

Thèse de doctorat / 2020

Université Panthéon - Assas

École doctorale d'économie, gestion, information et communication

Thèse de doctorat en informatique

soutenue le

Composantes géantes sur des flux de données



UNIVERSITÉ PARIS II
PANTHÉON-ASSAS

Achraf Lassoued

Sous la direction de Michel De Rougemont

Rapporteur :

Mme. Vilnat Anne, Professeur associé à l'Université Paris-Saclay

M. Laurent Dominique, Professeur émérite à l'Université Cergy-Pontoise

Membres du jury :

Mme. Chevalier Céline, Maître de conférences à l'Université Paris II

M. Spyrtatos Nicolas, Professeur émérite à l'Université Paris-Sud

...

1 **Avertissement**

- 2 L'université n'entend donner aucune approbation ni improbation aux opinions émises dans
- 3 cette thèse ; ces opinions doivent être considérées comme propres à leur auteur.

4 Remerciements

5 Ce manuscrit constitue la synthèse de trois ans de doctorat effectué au Centre de Re-
6 cherche en Économie et Droit (*CRED*) à l'Université Paris II Panthéon-Assas.

7
8 Je tiens à remercier et exprimer ma gratitude à mon directeur de thèse, le professeur
9 Michel de ROUGEMONT qui a rendu ce travail possible. Ses précieux conseils, avis et
10 suggestions ont été d'une valeur inestimable à toutes les étapes du travail.

11
12 J'adresse un remerciement particulier à Madame Isabelle Bril, Directrice de recherches
13 (LACITO-CNRS) et Directrice d'Etudes à l'EPHE (Typologie et Typologie des langues
14 austronésiennes), pour son rôle déterminant dans l'encadrement et le financement d'une
15 partie de ma thèse.

16
17 Je remercie sincèrement Madame Anne Vilnat, Professeur à l'Université Paris-Saclay,
18 ainsi que Monsieur Laurent Dominique, Professeur émérite à l'Université de Cergy Pan-
19 toise, pour l'intérêt qu'ils ont témoigné à l'égard de mon travail en acceptant la charge de
20 rapporteurs.

21
22 J'exprime ma reconnaissance à Madame Céline Chevalier, Maître de conférences à l'Uni-
23 versité Paris 2, et Monsieur Nicolas Spyratos, Professeur émérite à l'Université Paris-sud
24 pour leurs participation au jury en tant qu'examineurs.

25
26 Je tiens également à remercier mes collègues, camarades du lycée Saint Nicolas et l'EPF,
27 pour tout le soutien qu'ils m'ont apporté lors de la réalisation de cette thèse.

28
29 Enfin, je voudrais remercier ma famille pour tout leur amour et leurs encouragements.
30 Pour mes parents qui m'ont soutenu dans toutes mes activités, et ma soeur qui a été soli-
31 daire, encourageante et patiente, au cours de cette thèse.

32
33 Je vous remercie.

34 Achraf Lassoued
35 Université Panthéon-Assas

36 **Résumé :**

37 Nous étudions des données de nature diverse sous forme de flux, en particulier :

- 38 • Base de données,
- 39 • Réseaux sociaux,
- 40 • Données de texte.

41 Pour une base de données qui suit un schéma relationnel, un schéma d'analyse *OLAP*
42 (Online Analytical Processing) définit une des tables de la base de données comme une table
43 d'analyse. Nous supposons que les tuples de la table d'analyse arrivent sous forme d'un
44 flux. Nous étudions l'approximation des requêtes *OLAP*, en échantillonnant de manière
45 non uniforme les tuples du flux sans stocker les données d'analyse et donnons un modèle
46 de préférence dans ce cadre.

47 Dans le cas du réseau social *Twitter*, nous observons un flux de tweets qui contiennent
48 un tag donné et le transformons en un flux d'arêtes d'un graphe. Nous souhaitons étudier
49 l'existence des grands clusters dans le graphe ainsi obtenu. Nous proposons une méthode
50 d'échantillonnage uniforme qui va associer au graphe un sous-graphe aléatoire et étudions
51 les composantes géantes de ce sous-graphe aléatoire comme témoin des grands clusters du
52 graphe d'origine.

53 Pour un flux de texte, nous considérons les paires de mots dans une phrase lemmatisée
54 comme des arêtes d'un graphe où les nœuds sont les mots. Nous transformons le flux de
55 texte en flux d'arêtes. Nous échantillonnons les arêtes proportionnellement à la similarité
56 *Word2vec* des mots. Nous analysons ensuite les composantes géantes.

57 Nous étendons les vecteurs *Word2vec* en prenant en compte la morphologie d'une
58 langue, en particulier la structure des préfixes et des suffixes d'un mot. Les nouveaux vec-
59 teurs d'un mot ont 3 composantes : un vecteur pour le préfixe, un vecteur pour la racine
60 et un vecteur pour le suffixe. Nous définissons un vecteur probabiliste pour les phrases.

61 Sur les trois types de données, nous avons échantillonné selon des distributions précises.
62 Les résultats principaux de cette thèse montrent comment approcher les propriétés de ces
63 données avec ces échantillons.

64 *Descripteurs : Algorithmes de streaming, Approximation, Requêtes OLAP, Préférences, Cluste-*
65 *ring, Graphes dynamiques, TAL, Morphologie*

66 **Abstract :**

67 We study different types of data as a stream :

- 68 • Databases,
- 69 • Social Networks,
- 70 • Text data.

71 For a database that follows a relational schema, an *OLAP* (Online Analytical Proces-
72 sing) analysis schema defines one of the database tables as an analysis table. We suppose
73 that the tuples of the analysis table arrive as a stream. We study the approximation of
74 *OLAP* queries, by sampling with weights the tuples of the stream without storing the
75 analysis data. We give a *preference* model in this context.

76 For the social network *Twitter*, we observe a stream of tweets that contain any given
77 tag and transform it into a stream of edges of a graph. We study the existence of large
78 clusters in the generated graph. We propose a uniform sampling method that will associate
79 a random subgraph of the graph and study the giant components of this random subgraph
80 as a witness of the large clusters of the original graph.

81 For a stream of text, we consider the pairs of words in a lemmatized sentence as edges
82 of a graph where the nodes are the words. We transform the stream of text into a stream
83 of edges. We sample the edges proportionally to the *Word2vec* similarity of the words. We
84 then analyze the giant components.

85 We extend the classical *Word2vec* vectors, in order to take into account the morphology
86 of the words, i.e. the prefixes, roots and suffixes. A new vector has 3 components, one for
87 the prefix, one for the root and one for the suffix. We define a probabilistic vector associated
88 with the sentences.

89 With the three types of data, we sampled according to specific distributions. The main
90 results of this thesis show how the properties of these data can be approximated with these
91 samples.

92 *Keywords : Streaming algorithms, Approximation, OLAP queries, Preferences, Clustering, Dy-*
93 *namic graphs, NLP, Morphology*

94 Sommaire

95	Remerciements	5
96	1 Introduction	13
97	1.1 Comparaisons avec les autres techniques d'analyse	14
98	1.2 Organisation de la thèse	15
99	1.3 Publications	16
100	2 Préliminaires d'algorithmique	17
101	2.1 Échantillonnage des flux de données	17
102	2.1.1 Reservoir Sampling	18
103	2.1.1.1 Justification de l'algorithme	18
104	2.1.1.2 Exemple	19
105	2.1.2 Réservoir pondéré	21
106	2.1.3 Echantillonnage biaisé	23
107	2.2 Graphes	24
108	2.2.1 Définitions	24
109	2.3 Graphes aléatoires	25
110	2.3.1 Graphe aléatoire uniforme	25
111	2.3.2 Modèle Erdős-Renyi	25
112	2.3.3 Modèle de configuration	26
113	2.3.4 Composante géante	27
114	2.3.4.1 2-Core	28
115	2.3.5 Conclusion	29
116	3 Préliminaires d'analyse du texte	31
117	3.1 Prétraitement des données	31
118	3.1.1 Lemmatisation et Entités Nommées	31
119	3.1.1.1 Décomposition grammaticale classique	34
120	3.1.1.2 La décomposition par dépendances	35
121	3.2 Vectorisations des mots : <i>Word2vec</i>	36
122	3.3 Classification	39
123	3.3.1 IRaMuTeQ	39

124	3.3.2	LDA	39
125	3.4	Mécanismes d'Attention	40
126	3.4.1	Conclusion	41
127	4	Analyse OLAP sur des flux de données	43
128	4.1	Entrepôts de données et requêtes OLAP	44
129	4.2	Approximation à partir d'échantillons	45
130	4.2.1	Construction des échantillons à partir d'un entrepôt de données	48
131	4.2.2	Approximation des requêtes OLAP sur les données de streaming	49
132	4.3	Préférences pour l'analyse OLAP	49
133	4.3.1	Préférences comme distribution sur les <i>Mesures</i>	51
134	4.3.2	Conclusion	52
135	5	Analyse du réseau social Twitter en streaming	53
136	5.1	Récupération des données	53
137	5.1.1	API Search Twitter	54
138	5.1.2	API streaming Twitter	55
139	5.2	Caractéristiques des données Twitter	56
140	5.2.1	Structure	56
141	5.2.2	Contenus	56
142	5.3	Construction d'un graphe social à partir d'un arbre JSON	58
143	5.3.1	Twitter graphe d'un flux de données	59
144	5.4	Propriétés d'un graphe Twitter	60
145	5.4.1	Distribution des degrés	60
146	5.4.2	Propriétés à analyser	61
147	5.4.2.1	Communautés et sous-graphes denses	62
148	5.4.2.2	Graphes dynamiques dans un flux d'arêtes	64
149	5.4.2.3	Corrélation de contenu entre flux	65
150	5.4.3	Conclusion	66
151	6	Analyse de texte en streaming	67
152	6.1	Algorithmes de streaming	67
153	6.2	Représentation du texte en streaming	68
154	6.3	Stabilité	68
155	6.4	Classification	69
156	6.4.1	Distances entre les composantes géantes	69
157	6.4.2	Classification en k -classes	70
158	6.5	Phrases centrales et mécanismes d'attention	70
159	6.6	Préférences	70
160	6.7	Conclusion	73

161	7	Modèle de morphologie dans les langues naturelles	75
162	7.1	Langue avec une forte morphologie	76
163	7.2	Un modèle statistique pour la morphologie	76
164	7.2.1	Statistiques de base pour la langue <i>Amis</i>	77
165	7.2.2	Représentation vectorielle des préfixes, racines et suffixes	78
166	7.2.3	Distributions et vecteurs représentatifs	80
167	7.2.4	Vecteurs de phrases	82
168	7.3	Un aperçu syntaxique de la langue <i>Amis</i>	82
169	7.4	Classification des phrases à partir des préfixes et suffixes	84
170	7.4.1	Vecteur de phrases et mécanismes d'attention	86
171	7.5	Grammaires et statistiques	87
172	7.5.1	Les meilleurs arbres de dérivation	88
173	7.5.2	Comparaison avec d'autres modèles stochastiques	88
174	7.5.3	Conclusion	89
175	8	Recherche et explications	91
176	8.1	Corrélation de flux	91
177	8.2	Phylogénie entre flux	92
178	8.3	Distance entre mots, relative à des flux	93
179	8.4	Moteur de recherche par similarité	94
180	8.4.1	Algorithme de recherche	95
181	8.4.2	Justification de l'algorithme de recherche	95
182	8.4.3	Structures de données	96
183	8.5	Explications	97
184	8.6	Conclusion	98
185	9	Expériences	99
186	9.1	Flux Twitter	99
187	9.1.1	Expérience sur un flux	100
188	9.1.2	Expérience sur plusieurs flux	101
189	9.2	Analyse de texte	102
190	9.2.1	Échantillonnage uniforme ou pondéré	103
191	9.2.2	Classification	106
192	9.3	Décompositions hiérarchiques des composantes géantes	108
193	10	Conclusion	111
194		Bibliographie	113
195	A	Phrases de la langue Amis	123

1 Introduction

Des flux de données sont disponibles partout autour de nous, soit sous forme de données que nous recherchons (pages web, documents), soit sous forme de données envoyées par une plateforme numérique suite à un abonnement. Par exemple une vidéo envoyée par *Netflix* ou un flux envoyé par un réseau social comme *Twitter*.

Les réseaux sociaux font aujourd'hui partie de notre vie quotidienne après l'apparition de Facebook et Twitter. Les données issues de ces réseaux sont très différentes des données classiques que nous utilisons en exploration de données. Ces réseaux génèrent des flux de données volumineux qui permettent de développer des nouveaux services et de créer des nouvelles valeurs économiques. Dans ce cadre la règle des quatre "V" symbolise ce nouveau monde :

- *Volume* : la quantité volumineuse d'informations.
- *Variété* : l'hétérogénéité des formats, de types, et de qualité des informations.
- *Vélocité* : l'aspect dynamique et temporel des données.
- *Véracité* : le fait que les données contiennent beaucoup de bruit.

Dans cette thèse, on introduit des méthodes qui gèrent des flux volumineux, pour une variété de données, qui évoluent très vite dans le temps et qui ont beaucoup de bruit. Ces méthodes sont basées sur un échantillonnage des données selon différentes lois. La loi la plus simple est la loi uniforme. On analyse les types de données suivants :

- Entrepôts de données au sein d'une base de données.
- Graphes issus des réseaux sociaux .
- Texte.

Notre méthode d'analyse de données s'applique à ces trois types de données et suit les phases suivantes :

- Nous échantillons des flux de données avec des distributions particulières pour ne retenir qu'un nombre constant d'échantillons. Pour chaque élément du flux, nous décidons avec une certaine probabilité si l'élément est conservé ou pas. L'ensemble des échantillons est stocké dans un *Réservoir* qui est une structure aléatoire.
- L'analyse de données consiste à observer les composantes géantes du Réservoir.

- 226 • On introduit une distance entre ces composantes géantes qui permettra de classifier et
227 de rechercher les propriétés essentielles des données.

228 Supposons un flux de tuples t associés à un entrepôt de données qu'on cherche à ana-
229 lyser à l'aide de requêtes d'analyse OLAP (Online Analytical Processing) pour certaines
230 Mesures M de ces tuples. Nous montrons comment approximer les réponses aux requêtes
231 OLAP en échantillonnant les tuples selon une certaine loi qui dépend de la Mesure M . On
232 obtient alors des réponses approximatives en temps constant sans lire tout l'entrepôt de
233 données.

234

235 Les réseaux sociaux génèrent des flux de données, plus précisément des flux d'arêtes
236 d'un graphe. À partir de certains mots clés, Twitter génère tous les tweets qui contiennent
237 ces mots clés. Typiquement, nous recevons 10^3 tweets par minute. Une méthode classique
238 transforme le flux de tweets en un flux d'arêtes d'un graphe où les nœuds sont les tags.
239 L'analyse des grandes communautés du graphe Twitter est basée sur la détection de compo-
240 santes géantes dans un graphe aléatoire. Nous utilisons un échantillonnage de Réservoir et
241 dans ce cas, suivons la loi uniforme. Nous analysons les composantes géantes du Réservoir
242 de taille constante k .

243

244 Pour l'analyse de texte, on utilise la technique *Word2vec*, introduite dans [44, 51] qui
245 associe un vecteur de petite dimension à chaque mot. La similarité de deux mots est le
246 produit de deux vecteurs. Un abonnement à un flux RSS d'article d'actualité génère un
247 flux important de texte. Nous utilisons des méthodes classiques pour le prétraitement du
248 texte par exemple la lemmatisation. Nous considérons ensuite les paires de mots d'une
249 même phrase comme des arêtes d'un graphe où les noeuds sont des mots. Le flux de texte
250 est transformé ainsi en un flux d'arêtes. Nous échantillonnons les arêtes avec un poids
251 proportionnel à la similarité *Word2vec* des mots. Nous analysons ensuite les composantes
252 géantes qui deviennent stables dans ce modèle.

253

254 Nous étendons l'approche *Word2vec* en utilisant la morphologie d'une langue. Pour
255 un mot qui se décompose selon la structure *préfixe-racine-suffixe*, on peut généraliser la
256 conception de vecteur aux préfixes, aux racines et aux suffixes. Cette analyse a été faite
257 à partir d'une langue naturelle l'*Amis* et nous développons un modèle statistique dans ce
258 cadre.

259 1.1 Comparaisons avec les autres techniques d'analyse

260 Il existe de nombreuses méthodes pour analyser les bases de données. En général, les
261 données sont stockées et ne sont pas considérées comme des flux. L'approximation est ra-
262 rement un sujet central dans le cadre de l'analyse OLAP.

263

264 Dans les réseaux sociaux, une méthode classique pour trouver des communautés, c'est-à-
265 dire des sous-graphes denses, est d'utiliser des méthodes spectrales telles que la décomposition
266 par module. La complexité $O(n^3)$ de ces méthodes n'est pas adaptée aux données volumi-
267 neuses. D'autres méthodes [8, 29, 43, 25] se concentrent sur les solutions algorithmiques
268 pour approximer des sous-graphes denses. Un des buts est d'avoir un algorithme de flux
269 qui utilise une mémoire sous-linéaire pour détecter des sous-graphes denses. Il n'est pas
270 possible de trouver un tel algorithme, pour la notion de complexité dans le pire des cas.
271 Nous présentons cependant, un algorithme efficace pour les graphes qui satisfont une dis-
272 tribution de degré qui suit une loi de puissance.

273
274 Pour le texte, nous échantillons selon une distribution, analysons des composantes
275 géantes d'un graphe et introduisons une distance entre ces composantes. Nous pouvons
276 ensuite classifier les textes en utilisant l'algorithme classique de *k-means* où k est le
277 nombre de classes. Chacune des k classes est un sujet, avec une composante géante comme
278 représentant. Notre méthode est beaucoup plus simple que les techniques classiques utilisées
279 dans le traitement de texte comme IRaMuTeQ [52] ou LDA (Latent Dirichlet Allocation)
280 [10]. Les mécanismes d'attention [40, 60] suivent les mêmes idées.

281

282 1.2 Organisation de la thèse

283 Dans le chapitre 2 nous introduisons les préliminaires algorithmiques, les concepts de
284 théorie des graphes et en particulier les différents modèles de graphes aléatoires.

285

286 Dans le chapitre 3, nous introduisons les préliminaires d'analyse de texte. Nous allons
287 présenter les méthodes classiques utilisées dans le monde du traitement du langage naturel.

288

289 Dans le chapitre 4, nous étudions les données d'un entrepôt de données sous forme de
290 flux. Nous montrons comment approximer des requêtes OLAP dans ce cadre en suivant
291 une loi d'échantillonnage particulière, c'est un des résultats principaux de la thèse.

292

293 Dans le chapitre 5, nous présentons plusieurs méthodes pour récupérer les données
294 Twitter sous forme d'un flux depuis ce réseau social, puis nous transformons ces données
295 en un flux d'arêtes et échantillons ces arêtes avec une loi uniforme pour analyser les
296 composantes géantes du graphe Twitter.

297

298 Dans le chapitre 6, nous présentons des méthodes pour analyser du texte sous forme
299 d'un flux. Nous transformons chaque phrase en un flux d'arêtes et échantillons ces
300 arêtes selon une distribution *Word2vec*. Un des résultats importants de la thèse est de
301 montrer comment l'échantillonnage selon cette distribution est stable. Nous montrons en-
302 suite comment utiliser les composantes géantes pour classifier le texte.

303

304 Dans le chapitre 7, nous montrons que notre méthode peut se généraliser en étudiant
305 la morphologie des mots. Nous présentons un exemple d'analyse sur la langue *Amis*, une
306 langue austronésienne de Taiwan.

307

308 Dans le chapitre 8, nous présentons un moteur de recherche basé sur les distances entre
309 les composante géantes introduites dans les chapitres précédents.

310

311 Dans le chapitre 9, nous présentons les différentes expériences réalisées dans le cadre
312 de cette thèse.

313

314 1.3 Publications

315 Cette thèse est basée sur les publications suivantes :

- 316 ● A statistical model for morphology inspired by the Amis language, *Proceedings of the*
317 *Language, Ontology, Terminology and Knowledge Structures Workshop*, LOTKS 2017
318 [16] (version étendue de l'article soumise à un journal).
- 319 ● Topics as giant components in texts generated by a stream, *soumis à une conférence*.
- 320 ● Preferences as distributions on streaming data, *soumis à International Workshop on*
321 *Information Search, Integration and Personalization*, ISIP2020.

322 2 Préliminaires d'algorithmique

323 Dans ce chapitre nous introduisons des méthodes d'approximation en échantillonnant
324 un flux d'information en particulier pour un flux d'arêtes d'un graphe G . On s'intéresse
325 aux algorithmes qui détectent la présence de cluster dans ce graphe ou qui approximent
326 les clusters. Pour des flux d'arêtes, nous considérons des algorithmes probabilistes qui
327 échantillonnent les arêtes et génèrent des graphes aléatoires.

328 Nous considérons les composantes géantes du graphe aléatoire et interprétons chaque
329 composante géante comme le représentant d'un cluster. Dans le cas de l'analyse du texte,
330 un algorithme standard transforme le flux de texte en un flux d'arêtes. On réduit l'analyse
331 du texte à l'analyse du graphe aléatoire associé.

332 Nous distinguons :

- 333 • Les méthodes d'échantillonnage,
- 334 • Les graphes aléatoires et l'analyse des composantes géantes.

335 2.1 Échantillonnage des flux de données

336 Un panorama des différentes techniques d'échantillonnage sur les graphes est présenté
337 dans [34]. Nous considérons seulement l'échantillonnage d'arêtes dans un flux d'arêtes :

- 338 • Échantillonnage uniforme par Réservoir,
- 339 • Échantillonnage pondéré,
- 340 • Échantillonnage biaisé.

341 Les problèmes sous-jacents liés à l'échantillonnage de flux de données nécessitent l'ap-
342 plication de nombreuses idées et techniques issues de l'échantillonnage classique, mais ils
343 nécessitent également de nouvelles innovations significatives.

344 Notre discussion commence par l'obtention d'un échantillon d'arêtes, lorsque le nombre
345 d'échantillon souhaité est beaucoup plus petit que le nombre total d'éléments dans le flux.
346 Les algorithmes les plus simples consistent à échantillonner le flux de données avec une
347 distribution uniforme.

348 Pour certaines applications, il peut être souhaitable d'échantillonner selon une distribu-
349 tion biaisée vers des éléments plus récents. Dans les sections suivantes, nous allons discuter

350 des schémas d'échantillonnage uniforme, pondéré en supposant que chaque élément à un
351 poids et biaisé en favorisant les éléments les plus récents.

352 2.1.1 Reservoir Sampling

353 Le Reservoir Sampling est un algorithme probabiliste, introduit par Jeffrey Vitter[62].
354 Soit R un Réservoir contenant au plus k arêtes dans un flux de m arêtes e_1, e_2, \dots, e_m . Nous
355 remplissons d'abord le Réservoir avec e_1, e_2, \dots, e_k . Pour $i > k$:

- 356 • On garde e_i avec probabilité k/i .
- 357 • Si on garde e_i , on enlève une des arêtes du Réservoir (avec probabilité $1/k$) pour faire
358 de la place à e_i .

Algorithme 1 Approche classique du Reservoir Sampling

INPUT :

k // size of a Reservoir of size k ,

e_1, e_2, \dots, e_m // Stream

$i = 0$

for each edge arriving from the input stream **do**

$i = i+1$

if $i \leq k$ **then**

add the edge to the Reservoir

else

decide with the probability $\frac{k}{i}$ whether to accept the edge

if the edge is accepted **then**

replace a randomly selected edge in the Reservoir with the accepted edge

end if

end if

end for

359 2.1.1.1 Justification de l'algorithme

360 La propriété du Réservoir est que chaque arête e_i a une probabilité $\frac{k}{m}$ d'être dans le
361 Réservoir, c'est-à-dire uniforme. On peut le montrer par récurrence sur m . Supposons que
362 la propriété soit vraie à l'étape $m - 1$, toutes les arêtes ont la probabilité $\frac{k}{m-1}$. L'élément
363 m a une probabilité $\frac{k}{m}$ de faire partie du Réservoir. Pour les éléments qui en font déjà
364 partie, leur probabilité de rester se décompose en deux événements :

- 365 • On ne prend pas l'élément e_m , la probabilité c'est : $1 - \frac{k}{m}$.
- 366 • On prend l'élément e_m , avec probabilité $\frac{k}{m}$ et on survit à l'éjection avec $\frac{k-1}{k}$.

367 La probabilité de rester dans le Réservoir est donc la somme :

368

$$\frac{m-k}{m} + \frac{k}{m} \frac{k-1}{k} = \frac{m-1}{m}.$$

369

Pour être dans le Réservoir à l'itération m il faut que deux événements soient réalisés :

370

- Être dans le Réservoir à l'étape $m - 1$ et on sait que par récurrence la probabilité est :

371

$$\frac{k}{m-1}.$$

372

- Rester dans le Réservoir : la probabilité de rester est : $\frac{m-1}{m}$ d'après ce qu'on vient de voir.

373

374

La probabilité d'être dans le Réservoir à l'étape m est donc :

375

$$\frac{m-1}{m} \frac{k}{m-1} = \frac{k}{m}.$$

376

2.1.1.2 Exemple

377

L'espace probabiliste Ω est déterminé par les choix effectués à chaque étape de l'algorithme. L'algorithme s'applique à tous types de données. On peut présenter l'espace probabiliste sous forme d'un arbre qui reproduit les décisions de l'algorithme.

