-----------|
| 环保 | 207 | | | | |
| 环保部 | 9 | | | | |
| 环保法 | 6 | | | | |
| 环保局 | 3 | | | | |
| 环保督 | 1 | | | | |
| 环保舆 | 1 | | | | |

Affichage de 1 à 6 des 6 items (filtrage à partir des 11,748 items)

Préc. 1 Suiv.

figure(4)

On voit que ce cible “环保” (la protection environnementale) apparaît dans ce corpus 207 fois. Vu qu’il y a deux termes en mandarin qui désignent la même chose “la protection environnementale”, il vaut mieux de chercher et observer avec le deuxième term: “环境保护”. Le premier term “环保” est l’abréviation du deuxième “环境保护”.

Dictionnaire

Copy CSV Excel PDF Print

Recherche : 环境保护

| Item | Fq | Concordance | Ventilation | Carte | Sélection |
|---------------------------|----|-------------|-------------|-------|-----------|
| No matching records found | | | | | |

Affichage de 0 à 0 des 0 items (filtrage à partir des 11,748 items)

Préc. Suiv.

figure(5)

Comme on a dit que l’abréviation “环保” (la protection environnementale) est employée 207 fois dans notre corpus, alors le term en toute lettre qui dit la même chose “环境保护” est apparu aucune fois. En fait, cela est à cause du gros fichier qu’on avait préparé en amont pour le corpus. En raison de la particularité du mandarin, nous somme obligé de faire la segmentation du texte en chinois dans le but de faire l’analyse. Alors, le processus de segmentation a fait diviser le groupe de mots “protection environnementale” en “protection” et “environnemental”, en l’occurrence, “环境保护”

qui est divisé en “环境” et “保护”. C’est la raison que la fréquence de “protection environnementale” (“环境保护”) est zéro. Pour vérifier si ce terme “environnement” est employé ensemble avec “protection”, on peut se servir de la fonction “Coocs”.

Cooccurrences

Cette fonction permet de calculer les cooccurents du pôle source avec les contextes gauches et droites.

PÔLE SOURCE 环境

ement Trame Cadre SR/Patron Section **Coocs** Bi-Texte Dépendance Sélection Export

Calcul de cooccurents

NB TERME GAUCHE 10 **NB TERME DROITE** 10

Cooccurents* Cooccurents* sur partie sélectionnée

Réseau Cooccurents* Réseau Cooccurents* sur partie sélectionnée

STOPLISTE="GESTIONNAIRE DE SÉLECTION"

| Cooc | FqCooc | CoFreq | IndSP |
|------|--------|--------|-------|
| 保护 | 423 | 334 | ** |
| 生态 | 743 | 411 | ** |
| 污染 | 973 | 464 | ** |
| 部 | 65 | 72 | ** |
| 治理 | 322 | 163 | 47 |
| 抄报 | 25 | 35 | 43 |
| 手 | 28 | 35 | 39 |
| 问题 | 320 | 136 | 30 |
| 质量 | 147 | 80 | 27 |
| 保护法 | 15 | 21 | 25 |
| 改善 | 63 | 43 | 20 |
| 自净 | 12 | 15 | 16 |

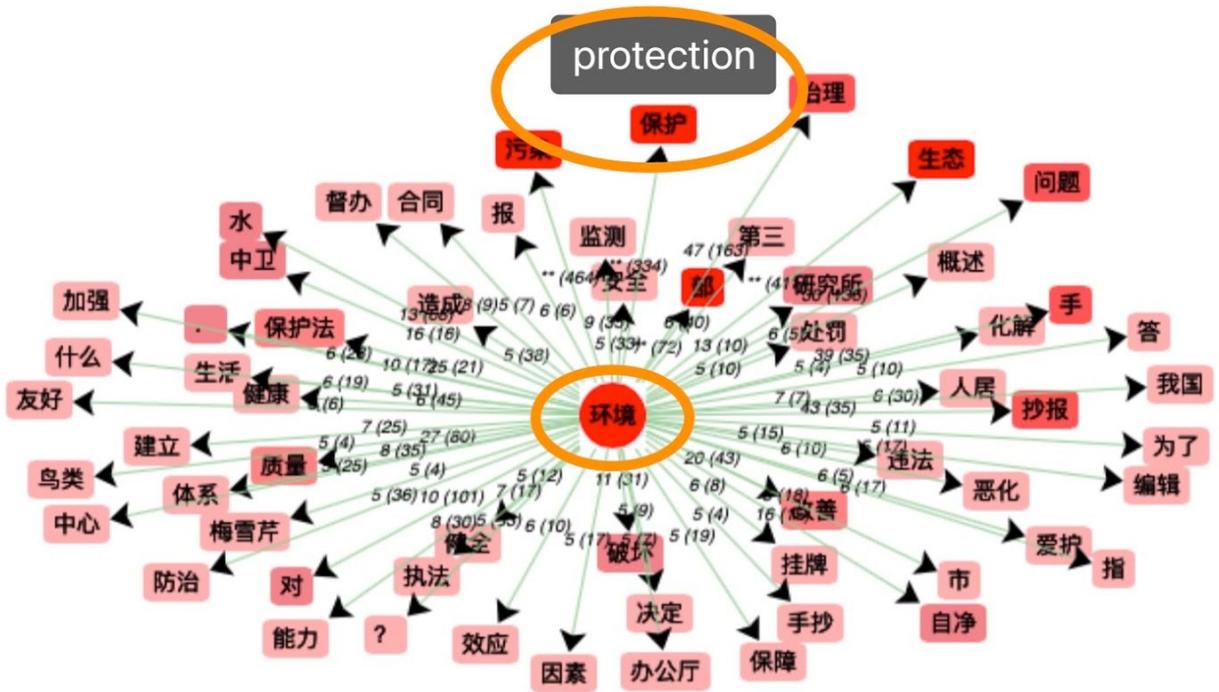
Couleur Noeud-Cooc : IndicSp<5<IndicSp<10<IndicSp<20<IndicSp<30<IndicSp<50<IndicSp

Arc Label : POLE-IndiceSpécif(Co-Freq)→COOC



figure(6)

On voit que quand on saisit “environnement” dans le pôle source, ensuite on tape le bouton “Cooocs”. Le terme le plus fréquent en cooccurrence avec “environnement” est “protection”.



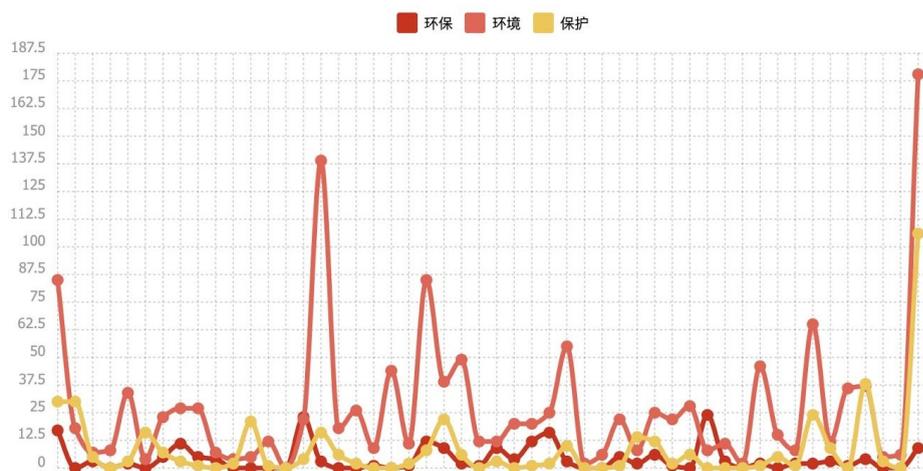
figure(7)