379

380

- Chaque feuille représente un événement, le résultat du tirage.

381

- Les branches représentent les différentes possibilités quand on rajoute un nouvel élément et la probabilité est indiquée sur chaque branche.

382

383

On va présenter le Réservoir quand on a un flux d'arêtes mais on peut aussi l'appliquer à un flux de valeurs. La figure 2.1 montre un exemple où le flux est composé des valeurs 1, 2, 3, 4, 5 avec $k = 3$. Le premier tirage a lieu quand on lit la valeur 4 et le Réservoir contient 1, 2, 3, la branche gauche correspond à la situation où on ne garde pas l'élément 4 avec probabilité $1/4$ et les trois autres branches correspondent au cas où on garde la valeur 4 et on enlève 1 sur la branche 2, on enlève 2 sur la branche 3, on enlève 3 sur la branche 4. On répète l'algorithme sur l'entrée 5, et l'arbre global est indiqué sur la figure 2.1.

389

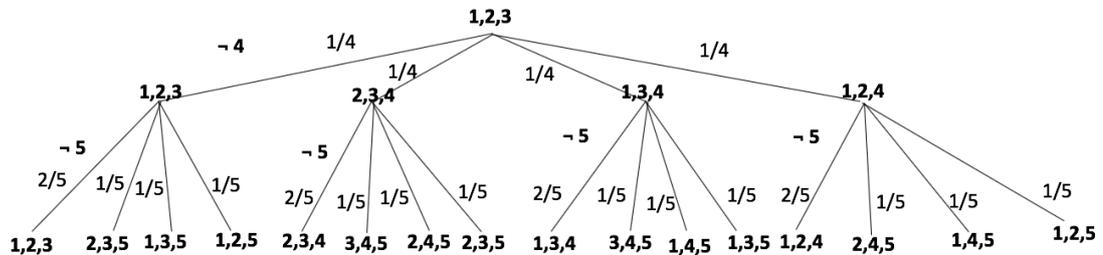


FIGURE 2.1: Espace probabiliste

La propriété essentielle de cette technique, est que tous les sous-ensembles sont équiprobables. Le nombre de sous-ensembles de 3 éléments parmi 1,2,3,4,5 est :

$$\binom{5}{3} = \frac{5!}{3!(5-3)!} = \frac{5 * 4 * 3 * 2}{3 * 2 * 2} = 10$$

390 On va montrer que la probabilité que le Réservoir contient un sous ensemble est égale
391 1/10 suivant deux explications.

Comment peut-on calculer dans ce tirage la probabilité d'avoir un sous ensemble S ? Pour un sous-ensemble S donné la probabilité qu'un élément donné soit dans le Réservoir est $\frac{k}{m}$. Maintenant si on prend un deuxième élément parmi le deuxième élément de S , la probabilité est $\frac{k-1}{m-1}$ et $\frac{k-2}{m-2}$ pour le troisième. Si on prend $k = 3$ et $m = 5$ la probabilité d'avoir un sous ensemble :

$$Prob(S) = \frac{3}{5} * \frac{2}{4} * \frac{1}{3} = \frac{3 * 2 * 1}{5 * 4 * 3} = \frac{6}{60} = \frac{1}{10}$$

392 La probabilité qu'un élément soit dans le Réservoir est 3/5, la probabilité qu'un deuxième
393 élément soit dans le Réservoir est 2/4 et la probabilité qu'un troisième élément soit dans
394 le Réservoir est 1/3. La probabilité d'avoir un sous-ensemble à trois éléments est bien le
395 produit c'est à dire 1/10, ce qui correspond aux probabilités observées dans la figure 2.1.

396 L'arbre de la figure 2.1 a 16 feuilles et chaque noeud a 4 successeurs. Parmi les 16
397 feuilles, 4 ont une probabilité de 1/10 et 12 autres ont une probabilité de 1/20. Il existe 4
398 sous-ensembles parmi les 10 qui ont une seule feuille et une probabilité de 1/10 et 6 autres
399 qui ont 2 feuilles, chacune de probabilité 1/20. Chaque sous-ensemble parmi les 10 a donc
400 bien la même probabilité de 1/10

401 Maintenant si on regarde le flux 1, 2, 3, 4, 5..... n avec $k = 3$, on obtient l'arbre de la
402 figure 2.2 qui montre le cas général.

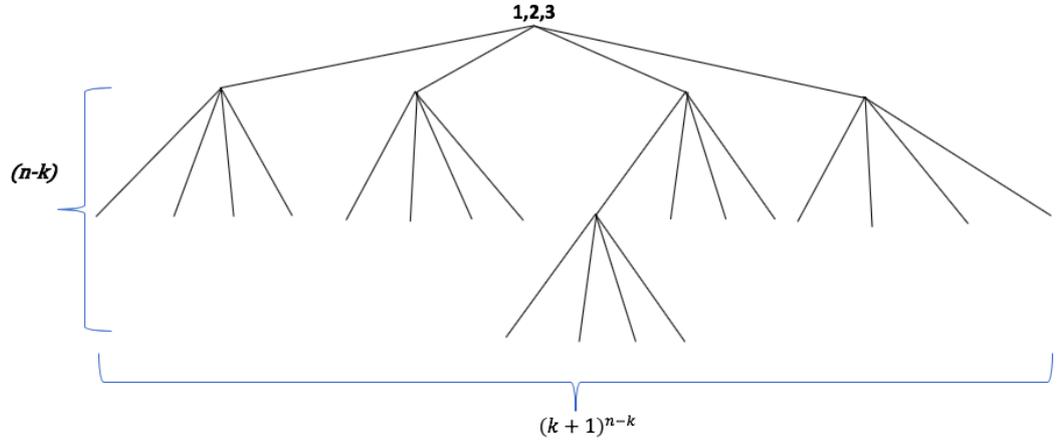


FIGURE 2.2: Espace probabiliste cas général

403 2.1.2 Réservoir pondéré

404 Dans le cas des arêtes e avec des poids w_e , l'échantillonnage du Réservoir pondéré
 405 conserve k arêtes avec la probabilité $\left(\frac{k * w_e}{\sum w_e}\right)$. Soit e_1, e_2, \dots, e_m un flux d'arêtes, soit w_1, w_2, \dots, w_m
 406 le poids qui correspond à chaque arête et $R_m = e_{i_1}, e_{i_2}, \dots, e_{i_k}$ le Réservoir à l'état m de taille
 407 k .

408

Lemme 1. Pour tout m la probabilité qu'une arête e_i soit dans le Réservoir R_m est

$$\frac{k * w_i}{\sum_{i=1, \dots, m} w_i}.$$

409 Comme dans le cas du Reservoir Sampling, on va montrer ce résultat par récurrence.
 410 On suppose que c'est vrai pour l'étape m et on veut démontrer que c'est vrai à l'étape
 411 $m + 1$.

Démonstration. Supposons que

$$Prob[e_i \in R_m] = \frac{k * w_i}{\sum_{i=1, \dots, m} w_i}$$

Algorithme 2 Approche classique du Reservoir pondéré

INPUT :

k // size of a Reservoir of size k ,

e_1, e_2, \dots, e_m // stream of edges,

w_1, w_2, \dots, w_m // weight of edges

$i = 0$

for each edge arriving from the input stream **do**

$i = i + 1$

if $i \leq k$ **then**

 add the edge to the Reservoir

else

$j = \text{random}(0, 1)$

 decide with the probability

$$\frac{K \cdot w_i}{\sum_{j \leq i} w_j}$$

 whether to insert the edge

if the edge is accepted **then**

 replace a randomly selected edge in the Reservoir with the accepted edge.

end if

end if

end for

en suivant l'algorithme, la probabilité que e_{m+1} soit dans R_{m+1} est :

$$Prob[e_{m+1} \in R_{m+1}] = \frac{k * w_{m+1}}{\sum_{i=1, \dots, m+1} w_i}$$

412 La probabilité qu'une arête soit dans le Réservoir à l'étape $m+1$, est la probabilité que
413 cette arête soit à l'étape m et qu'elle survive au tirage de l'arête e_{m+1} . La probabilité de
414 survivre est similaire à celle qu'on a calculée dans le cas du Réservoir uniforme avec deux
415 événements :

416 • C'est la probabilité que e_m ne soit pas prise : $1 - \frac{(k * w_{m+1})}{\sum_{i=1, \dots, m+1} w_i}$.

417 • La probabilité que e_m soit prise mais que l'arête ne soit pas sélectionnée pour sortir du
418 Réservoir : $\frac{k * w_{m+1}}{\sum_{i=1, \dots, m+1} w_i} * \frac{(k-1)}{k}$.

$$i \in [1, \dots, m], Prob[e_i \in R_{m+1}] = Prob[e_i \in R_m] * \left[1 - \frac{(k * w_{m+1})}{\sum_{i=1, \dots, m+1} w_i} + \frac{k * w_{m+1}}{\sum_{i=1, \dots, m+1} w_i} * \frac{(k-1)}{k} \right]$$

$$\begin{aligned} & \frac{k * w_i}{\sum_{i=1, \dots, m} w_i} * \left[\frac{\sum_{i=1, \dots, m+1} w_i - k * w_{m+1} + k * w_{m+1} - w_{m+1}}{\sum_{i=1, \dots, m+1} w_i} \right] \\ & = \\ & \frac{k * w_i * \sum_{i=1, \dots, m} w_i}{\sum_{i=1, \dots, m} w_i * \sum_{i=1, \dots, m+1} w_i} = \frac{k * w_i}{\sum_{i=1, \dots, m+1} w_i} \end{aligned}$$

419

□

420 2.1.3 Echantillonnage biaisé

421 En général, les échantillons biaisés reposent sur une fonction de biais temporel, afin de
422 favoriser les données récentes dans le flux. Le Réservoir biaisé présenté par Aggarwal [2]
423 fournit un biais exponentiel en faveur des éléments récemment observés.

424 La fonction de biais est donnée par $f(i, t)$ où i est le i ème élément de flux et t est un
425 horodatage futur.

$$f(i, t) = e^{-\lambda(t-i)} \quad (2.1)$$

426 Le paramètre λ définit le taux du biais, il est généralement compris entre $[0, 1]$

427 2.2 Graphes

428 Les définitions qui sont présentées dans cette section suivent largement celles de Biggs
429 [9] et de Bollobás [12].

430 2.2.1 Définitions

431 Un graphe G est constitué d'un ensemble de nœuds V et d'un ensemble d'arêtes E .
432 L'ensemble E est constitué de paires de nœuds (u_1, u_2) où $u_1, u_2 \in V$. Les arêtes peuvent
433 être dirigées ou non dirigées si le graphe est symétrique.

434
435 Un sous-graphe G_s peut être généré à partir de G soit en sélectionnant un sous-ensemble
436 V_s de V , soit par un sous-ensemble E_s de E . Un sous-graphe G_s généré à partir d'un
437 sous-ensemble V_s présente toutes les arêtes de G qui connectent des nœuds dans V_s . Un
438 sous-graphe G_s généré à partir d'un sous-ensemble E_s de E , présente les nœuds de G qui
439 existent dans E_s .

440
441 **Degré.** Le degré (ou valence) d'un nœud u est le nombre d'arêtes qui lui sont incidentes.
442 Étant donné un graphe symétrique $G = (V, E)$ la somme des degrés d'un nœud u est donnée
443 par l'équation suivante :

$$\sum_{u \in V} d(u) = 2|E| \quad (2.2)$$

444 Dans un graphe orienté, nous pouvons distinguer le degré d'un nœud u par son degré
445 entrant $d^-(u)$ défini comme le nombre d'arêtes dirigées vers le nœud u et le degré sortant
446 $d^+(u)$ défini comme le nombre d'arêtes sortants du nœud u .

447
448 **Densité.** La densité d'un graphe $G = (V, E)$ est définie par le rapport entre le nombre
449 d'arêtes divisées par le nombre d'arêtes possibles. La densité d'un graphe non orienté est
450 donc égale à :

$$D = \frac{2|E|}{|V| \cdot (|V| - 1)} \quad (2.3)$$

451 Pour un graphe orienté, la densité est égale à :

$$D = \frac{|E|}{|V| * (|V| - 1)} \quad (2.4)$$

452 Notons que le nombre maximal d'arêtes d'un graphe non orienté est :

$$\frac{|V| * (|V| - 1)}{2} \quad (2.5)$$

453 Un autre concept utile est le degré moyen ou la connectivité d'un graphe définie comme :

$$\frac{2|E|}{|V|}. \quad (2.6)$$

454 **Demi-arêtes.** Étant donné un graphe symétrique $G = (V, E)$ et une arête $e \in E$, on
 455 considère une opération de *coupure* qui définit $G'_e = (V \cup V', E')$, obtenu par la transfor-
 456 mation suivante. Si $e = (v_1, v_2) \in E$, on génère deux nouveaux noeuds v'_1, v'_2 dans V' et
 457 on remplace l'arête e par les deux arêtes symétriques (v_1, v'_1) et (v_2, v'_2) de E' .

458 $G'_e = (V \cup V', E')$ est défini par $V' = \{v'_1, v'_2\}$ et $E' = E \cup \{(v_1, v'_1), (v_2, v'_2)\} - \{(v_1, v_2)\}$.
 459 Une *demi-arête* est une arête entre un noeud de V et un noeud de V' . L'opération inverse
 460 consiste à reconnecter les demi-arêtes pour obtenir G à partir de G'_e .

461 2.3 Graphes aléatoires

462 Un graphe aléatoire est un graphe généré par un processus aléatoire, celui-ci peut être
 463 obtenu de différentes manières. Une des possibilités est de choisir un graphe parmi une
 464 collection de graphes possibles avec une certaine probabilité, une autre possibilité consiste
 465 à choisir une arête avec une certaine probabilité. Plusieurs modèles de graphes aléatoires
 466 possibles, dans ce chapitre on va mentionner que trois modèles et finalement on va utiliser
 467 le modèle Erdős-Renyi et le modèle de configuration.

468 2.3.1 Graphe aléatoire uniforme

469 **Definition 1.** *Un graphe aléatoire tiré selon $G(n, m)$ est un graphe obtenu en échantillonnant*
 470 *uniformément tous les graphes non isomorphes à n sommets et m arêtes.*

471 La probabilité de sélectionner un graphe \widehat{G} nécessite de déterminer la taille de l'ensemble
 472 de tous les graphes possibles.

473 2.3.2 Modèle Erdős-Renyi

474 Nous pouvons également générer un graphe aléatoire en décidant pour chaque arête si
 475 on la garde avec une probabilité p , en partant du graphe complet.

476 **Definition 2.** *Un graphe aléatoire \widehat{G} tiré selon $G(n, p)$ est obtenu en partant de l'ensemble*
 477 *des sommets $V = \{1, 2, 3...n\}$, et en reliant chaque paire de sommets i, j où $i \neq j$ par une*
 478 *arête avec probabilité p tel que $0 \leq p \leq 1$.*

479 Ce modèle est généralement appelé le modèle des graphes aléatoires Erdős-Renyi (ER),
 480 décrit par Erdős-Renyi dans deux documents de 1959 [27] et de 1960 [28].

481 2.3.3 Modèle de configuration

482 Les graphes de type Erdős-Renyi ont une distribution de degrés de type gaussienne.
 483 Dans des nombreux cas pratiques, par exemple dans les graphes issus des réseaux sociaux,
 484 les graphes ont une distribution de degré qui n'est pas gaussienne mais qui suit une loi de
 485 puissance. Nous introduisons donc le modèle de configuration qui est un modèle de graphes
 486 aléatoires qui intègre une distribution arbitraire de degrés.

487
 488 Une distribution de degré \mathcal{D} , est une séquence qui pour un n fixe détermine le nombre
 489 de nœuds de degré 1, de degré 2, de degré i tel que la somme de ces nombres est n . Étant
 490 donné cette séquence, le nombre d'arêtes m d'un graphe avec n nœuds est la somme des
 491 degrés divisée par 2.

492 **Definition 3.** *Étant donné une séquence de degrés arbitraire \mathcal{D} , un graphe aléatoire \widehat{G}*
 493 *tiré selon le modèle de configuration est un graphe aléatoire à n sommets, où on associe d_i*
 494 *demi-arêtes à chaque nœuds u_i selon \mathcal{D} . On choisit un appariement (matching) uniforme*
 495 *entre les demi-arêtes ; les demi-arêtes sont choisies de manière aléatoire, afin de former*
 496 *des arêtes.*

497 Une propriété fondamentale du modèle de configuration est que la probabilité d'avoir
 498 une arête entre les sommets i et j :

$$p_{ij} = \frac{d_i \cdot d_j}{2m} \quad (2.7)$$

499 Par exemple, pour une séquence de degrés $\mathcal{D}_{\neq} = (3, 2, 1)$, nous pouvons générer un
 500 graphe aléatoire à partir de cette séquence. En faisant un matching aléatoire uniforme sur
 501 les (demi-arêtes) attachées aux sommets, nous pouvons en déduire que chaque matching de
 502 ce type se produit avec une probabilité égale. Ce processus peut générer des multi-graphes,
 503 c'est-à-dire des graphes avec plusieurs arêtes parallèles. Cependant tous les graphes simples
 504 sont équiprobables. La figure 2.3 illustre le processus.

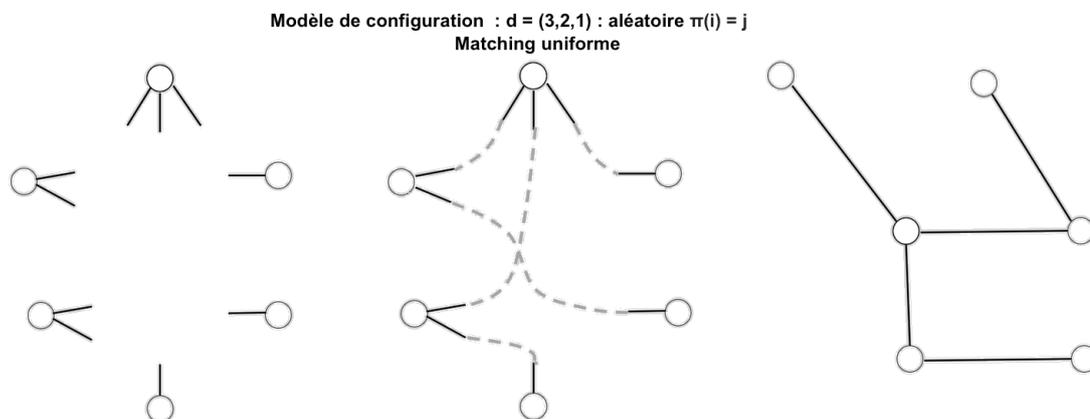


FIGURE 2.3: Exemple du modèle de configuration

505 2.3.4 Composante géante

506 Une composante géante est une composante connexe d'un graphe aléatoire donné, qui
507 contient une fraction constante des sommets du graphe.

Definition 4. Une famille de graphe $\{G_n : n \in N\}$ admet une composante géante si :

$$\exists \alpha, n_0 \quad \forall n \geq n_0 \quad \exists \text{ une composante connexe } C : |C| \geq \alpha \cdot n$$

508 Avec un graphe aléatoire $G(n, p)$, si la valeur de p est faible, alors on a une faible densité
509 d'arêtes et le graphe ne contient pas de composante géante. Une valeur plus élevée de p
510 implique une densité d'arête élevée et nous observons ainsi l'émergence de composantes
511 géantes dans le graphe.

512 Le résultat fondamental d'Erdős-Renyi est de réaliser que si $p > \frac{1}{n}$ alors on a une
513 composante géante et la transition est brutale. Pour tout ϵ , si $p = \frac{1}{n} - \epsilon$ on n'a pas de
514 composante géante, alors que pour $p = \frac{1}{n} + \epsilon$ on a une composante géante. Ce résultat est
515 vrai avec grande probabilité. Quand la transition est brutale, on dit qu'on a une *transition*
516 *de phase*.

517 Le modèle Erdős-Renyi prend un graphe complet et applique des tirages sur chacune
518 des arêtes. Si on applique Erdős-Renyi non pas au graphe complet mais sur une γ -clique,
519 on peut observer que la probabilité de transition est $\frac{1}{\gamma \cdot n}$, démontré dans [42].

520 **Lemme 2.** Si G est une γ -clique alors pour le modèle Erdős-Renyi, il existe une compo-
521 sante géante si $p \geq \frac{1}{\gamma \cdot n}$

522 Dans la figure 5.2, le nombre de noeuds $n = 100$, la taille de la composante géante
523 $|C| = 40$. Il suffit de prendre $\alpha = 0,4$.

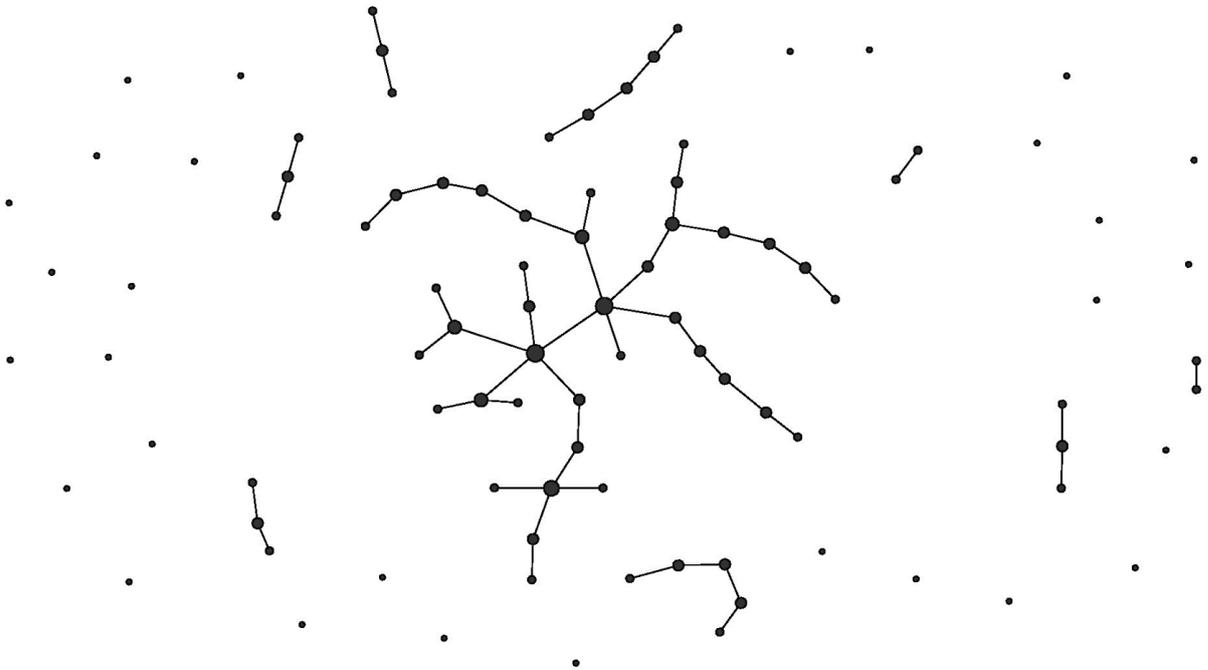


FIGURE 2.4: Exemple de Composante géante

524 Dans le modèle de configuration, nous cherchons une condition similaire pour observer
 525 une composante géante. Un résultat fondamental des travaux de Molloy-Reed [48] donne
 526 les conditions suffisantes pour l'existence de composante géante dans le modèle de configu-
 527 ration pour une distribution de degré D . Parmi ces conditions il faut que $\mathbb{E}[D^2] - 2\mathbb{E}[D] > 0$
 528 pour avoir une composante géante.

529 De nombreuses études dans le monde de la théorie des graphes se concentrent sur
 530 l'existence de composantes géantes ; ces recherches visent à trouver la transition de phase,
 531 en particulier lorsque nous avons une distribution arbitraire de degrés D . Dans notre cas,
 532 nous avons un Réservoir d'arêtes obtenu par la combinaison des deux modèles : nous
 533 avons un graphe G qui suit une distribution de degrés fixe (la loi de puissance), puis nous
 534 appliquons Erdős-Renyi pour obtenir un Réservoir. Sous quelles conditions, observons nous
 535 des composantes géantes dans le Réservoir ?

536 2.3.4.1 2-Core

537 Une approche classique de théorie des graphes est de construire un 2-Core à partir d'un
 538 graphe connexe. On obtient un graphe G_1 , si on répète l'élimination des noeuds de degré 1
 539 de G_1 , on obtient un nouveau graphe G_2 . Pour obtenir le 2-Core d'un graphe, on élimine
 540 tous les noeuds de degré 1 et on répète cette opération jusqu'à obtenir des noeuds de degré
 541 au moins 2. On peut généraliser cette notion au d -Core pour $d > 2$.

542 Nous allons utiliser la notion de 2-Core dans la section 9.3 pour approximer une com-
 543 posante géante. La figure 2.5 montre un graphe et son 2-Core est spécifié par le cercle, en
 544 éliminant les noeuds 7, 8, 9. À la première itération on élimine 7, 9, à la deuxième itération
 545 on élimine 8.

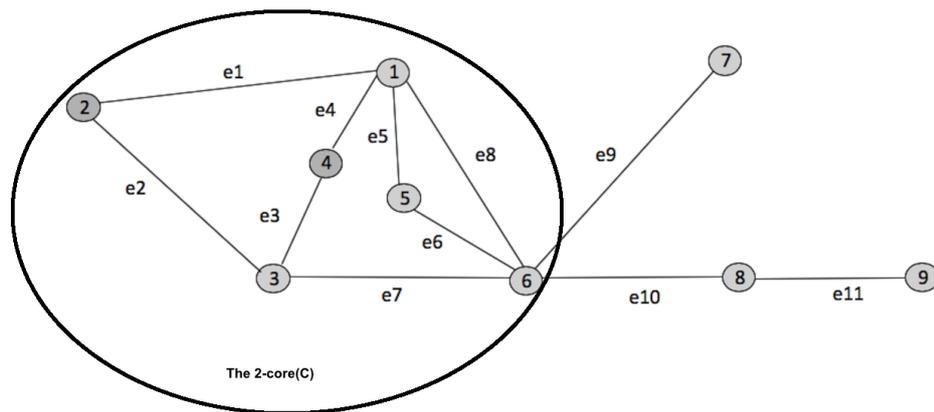


FIGURE 2.5: 2-Core d'une composante géante.

546 2.3.5 Conclusion

547 Dans ce chapitre, nous avons présenté différentes techniques d'échantillonnage pour
 548 les types de données considérés dans les chapitres suivants. Nous avons mentionné trois
 549 modèles de graphes aléatoires. Nous allons utiliser le Modèle Erdős-Renyi et le modèle de
 550 configuration.

551 Le *Réservoir* d'arêtes combine les deux modèles. Nous avons ensuite introduit la notion
 552 de composante géante et nous souhaitons savoir s'il existe des composantes géantes dans
 553 le Réservoir.

554 3 Préliminaires d'analyse du texte

555 Dans ce chapitre, on présente des techniques classiques d'analyse du texte :

- 556 • Le prétraitement qui va utiliser la décomposition grammaticale puis ensuite une recon-
557 naissance d'entité.
- 558 • La représentation de mots par vecteur de petite dimension *Word2vec*.
- 559 • La classification de documents.

560 Le prétraitement est une opération essentielle, basée sur l'apprentissage des entités.
561 La représentation d'un mot par vecteur suit l'analyse de la corrélation des occurrences de
562 deux mots dans la même phrase, ou le deuxième moment de la distribution des mots¹.

563 Étant donnée une matrice de corrélation (n, n) , ou toute matrice positive semi-définie
564 symétrique, on peut associer un vecteur v_i dimension n à valeur réelle pour $i = 1, \dots, n$, de
565 telle sorte que le produit scalaire $v_i \cdot v_j$ soit égal à $A(w_i, w_j)$. Si nous ne sélectionnons que les
566 grandes valeurs propres de A , nous pouvons obtenir des vecteurs de plus petite dimension
567 tels que $v_i \cdot v_j \simeq A(w_i, w_j)$.

568 La classification de documents est une opération classique. On va décrire deux méthodes
569 pour réaliser des ce type de tâches.

570 3.1 Prétraitement des données

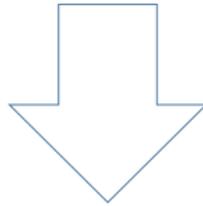
571 3.1.1 Lemmatisation et Entités Nommées

572 L'objectif de la lemmatisation est de réduire les formes flexionnelles et parfois les formes
573 dérivées d'un mot à une forme de base commune, la racine, en utilisant une décomposition
574 grammaticale. La lemmatisation consiste en une analyse morphologique des mots, visant à
575 supprimer les terminaisons inflexibles et à trouver la forme de base d'un mot. Celle-ci est
576 aussi appelée lemme qu'on peut aussi trouver dans un dictionnaire. La figure 3.1 montre
577 la lemmatisation de la phrase :

578 *When I arrived in New York, around June 10th 2018, I lived across the First National*
579 *Bank of Chicago.*

1. Le deuxième moment est la distribution des paires de mots, le troisième moment est la distribution des triplets de mots et le k -th moment est la distribution des k mots dans la même phrase.

When I arrived in New York, around June 10th 2018, I lived across the First National Bank of Chicago.



When I arrive in New-York, around June 10th 2018, I live across the First-National-Bank of Chicago.

FIGURE 3.1: Phrase lemmatisée.

580 La reconnaissance d'entités nommées (NER) est une sous-tâche d'extraction d'informa-
581 tions qui cherche à localiser et à classer les entités nommées mentionnées dans un texte dans
582 des catégories prédéfinies telles que les noms de personnes, les organisations, les lieux, etc
583 [49]. Le NER agit non seulement comme un outil autonome d'extraction d'informations
584 (IE), mais joue également un rôle essentiel dans diverses applications de traitement du
585 langage naturel (NLP) telles que le *compréhension de texte* [66], les *systèmes de question-*
586 *réponse* [24], la *traduction automatique* [5], etc.

587 La classification consiste à attribuer une étiquette telle que "nom", "lieu", "date",
588 "email", etc. Si le NER désigne au départ l'extraction de noms propres, de noms de lieux
589 et de noms d'organisations, ce concept s'est étendu à d'autres entités telles que la date,
590 l'email, le montant d'argent, etc.

591 Il existe trois grandes approches des NER :

- 592 ● basée sur un lexique
- 593 ● basée sur des règles grammaticales
- 594 ● basée sur l'apprentissage machine.

595 Un système (NER) peut combiner plusieurs de ces techniques [37]. L'approche grammati-
596 cale est utilisée dans l'étiquetage morphosyntaxique où Part-of-Speech (POS) Tagging ou
597 on essaye d'attribuer une étiquette à chaque mot d'une phrase mentionnant la fonctionna-
598 lité grammaticale d'un mot (Nom propre, adjectif, déterminant...).

599 Le tableau 3.1 présente les principaux outils utilisés pour le marquage des (NER), le
600 *SpaCy* et le *NLTK* sont basés sur des approches statistiques. Alors que le NER du *CoreNLP*
601 NER de Stanford utilise un processus d'apprentissage à l'aide d'un réseau de neurones qui
602 nécessite une grande masse de données.

603 Prenons l'exemple suivant : " *When I arrived in New York, around June 10th 2018,*
604 *I lived across the First National Bank of Chicago.*". La figure 3.2 montre la classification

Library Name
SpaCy
CoreNLP
NLTK

TABLE 3.1: Outils pour la reconnaissance des entités

605 des entités nommées de cet exemple avec l’outil de Stanford et l’utilisation de la librairie
606 SpaCy.

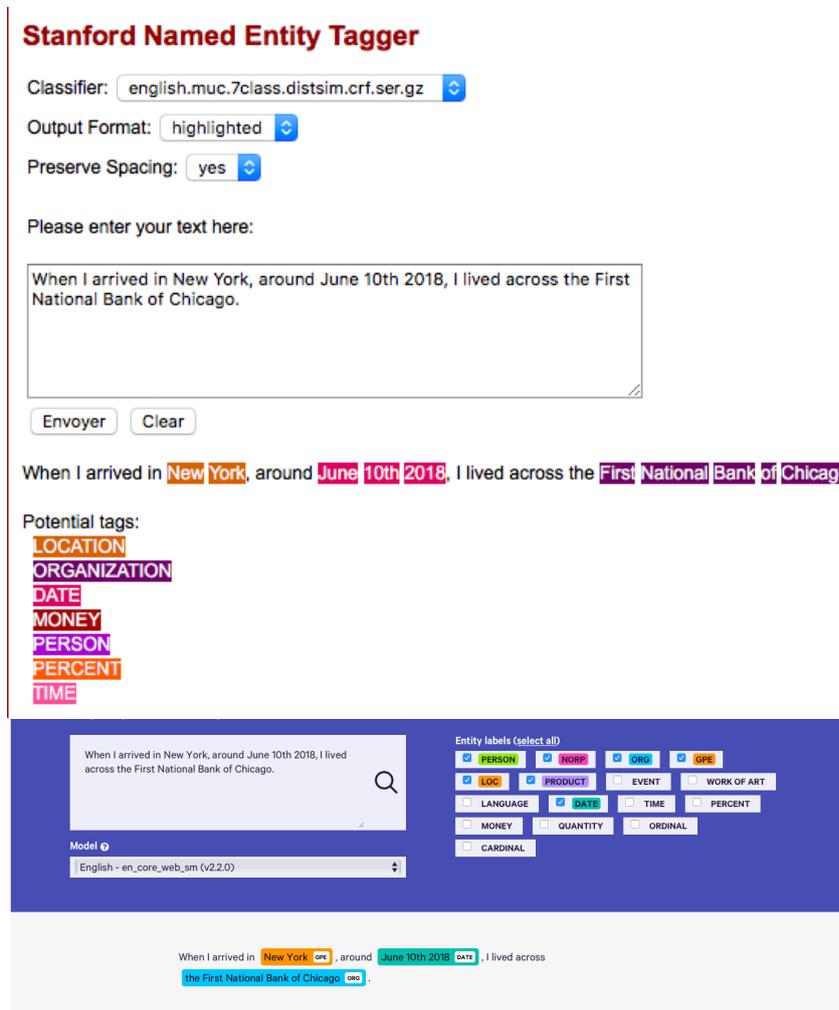


FIGURE 3.2: Outil NER de Stanford avec la librairie SpaCy.

607 Ces opérations ne sont pas parfaites : en particulier certains systèmes vont reconnaître
608 *National Bank of Chicago* comme entité et considérer *First* comme un adjectif.

609 3.1.1.1 Décomposition grammaticale classique

610 Une des représentations grammaticales classiques est celle des grammaires hors-contexte
611 (CFG, Context-Free Grammars) largement utilisée pour la description des langages de
612 programmation.

613 La figure 3.3 montre une grammaire hors-contexte très simple. Dans ce cas, l'ensemble
614 de non-terminaux spécifie quelques catégories syntaxiques de base : par exemple *S* signifie
615 "phrase", *NP* "expression nominale", *VP* "expression verbale", etc. Σ contient l'ensemble

$N = \{S, NP, VP, PP, DT, Vi, Vt, NN, IN\}$
 $S = S$
 $\Sigma = \{\text{sleeps, saw, man, woman, dog, telescope, the, with, in}\}$

$R =$

S	\rightarrow	NP	VP
VP	\rightarrow	Vi	
VP	\rightarrow	Vt	NP
VP	\rightarrow	VP	PP
NP	\rightarrow	DT	NN
NP	\rightarrow	NP	PP
PP	\rightarrow	IN	NP

Vi	\rightarrow	sleeps
Vt	\rightarrow	saw
NN	\rightarrow	man
NN	\rightarrow	woman
NN	\rightarrow	telescope
NN	\rightarrow	dog
DT	\rightarrow	the
IN	\rightarrow	with
IN	\rightarrow	in

FIGURE 3.3: Un simple exemple d'une Grammaires hors-contexte, l'ensemble des non-terminaux N contient un ensemble de base de catégories syntaxiques : S = Sentence, VP = Verb Phrase, NP = Noun Phrase, PP = Prepositional Phrase, DT = Determiner, Vi = Intransitive Verb, Vt = Transitive Verb, NN = Noun, IN = Preposition. Σ est l'ensemble des mots possibles dans la langue.

616 des mots du vocabulaire. Le symbole de départ dans cette grammaire est S : cela spécifie
 617 que chaque arbre d'analyse a S comme racine.

618 Une décomposition est un arbre syntaxique qui est le témoin du fait que la phrase
 619 est grammaticalement correcte. Dans le cas d'une grammaire context-free on l'obtient
 620 avec l'algorithme *CKY* (Cocke–Kasami–Younger) [36] qui permet de trouver l'arbre de
 621 décomposition en $O(n^3)$ en utilisant la programmation dynamique.

622 3.1.1.2 La décomposition par dépendances

623 La décomposition par dépendances [57] procède différemment de l'algorithme *CKY*
 624 et décompose la phrase en deux parties : la partie gauche qui est la pile, la partie droite
 625 qui est le buffer. Le but de la décomposition par dépendances est de trouver un arbre de
 626 décomposition. À chaque étape, l'algorithme de décomposition lit un nouveau mot et doit
 627 décider une des trois actions possibles :

- 628 • Empiler le mot.
- 629 • Trouver les dépendances [57] à l'intérieur des mots de la pile .

- 630 • Certains mots vont devenir des mots importants qu'on va rattacher à la racine, par-
631 exemple les verbes.

632 L'algorithme de décomposition apprend à partir d'exemples l'action de décomposition
633 pour un mot donné. Un exemple est un arbre de décomposition associé à une phrase avec
634 des différentes dépendances, comme l'indique le treebanks de la figure 3.4. Le treebank
635 décrit des milliers d'exemples de cette nature.

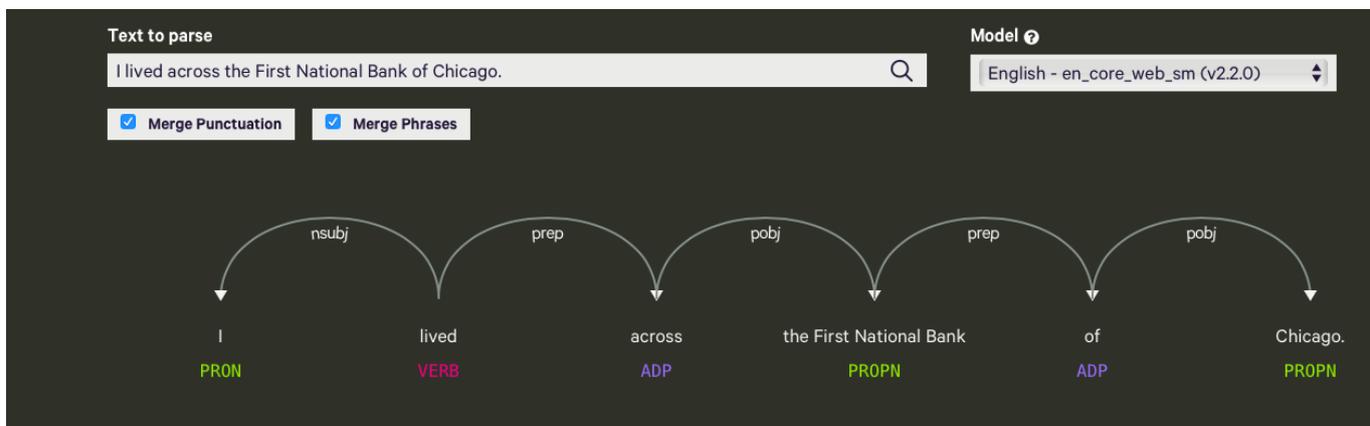


FIGURE 3.4: Treebank

636 3.2 Vectorisations des mots : *Word2vec*

637 Les représentations vectorielles des mots (*word embeddings*) essaient de saisir les rela-
638 tions entre les mots sous forme de distance ou d'angle, et ont de nombreuses applications en
639 traitement du langage naturel et en apprentissage machine. Les méthodes les plus simples
640 utilisent des vecteurs de mots obtenus directement à partir de la matrice de co-occurrence.
641 Les heuristiques de repondération sont connues pour améliorer ces méthodes, tout comme
642 la réduction de dimension [23]. Certaines méthodes de repondération sont non linéaires,
643 notamment la prise en compte de la racine carrée des nombres de co-occurrence [53], ou le
644 logarithme, ou les informations ponctuelles *PMI* (Pointwise Mutual Information) qui s'y
645 rapportent [20].

646 Les réseaux neuronaux [[55], [54], [7], [21]] proposent une autre façon de construire
647 des *embeddings*. Le vecteur associé à un mot est simplement la représentation interne du
648 réseau neuronal lorsque l'entrée est le mot donné. Cette méthode est non linéaire et non
649 convexe. Elle a été popularisée par *Word2vec*, une famille de modèles basés sur l'énergie
650 [[44], [45], [46]], et par une approche de factorisation de matrice appelée *GloVe* [51]. Le
651 premier article a également montré comment résoudre des analogies en utilisant l'algèbre
652 linéaire sur les *embeddings*. Les expériences et la théorie ont été utilisées pour suggérer

653 que ces nouvelles méthodes sont liées aux anciens modèles basés sur le PMI, mais avec de
654 nouveaux hyperparamètres et/ou des méthodes de repondération des termes [39].

655 Les chercheurs en traitement du langage utilisent des approximations de faible di-
656 mension. Les méthodes *Word2vec* et *Glove* au départ, sont des méthodes expérimentales.
657 L'article [4] fournit une justification théorique aux modèles non linéaires comme *PMI*,
658 *Word2vec* et *GloVe*, ainsi qu'à certains choix d'hyperparamètres.

659 Étant donnée une matrice de corrélation (n, n) , ou toute matrice positive semi-définie
660 symétrique, on peut associer un vecteur v_i de dimension n pour $i = 1, \dots, n$, de telle sorte
661 que le produit scalaire $v_i \cdot v_j$ est égal à $A(w_i, w_j)$. Si nous ne sélectionnons que les grandes
662 valeurs propres de A , nous pouvons obtenir des vecteurs de plus petite dimension tels que
663 $v_i \cdot v_j \simeq A(w_i, w_j)$. On peut considérer les coefficients de la matrice de corrélation comme
664 les probabilités d'observer une paire de mots dans une phrase tirée aléatoirement dans ce
665 document puis ensuite en prenant deux mots tirés aléatoirement dans cette phrase. Les
666 valeurs diagonales de la matrice donnent le premier moment de cette distribution.

667 Pour $n = 10^4$, cette méthode n'est pas réaliste et on doit utiliser des méthodes plus
668 efficaces telles que *Word2vec*, qui est utilisée pour obtenir des vecteurs de dimension
669 300 lorsque le dictionnaire comporte $n = 10^4$ mots, c'est-à-dire réduire la dimension.
670 Considérons un corpus de 10000 mots, et dans ce corpus les deux mots suivants (ants,
671 abandon). Chaque mot est codé par son vecteur one-hot²(ici des matrices de colonnes de
10 000 lignes) comme le montre la figure 3.5.

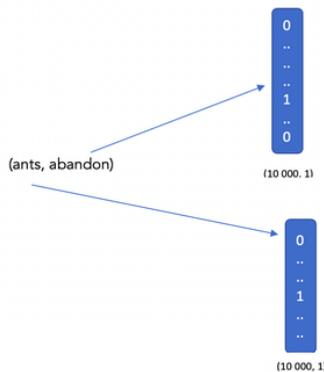


FIGURE 3.5: Vecteur one hot

672 Un réseau de neurones récurrents RNN observe des fenêtres de M mots dans une fenêtre
673 coulissante. Les variables latentes du réseau définissent la matrice de la figure 3.6.
674

2.

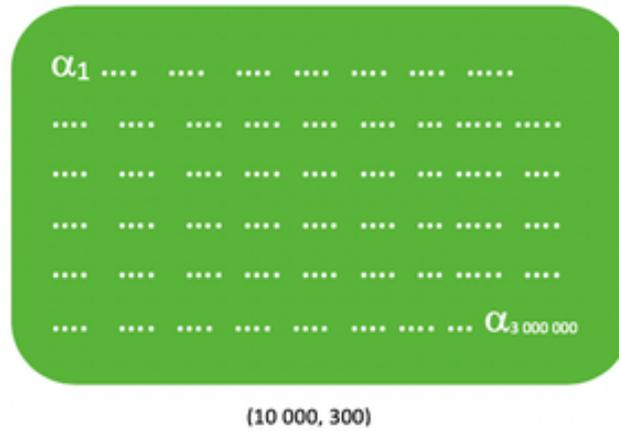


FIGURE 3.6: Matrice 10 000 lignes et 300 colonnes

675 Cependant, après la transformation du couple de mots en un vecteur one hot, la figure 3.7 montre comment la couche cachée sera transformée. La similarité entre deux mots est

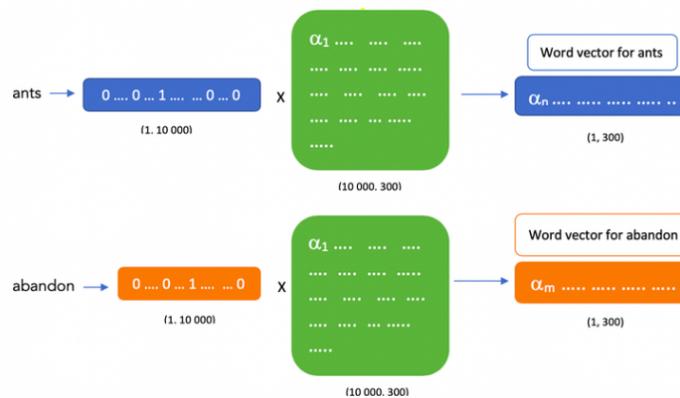


FIGURE 3.7: Matrice 10 000 lignes et 300 colonnes

676
677 le produit scalaire des vecteurs de cette matrice de la figure 3.8.

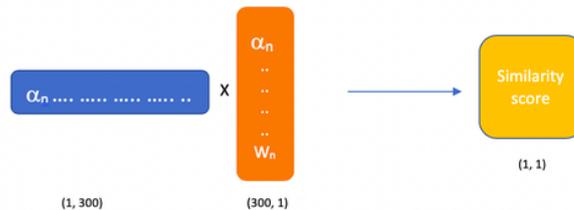


FIGURE 3.8: Vecteur de similarité

678 3.3 Classification

679 La classification de documents consiste à définir un nombre fixe de classes, en général
680 un paramètre k de manière que chaque document soit proche d'une des classes.

681 3.3.1 IRaMuTeQ

682 IRaMuTeQ [52](Interface de R pour les Analyses Multidimensionnelles de Textes et
683 de Questionnaires) présente des analyses statistiques des matrices associées à un ensemble
684 de documents. Soit M la matrice avec i lignes, où i est l'index du i -ème documents et j
685 colonnes où j est l'index du j -ème mots, $M(i, j)$ est le nombre d'occurrences du mot j dans
686 le document i . IRaMuTeQ permet de faire différents types d'analyse : statistiques tex-
687 tuelles classiques, spécificités des groupes, classification hiérarchique descendante, analyse
688 de similarité de mots et de documents. La méthode IRaMuTeQ [52] permet des classifi-
689 cations hiérarchiques basées sur l'analyse en composantes principales (PCA) qui remonte
690 aux années 60s. En pratique on peut utiliser la décomposition SVD (Singular Value De-
691 composition) d'une matrice (n, n) qui a cependant une complexité temporelle en $O(n^3)$.

692 L'inconvénient de cette méthode est le passage à l'échelle quand n est très grand. Pour
693 l'analyse de telles données, d'autres méthodes sont apparues telles que *LDA*.

694 3.3.2 LDA

695 Un modèle thématique, par exemple *LDA* (Latent Dirichlet Allocation) [10] apprend
696 une représentation de faible dimension des documents comme un mélange de variables
697 latentes appelées sujets (*Topics*). Les sujets sont des distributions multinomiales sur un
698 vocabulaire prédéfini de mots. Puisque les mots sont observés, les distributions des do-
699 cuments et des sujets, respectivement θ et ϕ sont conditionnellement indépendantes. En
700 outre, les documents ne sont pas directement liés aux mots. Cette relation est plutôt régie

Symbol	DESCRIPTION
D	total number of documents
k	number of topics
V	total number of unique words
w	the unique word associated with the i th token
Z	the topic associated with w

TABLE 3.2: Notation utilisée dans le cadre de LDA

701 par des variables latentes supplémentaires, z , introduites pour représenter la responsabi-
 702 lité d'un sujet particulier dans l'utilisation de ce mot dans le document, c'est-à-dire le(s)
 703 sujet(s) sur lequel (lesquels) le document est focalisé. Un glossaire des notations utilisées
 704 dans le document est donné dans le tableau 3.2.

705 Pour compléter son processus de génération des documents, le modèle *LDA* considère
 706 pour les distributions de documents sur les sujets et les distributions de sujets sur les mots
 707 des distributions de *Dirichlet priors sparse* c'est-à-dire chaque document a peu de topics
 708 et le support de distribution est petit. Le *LDA* est basé sur l'hypothèse de l'échangeabilité
 709 des mots dans un document et des documents dans un corpus.

710 Dans le cadre de *LDA*, les documents sont supposés être générés de manière aléatoire
 711 à partir d'un ou plusieurs sujets, chacun d'entre eux étant une distribution de mots. Les
 712 thèmes sont considérés comme des variables latentes, et *LDA* s'exécute en déduisant les
 713 thèmes des documents par un processus de *Dirichlet*. L'algorithme échantillonne les docu-
 714 ments à plusieurs reprises et modifie les thèmes pour mieux les adapter jusqu'à atteindre
 715 une convergence spécifiée.

716 3.4 Mécanismes d'Attention

717 Les mécanismes d'attention [40, 60] prévoient pour un mot w_i dans une phrase, la
 718 distribution des mots les plus corrélés. La corrélation d'un mot w_j avec w_i est approxima-
 719 tivement la valeur $v(w_i).v(w_j)$. On peut alors calculer $\{v(w_i).v(w_j) : j \neq i\}$ pour un mot w_i
 720 fixe et normaliser les valeurs pour obtenir une distribution. Ce mécanisme est principale-
 721 ment utilisé dans les transformateurs [63] pour la traduction automatique. Nous montrons
 722 dans la section 6.5 que le Réservoir pondéré peut reconstruire les distributions d'attention
 723 les plus significatives. Nous définissons les *distributions les plus significatives* comme les
 724 distributions des nœuds centraux dans un cluster.

725 L'analyse de l'attention de la phrase :

726 *On September 11th, 2001, the attack on the World Trade Center had a huge impact on*
 727 *American politics.*

728 pour les entités *September 11th, 2001* et *the World Trade Center* est donnée dans la figure
 729 3.9.