Cadre - Ventilation

1, Fréquence absolue



Graphiques de ventilation (annotation : 1) sur la partition dump : Fréquence absolue (haut) / Fréquence relative (centre) / Spécificité (bas)
Clic sur item dans légende pour masquer/afficher courbe de l'item visé



figure(8)

Le bouton de la ventilation nous permet d'apercevoir la répartition du mot en partie. La ventilation nous représente par les graphiques sur la fréquence absolue qui compte le nombre de fois que le pôle apparaît dans chaque partie, la fréquence relative qui illustre le nombre de fois que le pôle apparaît au total des mots dans la partie. En l'occurrence, le cible clé “环境” (environnement) apparaît dans des parties du corpus de façon instable par l'observation ces graphiques de ventilation. Par exemple, le cible clé “环境” représente 85 fois dans la première partie du corpus correspondant au texte “1-1.txt” par le graphique de fréquence absolue. Egalement, il existe des partie qui montre très peu de fréquence de termes, par exemple, le texte “1-19.txt” dont “环境” (environnement) apparaît juste 4 fois dans la partie. Alors on voit qu'il y a un pic qui représente le plus nombre de fréquence sur le dernier text correspondant au texte numéro neuf “1-9.txt”. Dans

ce texte, le cible clé “环境” représente 179 fois comme la fréquence absolue. Je vérifie les valeurs des fréquences absolues respectivement et je les montre en-dessous(figure(9)):

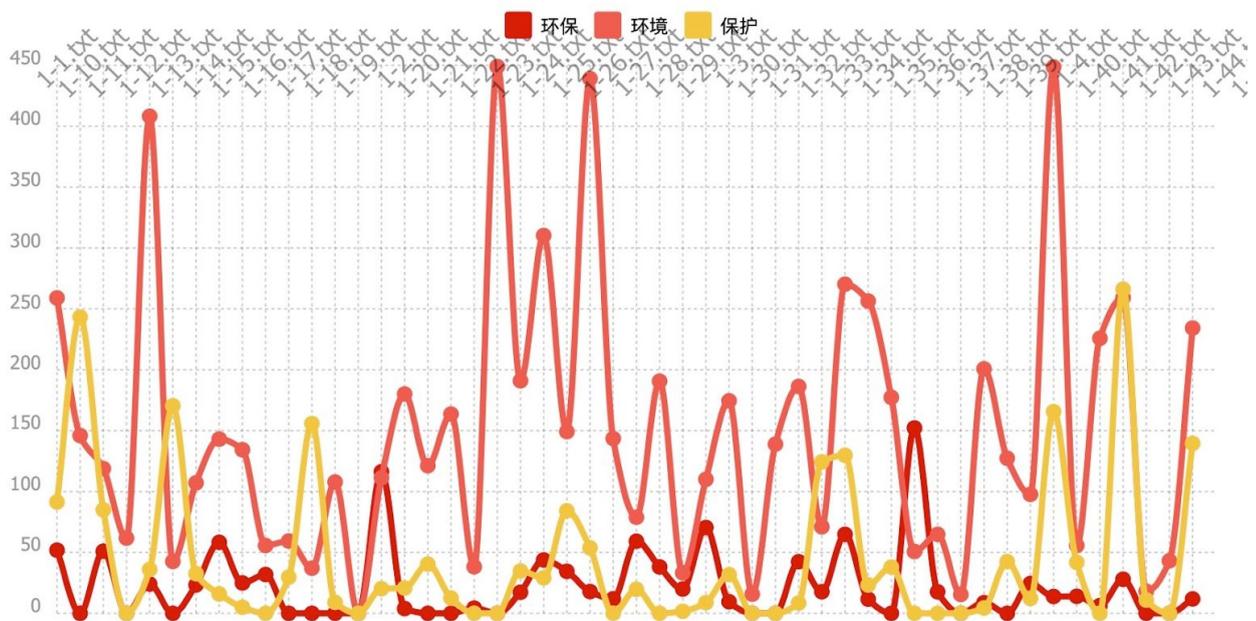




figure(9)

Enfin, même si ce sont tous les textes ayant le même thématique “la protection environnementale”, ils ont différents fréquences absolues du terme clé. On peut déduire que le nombre de fois d’apparition du terme clé n’est pas un facteur très important pour un thématique d’un texte. Ici, la thématique reste la même, alors le nombre de fois d’apparition du clé varie entre 2 fois et 179 fois.

2, Fréquence relative



figure(10)

En ce qui concerne la fréquence relative, on voit que les fréquence du cible “环境” (l’environnement) est relativement beaucoup plus hautes que le cible “环保” (la protection environnementale) et “环境保护” (la protection environnementale). Ayant recours à la fonction de “Coocs”, on peut chercher la raison de cette haute fréquence relative de “环境”.

PÔLE SOURCE 环境

Cadre SR/Patron Section **Coocs** Bi-Texte Dépendance Sélé

Calcul de cooccurrents

NB TERME GAUCHE 10 **NB TERME DROITE** 10

Cooccurrents* Cooccurrents* sur partie sélectionnée

Réseau Cooccurrents* Réseau Cooccurrents* sur partie sélectionnée

STOPLISTE="GESTIONNAIRE DE SÉLECTION"

(clic sur mot : contextes)

| Cooc | FqCooc | CoFreq | IndSP |
|------|--------|--------|-------|
| 保护 | 423 | 334 | ** |
| 生态 | 743 | 11 | ** |
| 污染 | 973 | 464 | ** |
| 部 | 65 | 72 | ** |
| 治理 | 322 | 163 | 47 |
| 抄报 | 25 | 35 | 43 |
| 手 | 28 | 35 | 39 |
| 问题 | 320 | 136 | 30 |
| 质量 | 147 | 80 | 27 |
| 保护法 | 15 | 21 | 25 |
| 改善 | 63 | 43 | 20 |

figure(11)

PÔLE SOURCE 环保

Cadre SR/Patron Section **Coocs** Bi-Texte Dépendance Sélection

Calcul de cooccurrents

NB TERME GAUCHE 10 **NB TERME DROITE** 10

Cooccurrents* Cooccurrents* sur partie sélectionnée

Réseau Cooccurrents* Réseau Cooccurrents* sur partie sélectionnée

STOPLISTE="GESTIONNAIRE DE SÉLECTION"

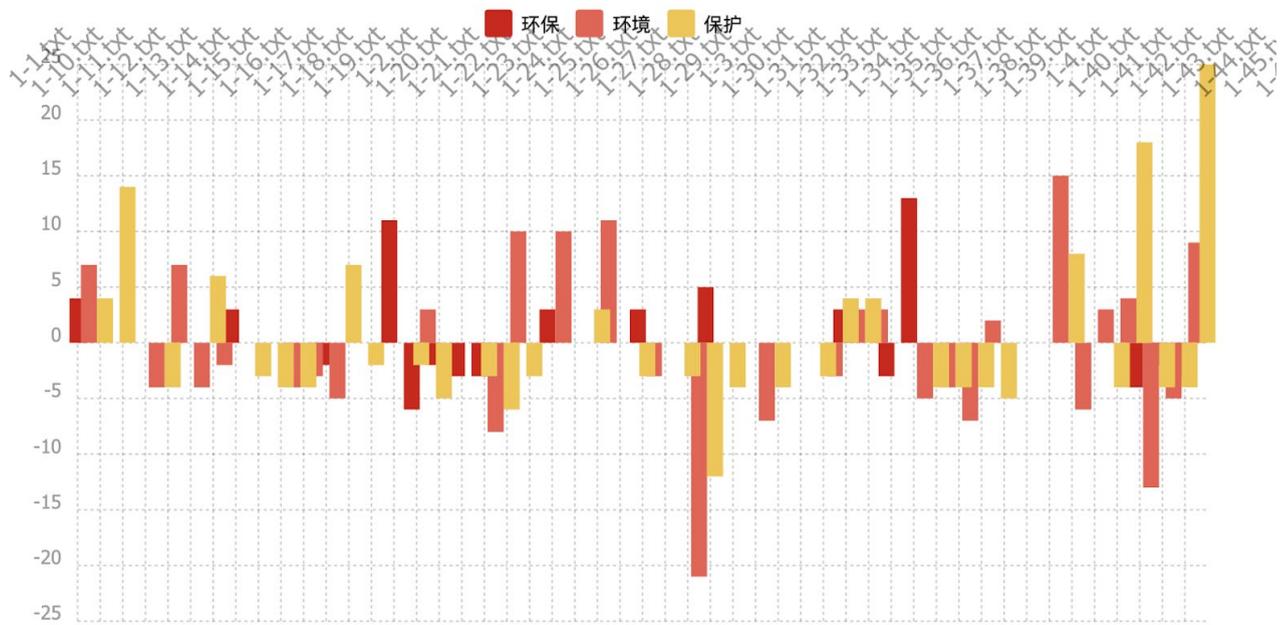
(clic sur mot : contextes)