730 3.4.1 Conclusion

731 Dans ce chapitre, nous avons présenté des techniques utilisées pour le prétraitement
732 de données que nous allons ensuite utiliser comme des boites noires. Nous utiliserons la
733 méthode *Word2vec* dans le chapitre ??, pour analyser des données de texte en streaming.
734 Nous avons présenté deux méthodes classiques de classification pour les comparer avec
735 notre méthode dans le chapitre ?. Les mécanismes d'Attention sont directement liés aux
736 méthodes d'échantillonnages que nous allons utiliser.

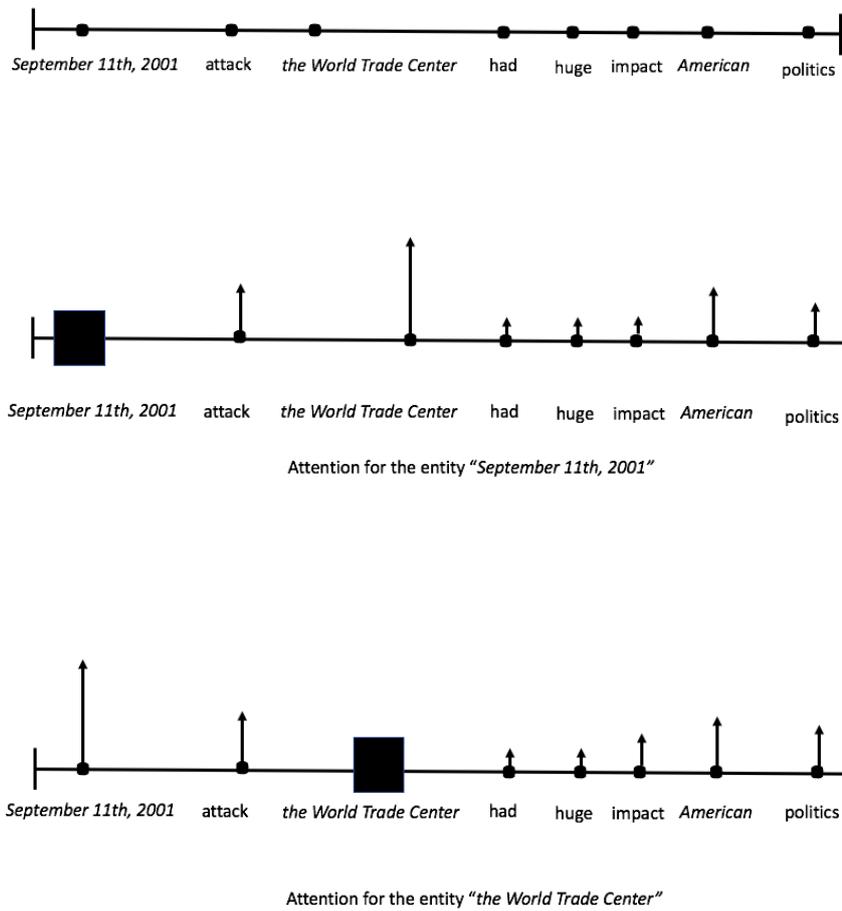


FIGURE 3.9: L'attention pour un mot spécifique.

737 4 Analyse OLAP sur des flux de 738 données

739 Un entrepôt de données selon Bill Inmon [35], est une collection non volatile et variable
740 dans le temps, de données orientées vers un sujet et intégrées à partir de sources de données.
741 Chacun des termes est défini comme suit :

- 742 • Orienté sujet. Dans l'entrepôt de données, les données sont organisées par thème. Les
743 données propres à un thème, les ventes par exemple, sont rapatriées de différentes bases.
- 744 • Intégré. Les données proviennent de sources hétérogènes utilisant chacune un type de
745 format. Elles sont modifiées avant d'être intégrées dans l'entrepôt.
- 746 • Non volatile. Les données ne disparaissent pas et ne changent pas au fil des traitements,
747 au fil du temps.
- 748 • Variable dans le temps. Les données non volatiles sont horodatées. On peut ainsi visualiser
749 l'évolution dans le temps d'une valeur donnée. Le degré de détail de l'archivage est
750 relatif à la nature des données. Toutes les données ne méritent pas d'être archivées.

751 On peut considérer l'entrepôt de données avec deux points de vue différents :

- 752 • Un flux de tuples $t_1, t_2, \dots, t_n, \dots$, et on ne stocke pas tous les tuples.
- 753 • On stocke tous les tuples, pour former une relation.

754 Si on suit la première approche, on s'intéressera à l'approximation des requêtes d'analyse
755 et le résultat ne sera qu'approché. Si on suit la deuxième approche, l'entrepôt de données
756 est une table d'une base de données qui suit un modèle relationnel et permet de réaliser des
757 analyses. L'entrepôt de données est alors considéré comme une table où on peut seulement
758 ajouter des tuples.

759 Un schéma d'analyse ou un schéma OLAP spécifie des dimensions d'analyse et des
760 *Mesures* associées aux tuples de l'entrepôt.

761 Dans l'analyse OLAP classique, nous stockons l'ensemble de l'entrepôt de données et
762 une requête OLAP est définie par un filtre, d dimensions et une *Mesure* parmi un ensemble
763 fini des *Mesures* M_1, \dots, M_l . Les *préférences* sont des nouvelles *Mesures* considérées comme
764 une distribution sur M_1, \dots, M_l .

765 La première section introduit les principaux concepts : entrepôts de données, schémas
766 OLAP. La deuxième section détaille l'approximation des requêtes OLAP en streaming. La
767 troisième section étudie la gestion des préférences dans ce cadre.

768 4.1 Entrepôts de données et requêtes OLAP

769 Considérons un schéma relationnel avec trois relations (Achat, Clients, Produits) comme
770 indiqué sur la figure 4.1. Un schéma OLAP 4.2 ou schéma étoile pour ce schéma relationnel,
771 spécifie la table Achat comme table d'analyse et décrit un arbre où chaque nœud est
772 étiqueté soit par un nom de table, soit par un ensemble d'attributs. La racine est étiquetée
773 par la table d'analyse qui peut aussi être représentée par la liste de tous ses attributs. Une
774 arête de l'arbre existe s'il existe une dépendance fonctionnelle entre les attributs du nœud
775 d'origine et les attributs du nœud d'extrémité pour les tables du schéma relationnel 4.1.
776 Les *Mesures* sont des attributs numériques de la table Achat. Dans l'exemple de la figure
777 4.1, les trois attributs numériques sont : *Speed*, *Quantity*, *Price*.

Achat(ID_c, ID_p, Day, Speed, Quantity, Price)
Clients(ID_c, name, age, gender, csp)
Produits(ID_p, article, type, artiste, price)

FIGURE 4.1: Schéma relationnel

778 Un Entrepôt de données, noté I , est une table particulière stockant des tuples t avec
779 les attributs $A_1, \dots, A_m, M_1, \dots, M_l$, certains A_i étant des clés étrangères d'autres tables, et
780 M_1, \dots, M_l les *Mesures*. Des tables de dimensions fournissent d'autres attributs liés à la table
781 d'analyse par des clés.

782 Une requête OLAP pour un schéma S est déterminée par :

- 783 • Une condition de filtre qui sélectionne un sous-ensemble des tuples de l'entrepôt.
- 784 • Une *Mesure* et un opérateur d'agrégation (COUNT, SUM, AVG, ...).
- 785 • d dimensions qui sont représentées par des nœuds de l'arbre. Une dimension peut aussi
786 être appelée dans certains contextes comme un classificateur.

Une *Mesure* est un attribut numérique, et on considère dans la suite l'opérateur d'agrégation comme la somme (SUM). La réponse à une requête Q est un tableau multidimensionnel, selon les dimensions C_1, \dots, C_d et la *Mesure* $t.M$. Chaque tuple c_1, \dots, c_d, m de la réponse, où $c_i \in C_i$, est tel que :

$$m = \frac{\sum_{t:t.C_1=c_1, \dots, t.C_d=c_d} t.M}{\sum_{t \in I} t.M}.$$

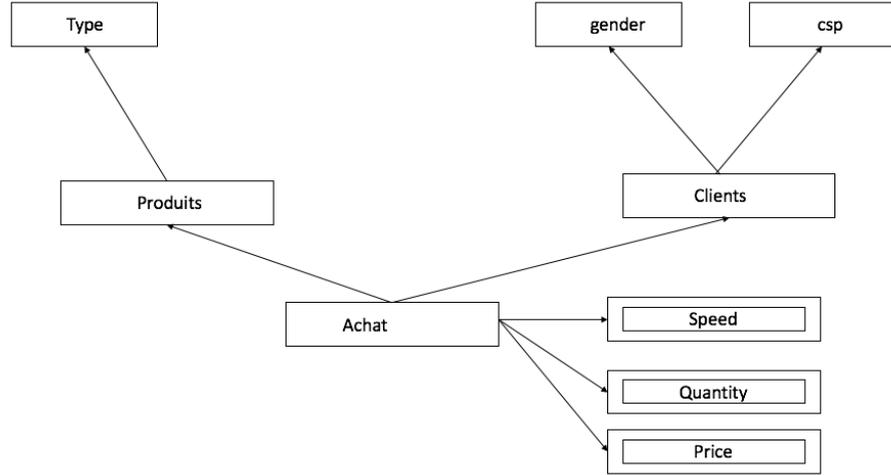


FIGURE 4.2: Schéma OLAP.

787 Nous considérons m , la valeur relative de la réponse alors que le numérateur est la valeur
 788 absolue de la réponse et nous écrivons Q_C^I comme la distribution ou *vecteur de densité*
 789 pour la réponse à Q sur les dimensions C_1, \dots, C_d et sur l'entrepôt de données I . Prenons
 790 par exemple Q_{gender} et $Q_{gender,type}$ comme des requêtes sans filtre, la dimension *gender* pour
 791 la première et (*gender, type*) pour la deuxième, la *Mesure* Prix et l'agrégation SUM.

792 La réponse à la requête Q est indiquée à la figure 4.3, considérée comme une distribution.
 793 *Gender* est un attribut dont les valeurs possibles sont $\{homme, femme\}$, et *type* est un
 794 attribut dont les valeurs possibles sont $\{pop, rock, classique\}$.

795 4.2 Approximation à partir d'échantillons

796 Dans notre contexte, nous approximations une distribution Q . Il y a généralement deux
 797 paramètres $0 \leq \varepsilon, \delta \leq 1$ pour l'approximation des algorithmes randomisés, où ε est l'erreur,
 798 et $1 - \delta$ la confiance. Considérons la réponse Q à une requête OLAP comme une distribution
 799 sur un support fini, et supposons qu'un algorithme probabiliste A calcule \hat{Q} . La norme L_1
 800 appliquée à $|Q - \hat{Q}|_1 = \sum_i |Q[i] - \hat{Q}[i]|$. L'algorithme (ε, δ) -approxime Q si :

$$Prob_{\Omega}[|Q - \hat{Q}| \leq \varepsilon] \geq 1 - \delta$$

801 L'espace probabiliste Ω est déterminé par les choix de l'algorithme A . Une approche
 802 classique consiste à prendre k échantillons aléatoires, où la probabilité d'un échantillon t
 803 est proportionnelle à sa *Mesure t.M*. La réponse $\hat{Q}[i]$ est la proportion d'échantillons où
 804 la valeur de l'attribut est i . Elle peut également être considérée comme la même requête

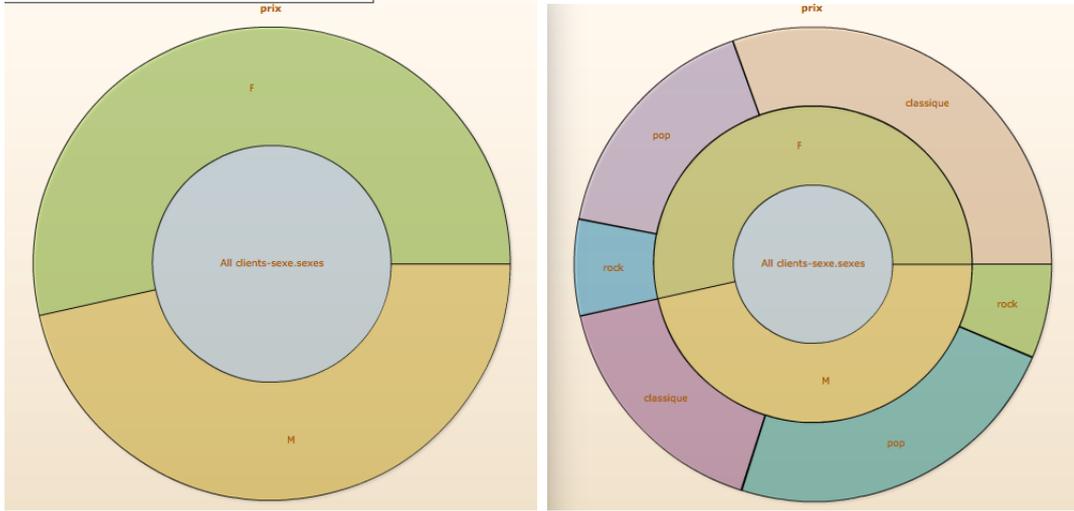


FIGURE 4.3: Une requête OLAP de dimension 1 et une requête de dimension deux.

805 OLAP sur un nouveau entrepôt de données construit avec k échantillons où la *Mesure*
806 originale M est remplacée par la constante 1.

807 Pour un entrepôt de données I , soit \widehat{I}_M un ensemble d'échantillons obtenus à l'aide de
808 l'échantillonnage pondéré sur une *Mesure* M , comme représenté dans la figure 4.5.

Algorithme 3 Algorithme pour calculer \widehat{Q}_C

Entrée : I

- Construire \widehat{I}_M ,
 - Pour toute valeur $c(C)$, la valeur $\widehat{Q}_{C=c}$ est le nombre de tuples \widehat{I}_M tels que $C = c$
 - \widehat{Q}_C est la distribution ainsi obtenue pour tout c .
-

809 Remarquons que \widehat{Q}_C est la réponse à la même requête OLAP sur \widehat{I}_M , où la nouvelle
810 mesure est une constante 1.

Soit $|C|$ le nombre de valeurs possibles c de l'attribut C , c'est-à-dire le support de C que nous supposons fini, et soit $1_{C=c}(t)$ la fonction caractéristique d'un tuple t pour le critère $C = c$,

$$1_{C=c}(t) = \begin{cases} 1 & \text{si } t.C = c \\ 0 & \text{sinon.} \end{cases}$$

811 Par hypothèse, le Réservoir R de taille k est exactement \widehat{I}_M et

$$Prob[t \in \widehat{I}_M] = \frac{k \cdot t.M}{T}$$

812 comme on a vu dans l'étude du chapitre 2.1.2.

813 Comme il a été montré dans [22], \widehat{Q}_C est une bonne approximation de Q_C . Pour chaque
814 valeur possible c de l'attribut C , l'espérance de $\widehat{Q}_{C=c}$ est le nombre attendu de tuples t tel
815 que $t.C = c$, est précisément $Q_{C=c}$.

816 **Lemme 3.** *L'espérance de \widehat{Q} , écrite $\mathbb{E}(\widehat{Q})$ est la valeur exacte $Q_{C=c}$.*

Démonstration. Supposons que nous ayons qu'une seule dimension C , pour simplifier. Considérons la densité pour $C = c$ et son espérance :

$$\begin{aligned} \mathbb{E}(\widehat{Q}_{C=c}) &= \frac{\sum_t \text{Prob}[t \in \widehat{I}_M] \cdot 1_{C=c}(t)}{k} = \\ &= \frac{\sum_{t: t.C=c} \text{Prob}[t \in \widehat{I}_M]}{k} = \\ &= \frac{\sum_{t: t.C=c} \frac{k.t.M}{T}}{k} = \frac{\sum_{t: t.C=c} t.M}{T} = Q_{C=c} \end{aligned}$$

817

□

Pour la distribution Q_C , comme vecteur, nous avons $\mathbb{E}(\widehat{Q}_C) = Q_C$. L'un des objectifs est d'obtenir un algorithme d'approximation A qui produit une distribution $A(Q_C) = \widehat{Q}_C$ de telle sorte que

$$\text{Prob}[\| Q_C - \widehat{Q}_C \|_1 \leq \varepsilon] \geq 1 - \delta$$

818 où nous prenons la norme L_1 entre les deux vecteurs. Elle peut être obtenue si nous prenons
819 plusieurs Réservoirs indépendants et prenons les valeurs moyennes de chaque composante.
820 La borne de *Chernoff-Hoeffding* [33] permettra d'obtenir l'inégalité ci-dessus.

Soit X_1, \dots, X_n des variables aléatoires indépendantes délimitées par l'intervalle : $0 \leq X_i \leq 1$. La moyenne empirique de ces variables est :

$$\bar{X} = \frac{1}{n}(X_1 + \dots + X_n).$$

L'une des inégalités de *Hoeffding* est :

$$P(\bar{X} - \mathbb{E}[\bar{X}] \geq t) \leq e^{-2nt^2}$$

821 La déviation à la moyenne décroît exponentiellement en fonction du nombre de tirages n .
822 Avec un seul Réservoir, l'algorithme A qui produit \widehat{Q}_C peut aussi être analysé avec
823 l'inégalité de Markov¹, mais l'approximation n'est pas aussi bonne :

824 **Lemme 4.** *Pour toutes les valeurs c , il existe ε, δ tel que $\text{Prob}[\| Q_{C=c} - \widehat{Q}_{C=c} \|_1 \leq \varepsilon] \geq$
825 $1 - \delta$.*

1. Pour $a \geq 0$, $\text{Prob}[\widehat{X} \geq a \cdot \mathbb{E}(X)] \leq 1/a$

Démonstration. Pour $a \geq 1$, $Prob[\widehat{Q}_{C=c} \geq a \cdot \mathbb{E}(\widehat{Q}_{C=c})] \leq 1/a$ qui peut également être réécrite comme :

$$Prob[|\widehat{Q}_{C=c} - Q_{C=c}| \geq (a-1) \cdot Q_{C=c}] \leq 1/a$$

Nous fixons $\varepsilon = (a-1) \cdot Q_{C=c}$, donc $a = 1 + \varepsilon/Q_{C=c}$ et $\delta = 1/a \simeq 1 - \varepsilon/Q_{C=c}$:

$$Prob[|\widehat{Q}_{C=c} - Q_{C=c}| \geq \varepsilon] \leq \delta$$

$$Prob[|\widehat{Q}_{C=c} - Q_{C=c}| \leq \varepsilon] \geq 1 - \delta$$

826

□

827 4.2.1 Construction des échantillons à partir d'un entrepôt de données

828 Un des buts, est d'approximer une requête OLAP à partir des échantillons pris avec
829 une probabilité proportionnelle à une *Mesure*, comme par exemple *Speed*, qui représente la
830 vitesse de la transaction ou l'inverse du temps de transaction pour l'entrepôt de données
831 Achat. Comment réaliser un tel tirage ?

832 Nous modifions le schéma de la table Achat en ajoutant une colonne qui représente
833 le poids cumulatif de la *Mesure Speed*. Le nouvel attribut cumulatif est *Cumulative speed*
834 comme le montre la figure 4.4 en partant de la valeur cumulative de 0.

835 Dans l'exemple de la figure 4.4, la valeur maximum de l'attribut *Cumulative speed* est
836 de 5200. La procédure est alors la suivante :

- 837 • tirer une valeur uniforme λ entre 0 et 5200.
- 838 • trouver un tuple t tel que sa valeur cumulative est strictement plus grande ou égale à λ
839 et le tuple précédent à une valeur cumulative plus petite que λ .

840 Dans ce cas, la probabilité de choisir tuple t est proportionnelle à sa *Mesure Speed*. Nous
841 procédons par dichotomie, en comparant λ à la valeur cumulative du tuple en position $n/2$.
842 Si λ est supérieur, on isole la partie la plus récente et on recommence. Si λ est inférieur,
843 on sélectionne la partie la plus ancienne et on recommence.

Datawarehouse

ID_C	ID_P	Day	Speed	Cumulative speed
10	100	Monday	1000	1000
20	200	Monday	300	1300
10	300	Tuesday	800	2100
20	100	Tuesday	500	2700
30	200	Tuesday	1000	3700
10	100	Sunday	700	4400
20	200	Sunday	100	4500
30	300	Sunday	700	5200

FIGURE 4.4: Entrepôt de données.

844 Pour approximer une requête OLAP, il faut appliquer l'algorithme ?? décrit dans la
845 section 4.2.

846 4.2.2 Approximation des requêtes OLAP sur les données de streaming

847 Supposons que les tuples de l'entrepôt de données arrivent comme un flux $t_1, t_2, \dots, t_i, \dots$
848 Nous appliquons la méthode d'échantillonnage pondéré, détaillée dans le chapitre 2.1.2.

849 La principale propriété est que chaque tuple t_i est pris proportionnellement à sa *Mesure*.
850 C'est un argument simple classique que nous rappelons dans la section 2.1.2. Une fois le
851 Réservoir construit, on peut approximer des requêtes OLAP comme dans la section 4.2,
852 sans stocker tous les tuples.

853 4.3 Préférences pour l'analyse OLAP

854 Pour l'analyse de données, les préférences englobent une grande variété de possibilités.
855 Elles peuvent être définies comme des conditions logiques, telles que la géolocalisation, mais
856 peuvent aussi être des profils globaux représentés comme une distribution sur un support
857 fini. Un profil Netflix est une telle distribution sur 250 catégories. Il s'agit généralement
858 d'un vecteur sparse α de dimensions 250 tel que $\sum_i \alpha_i = 1$. Une préférence générale est une
859 combinaison de ces deux composantes : les contraintes logiques et des profils de distribution.
860 Ces deux composantes ont des impacts très différents sur l'analyse des données.

Supposons un schéma OLAP, avec plusieurs *Mesures* M_1, M_2, M_3 comme dans la figure 4.2. Une préférence pour le schéma OLAP est une distribution sur le support M_1, M_2, M_3 .

La figure 4.5 décrit l'entrepôt de données Achat avec deux clés étrangères ID_C pour la table $Clients$ et ID_P pour la table $Produits$. Il y a 3 *Mesures* et une préférence est une distribution telle que :

$$M_\alpha = \alpha \cdot \text{Speed} + \beta \cdot \text{Quantity} + (1 - \alpha - \beta) \cdot \text{Price}$$

861 pour certains $0 \leq \alpha, \beta \leq 1$.

862 Un profil de préférence est une distribution sur les *Mesures* possibles du tuple t . Ima-
863 ginons le scénario dans lequel nous avons un flux de tuples t et nous voulons analyser
864 rapidement les données avec des requêtes OLAP et la nouvelle *Mesure* M_α . Nous construi-
865 sons autant de Réservoirs que de *Mesures*. Le profil est inconnu lorsque nous prélevons
866 les échantillons. Nous pouvons construire un nouveau Réservoir à partir des échantillons
867 existants et l'utiliser pour approximer les requêtes OLAP pour la nouvelle *Mesure*.

868 Pour une taille du Réservoir k donnée et le schéma OLAP de la figure 4.2, nous construi-
869 sons trois Réservoirs indépendants, R_1 , R_2 et R_3 avec k échantillons pour chacune des *Me-*
870 *sures*. La construction suit l'échantillonnage k -pondéré, comme le montre la figure 4.5. Si
871 une préférence M_α est définie à tout moment, *après avoir prélevé les échantillons*, nous
872 pouvons construire un nouveau Réservoir R_α de taille k . Nous montrons que l'évaluation
873 de la requête Q sur R_α est une approximation de Q , comme dans la section 4.2.2.

874

875 **Construction d'un Réservoir R_α**

- 876 • Prenons $\alpha \cdot k$ échantillons uniformes de R_1 ,
- 877 • Prenons $\beta \cdot k$ échantillons uniformes de R_2 ,
- 878 • Prenons $(1 - \alpha - \beta) \cdot k$ uniform samples from R_3 .

Nous montrons que l'évaluation de la requête Q sur R_α est une approximation de Q .
Nous supposons pour simplifier que $p = 2$. Une préférence s'écrit comme :

$$M_\alpha = \alpha \cdot M_1 + (1 - \alpha) \cdot M_2$$

879 pour $0 \leq \alpha \leq 1$.

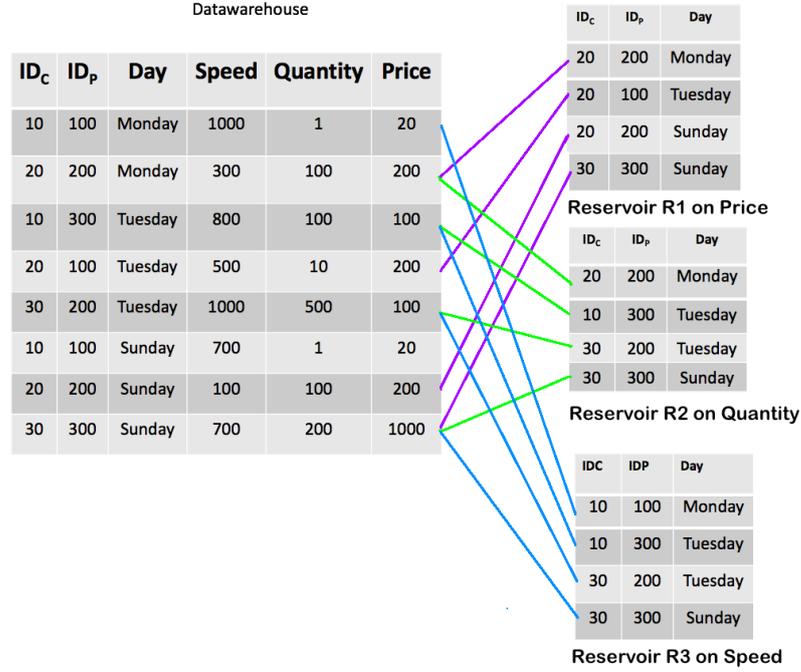


FIGURE 4.5: Échantillons provenant d'un entrepôt de données pour les trois *Mesures*.