| Cooc | FqCooc | CoFreq | IndSP |
|------|--------|--------|-------|
| 产业 | 93 | 29 | 25 |
| 舆情 | 23 | 17 | 24 |
| 举报 | 82 | 23 | 19 |
| 督查 | 37 | 11 | 10 |
| 装备 | 30 | 10 | 10 |
| 节能 | 40 | 10 | 9 |
| 轮 | 13 | 7 | 9 |
| 扫描 | 11 | 6 | 8 |
| 中央 | 53 | 11 | 8 |
| 垂改 | 6 | 5 | 8 |
| 100 | 15 | 6 | 7 |

figure(12)

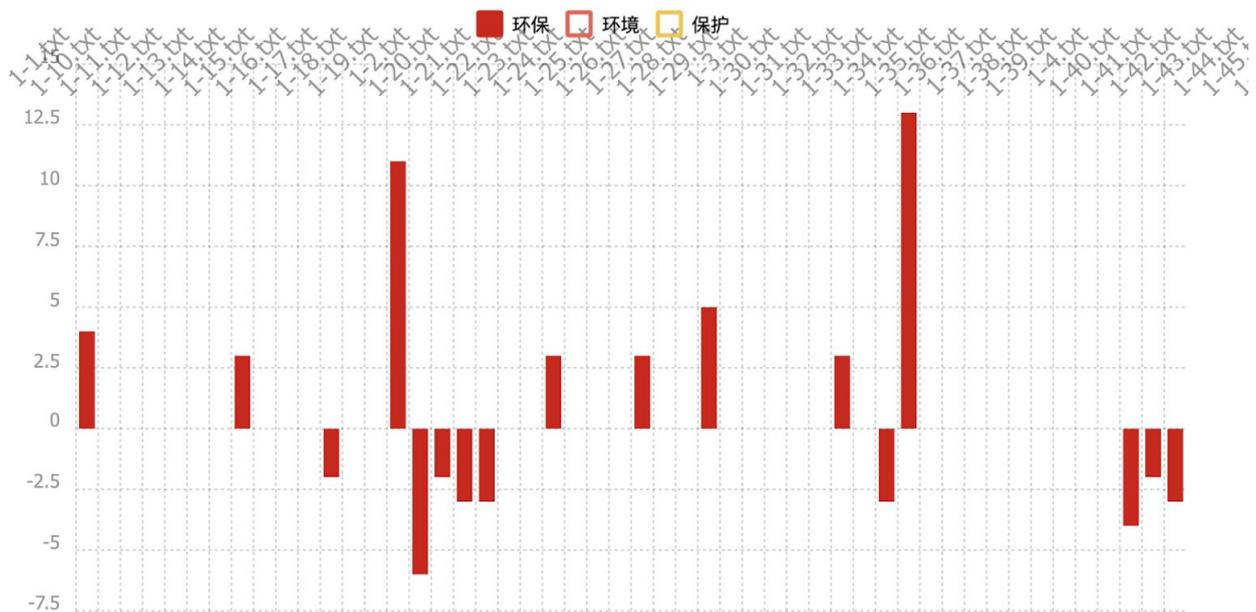
On voit que le cible “环境” est co-présenté très souvent avec d’autres termes. Par exemple, “la pollution de l’environnement”, “l’environnement écologique” etc. La fréquence de cooccurrence du terme “环境”(l’environnement) avec le terme “l’écologie” est 743 et celle avec le terme “la pollution” est 973. Alors on peut aussi observer les cooccurrences du cible “环保”(la protection environnementale), ce cible co-présente avec les termes comme “l’industrie”, “la situation” etc. Mais les fréquence de cooccurrence avec ces termes ne sont pas autant hautes que celles entre le cible “环境” et d’autres termes. C’est pour ça que les fréquences du cible “环境” (l’environnement) est relativement beaucoup plus hautes que les deux autres termes.

3, Spécificité



figure(13)

En ce qui concerne la spécificité, cette fonction sert à illustrer les comparaisons entre les parties. Avec la spécificité, nous nous attendons à avoir le nombre de fois équivalent que le pôle source apparaît pour chaque partie. S'il est sur-représenté, cela voudrait dire que le cible est particulièrement présent dans une partie par rapport aux autres parties; Par contre, s'il est en sous-représentation, celui-ci est particulièrement absent dans une partie comparé aux autres. En l'occurrence, on voit que ces trois termes apparaissent dans les textes de manière irrégulière. Si on veut juste observer les distributions de cible "环保" (la protection environnementale), on décoche les autres deux cases. Voilà ce que nous obtenons:



figure(14)

On constate que dans le texte numéro 15 et le texte numéro 38, le cible “环保” (la protection environnementale) sont particulièrement présent, donc sur-représenté dans ces deux parties par rapport aux autres. Egalement, il existe des parties dont le cible source est sous-représenté.

Barycentre temporel / Coefficient Von Neumann

Vu que ce corpus n’est pas récolté et structuré de façon chronologique, on ne va pas appliquer la fonction de Barycentre temporel et Coefficient Von Neumann. Parce que le bouton BT sert à constater à partir de quel moment ou une zone temporel, le cible est présent de la même quantité à gauche et à droite sur la ligne chronologique.

Sections -- Cooccurrents

On va regarder l’opération “Cooccurrents” sur la Carte des Sections.



(clic sur mot : contextes)