880 4.3.1 Préférences comme distribution sur les *Mesures*

881 Soit Q_C^1 la distribution obtenue par l'analyse sur la dimension C pour la première
 882 *Mesure* M_1 et \widehat{Q}_C^1 l'approximation sur le Réservoir R_1 . De la même façon, Q_C^2 est la
 883 distribution obtenue par l'analyse sur la dimension C pour la deuxième *Mesure* M_2 et \widehat{Q}_C^2
 884 est l'approximation sur le Réservoir R_2 .

885 Considérons une nouvelle *Mesure* M_α , définie à tout moment, et le Réservoir R_α défini
 886 dans la section 4.3. Soit \widehat{Q}_C^α la réponse de Q_C sur R_α et montrons le lien entre ces deux
 887 valeurs.

888 **Theorem 1.** Soit Q_C^α la réponse de Q sur la dimension C pour la *Mesure* M_α . L'espérance
 889 de \widehat{Q}_C^α est Q_C^α .

890 *Démonstration.* Evaluons $\mathbb{E}(\widehat{Q}_{C=c}^\alpha)$, l'espérance de la densité pour $C = c$ sur le Réservoir
 891 R_α . C'est le nombre d'échantillons prévu avec $C = c$ divisé par k le nombre total d'échantillons.
 892 Le nombre attendu d'échantillons est :

$$\sum_t 1_{C=c}(t) \cdot Prob[t \in R_\alpha] = \sum_{t:t.C=c} \left(\alpha \cdot \frac{k.t.M_1}{T} + (1 - \alpha) \cdot \frac{k.t.M_2}{T} \right)$$

893 car chaque t tel que $C = c$ dans R_α est pris avec probabilité $(\alpha \cdot \frac{k.t.M_1}{T} + (1 - \alpha) \cdot \frac{k.t.M_2}{T})$.

894 Par conséquent :

$$\mathbb{E}(\widehat{Q}_{C=c}^\alpha) = \frac{\sum_{t:C=c} (\alpha \cdot \frac{k.t.M_1}{T} + (1 - \alpha) \cdot \frac{k.t.M_2}{T})}{k}$$

$$\begin{aligned} \mathbb{E}(\widehat{Q}_{C=c}^\alpha) &= \frac{\sum_{t:C=c} (\alpha.t.M_1 + (1 - \alpha)t.M_2)}{T} \\ &= Q_{C=c}^\alpha \end{aligned}$$

895 L'espérance de la densité $\widehat{Q}_{C=c}^\alpha$ est précisément $Q_{C=c}^\alpha$. □

896 **Theorem 2.** *Il existe (ε, δ) tel que Q_C^α sur la dimension C peut être (ε, δ) -approximé par*
897 \widehat{Q}_C^α .

898 On peut approximer $Q_{C=c}^\alpha$ en prenant plusieurs Réservoirs indépendants et appliquer
899 Chernoff-Hoeffding. Avec l'inégalité de Markov nous obtenons une approximation plus
900 faible, comme dans lemma 4.

901 4.3.2 Conclusion

902 Dans cette partie, nous avons montré comment approximer les requêtes OLAP, à partir
903 d'un entrepôt de données en streaming, à l'aide d'un tirage proportionnel aux *Mesures*.
904 Nous avons introduit une notion de préférence comme une distribution sur les *Mesures*.
905 Nous avons également montré, que nous pouvons généraliser l'approximation sur toute
906 préférence.

907 5 Analyse du réseau social Twitter 908 en streaming

909 De nombreux réseaux sociaux partagent des flux de données avec des utilisateurs, par
910 exemple Twitter, un réseau social de type microblogging créé en 2006. Les utilisateurs
911 publient des messages appelés *tweets* qui sont partagés selon certaines règles simples. Il a
912 été conçu à l'origine comme un service basé sur le SMS où les messages étaient limités à
913 160 caractères. Aujourd'hui, les *tweets* sont limités à 280 caractères, laissant 20 caractères
914 pour le nom de l'utilisateur.

915 L'étude des Tweets est un sujet classique en particulier pour *l'analyse de sentiments*.
916 On cherche à savoir pour un sujet donné si le Tweet est positif ou négatif. L'analyse de
917 sentiments est présentée dans [38, 30].

918 Dans la première section, nous expliquons comment récupérer les données par des
919 API. Dans la deuxième section nous décrivons le type de données qu'on peut récupérer.
920 Dans la troisième section, nous décrivons comment le flux de tweet est transformé en un
921 flux d'arêtes d'un graphe. Dans la quatrième section, nous montrons comment décider de
922 manière approchée certaines propriétés de graphes Twitter.

923 5.1 Récupération des données

924 Les tweets publiés par des comptes publics peuvent être récupérés à l'aide de l'une des
925 deux API Twitter suivantes :

- 926 1. API Search REST : Renvoie une collection de Tweets correspondant à une requête
927 spécifiée.
- 928 2. API Streaming : Renvoie un flux de Tweets correspondant à un filtre spécifié. Dans
929 notre cas, nous utilisons l'API Streaming.

930 Pour utiliser les API Twitter, il est nécessaire d'avoir un compte développeur¹. Une
931 fois réalisé, il nous faut créer une application pour pouvoir obtenir des clés et jetons afin
932 de s'authentifier sur ces API.

1. <https://apps.twitter.com/>

933 5.1.1 API Search Twitter

934 L'API search permet de retourner des tweets qui répondent à une requête q . Cette
935 requête q peut être composée à partir de plusieurs paramètres tels que :

Paramètres	Obligatoire	Description	Exemple
q	oui	Renvoie les tweets concernant des termes contenus dans le message du tweet.	France
geocode	non	Renvoie les tweets d'utilisateurs situés dans un rayon donné de la latitude / longitude donnée.	37.781157 -122.398720 1mi
lang	non	Renvoie les tweets écrits dans une langue spécifique	fr
result_type	non	Renvoie le type de tweets par exemple les plus populaires (les plus retweetés), les plus récents ou mixtes (mélange entre les plus populaires et les plus récents).	mixed
count	non	Nombre de tweets retournés	100
until	non	Retourne les tweets créés avant une date donnée	2019-01-10
since_id	non	Renvoie les résultats avec un identifiant supérieur (c'est-à-dire plus récent) à l'identifiant spécifié.	12345
max_id	non	Retourne les résultats avec un ID inférieur à (ou antérieur à) ou égal à l'ID spécifié.	54321
include_entities	non	Les entités ne seront pas incluses si elles sont définies sur false.	false

936 Par exemple, si nous voulons recevoir les tweets les plus récents contenant le terme
937 "France" écrits en français, nous devons réaliser la requête suivante :

938

```
939 $curl --request GET
940 --url 'https://api.twitter.com/1.1/search/tweets.json?q=france$lang=fr&result_type=recent'
941 --header 'authorization:_OAuth_oauth_consumer_key="consumer-key-for-app",
942 _oauth_nonce="generated-nonce",_oauth_signature="generated-signature",
943 _oauth_signature_method="HMAC-SHA1",_oauth_timestamp="generated-timestamp",
944 _oauth_token="access-token-for-authed-user",_oauth_version="1.0"
```

945 Nous présentons la structure et le format de la réponse à la requête q dans la section
946 suivante.

947 La search API présente des limites :

- 948 1. Le nombre de tweets retournés par requête ne peut pas dépasser 1500 tweets.
- 949 2. Il est impossible de retrouver les tweets qui ont été envoyés il y a plus d'une semaine.

950 Ci-dessous nous présentons un exemple d'utilisation en Python² de l'API search de
951 Twitter en utilisant la librairie Tweepy³

952

2. Python est un langage de programmation interprété

3. Tweepy est une librairie écrite en python, cette librairie permet à Python de communiquer avec la plateforme Twitter et d'utiliser son API.

```

import tweepy

# consumer keys and access tokens, used for OAuth
consumer_key = 'YOUR-CONSUMER-KEY'
consumer_secret = 'YOUR-CONSUMER-SECRET'
access_token = 'YOUR-ACCESS-TOKEN'
access_token_secret = 'YOUR-ACCESS-TOKEN-SECRET'

# OAuth process, using the keys and tokens
auth = tweepy.OAuthHandler(consumer_key, consumer_secret)
auth.set_access_token(access_token, access_token_secret)

# creation of the actual interface, using authentication
953 api = tweepy.API(auth)

# search recent tweets containing text "France"
# search french language tweets
# limit response 5 tweets
search = tweepy.Cursor(api.search, q="France", result_type="recent", lang="fr").items(5)

# print the searched tweets
for item in search:
    print (item.text)

```

Code 5.1: Exemple d'utilisation de l'API search en Python

954 5.1.2 API streaming Twitter

955 Les requêtes de flux s'exécutent en continu. Les différentes API de flux de données four-
956 nies par Twitter permettent aux développeurs d'avoir accès aux données de flux globales
957 de Twitter. Twitter propose plusieurs points d'accès au flux de données :

- 958 1. Flux publics : Il donne accès aux flux de données accessibles au public sur Twitter.
959 Généralement, ils sont utilisés pour suivre des sujets spécifiques, individuels et pour
960 l'exploration de données.
- 961 2. Flux d'utilisateurs : Ils fournissent un accès aux données Twitter relatives à un
962 utilisateur particulier. Vous avez accès aux tweets, aux abonnés et aux amis d'un
963 utilisateur que vous souhaitez.
- 964 3. Flux de site : Les flux de site sont utilisés par les serveurs qui se connectent à Twitter
965 en tant que serveur proxy⁴.

Paramètre	Description
Follow	Liste des utilisateurs auxquels renvoyer les statuts dans le flux.
Track	Mots-clés à suivre.
Locations	Spécifie une zone géographique à suivre.
Delimited	Spécifie si les messages doivent être séparés par une longueur.
Stall_warnings	Spécifie si les avertissements de décrochage doivent être livrés.

TABLE 5.1: Paramètres pour le point d'accès "POST statuses/filter"

4. De nombreux sites web utilisent ce point d'accès pour afficher les tweets dans une page web.

966 Pour l'analyse en temps réel des données de Twitter, les flux publics constituent une
967 bonne source d'API.

968 5.2 Caractéristiques des données Twitter

969 Toutes les API Twitter qui renvoient des tweets fournissent ces données codées à l'aide
970 de la notation JSON⁵. Le format JSON est basé sur des paires clé-valeur, avec des attributs
971 nommés et des valeurs associées. Ces attributs et leur état sont utilisés pour décrire les
972 objets.

973 5.2.1 Structure

974 La figure ci-dessous est un exemple d'un tweet émis par l'API Twitter. Les informations
975 fournies contiennent des données relatives au tweet : date de création du tweet, texte du
976 tweet, adresses URL ou mentions d'utilisateur, etc. Des informations sur Retweet sont
977 également fournies avec les informations d'origines. Les informations concernant le créateur
978 de la publication sont également associées aux informations de tweet.

```
"created_at": "Fri Jan 23 23:57:36 +0000 2015",
"id": 558775589612437504,
"id_str": "558775589612437504",
"text": "Thanks, polar vortex: Attendance dips at major Chicago museums in 2014 http://t.co/jQLEurk9s http://t.co/3bss2nGemx",
"source": "\u003ca href=\"http://twitter.com\" rel=\"nofollow\"\u003eTwitter Web Client\u003c/a\u003e",
"truncated": false,
"in_reply_to_status_id": null,
"in_reply_to_status_id_str": null,
"in_reply_to_user_id": null,
"in_reply_to_user_id_str": null,
"in_reply_to_screen_name": null,
"user": {
  "id": 7313362,
  "id_str": "7313362",
  "name": "Chicago Tribune",
  "screen_name": "chicagotribune",
  "location": "Chicago, IL",
  "url": "http://www.chicagotribune.com/",
  "description": "Chicago Tribune news, features and so much more live from our newsroom. A part of your life since 1847.",
  "protected": false,
  "verified": true,
  "followers_count": 321548,
  "friends_count": 523,
  "listed_count": 8074,
  "favourites_count": 34,
  "statuses_count": 47367,
  "created_at": "Sat Jul 07 14:10:07 +0000 2007",
  "utc_offset": -21600,
  "time_zone": "Central Time (US & Canada)",
  "geo_enabled": false,
  "lang": "en",
```

FIGURE 5.1: Exemple de Twitter JSON

979 5.2.2 Contenus

980 Sur Twitter, nous recevons de nombreux objets, les Tweets et leurs auteurs. Ces objets
981 encapsulent tous les attributs principaux qui décrivent l'objet. Chaque Tweet comporte

5. JavaScript Object Notation (JSON) est un format de données textuelles dérivé de la notation des objets du langage JavaScript.

982 un auteur, un message, un identifiant unique, un horodatage de publication et parfois
 983 des métadonnées géographiques partagées par l'utilisateur. Chaque utilisateur a un nom
 984 Twitter, un identifiant, un nombre d'abonnés et le plus souvent une biographie de compte.

985 Avec chaque Tweet, nous générons également des objets 'entité', qui sont des tableaux
 986 de contenus de Tweet communs tels que des hashtags, des mentions, des médias et des
 987 liens. S'il existe des liens, la charge JSON peut également fournir des métadonnées telles
 988 que l'URL entièrement déroulée, ainsi que le titre et la description de la page Web. Ainsi,
 989 en plus du contenu du texte, un Tweet peut avoir plus de 150 attributs associés. Reflétant
 990 la hiérarchie JSON ci-dessus, voici une descriptions supplémentaires de ces objets :

- 991 1. Tweet : Également appelé objet 'Statut', possède de nombreux attributs de niveau
 992 racine, parent de d'autres objets.
 - 993 (a) Utilisateur : Métadonnées au niveau du compte Twitter. Inclut tous les enri-
 994 chissements disponibles au niveau du compte, tels que Profile geo.
 - 995 (b) Entités : Contient jusqu'à quatre photos natives, ou une vidéo ou un fichier GIF
 996 animé.
 - 997 (c) Entités étendues : Contient les tableaux d'objets de #hashtags, @mentions, \$
 998 symbol, URL et media.
 - 999 (d) Lieux : Parent à l'objet 'coordonnées'.

1000 Lors de l'acquisition de données Tweet, l'objet principal est l'objet Tweet, qui est un
 1001 objet parent de plusieurs objets enfants. Par exemple, tous les Tweets incluent un objet
 1002 Utilisateur qui décrit l'auteur du Tweet. Si le tweet est géolocalisé, un objet "lieu" est
 1003 inclus. Chaque Tweet comprend un objet "entités" qui encapsule des tableaux de hashtags,
 1004 de mentions utilisateur, d'URL et de médias natifs. Si le Tweet contient un média "attaché"
 1005 ou "natif" (photos, vidéo, GIF animé), il y aura un objet "extended_entities".

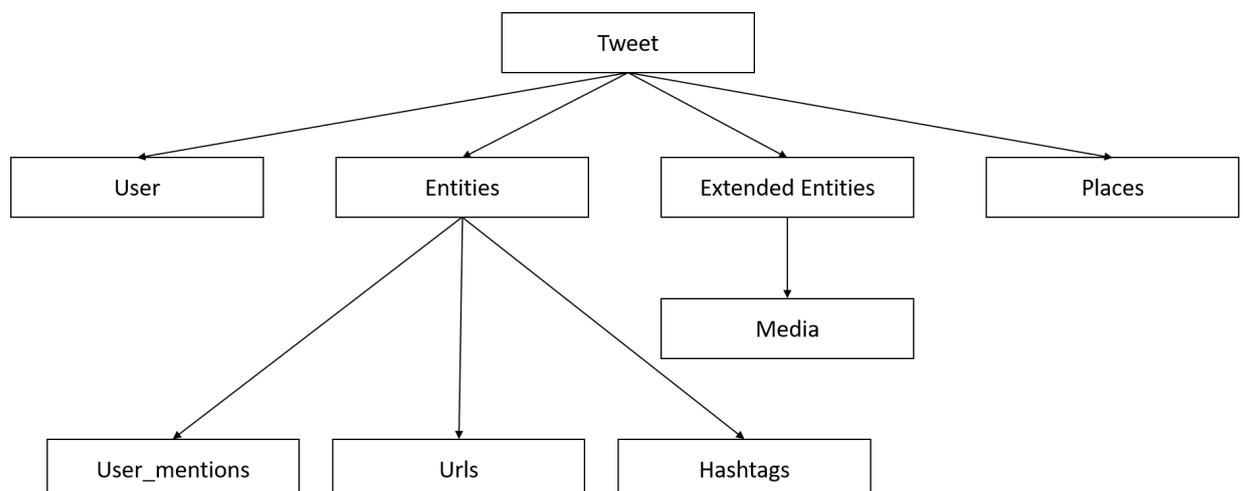


FIGURE 5.2: Exemple d'arborescence Twitter

1006 5.3 Construction d'un graphe social à partir d'un arbre JSON

1007 Nous distinguons deux grands types de graphe possibles sur Twitter : le graphe clas-
1008 sique d'utilisateur où les arêtes déterminent les relations entre les utilisateurs 5.4. Un
1009 autre graphe possible est déterminé par des événements. Nous définissons ici un événement
1010 comme étant l'activité générée par l'émission de tweets. Lors de l'émission d'un tweet, un
1011 certain nombre d'informations sont produites. Ces informations peuvent être également
1012 représentées sous forme de graphe.

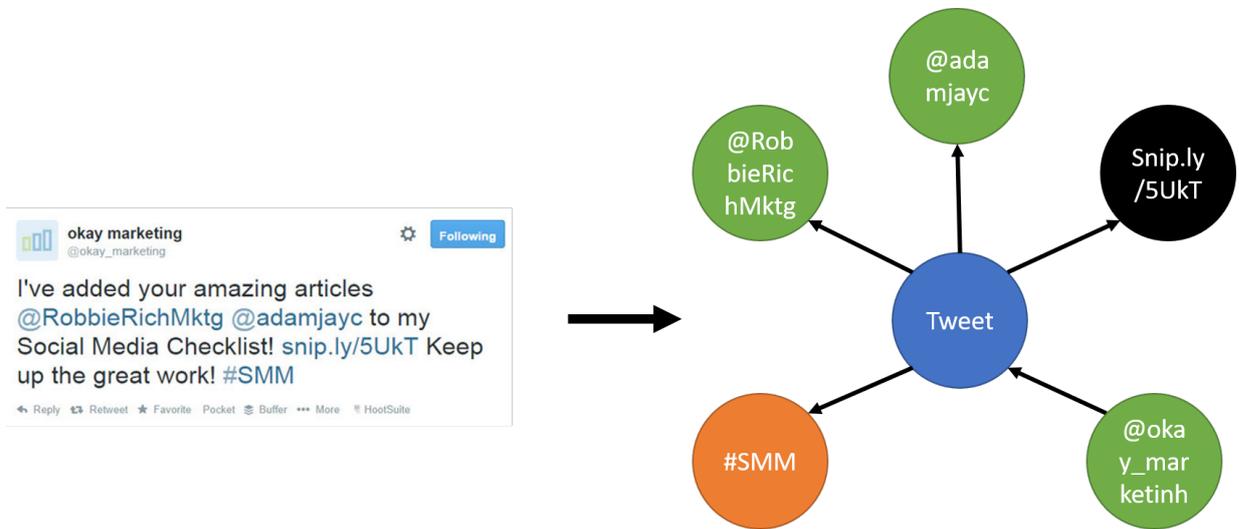


FIGURE 5.3: D'un tweet à un graphe

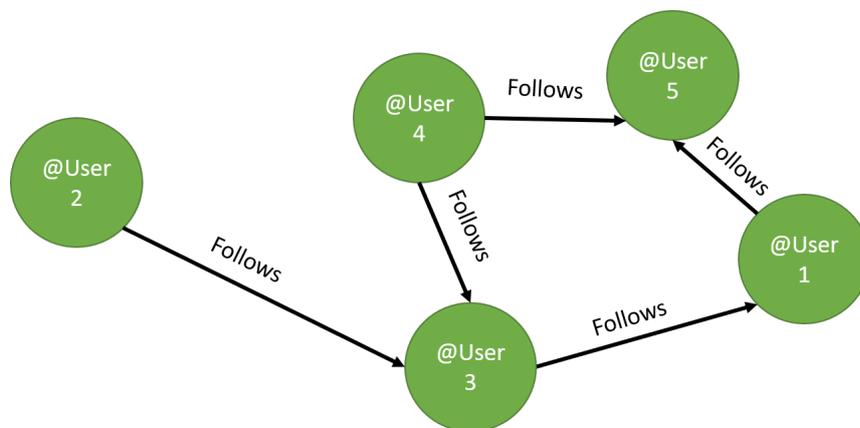


FIGURE 5.4: Exemple de graphe non symétrique d'utilisateurs sur Twitter

1013 Dans la suite, nous considérons les graphes de tweets générés sur un sujet donné, défini
1014 par un ou plusieurs tags.

1015 5.3.1 Twitter graphe d'un flux de données

1016 Dans les réseaux sociaux et le crowdsourcing, nous observons d'importants flux de
1017 données en ligne, le plus souvent des arêtes e_1, \dots, e_m d'un graphe. Compte tenu d'un
1018 ensemble de tags tels que $\{\#Ethereum, \#Bitcoin\}$, ou $\{\#Amazon\}$, Twitter fournit un
1019 flux de tweets représenté par des arbres JSON dont le contenu C (le texte du tweet)
1020 contient au moins un de ces tags. Le *Graphe Twitter* du flux est le graphe $G = (V, E)$ avec
1021 des arêtes multiples E où V est l'ensemble des tags $\#x$ ou $@y$. Pour chaque tweet envoyé
1022 par $@y$ qui contient des tags, nous ajoutons les arêtes $(@y, \#x)$ et $(@y, @z)$ dans E .

1023 Dans notre approche, nous considérons l'hypergraphe où nous ajoutons le contenu C
1024 à chaque arête. Nous avons donc les hyperarêtes $(@y, \#x, C)$ et $(@y, @z, C)$. Les URL
1025 qui apparaissent dans le tweet peuvent également être considérées comme des nœuds mais
1026 nous les ignorons pour des raisons de simplicité. Un flux de tweets est ensuite transformé
1027 en un flux d'arêtes e_1, \dots, e_m, \dots , bien que chaque arête soit une hyperarête, qui stocke
1028 également un timestamp.

1029 Les réseaux sociaux tels que Twitter évoluent de manière dynamique, et des sous-
1030 graphes denses apparaissent et disparaissent au fil du temps à mesure que l'intérêt pour
1031 des événements particuliers augmente et disparaît.

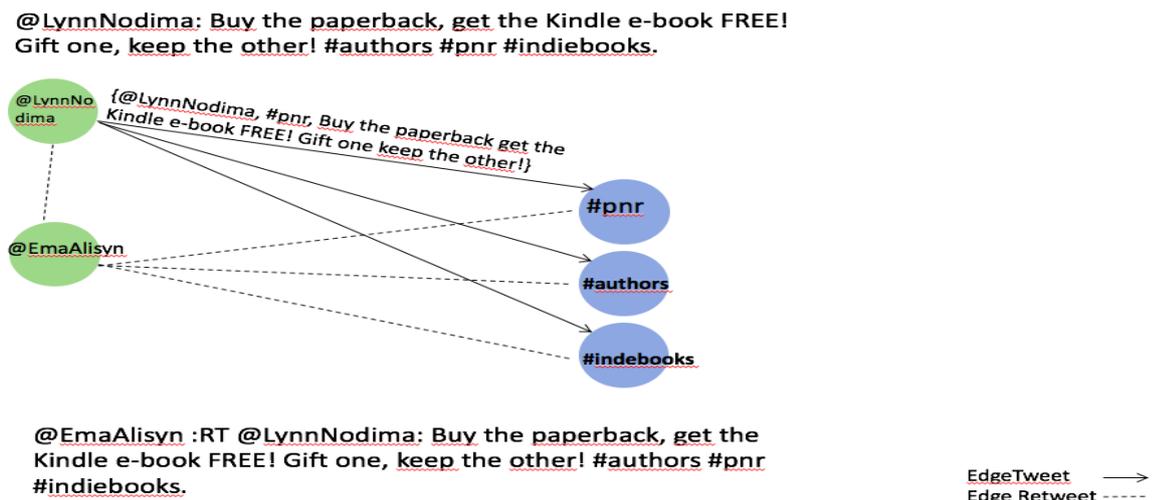


FIGURE 5.5: Graphe Twitter

1032 5.4 Propriétés d'un graphe Twitter

1033 L'analyse des réseaux sociaux utilise plusieurs concepts issus de la théorie des graphes
 1034 qui fournit des outils pour représenter formellement les réseaux sociaux et quantifier leurs
 1035 propriétés structurelles. Les graphes des réseaux sociaux satisfont certaines propriétés
 1036 générales, comme une distribution particulière de degrés présentée dans la section 5.4.1
 1037 et un petit diamètre. On étudie dans la section 5.4.2 des solutions algorithmiques efficaces
 1038 pour ces classes de graphes.

1039 5.4.1 Distribution des degrés

Dans un graphe non-orienté $G = (V, E)$, la distribution des degrés est définie par :

$$P(k) = \frac{|\{u_i \in V : \text{deg } u_i = k\}|}{|V|}, \quad (5.1)$$

1040 $P(k)$ donne la fraction (ou le pourcentage) des nœuds dans V , dont le degré est
 1041 $k \in \{0, 1, \dots, n-1\}$, $n = |V|$. Par conséquent, $P(k)$ peut être interprété comme la proba-
 1042 bilité pour un nœud (tiré au hasard V) d'avoir exactement k arêtes.

1043 Dans les graphes aléatoires Erdős-Rényi, il a été montré dans [3] que $P(k)$ suit une
 1044 distribution de Poisson dont le pic est situé à $\langle k \rangle$ ($\langle \cdot \rangle$ désigne la valeur d'attente).

1045 Une observation générale des réseaux sociaux est de constater que tous les graphes
 1046 suivent une distribution de degrés comme le montre la figure 5.6. La figure 5.7, illustre
 1047 la distribution en loi de puissance des degrés dans un graphe Twitter. Dans une loi de
 1048 puissance, la probabilité pour un nœud donné d'avoir k arêtes est définie comme suit :
 1049
 1050

$$P(k) \sim k^{-\lambda} \quad (5.2)$$

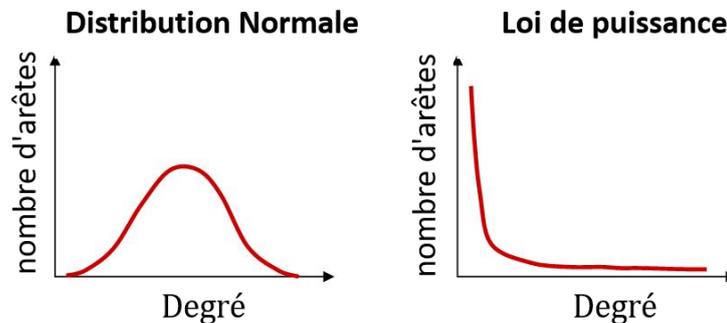


FIGURE 5.6: Distributions de degrés

Degree Report

Results:
Average Degree: 1,631

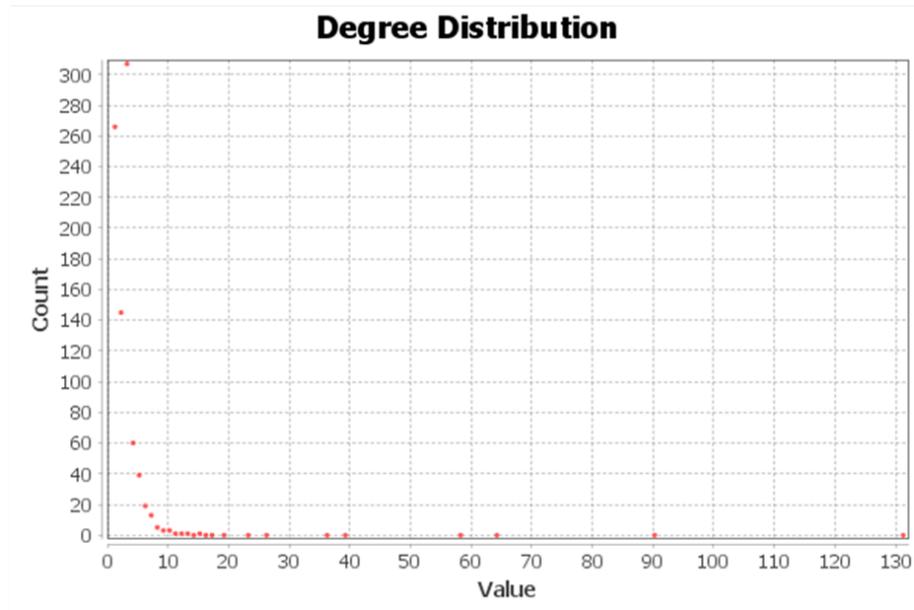


FIGURE 5.7: Exemple de la distribution de degrés dans un graphe Twitter

1051 Notre modèle de graphe a une distribution de degrés qui suit une loi de puissance, telle
1052 que la loi de Zipf où :

$$Prob[d(i) = j] = \frac{c}{j^2} \text{ avec } c = \frac{6}{\pi^2} \quad (5.3)$$

1053 car :

$$\sum_{i \geq 1}^{\infty} \frac{1}{i^2} = \frac{\pi^2}{6}$$

1054 5.4.2 Propriétés à analyser

1055 Pour un graphe social G , il existe de nombreuses propriétés qui peuvent ou ne peuvent
1056 pas être satisfaites. Par-exemple l'existence de grands sous-graphes denses, une variante
1057 du problème *Maxclique*.

1058 Si on suit plusieurs flux, par exemple deux flux, on obtient deux graphes sociaux G_1
1059 et G_2 . Comment peut-on étudier leur corrélation ? Nous étudions cette question dans la
1060 section 5.4.2.3.

1061 5.4.2.1 Communautés et sous-graphes denses

1062 Soit $S \subseteq V$ et $E(S)$ l'ensemble des arêtes internes, c'est-à-dire les arêtes $e = (u, v)$ où
1063 $u, v \in S$.

1064 **Definition 5.** Un γ -cluster S est un sous-graphe qui dépend d'un paramètre $\gamma \leq 1$ tel que
1065 $|E(S)| \geq \gamma \cdot |S| \cdot |S - 1|/2$.

1066 Un sous-ensemble S est (δ) -grand, si la taille de S est plus grande que $\delta\sqrt{n}$. Le problème
1067 (γ, δ) -grand sous-graphe dense consiste à décider s'il existe un γ -cluster S qui est (δ) -grand.

1068 Si $\gamma = 1$, l'ensemble S est une clique. En général ce problème généralise *MaxClique*, un
1069 problème NP-dur, qui est difficile à approcher [31].

1070 **Definition 6.** *MaxClique* est un problème d'optimisation, qui consiste à trouver la taille
1071 de la plus grande clique d'un graphe.

1072

1073 On peut aussi considérer *MaxClique* comme un problème de décision. L'entrée est un
1074 graphe G_n et un entier $\lambda \leq n$. On décide s'il existe une clique de taille plus grande que λ .

1075

1076 La densité classique de S est le rapport $\rho = |E[S]|/|S|$. On peut vouloir trouver des sous-
1077 graphes avec des nœuds S qui maximisent ρ . Dans le cas d'un flux d'arêtes, l'approximation
1078 des sous-graphes denses est bien étudiée dans [8, 26, 29, 43] et une limite inférieure de
1079 l'espace $\Omega(n)$ est connue [6].

1080 Les deux versions du problème sont difficiles à approximer, et nous allons modifier la
1081 notion d'approximation pour la rendre compatible avec la classe de graphes sociaux, en
1082 utilisant la première version de graphes denses. Une famille de graphes G_n a un (γ, δ) -
1083 cluster s'il existe (γ, δ) tel que pour tout n assez grand $\exists S$ tel que S est un (γ, δ) -cluster
1084 pour G_n .

1085 Soit μ la distribution uniforme des graphes qui suivent une loi de degrés comme une
1086 loi de puissance. Soit un langage L associé à un problème de décision, c'est-à-dire L est
1087 l'ensemble de toutes les instances positives.

1088 **Definition 1.** Un algorithme randomisé A est asymétrique pour la distribution μ pour
1089 décider le langage L s'il satisfait les deux conditions :

- 1090 — Pour tout $x \in L$, $\text{Prob}_{\Omega}[A(x) \text{ accepte}] \geq 1 - \varepsilon$
- 1091 — Si $x \notin L$ est tiré selon μ , $\text{Prob}_{\mu \times \Omega}[A(x) \text{ rejette}] \geq 1 - \varepsilon$

1092 L'approximation n'est plus pour le cas le plus défavorable. La deuxième condition
1093 spécifie une probabilité de rejet pour une entrée négative tirée selon μ .

1094 L'algorithme A va considérer la taille des composantes géantes du Réservoir, pour une
1095 taille de k bien choisie et décider la propriété P : Il existe un (γ, δ) -cluster ?

1096

Algorithme A1

Entrée : un flux de m arêtes d'un graphe G .

Sortie : Accepter s'il existe un (γ, δ) -cluster, sinon Rejeter. • Maintenir un Réservoir R de taille $k = c \cdot \sqrt{n} \log n / (4 * \gamma \delta)$.

• Soit C la plus grande composante connexe de R .

• Accepter si $|C| \geq \Theta(n^{1/8} \log^2 n)$, sinon Rejeter.

1097 **Détection de communautés.** Soit C la plus grande composante connexe d'un Réservoir
1098 de taille $k = (\alpha \sqrt{n} \log n)$. Considérons l'algorithme A1, qui dépend de γ, δ où $\alpha = 1/\gamma \delta$
1099 est un paramètre auxiliaire.

1100 La borne $n^{1/8} \log^2 n$ est une application directe du théorème 5.4.2.1 de Molloy-Reed.
1101 On va montrer que A1 est un algorithme asymétrique pour la distribution μ . La première
1102 condition suit le résultat suivant :

1103 **Théorème 1.** Si G a $O(n \cdot \log n)$ arêtes et un (γ, δ) -cluster S , l'algorithme A1 Accepte
1104 presque sûrement.

Démonstration. Pour un graphe G avec $cn \cdot \log n / 4$ arêtes, comme ceux dont la distribution de degrés est une loi de puissance, considérons le taux d'échantillonnage :

$$\frac{k}{m} = \frac{c \cdot \sqrt{n} \log n / 4}{\gamma \cdot \delta \cdot cn \cdot \log n / 4} = \frac{1}{\gamma \cdot \delta \cdot \sqrt{n}}$$

1105 Celui-ci est indépendant des arêtes. On a donc dans le Réservoir un graphe tiré selon
1106 Erdos-Renyi avec une probabilité de $\frac{k}{m}$.

Puisque

$$|S| > \delta \cdot \sqrt{n}$$

$$\frac{k}{m} = \frac{1}{\gamma \cdot \delta \cdot \sqrt{n}} > \frac{1}{\gamma \cdot |S|}$$

$$\frac{k}{m} = \frac{1}{\gamma \cdot \delta \cdot \sqrt{n}} > \frac{(1 + \epsilon)}{\gamma \cdot |S|}$$

D'après le lemme 2, il existe une composante géante si

$$\frac{k}{m} > \frac{(1 + \epsilon)}{\gamma \cdot |S|}$$

1107 L'algorithme A1 accepte donc presque sûrement. □

1108 **Non Détection.** La deuxième condition à vérifier est plus difficile. L'étude est basée
1109 sur les graphes aléatoires satisfaisant une distribution de degrés, comme loi de puissance. Le
1110 résultat central de Molloy-Reed [48] donne une condition suffisante pour la non existence

1111 de composantes géantes et plus précisément une borne sur la taille de la plus grande
1112 composante connexe.

1113 Soit $d_i(n)$ le nombre de noeuds de degré i pour un graphe de taille n . La distribution
1114 $(d_i(n))_{i,n}$ est réalisable si :

- 1115 1. pour tout n , il existe un graphe de taille n dont la distribution de degré est $(d_i(n))_i$,
- 1116 2. pour tout i , $\ell_i = \lim_{n \rightarrow \infty} d_i(n)/n$ existe.

1117 Considérons une distribution réalisable telle que :

- 1118 1. $Q(\mathcal{D}) = \sum_i (i^2 - 2i)\ell_i$ est plus petit qu'une constante inférieure à 0,
- 1119 2. le degré maximum est inférieur $n^{1/9}$,
- 1120 3. le degré moyen est $O(1)$.
- 1121 4. la convergence de $i(i - 2)d_i(n)/n$ vers sa limite ℓ_i est uniforme.

1122 Soit μ la distribution des graphes aléatoires qui suivent les $d_i(n)$ dans le modèle de
1123 configuration. [48] montre que presque surement, la taille de la plus grande composante
1124 connexe est au plus $Bn^{1/4}$ pour une constante B qui dépend de $Q(\mathcal{D})$. D'autre part aucune
1125 composante connexe de G a plus d'un cycle et il y a au plus $2Bn^{1/4}$ cycles.

1126 Pour une instance aléatoire selon μ , la taille de la plus grande composante connexe est
1127 inférieure au seuil de l'algorithme A1. Notez que la variable n de Molloy-Reed s'applique
1128 au nombre de noeuds du graphe du Reservoir, environ \sqrt{n} . D'où le seuil $n^{1/8} \log^2 n$, utilisé
1129 par l'algorithme A1.

1130 **Théorème 2.** [47] Soit G un graphe aléatoire tiré uniformément selon la loi μ . l'algorithme
1131 A1 Rejette presque surement et est correct.

1132 On peut trouver des entrées adversaires dans ce cas sont des graphes, qui n'ont pas
1133 de (γ, δ) -cluster pour lesquelles l'algorithme se trompe dans le sens il va dire il existe un
1134 (γ, δ) -cluster avec une certaine probabilité. Ces entrées ne suivent pas la loi de puissance.
1135 L'approximation est donc dans ce sens précis.

1136 Ce résultat se généralise aux graphes dynamiques générés par des fenêtres coulissantes.
1137 On détectera les composantes géantes et on pourra suivre leurs évolutions dans le temps.

1138 5.4.2.2 Graphes dynamiques dans un flux d'arêtes

1139 Nous suivons le modèle des fenêtres glissantes, comme le montre la figure 5.8, défini par
1140 un intervalle de temps fixe τ . Si le débit (le nombre d'arêtes par unité de temps) du flux
1141 est fixe, chaque fenêtre admet le même nombre d'arêtes. Ce n'est pas le cas en pratique,
1142 car le taux fluctue d'un facteur de 2 à tout moment.

1143
1144 Chaque fenêtre de longueur τ se termine aux instants $t_1, t_2, \dots, t_i, \dots$. Chaque $t_i =$
1145 $\tau + \lambda \cdot (i - 1)$ pour $i > 1$ et $\lambda < \tau$ détermine une fenêtre de longueur τ et un graphe
1146 G_i défini par les arêtes dans la fenêtre ou dans un intervalle de temps $[t_i\tau, t_i]$. Le nombre

1147 d'arêtes dans une fenêtre peut augmenter ou diminuer et reflète la vitesse croissante ou
 1148 décroissante d'un flux. Les fenêtres consécutives se chevauchent dans un facteur τ/λ , en-
 1149 viron 50% dans les expériences. En pratique, $\tau = 60$ minutes et $\lambda = 30$ minutes.

1150

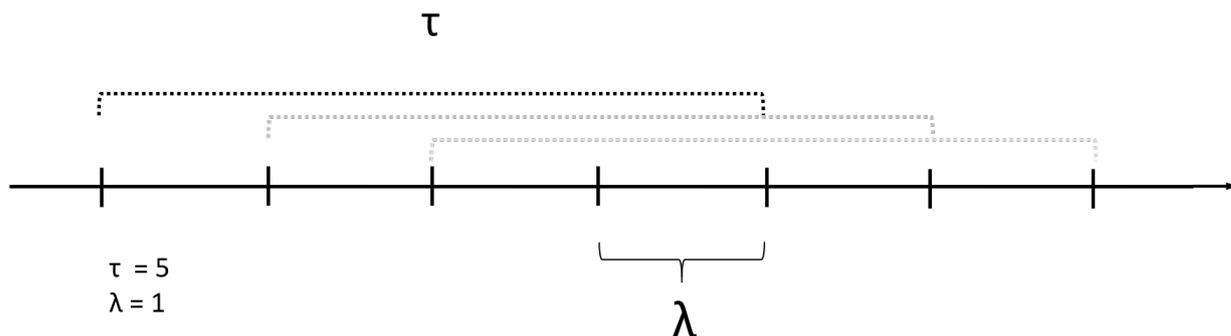


FIGURE 5.8: Step reservoir sampling.

1151 5.4.2.3 Corrélation de contenu entre flux

La corrélation classique, aussi appelée corrélation de Pearson $p(X, Y)$ de deux variables aléatoires X, Y de moyenne μ et l'écart-type σ est :

$$\frac{\mathbb{E}[(X - \mu_X)(Y - \mu_Y)]}{\sigma_X \cdot \sigma_Y} \quad (5.4)$$

1152 Comment peut-on l'étendre aux graphes? On pourrait choisir certains paramètres des
 1153 graphes et prendre les corrélations entre ces paramètres. Les graphes sociaux ont cependant
 1154 des statistiques très similaires et la corrélation serait importante. Le noyau de l'information
 1155 semble se cacher dans la structure des clusters.

1156 Soit deux graphes G_1 et G_2 issus de deux flux distincts. Une première approche de leur
 1157 corrélation de contenu serait de considérer la similarité Jaccard⁶ $J(V_1, V_2)$ sur les domaines
 1158 des deux graphes. Cette définition a plusieurs inconvénients : elle est indépendante des
 1159 structures des graphes, elle est très sensible au bruit, et n'est pas bien adaptée lorsque les
 1160 tailles sont très différentes. On doit également stocker l'ensemble des arêtes de l'ensemble
 1161 des graphes.

1162 Nous proposons plutôt l'approche suivante : nous appliquons la similarité de Jaccard
 1163 uniquement aux clusters des graphes. Soit $C_i = \bigcup_j C_{i,j}$ l'ensemble des clusters du graphe
 1164 G_i , pour $i = 1$ ou 2 .

6. La similarité Jaccard ou Index entre deux ensembles A et B est $J(A, B) = |A \cap B|/|A \cup B|$. La distance Jaccard est de $1 - J(A, B)$.

1165 **Definition 7.** *La corrélation de deux graphes est le coefficient $\rho = J(C_1, C_2)$.*

1166 Cette approche exploite la structure des graphes, est insensible au bruit et s'adapte bien
1167 lorsque les graphes ont des tailles différentes. Dans la suite, on considère pour simplifier
1168 que les graphes ont un seul cluster. Il est alors possible d'approximer la corrélation en
1169 utilisant les composantes géantes des graphes. Il suffit de considérer le 2-core de chaque
1170 composante géante. En effet, comme le montre [47], le 2-core de la composante géante du
1171 Réservoir est une bonne approximation, au sens ensembliste, de chaque C_i .

Algorithme A2

Entrée : deux flux d'arêtes définissant les graphes G_1, G_2 .

Sortie : une approximation de la corrélation ρ .

- Maintenir deux Réservoirs R_1, R_2 de taille $k = \Theta(c.\alpha\sqrt{n}\log n/4)$.
 - Soit \widehat{C}_1 la plus grande composante connexe de R_1 et \widehat{C}_2 la plus grande composante connexe de R_2 .
 - Soit $\widehat{C}'_1 = 2\text{-core}(\widehat{C}_1)$ et $\widehat{C}'_2 = 2\text{-core}(\widehat{C}_2)$.
 - Renvoyer $J(\widehat{C}'_1, \widehat{C}'_2)$
-

1172 On peut donc approximer chaque C_i par \widehat{C}'_i et ainsi la similarité de Jaccard.

1173 **Théorème 3.** *L'algorithme A2 approxime le coefficient $\rho = J(C_1, C_2)$.*

1174 Pour chaque fenêtre on fait cette analyse. Dans le cas des fenêtres dynamiques la fonc-
1175 tion de corrélation dépendra du temps et donnera une approximation sur chaque fenêtre.

1176 5.4.3 Conclusion

1177 Nous avons montré comment analyser les graphes Twitter. On cherche en particulier
1178 à décider s'il existe de grands clusters. Notre contribution est une implémentation de la
1179 méthode présentée dans [42] dans le contexte collaboratif Python de Google. L'environne-
1180 ment Colab, permet d'implémenter les algorithmes écrit en Python, dans un contexte logi-
1181 ciel des serveurs de Google qui lisent directement les données de Twitter. Ce développement
1182 logiciel a servi de motivation pour imaginer des échantillonnages sur des distributions plus
1183 sophistiquées dans le cadre de l'analyse du texte.

1184 6 Analyse de texte en streaming

1185 Le but du chapitre est de présenter un modèle statistique construit à partir des données
1186 de texte en streaming, sans stocker l'ensemble du texte.

1187 Dans le cas d'un texte classique, chaque phrase génère des arêtes où les nœuds sont les
1188 mots et les arêtes relient deux mots de la phrase, éventuellement deux mots consécutifs.
1189 Nous étudions les grandes composantes de ces graphes dans des fenêtres coulissantes en
1190 échantillonnant un nombre fixe d'arêtes avec deux distributions différentes, la distribu-
1191 tion uniforme ou une distribution pondérée. Ces distributions définissent des sous-graphes
1192 aléatoires du graphe d'origine. L'analyse est basée sur l'étude des composantes géantes de
1193 ces sous-graphes aléatoires.

1194 Pour un texte standard, un flux RSS d'articles d'actualité peut générer un flux impor-
1195 tant de phrases. Le graphe généré peut donc être très grand. Nous observons que pour
1196 l'échantillonnage uniforme des arêtes, les composantes géantes ne sont pas stables. Nous
1197 observons aussi que si nous échantillonnons les arêtes proportionnellement à la similarité
1198 des mots, donnée par *Word2vec* [44], les composantes géantes deviennent stables. Nous
1199 pouvons également classifier le texte en introduisant une distance naturelle entre les com-
1200 posantes géantes et utiliser l'algorithme classique *k-means*.

1201 La première section décrit les algorithmes de streaming adapté à ce contexte. La
1202 deuxième section détaille l'approche du texte en streaming. La troisième section intro-
1203 duit la stabilité. La quatrième section détaille la classification. La cinquième section étudie
1204 la gestion des préférences dans ce cadre.

1205 6.1 Algorithmes de streaming

1206 Les algorithmes de streaming introduits dans le chapitre 2 lisent les données sous
1207 forme d'un flux de taille n et maintiennent une petite mémoire : si possible constante
1208 indépendante de n , ou $poly(\log)n$ ou au moins sous-linéaire $o(n)$ dans la taille du flux. Les
1209 algorithmes classiques peuvent considérer un flux de valeurs numériques $x_i \in \{1, 2, \dots, n\}$,
1210 de mots sur un alphabet Σ , ou d'arêtes $e_i = (v_j, v_k)$ d'un graphe $G = (V, E)$ où $v_j, v_k \in V$.

1211 Une technique importante appelée *Reservoir sampling* introduit dans la section 2.1.1
1212 conserve k éléments du flux avec une *distribution uniforme*. Dans un flux de longueur m ,
1213 chaque élément a une probabilité k/m d'être dans le Réservoir. Dans le cas d'un graphe

1214 où les arêtes on des poids, le Réservoir conserve k éléments du flux avec une *distribution*
 1215 *non uniforme*. Chaque élément e_i de poids w_i d'un flux à une probabilité $k \cdot w_i / \sum_i w_i$
 1216 d'être dans le Réservoir. Une autre situation est de lire des arêtes et d'avoir un oracle qui
 1217 va donner le poids de chaque arête. Dans notre cadre, à chaque fois nous lisons une arête,
 1218 nous faisons appel à la fonction *Word2vec* pour calculer la similarité entre deux mots. Si
 1219 k est sous-linéaire, par-exemple $O(\sqrt{m})$, on obtient un algorithme d'espace sous-linéaire.

1220 6.2 Représentation du texte en streaming

Supposons que nous lisons un texte sous forme de flux : $w_1, w_2, \dots, w_n, \dots$ où w_i est un mot d'un vocabulaire important (10^4 mots environ). Pour chaque phrase, nous utilisons une étape de prétraitement, la *lemmatisation*, qui élimine les mots très fréquents (articles, verbes modaux), reconnaît les Entités (par exemple "World Trade Center") et normalise chaque mot. Nous échantillons ensuite deux mots w_i, w_j de la même phrase avec le poids :

$$v(w_i) \cdot v(w_j)$$

1221 où v est la fonction *Word2vec* qui transforme un mot en un vecteur de dimension 300.
 1222 La fonction *Word2vec* peut-être donnée comme un oracle ou peut-être construite tout
 1223 en lisant le flux du texte. Nous pouvons alors associer un Réservoir de taille k , en utili-
 1224 sant l'échantillonnage pondéré, et étudier les composantes géantes du graphe défini par le
 1225 Réservoir.

1226 Le produit $v(w_i) \cdot v(w_j)$ est aussi appelé la *Similarité* de (w_i, w_j) .

1227

1228 6.3 Stabilité

Considérons deux expériences sur le même flux (ou document D) avec deux Réservoirs *indépendants*. Soit V_1 (resp. V_2) l'ensemble des noeuds des composantes géantes de C_1 (resp. C_2) comme dans la section 2.3.4. Pour un flux D et un Réservoir de taille k , soit la variable aléatoire $\rho(D, k)$ définie comme suit :

$$\rho(D, k) = \frac{|V_1 \cap V_2|}{|V_1|}$$

1229 La variable aléatoire dépend de deux expériences et on peut aussi prendre son espérance
 1230 $\mathbb{E}(\rho(D, k))$. Pour les flux Twitter, l'échantillonnage uniforme offre une bonne stabilité,
 1231 comme le suggère l'analyse théorique des composantes géantes. Au contraire, pour les
 1232 textes, l'échantillonnage uniforme est instable comme on le montre dans la partie expérimentale
 1233 du chapitre 9.

1234 6.4 Classification

1235 Nous introduisons d’abord une distance naturelle entre des composantes géantes, qui
1236 peut être généralisée à des ensembles de composantes. Nous pouvons ensuite appliquer
1237 l’algorithme *k-means*¹

1238 L’analyse du texte classique cherche à reconnaître le contenu d’un texte en particulier
1239 quand on a *k* possibilité. Un problème classique dans l’analyse du langage comme dans
1240 le chapitre 3 est de pouvoir détecter si le contenu du texte fait référence à une classe
1241 particulière.