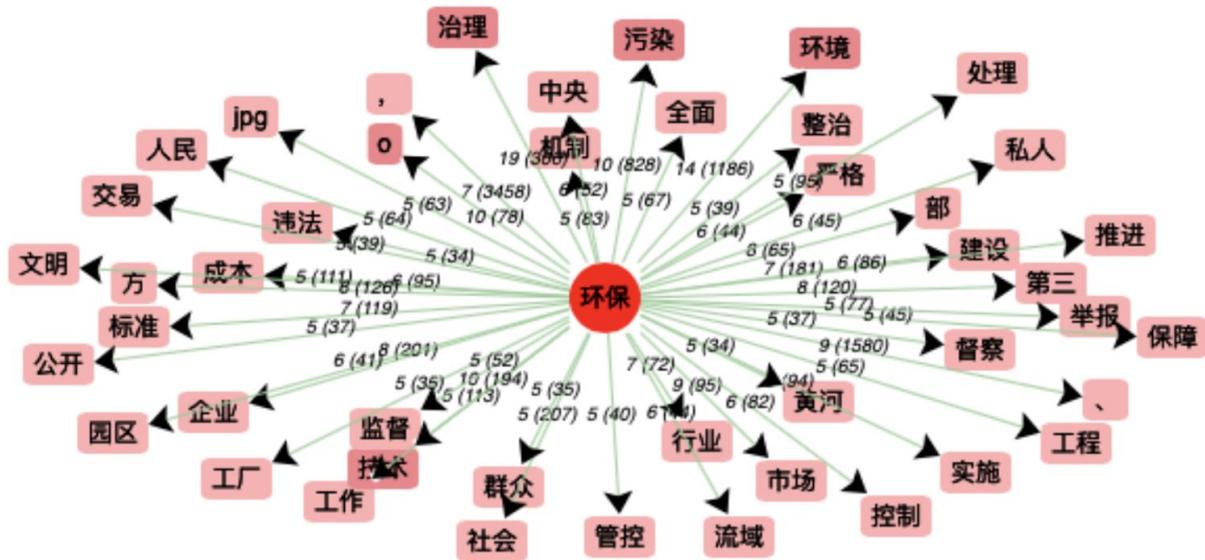
| Cooc | FqCooc | CoFreq | IndSP |
|------|--------|--------|-------|
| 治理 | 322 | 306 | 19 |
| 环境 | 1389 | 1186 | 14 |
| 污染 | 973 | 828 | 10 |
| 技术 | 208 | 194 | 10 |
| o | 78 | 78 | 10 |
| , | 1916 | 1580 | 9 |
| 市场 | 97 | 95 | 9 |
| 第三 | 127 | 120 | 8 |
| 企业 | 220 | 201 | 8 |
| 部 | 65 | 65 | 8 |
| , | 4300 | 3458 | 7 |

Couleur Noeud-Cooc : IndicSp<5<IndicSp<10<IndicSp<20<IndicSp<30<IndicSp<50<IndicSp
Arc Label : POLE-IndiceSpécif(Co-Freq)→COOC



figure(15)

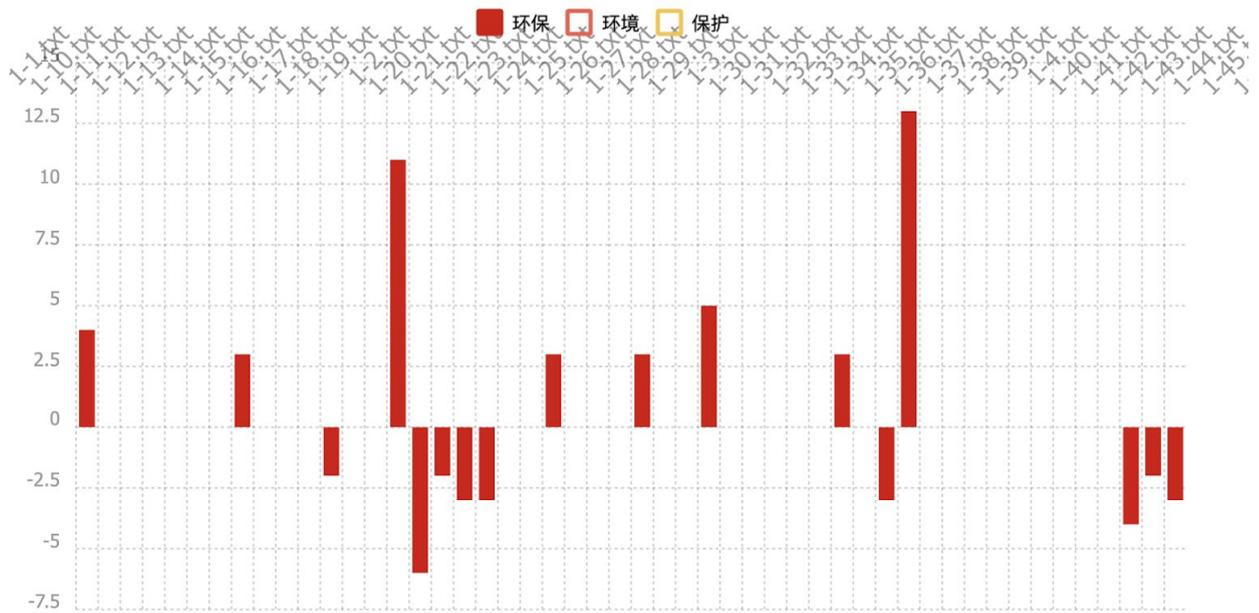
En saisissant le cible “环保” (la version raccourcie de “la protection environnementale”), on observe qu’il est co-présenté avec divers groupes de mots. Par exemple, “la technique”, “l’aménagement”, “la pollution”, “le marché”, “l’entreprise” etc. Grâce à cette fonction, nous obtenons deux représentations de cooccurrences du cible source. Premièrement, c’est le tableau contenant quatre colonnes qui sont le mot cooccurrent, la fréquence, la co-fréquence de deux mots ou deux groupes de mots et l’indice de spécificité. Deuxièmement, comme l’image montrée en-dessous, il illustre en une seule représentation tous les cooccurrents du cible source avec la co-fréquence et l’indice de spécificité.