1242 6.4.1 Distances entre les composantes géantes

1243 Considérons deux composantes géantes $C_1 = (V_1, E_1)$ à l’instant t_1 et $C_2 = (V_2, E_2)$ à
1244 l’instant $t_2 \geq t_1$ sous forme de graphes étiquetés. Il existe plusieurs distances possibles :

- 1245 • La *distance de Jaccard* $\text{dist}_J(C_1, C_2) = 1 - J(V_1, V_2)$ où $J(V_1, V_2)^2$ est la similarité de
1246 Jaccard entre les domaines V_1 et V_2 .
- 1247 • La *distance d’édition* prend en compte les arêtes. Une *édition* est la suppression ou l’in-
1248 sersion d’une arête avec des étiquettes et $\text{dist}_E(C_1, C_2)$ est le nombre minimum d’éditions
1249 pour transformer C_1 en C_2 , divisé par $|E_1| + |E_2|$.
- La distance d’édition amortie tient compte de la différence de temps $|t_2 - t_1|$:

$$\text{dist}_A(C_1, C_2) = \text{dist}_E(C_1, C_2) + \alpha \cdot |t_2 - t_1|$$

1250 pour une constante α .

Chaque fenêtre coulissante peut avoir plusieurs composantes géantes de tailles différentes. Dans ce cas, nous écrivons C comme une distribution sur les composantes avec un poids proportionnel à leurs tailles. Par exemple, si C a deux composantes C_1 et C_2 , nous écrivons :

$$C = \frac{|C_1|}{|C_1| + |C_2|} \cdot C_1 + \frac{|C_2|}{|C_1| + |C_2|} \cdot C_2.$$

1251 On peut donc écrire :

1. Étant donné un ensemble de points (x_1, x_2, \dots, x_n) , on cherche à partitionner les n points en k ensembles disjoints S_1, S_2, \dots, S_k en minimisant la distance entre les points à l’intérieur de chaque partition :

$$\arg \min_{\mathbf{S}} \sum_{i=1}^k \sum_{\mathbf{x}_j \in S_i} \|\mathbf{x}_j - \boldsymbol{\mu}_i\|^2$$

où $\boldsymbol{\mu}_i$ est le barycentre des points dans S_i . On procède par itération à partir de $\boldsymbol{\mu}_{1..i}$ pour $i = 1, \dots, k$ aléatoires $\in (x_1, x_2, \dots, x_n)$. On étiquète tous les points qui sont les plus proches de $\boldsymbol{\mu}_{1..i}$ par i . On calcule ensuite les nouveaux barycentres des points étiquetés par i , soit $\boldsymbol{\mu}_{2..i}$. On réitère cette opération jusqu’à ce que les barycentres ne changent plus.

2. La similarité de Jaccard ou l’index entre deux ensembles A et B est $J(A, B) = |A \cap B| / |A \cup B|$.

$$C = \lambda \cdot C_1 + (1 - \lambda) \cdot C_2$$

1252 Nous pouvons alors généraliser les distances entre les composantes aux distances entre les
1253 distributions des composantes. Dans la section suivante, nous n'utilisons que la distance
1254 de Jaccard.

1255 6.4.2 Classification en k -classes

1256 Une séquence de fenêtres coulissantes génère une séquence de composantes géantes
1257 C_1, \dots, C_n avec toutes les distances entre les paires. Supposons pour simplifier que nous
1258 gardons la composante géante la plus grande. Nous pouvons ensuite les regrouper en k
1259 classes, avec l'algorithme classique k -means. Chaque classe i a un représentant C_i . Pour
1260 une nouvelle composante géante C , il suffit de vérifier le $\text{Min}_i \text{dist}(C_i, C)$ pour classifier
1261 C .

1262 6.5 Phrases centrales et mécanismes d'attention

1263 Chaque composante géante peut être analysée, en partant du nœud de degré maximal,
1264 le *centre*, et ses arêtes adjacentes, les *arêtes centrales*. Pour le 2-core de la figure 6.3, nous
1265 commençons par le nœud central 1 de degré maximum et ses arêtes centrales adjacentes
1266 e_1, e_4, e_5, e_8 . Chaque arête est de la forme (u, v, texte) et nous pouvons analyser la répartition
1267 de l'attention vue dans la section 3.4³ [60] des mots u et v dans chaque phrase *texte*.

1268 Par exemple l'arête e_8 est :
1269 $(\text{skate}, \text{jump}, \text{"A boy is jumping on skateboard in the middle of a red bridge."})$
1270 et l'arête e_4 est :

1271 $(\text{skate}, \text{sidewalk}, \text{"The boy skates down the sidewalk."})$
1272 L'analyse de l'attention de la première phrase pour les mots *skate* et *jump* est donnée dans
1273 la figure 6.1. Une composante géante fournit un *nœud central* et des *phrases centrales*. On
1274 peut analyser la composante géante de manière récursive, en prenant parmi les nœuds
1275 explorés celui qui a le plus grand degré et itérer la construction. À l'étape suivante, le
1276 nœud 6 serait exploré avec trois nouvelles arêtes.

1277 6.6 Préférences

Soit $S_1 = \{v_1(w) : w \in T_1\}$ l'ensemble des vecteurs associés aux mots w dans le processus *Word2vec* associé à un texte T_1 , et $S_2 = \{v_2(w) : w \in T_2\}$ l'ensemble similaire pour un texte T_2 . Les vecteurs peuvent être très différents et reflètent les différents sujets,

3. Pour un mot u , l'*attention de la phrase* est la distribution sur les autres mots avec un poids proportionnel à sa similarité *Word2vec*.

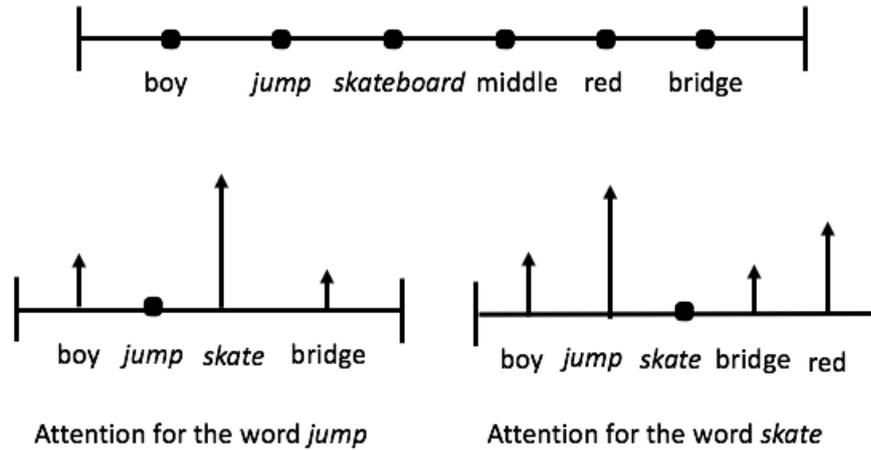


FIGURE 6.1: L'analyse de la phrase : *A boy is jumping on skateboard in the middle of a red bridge.*

comme le montre la figure 4.5. Une *préférence pure* est un tel ensemble de vecteurs construit à partir d'un texte de référence T . Une *préférence mixte* est une répartition sur S_i , par exemple $\lambda.S_1 + (1 - \lambda).S_2$ pour 2 ensembles purs, i.e. l'ensemble des vecteurs :

$$S_\lambda = \{\lambda.v_1(w) + (1 - \lambda).v_2(w) : w \in T_1 \cup T_2\}$$

1278 si un mot $w \notin T$, nous supposons que $v(w) = 0$.

1279 Nous montrons comment analyser un nouveau texte T comme un flux, en accordant
 1280 une préférence mixte arbitraire. Comme dans l'exemple précédent d'un flux de données
 1281 structurées, la distribution peut ne pas être connue à l'avance.

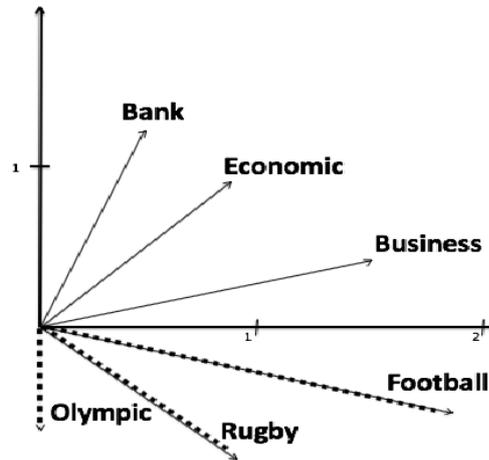


FIGURE 6.2: Deux ensembles de vecteurs dans 2 dimensions

1282 Supposons que nous analysions un nouveau texte T comme un flux, en lui donnant une
 1283 préférence mixte arbitraire, comme dans la figure 6.3. Nous maintenons deux Réservoirs R_1
 1284 pour les poids $v_1(w_i).v_1(w_j)$ et R_2 pour les poids $v_2(w_i).v_2(w_j)$. Nous pouvons construire
 1285 un nouveau Réservoir R_α , comme dans la section 4.3.

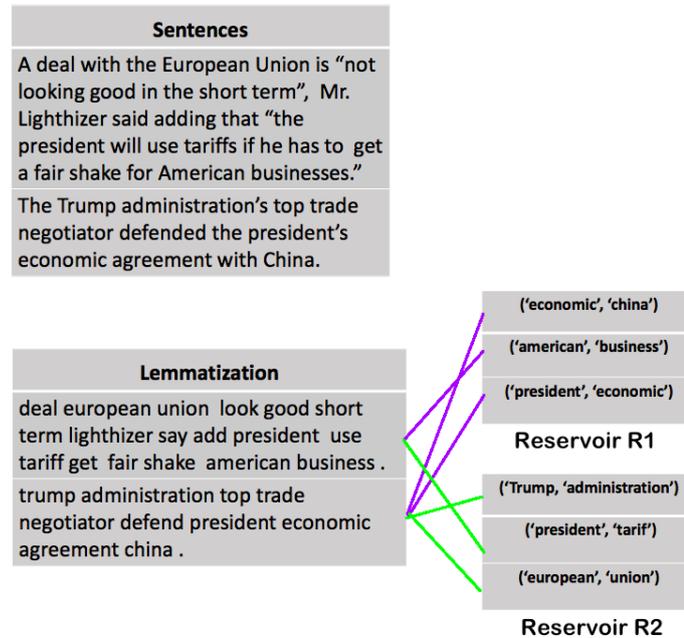


FIGURE 6.3: Deux Réservoirs pour un texte T

1286 **Theorem 3.** *Le Réservoir R_α échantillonne des paires de mots (w_1, w_2) dans un flux, avec*
 1287 *un poids proportionnel à $\alpha \cdot (v_1, v_2)_1 + (1 - \alpha) \cdot (v_1, v_2)_2$.*

1288 Nous étudions les composantes géantes du graphe associées aux arêtes du Réservoir R_α .
 1289 La figure 6.4 montre une composante géante typique et son 2-core. Le 2-core d'un Texte T
 1290 est la structure centrale utilisée pour classifier le Texte, ou pour associer une thématique.
 1291 Les distances entre les 2-core sont introduites dans la section 6.4.1.

1292 6.7 Conclusion

1293 Dans ce chapitre nous proposons une approche de l'analyse du texte en streaming,
 1294 sans stocker tout le texte. Nous avons construit un graphe aléatoire dont les arêtes sont
 1295 des échantillons de paires de mots d'une même phrase, choisis avec une distribution non
 1296 uniforme. Chaque composante géante définit un sujet, qui peut être décrit par des mots
 1297 centraux et des phrases centrales.

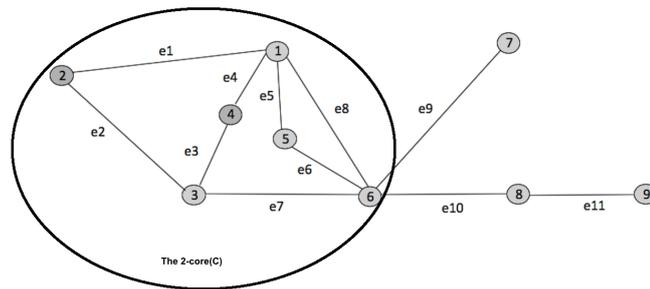


FIGURE 6.4: La composante géante du Réservoir R_α

1298 7 Modèle de morphologie dans les 1299 langues naturelles

1300 Les mots ont une structure interne, appelée aussi morphologie. Par exemple, le mot
1301 *préexistant* a un préfixe *pre-*, une racine *exist* et un suffixe *-ant*. Le suffixe *-ant* décrit
1302 l'information sur la catégorie syntaxique du mot. Dans ce cas, nous écrivons *pré-exist-ant*
1303 pour distinguer ces trois composantes. Considérons les verbes :

1304 *ap-porter, com-porter, col-porter, dé-porter, em-porter.*

1305 qui sont construits à partir de la racine *porter* avec différents préfixes qui changent le sens
1306 du mot. Si on associe un vecteur à chaque mot, quel est le rapport entre les vecteurs de
1307 deux mots qui ont la même racine mais deux préfixes différents par exemple *ap-porter*,
1308 *com-porter*?

1309 Un sujet important est d'étendre les vecteurs des mots aux phrases et plus généralement
1310 aux paragraphes et aux textes. Nous construisons un vecteur *Word2vec* structuré, associé
1311 à chacun des mots mais qui va être composé de trois parties, une partie préfixe, une
1312 partie racine et une partie suffixe. Tous les vecteurs de l'exemple ci-dessus vont partager la
1313 même racine. Dans le chapitre précédent, on a considéré les vecteurs *Word2vec* comme des
1314 vecteurs donnés par des boîtes noires. On souhaite améliorer ces vecteurs, en les structurant
1315 et ce chapitre explique comment réaliser cette opération.

1316 Dans notre modèle nous construisons des statistiques morphologiques qui vont per-
1317 mettre d'organiser nos vecteurs de manière structurée. Ce chapitre est basé sur l'article
1318 [16].

1319 Des nouvelles approches essaient de généraliser les vecteurs des mots aux phrases et
1320 plus généralement aux paragraphes et aux textes. Dans notre cadre, nous étendrons des vec-
1321 teurs aux phrases avec une approche complètement différente. Nous associons de nouveaux
1322 vecteurs à chaque phrase basés sur la structure des préfixes et suffixes qui capturent une
1323 information sémantique comme le Temps, la Voix, l'Humeur, la Force d'élocution. Cette
1324 approche, est particulièrement adaptée à certaines langues qui ont une forte morphologie
1325 comme l'*Amis*.

1326 Dans la première section, nous présentons la langue *Amis*. Dans la deuxième section,
1327 nous présentons notre modèle statistique pour saisir la morphologie d'une langue naturelle
1328 et l'appliquer à l'*Amis*. Dans la troisième section, nous donnons un aperçu syntaxique

1329 de *Amis*. Dans la quatrième section, nous étudions comment prédire approximativement
1330 la classe syntaxique d'une phrase simple. Dans la cinquième section, nous décrivons quel
1331 arbre de dérivation est préférable pour une grammaire donnée.

1332 7.1 Langue avec une forte morphologie

1333 *Amis* est l'une des quinze langues austronésiennes survivantes parlées à Taïwan. *Amis*
1334 est parlée le long de la côte orientale de Taïwan et compte quatre dialectes principaux
1335 présentant des différences significatives dans leur phonologie, leur lexique et leurs propriétés
1336 morphosyntaxiques. L'analyse porte sur les *Amis* du Nord : les données ont été recueillies
1337 par Isabelle Brill lors du travail de terrain [14] . Quelques phrases simples sont présentées
1338 en annexe A, où les affixes, les infixes et les clitiques¹ sont indiqués par des symboles
1339 spéciaux. Une propriété fondamentale de l'*Amis* est que les racines² sont le plus souvent
1340 sous-spécifiées et catégoriellement neutres [15] ; elles sont entièrement catégorisées (comme
1341 des noms, des verbes, des modificateurs, etc.) après avoir été dérivées et infléchies sous
1342 forme de mots morphosyntaxiques et projetées dans une phrase. Une autre propriété est
1343 que les préfixes, les suffixes, les affixes, contiennent énormément d'informations sémantiques
1344 et syntaxiques. Les bases nominales de noms sont signalées par l'article *u* ou par des
1345 démonstratifs. Les bases verbales ont des affixes de voix, les principaux étant la voix Acteur
1346 (Actor Voice) *mi-* (AV), la voix Patient (Undergoer Voice) *ma-* (UV), la voix Passive *-en*,
1347 la voix Locative *-an*. [19]. Un schéma syntaxique est donné dans la section 7.3.

1348 L'*Amis* est une langue austronésienne, considérée comme une protolange dans cette
1349 famille de 3000 langues parlées dans le Pacifique. De nombreuses études [14, 1] décrivent
1350 ces langues orales, leurs évolutions dans le temps et les grandes sous familles. Aujourd'hui,
1351 l'*Amis* est écrit en caractères chinois mais conserve sa structure linguistique unique. C'est
1352 un peu le *latin* des langues austronésiennes comme l'Indonésien et le Tagalog (langue parlée
1353 aux Philippines).

1354 7.2 Un modèle statistique pour la morphologie

1355 A partir de certains textes, nous pouvons analyser la distribution de fréquence des af-
1356 fixés. Avec une racine, la même analyse peut être effectuée. Nous appelons ces distributions
1357 statistiques les *statistiques morphologiques* de la langue. Il existe des approches similaires
1358 telles que [17] et [11] pour enrichir l'encastrement des mots avec des informations sur les
1359 caractères et les sous-mots. Dans notre cas, nous nous concentrons uniquement sur certains
1360 préfixes et suffixes. Nous considérons le premier et le deuxième moment des *statistiques*
1361 *morphologiques*. Nous pouvons alors déterminer quel préfixe est le plus probable dans un

1. le *clitique* est un affixe à l'intérieur d'un mot.

2. la *racine* est un mot atomique sans affixes. Les affixes sont soit inflexibles (c'est-à-dire qu'ils expriment une fonction sémantique ou syntaxique), soit dérivés (c'est-à-dire qu'ils créent différentes catégories).

1362 mot manquant d'une phrase, quel suffixe est peu probable étant donné un préfixe et une
1363 phrase, et la structure syntaxique des phrases simples. Nous considérons que ces statis-
1364 tiques sont très utiles pour saisir approximativement certains paramètres sémantiques et
1365 syntaxiques clés. *Amis*, une langue naturelle où la morphologie est importante, se prête
1366 bien à cette analyse.

1367 Nous avons construit un outil pour représenter la morphologie statistique de la langue
1368 *Amis*, étant donné un ensemble de textes où chaque mot a été décomposé en composants
1369 (c'est-à-dire préfixe, infixes, racine et suffixe).

- 1370 ● Nous analysons la distribution globale des racines et des affixes, c'est-à-dire les occur-
1371 rences les plus fréquentes.
- 1372 ● Étant donnée une racine (ou un préfixe, ou un suffixe), nous obtenons la distribution
1373 des paires (préfixe,suffixe) de cette racine, et de manière générale nous analysons la
1374 distribution de paire préfixes, ou de paire de suffixes.

1375 Compte tenu des nombreux textes oraux transcrits, nous avons d'abord construit un
1376 outil qui construit la distribution des préfixes, des suffixes et des racines, c'est-à-dire le
1377 nombre d'occurrences. À partir d'une racine, nous pouvons afficher la répartition de ses af-
1378 fixés. De même, nous pouvons donner un préfixe (respectivement un suffixe) et représenter
1379 la distribution des racines et des suffixes (resp. préfixes). Nous considérons alors les distribu-
1380 tions de préfixes, suffixes et racines au second moment et construisons leurs représentations
1381 vectorielles. Notre modèle statistique combine ces trois représentations pour chaque mot
1382 composé de préfixes, d'une racine et d'un suffixe. Il introduit également un vecteur pour
1383 les phrases.

1384 7.2.1 Statistiques de base pour la langue *Amis*

1385 La répartition de tous les préfixes et suffixes, étant donné 70 textes d'*Amis* avec plus de
1386 4000 mots, est donnée dans la figure 7.1. Tous les graphiques utilisent des valeurs absolues.

1387 L'outil *Morphix* fournit une interface où une racine (resp. un préfixe ou un suffixe)
1388 peut être sélectionnée ; la distribution des préfixes et des suffixes pour une racine donnée
1389 est affichée graphiquement, comme dans la figure 7.2.

1390 Étant donnée la distribution des (préfixes;suffixes)³ de la figure 7.2, qui à ce stade ne
1391 sont pas ordonnées selon les formes et fonctions des racines, nous obtenons par projection
1392 la distribution des préfixes et suffixes de la figure 7.2 pour cette racine spécifique.

3. Un mot peut avoir plusieurs préfixes et suffixes. Dans la figure 7.2, les paires les plus fréquentes (préfixes;suffixes) sont (*ma-*;), c'est-à-dire le préfixe *ma-* sans suffixe, (*ka-*;), c'est-à-dire le préfixe *ka-* sans suffixe, (*pa-se-*;), c'est-à-dire les deux préfixes *pa-* et *se-* sans suffixe, et (*ma;ay*), c'est-à-dire le préfixe *ma-* avec le suffixe *-ay*.

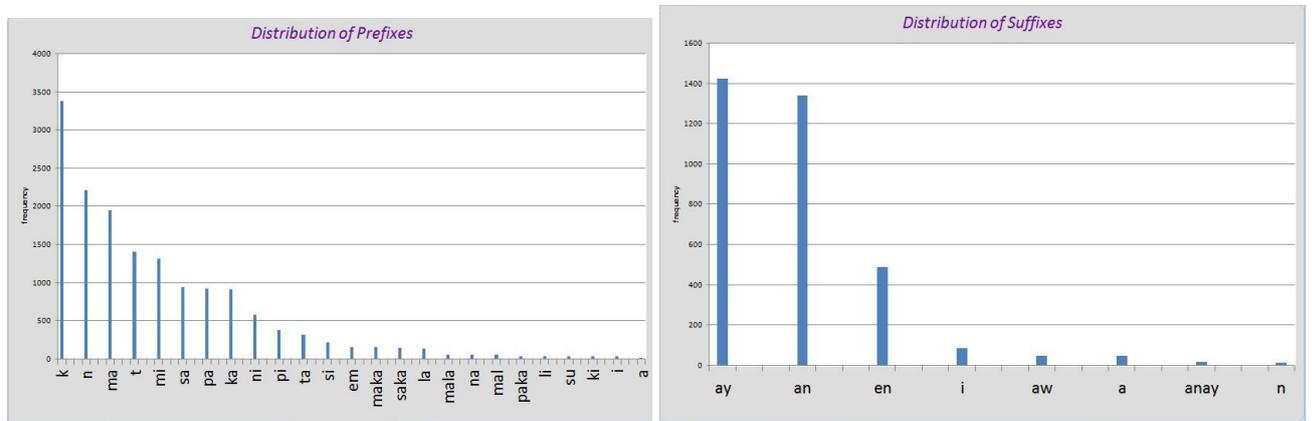


FIGURE 7.1: Préfixes et suffixes les plus fréquents.

1393 7.2.2 Représentation vectorielle des préfixes, racines et suffixes

1394 Étant donné une liste de préfixes n , la matrice de corrélation $M(n, n)$ mesure le nombre
 1395 de cooccurrences des préfixes dans une même phrase : $M(i, j)$ est le nombre d'occurrences
 1396 des préfixes i et j dans une même phrase et $M(i, i)$ est le nombre d'occurrences du préfixe
 1397 i . Comme M est une matrice semi-définie positive, il existe une matrice $U(n, n)$ telle que
 1398 $M = U.U^t$. Nous pouvons interpréter chaque préfixe i comme un vecteur v_i de dimension n
 1399 tel que $v_i.v_j = M(v_i, v_j)$. Une méthode classique *PCA* (Principal Component Analysis)
 1400 permet de réduire la dimension de ces vecteurs, le long de leurs composantes principales
 1401 définies par les grandes valeurs propres telles que $v_i.v_j \simeq M(v_i, v_j)$. La SVD (Singular
 1402 Value decomposition) produit les vecteurs propres (matrice S) et les valeurs propres λ_i
 1403 (matrice diagonale V). Si l'on projette les vecteurs propres sur les dimensions définies par
 1404 les grandes valeurs propres, nous réduisons la dimension.

1405 Considérons les 4 phrases d'*Amis* suivantes, structurées avec la racine *padang* 'help,
 1406 support'⁴ :

1407

1408 1. *Mi-padang k-u tumuk t-u suwal n-ira tatakulaq.*

1409 AV-help NOM-ART chief OBL-ART word GEN-that frog⁵

1410 The tumuk supported the words of the frog.

1411 2. *Isu Kungcu, yu ira k-u pa-padang-an,...*

1412 You Princess when exist NOM-ART RED-help-NMZ

4. La première ligne est le texte original où les préfixes et les suffixes sont identifiés. La deuxième ligne est l'analyse morphologique avec des étiquettes telles que AV, OBL,... décrites ci-dessous. La troisième ligne est la traduction.

5. Abbreviations : AV Actor Voice ; ART article ; CV conveyance voice ; GEN genitive ; IMP imperative ; INST.V instrumental voice ; LOC locative ; LV locative voice ; NFIN non-finite ; NOM nominative ; NMZ nominaliser ; OBL oblique ; PFV perfect ; PROH prohibitive ; RED reduplication ; UV undergoer voice.