figure(16)

En conclusion, cet outil d'analyse textométrique de données iTrameur nous aide à explorer le corpus de façon intuitive, générale et précise. Grâce au outil iTrameur dont l'objectif est définir et compter des unités dans les partie de corpus et contraster entre ces parties afin de découvrir éventuellement des phénomènes textuelles. Par exemple, avec la fonction "Dictionnaire", nous pouvons regarder tous les mots avec leurs fréquences et d'autres valeurs textométriques du corpus. En l'occurrence, notre cible clé "环保" (la protection environnementale) apparaît 207 fois. En plus, il est co-présenté avec divers groupes des mots, comme "la technique", "l'aménagement". Par le bouton "Ventilation", on observe toutes les distributions du cible source dans chaque partie du corpus avec la fréquence absolue, la fréquence relative et la spécificité. En général, la répartition du cible "环保" (la protection environnementale) n'est pas équilibré dans toutes les parties du corpus. On peut considérer que ça fait partie de la rédaction textuelle en chinois: les fréquences d'un terme n'est pas forcément liés avec le sujet d'un texte. Par exemple, le sujet "la protection environnementale" peut être sur-représenté ou sous-représenté dans un texte, alors cela n'empêche pas qu'il est toujours le sujet de textes. Dans certains textes dont le cible "环保" (la protection environnementale) est sous-représenté, par exemple, par la figure dessous (figure 17), on constate que notre groupe de mots dans le texte 22 (1-22.txt) est sur-représenté et le texte (1-23.txt) dont le cible est sous-représenté. Pour découvrir ce qui se passe dans le texte 23 dont le sujet est "la protection environnementale" et en même temps, ce sujet n'apparaît pas autant de fois que les autres textes, on peut se servir de bouton 'Parties' dans le Paramètre et on choisit le texte '1-23.txt', ensuite, on tape le mot 'l'environnement' (环境) dans la case 'Pôle Source' pour regarder ces cooccurrences en cliquant le bouton 'Cooccurrnents sur partie sélectionnée' dans 'Calcul de cooccurrents' (figure 18). Donc, on constate que dans le texte 1-23, le mot le plus occurrent avec 'l'environnement' ou 'environnemental(e)' n'est pas 'la protection', mais c'est 'la pollution/contamination' qui apparaît 170 fois et 'la protection' représente juste 16 fois dans le texte. En fait, le contexte est autour de la pollution environnementale, alors le sujet est la protection environnementale. Ce genre de phénomène peut expliquer la sous-représentation du terme quand il

est sujet. Il existe aussi une autre explication qui est la vitalité et l'art es mots. On peut toujours trouver des termes similaires ou des synonymes du sujet au lieu d'utiliser le même terme dans tout le contexte. Cela peut aussi expliquer la sous-représentation du sujet dans un texte.



figure(17)