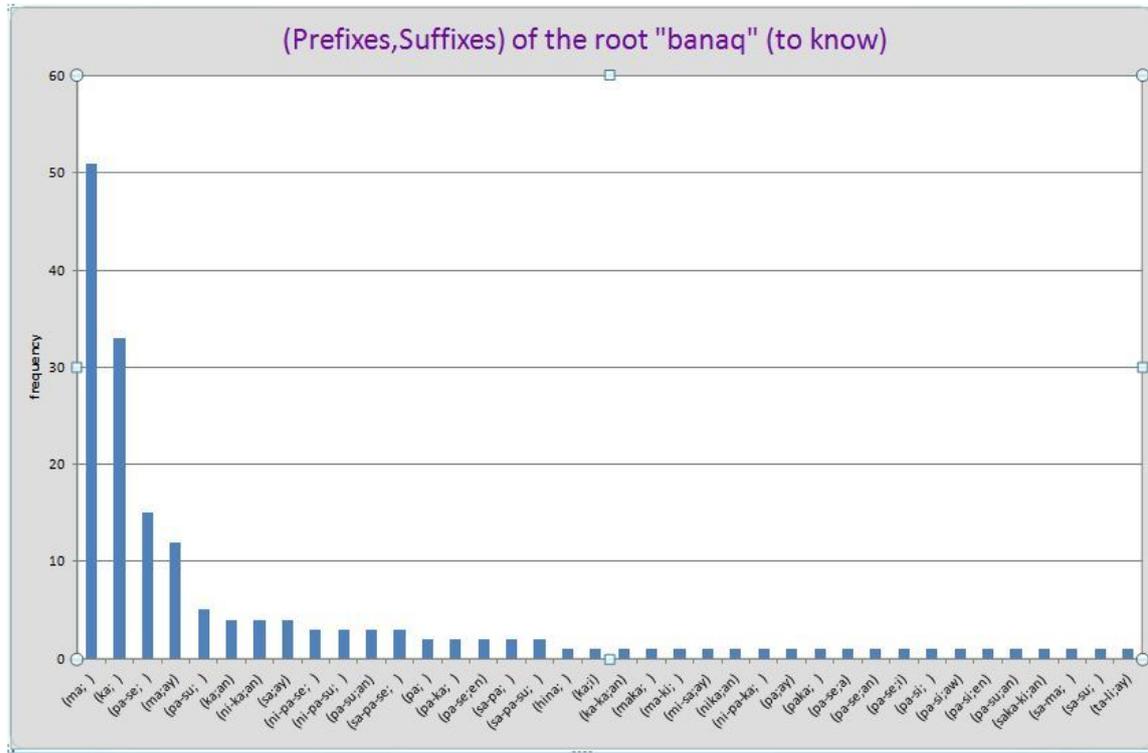


FIGURE 7.2: Les plus fréquents (préfixes ;suffixes) de la racine *banaq*. (' know').

1413 You Princess, when (you) had some help,...

1414 3. *Sulinay mi-padang k-u taw,...*

1415 Indeed AV-help NOM-ART people

1416 Indeed when people help,...

1417 4. *Aka-a ka-pawan t-u ni-padang-an n-u taw.*

1418 PROH-IMP NFIN-forget OBL-ART PFV.NMZ-help-NMZ GEN-ART people

1419 Then, you mustn't forget people's help.

1420

1421 Dans ces quatre phrases, il y a sept préfixes dans l'ordre : *k,ka,n,ni,mi,pa,t* et des suffixes *-an,*

1422 *-a*. La matrice M_p pour ces préfixes est :

$$M_p = \begin{matrix} & k & ka & n & ni & mi & pa & t \\ \begin{matrix} k \\ ka \\ n \\ ni \\ mi \\ pa \\ t \end{matrix} & \begin{pmatrix} 2 & 0 & 1 & 0 & 2 & 1 & 1 \\ 0 & 1 & 1 & 1 & 0 & 0 & 1 \\ 1 & 1 & 2 & 1 & 1 & 0 & 2 \\ 0 & 1 & 1 & 1 & 0 & 0 & 1 \\ 2 & 0 & 1 & 0 & 1 & 0 & 1 \\ 1 & 0 & 0 & 0 & 0 & 1 & 0 \\ 1 & 1 & 2 & 1 & 1 & 0 & 2 \end{pmatrix} \end{matrix}$$

1423 La première ligne de la matrice indique 2 occurrences de $k-$, 1 occurrence de la paire ($k-$, $n-$)
 1424 (première phrase), 2 occurrences de la paire ($k-$, $mi-$) (première et troisième phrases), 1 occurrence
 1425 de la paire ($k-$, $pa-$) (deuxième phrase) et 1 occurrence de la paire ($k-$, $t-$) (première phrase). Les
 1426 grandes valeurs propres de M_p sont de 2, 5 et 1, 7. Deux autres valeurs propres sont proches de
 1427 1 et les trois autres sont proches de 0. Si l'on décompose les vecteurs⁶ sur les deux plus grands
 1428 vecteurs propres, on obtient 7 vecteurs de dimension 2, un pour chaque préfixe dans la matrice
 1429 ci-dessous et dans la figure 7.3.

$$B = \begin{matrix} & & x_1 & & x_2 \\ & k & \left(\begin{array}{cc} -1.07140232 & 1.05796864 \\ -0.61507482 & -0.69205651 \\ -1.35624875 & -0.37685711 \\ -0.61507482 & -0.69205651 \\ -0.89917084 & 0.67091885 \\ -0.19923754 & 0.51352862 \\ -1.34624875 & -0.36685711 \end{array} \right) \\ & ka & & & \\ & n & & & \\ & ni & & & \\ & mi & & & \\ & pa & & & \\ & t & & & \end{matrix}$$

1430 $B * B^t$ est approximativement M_p . Le premier vecteur pour $k-$ a les coordonnées $-1.07, 1.06$.
 1431 Nous représentons graphiquement les 7 préfixes dans la figure 7.3. Une approche similaire peut
 1432 être suivie pour les suffixes et pour les racines. La figure 7.3 peut être utilisée pour prédire, à
 1433 partir d'un préfixe v , le préfixe suivant le plus probable v_{next} . C'est le vecteur v' qui maximise le
 1434 produit scalaire $|v.v'|$. Par exemple, étant donné le vecteur pour le préfixe $mi-$, le préfixe suivant
 1435 le plus probable est $k-$, dans ces 4 phrases. Pour les 95 phrases simples du corpus, les vecteurs
 1436 seraient légèrement différents.

1437 7.2.3 Distributions et vecteurs représentatifs

Soit δ la distribution des mots les plus fréquents, δ_P la distribution des préfixes les plus fréquents (resp. δ_R la distribution des racines) et soit π_p la projection qui associe le préfixe d'un mot. Par exemple, $\pi_p(mi-padang)=mi-$, un des principaux préfixes. De même, π_r associe la racine d'un mot, $\pi_r(mi-padang)=padang$. Ces distributions sont liées, principalement par des projections. Soit $\pi_p(\delta)$ la distribution telle que

$$\pi_p(\delta)(p) = \sum_{p \text{ est un préfixe de } w} \delta(w)$$

1438 Puis $\delta_P = \pi_p(\delta)$ et $\delta_R = \pi_r(\delta)$. La matrice de corrélation M_p des préfixes est également la
 1439 projection de la matrice de corrélation M des mots, c'est-à-dire $M_p = \pi_p(M)$.

1440 Pour chaque matrice de corrélation M_p, M_r, M_s des préfixes, racines et suffixes, on applique
 1441 la réduction de dimension et on obtient des vecteurs v_p de dimension n_p des préfixes, v_r de
 1442 dimension n_r pour les racines et v_s de dimension n_s pour les suffixes. La réduction importante
 1443 concerne M_r , dont la dimension est aussi grande que la taille du vocabulaire. Pour M_p, M_s , la
 1444 dimension est petite et la réduction est moins importante. Dans l'exemple précédent, la dimension

6. On utilise un package d'algèbre linéaire en Python pour obtenir la décomposition SVD et les projections.

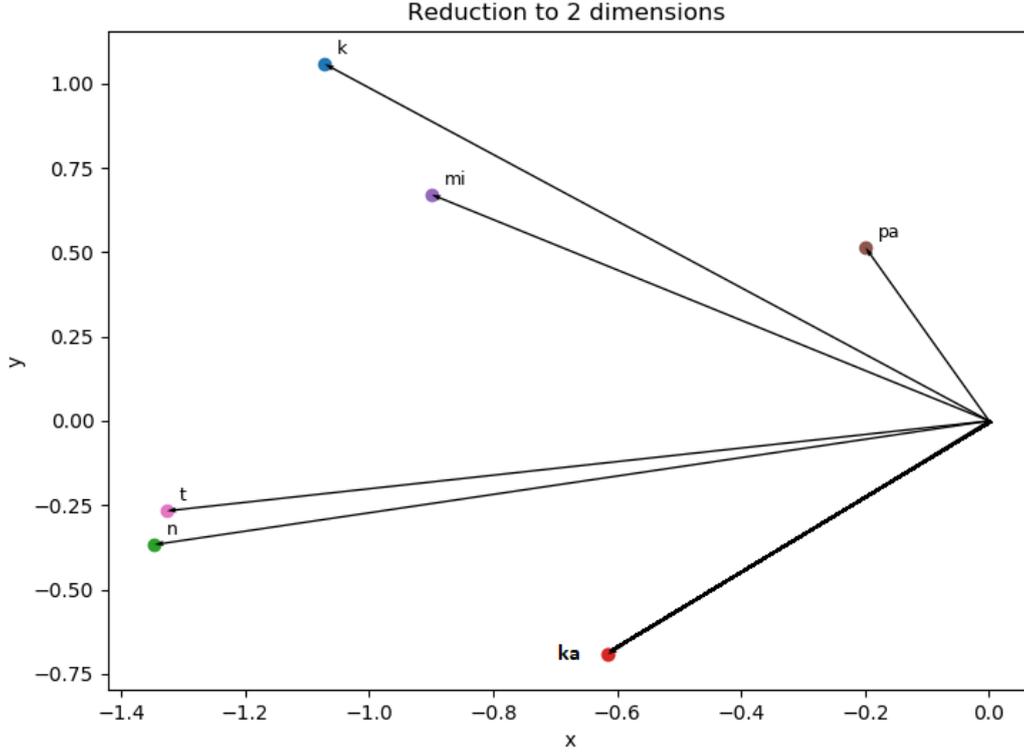


FIGURE 7.3: Les vecteurs des 7 préfixes les plus fréquents k -, ka -, n -, ni -, mi -, pa -, t - en deux dimensions.

1445 n_p des préfixes est 7 sans réduction de dimension et de 2 avec la réduction. Considérons un mot
 1446 $w=pre-root-suf$ avec un préfixe et un suffixe. Soit $v_p(pre)$ le vecteur du préfixe pre , $v_r(root)$ le
 1447 vecteur de la racine $root$ et $v_s(suf)$ le vecteur du suffixe suf . Le vecteur $v(w)$ est l'union des trois
 1448 vecteurs :

$$v(w) = \begin{pmatrix} v_p(pre) \\ v_r(root) \\ v_s(suf) \end{pmatrix}$$

Si on souhaite représenter des mots à plusieurs préfixes, on énumère les préfixes dans le vecteur $v(w)$. Pour deux mots w_i, w_j avec un préfixe, on définit :

$$\widetilde{M}(w_i, w_j) = M_p(pre_i, pre_j) + M_r(root_i, root_j) + M_p(suf_i, suf_j)$$

1449 soit la somme des corrélations des préfixes, racines et suffixes. Le point fondamental est que
 1450 pour deux mots quelconques w_i, w_j , le produit $v(w_i).v(w_j) \simeq \widetilde{M}(w_i, w_j)$. En effet, $v(w_i).v(w_j) =$
 1451 $v_p(pre_i).v_p(pre_j) + v_r(root_i).v_r(root_j) + v_s(suf_i).v_s(suf_j)$. Le produit $v_p(pre_i).v_p(pre_j)$ correspond approxi-

1452 mativement à $M_p(pre_i, pre_j)$ et de même pour les racines et les suffixes. D'où $v(w_i).v(w_j) \simeq$
1453 $\widetilde{M}(w_i, w_j)$. Pour les mots comportant plusieurs préfixes, cette approximation peut être généralisée.

1454 Notez que $\widetilde{M}(w_i, w_j)$ peut être très différent de $M(w_i, w_j)$. Il est possible que $M(w_i, w_j) = 0$,
1455 mais que ses préfixes, suffixes et racines aient de fortes corrélations, d'où $\widetilde{M}(w_i, w_j)$ peut être
1456 grand. Une théorie riche de ces vecteurs structurés peut être élaborée à l'aide de corrélations
1457 croisées, que nous n'utilisons pas pour l'instant.

1458 7.2.4 Vecteurs de phrases

1459 Un sujet important est d'étendre les vecteurs des mots aux phrases et plus généralement aux
1460 paragraphes et aux textes. Les réseaux de neurones récurrents [32, 65] ou d'autres approches
1461 basées sur l'apprentissage [59, 58] construisent un tel vecteur de variables latentes. Nous suivons
1462 une approche différente car nous voulons capturer des propriétés plus générales des phrases à l'aide
1463 des préfixes et suffixes. Définissons le vecteur probabiliste s d'une phrase comme un vecteur de
1464 dimension 5 dont les composantes sont des distributions sur certains domaines finis :

- 1465 • Valence : $\{0, 1, 2, 3\}$,
- 1466 • Voix : $\{AV, UV, LV, INST.V\}$,
- 1467 • Temps : $\{Present, Past, Future\}$,
- 1468 • Mood : $\{Indicative, Imperative, Hortative, Subjunctive\}$,
- 1469 • Force d'élocution : $\{Declarative, Negative, Exclamative\}$,

1470 La première composante s^1 , la *Valence*, est le nombre d'arguments du verbe. Nous la considérons
1471 comme une distribution sur le domaine $0, 1, 2, 3$ qui a 4 valeurs comme support. De même pour
1472 les autres composantes. Le vecteur s de dimension 5 a une taille globale de $18 = 4 + 4 + 3 + 4 + 3$.
1473 D'autres dimensions pourraient être utilisées, mais nous gardons ce type de vecteur pour simpli-
1474 fier. Si la troisième composante, une distribution sur le support $\{Present, Past, Future\}$ est le
1475 vecteur $[0, 1, 0]$, elle indique un Past (avec une probabilité de 1). Si la composante était $[\frac{1}{3}, \frac{1}{3}, \frac{1}{3}]$,
1476 elle indiquerait une répartition uniforme sur les trois temps.

1477 La figure 7.4 montre un vecteur sur la phrase :

1478 *5. mi-padang t-u suwal n-ira tatakulaq*

1479 '(he) supports the words of the frog.'

1480

Nous lisons la phrase w_1, w_2, \dots, w_n , et le vecteur $v_i = v(w_i)$ associé avec chaque mot w_i , défini dans la section précédente. Soit :

$$s_i = F(s_{i-1}, v_i)$$

1481 avec s_0 un état initial (vecteur zéro) et F une fonction. Nous construisons F par cas suivant les
1482 préfixes et les suffixes. Nous décrivons des règles syntaxiques plus avancées de l'*Amis* dans la
1483 section 7.3.

1484 7.3 Un aperçu syntaxique de la langue *Amis*

1485 L'ordre de base de l'*Amis* commence par le verbe. Les arguments sont marqués par les cas :
1486 nominatif (k -), génitif (n -), oblique (t -) [19]. Les affixes de voix (AV) *mi*-, (NAV) *ma*-, (UV) *ma*-,

Valence : {0,1,2,3}	[0,0,1,0]
Voix : {AV,UV,LV,INST.V}	[1,0,0,0]
Temps : {Present,Past,Fut}	[.9,.05,.05]
Mood : {Ind,Imp,Hort,Sub}	[1,0,0]
Force d'élocution : {Decl,Neg,Exclam}	[1,0,0]

FIGURE 7.4: Vecteur d'une phrase

1487 identifient également des classes de verbes ; l'appartenance à une classe peut être exclusive ou
 1488 permettre des alternances de voix avec une sémantique différente. Les verbes AV *mi-* désignent des
 1489 activités. Les verbes UV *ma-* désignent des réalisations accomplies par un agent sur des patients
 1490 spécifiques et totalement affectés ; les verbes NAV *ma-* comprennent des états, des propriétés. Le
 1491 système de voix est basé sur la co-sélection d'un argument nominatif (le sujet), et d'un affixe
 1492 de voix dont la sémantique correspond à celle du sujet nominatif. Les voix AV *mi-*, NAV *ma-*
 1493 et UV *ma-* sont limitées à des phrases déclaratives. Les phrases non déclaratives (telles que les
 1494 phrases négatives, impératives, hortatives) ont différentes formes. Comparons *ma-butiq cira* '(s)he
 1495 is asleep/sleeping' et *ka-butiq!* 'go to sleep!'.

1496 **Transitivité et alignement**⁷. Les verbes AV *mi-* et les verbes NAV *ma-* (Non-Actor Voice)
 1497 ont un alignement de type "accusatif" avec un argument marqué par *t-* comme dans (6) et (8). Le
 1498 sujet des verbes *mi-* est un acteur, tandis que le sujet des verbes "NAV" *ma-* est un non-acteur
 1499 (c'est-à-dire un thème ou une expérience, le siège d'une propriété ou d'un état). D'autre part,
 1500 les verbes UV *ma-* transitifs ont un alignement ergatif avec un sujet patient et un agent génitif
 1501 comme dans (7). L'ordre des mots (V-S-0 ou V-S-A) est donné pour chaque exemple :

1502

1503

1504 6. *Mi-melaw k-u wawa t-u tilibi.*
 1505 AV-look NOM-ART child OBL-ART TV V-S-O

1506 'The child is watching TV.'

1507 7. *Ma-melaw n-uhni k-u teker.*
 1508 UV-look GEN-3pl NOM-ART trap V-A-S

1509 'They saw the trap.'(lit. the trap was seen by him)

1510 8. *Ma-hemek k-aku t-u babainay. (*mi-)*
 1511 NAV-admire NOM-1sg OBL-ART boy V-S-O

7. L'alignement fait référence au type de marquage des arguments du verbe. Dans les langues accusatives, le sujet nominatif est marqué différemment de l'objet accusatif. Dans les langues ergatives, l'agent des verbes transitifs est marqué comme ergatif, tandis que le sujet des verbes intransitifs et le patient des verbes transitifs sont marqués de la même façon comme nominatif/absolutif.

1512 'I admire the boy.'
 1513 Toutes les autres voix s'alignent de manière ergative, comme UV *ma-*, donc UV *-en*, LV *-an*,
 1514 c'est-à-dire qu'elles ont un sujet nominatif correspondant à la sémantique de l'affixe de voix
 1515 (c'est-à-dire un patient, un lieu), et un agent génitif (si exprimé).

1516 7.4 Classification des phrases à partir des préfixes et suffixes

1517 Les préfixes et suffixes apparaissant dans une phrase fournissent certaines informations syn-
 1518 taxiques. Ici, nous limitons l'analyse à 95 phrases simples, principalement sous forme de verbes
 1519 AV et UV, données en annexe A. Nous codons chaque phrase comme dans la figure 7.5 ci-dessous,
 1520 avec ses préfixes, les racines et les suffixes. Nous demandons ensuite s'il est possible de récupérer
 1521 la syntaxe d'une phrase en observant uniquement les préfixes et les suffixes. Les racines lexicales
 1522 des mots sont remplacées par un "—" comme dans la figure 7.5.

Sentences	Sentences Encoding
Ma-pa-tangasa n-umisu k-aku iri.	Ma-pa— n— k— —.
Adihay k-u ni-urung t-u kidudung.	— k— ni— t— —.

FIGURE 7.5: Les phrases et leurs encodages.

1523

1524 Les *simples clauses* considérées appartiennent à l'un des modèles syntaxiques suivants, d'une
 1525 valeur de 10 :

1526 On considère 10 types de *clauses simples* :

1527

1528 . *V-A-S-O.* Verb-Agent-Subject-Object

1529 . *V-A-S.* Verb-Agent-Subject

1530 . *V-A-O.* Verb-Agent-Object

1531 . *V-A.* ...

1532 . *V-S-O.*

1533 . *V-S-A.*

1534 . *V-O-S.*

1535 . *V-S.*

1536 . *V-O.*

1537 . *V.*

1538

1539 Toutes les 95 phrases simples ont été étiquetées comme appartenant à l'une des 10 classes.

1540 Nous construisons maintenant un arbre de décision simple qui lit le codage d'une phrase et prédit
 1541 approximativement sa classe. L'arbre de décision présentés dans la figure 7.7 décide d'abord si

1542 la phrase contient un préfixe *ma-* ou un préfixe *mi-* ou un préfixe *k-* sans préfixe *ma-* ou *mi-*.

1543 L'arbre détecte ensuite la présence ou l'absence des préfixes de cas *n-*, *k-* et *t-*. Une branche

1544 avec l'étiquette *n-* suppose la présence du préfixe *n-* et avec l'étiquette $\neg n-$ suppose l'absence du

1545 préfixe *n-*. Chaque feuille de l'arbre est étiquetée avec l'une des classes syntaxiques.

1546 Les statistiques qui représentent la fréquence de ces 10 classes différentes parmi le corpus de
1547 95 phrases, sont données dans la figure 7.6.

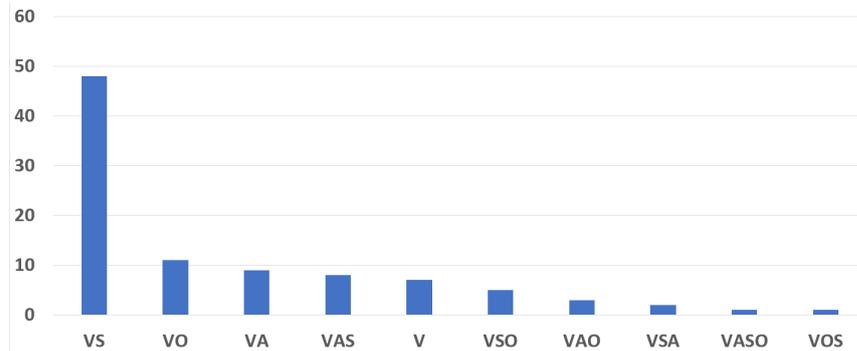


FIGURE 7.6: Fréquences des différentes classes syntaxiques existantes dans le corpus.

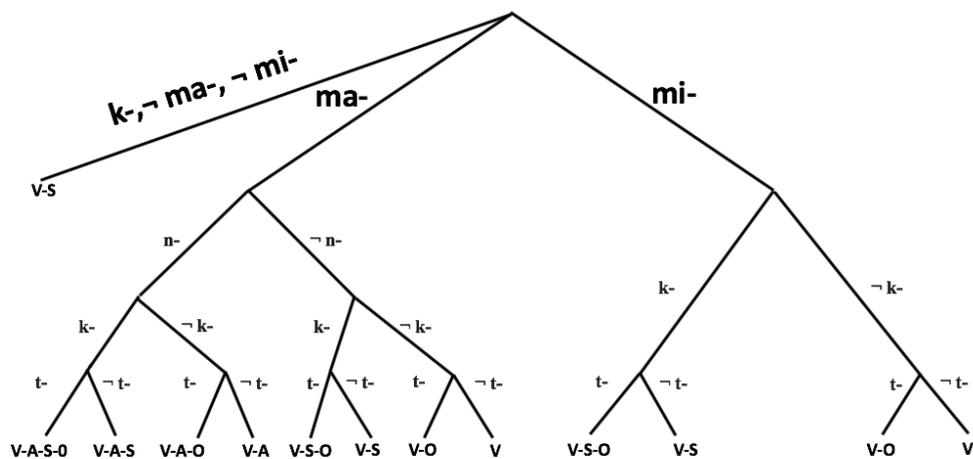


FIGURE 7.7: Arbre de décision pour prédire la classe syntaxique.

1548 Dans les phrases de l'annexe A, il y a 95 phrases affirmatives simples et 4 phrases négatives.
1549 Parmi les 95 phrases affirmatives, 22 phrases n'affichent pas les préfixes *ma-* ou *mi-* ou *k-*, c'est
1550 pourquoi 73 phrases ont été classées par l'arbre de décision. La prédiction est considérée comme
1551 *moyenne* si la classe prédite est à une distance de 1 pour la distance d'édition avec la classe exacte,
1552 *mauvaise* si la classe prédite est à une distance supérieure à 1 de la classe exacte. Par exemple,
1553 *V-S-O* est une prédiction *moyenne* de la classe *V-A-S-O* car elle se trouve à une distance de 1.
1554 *V-S* est une mauvaise prédiction de *V-A-S-O* car elle se trouve à une distance de 2. Nous avons
1555 obtenu une bonne prédiction dans 25 cas, une prédiction moyenne dans 42 cas et une mauvaise
1556 prédiction dans 6 cas. Les statistiques de prédiction sur les 95 phrases sont données dans la figure
1557 7.8.

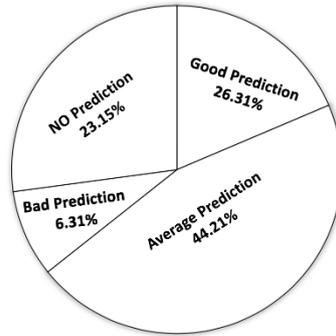


FIGURE 7.8: Statistiques des prédictions.

1558 Ce simple arbre de décision classe 77% des 95 phrases présentées en annexe A. La plupart du
1559 temps, la classification est de qualité moyenne : elle ne donne qu’une approximation de la classe.

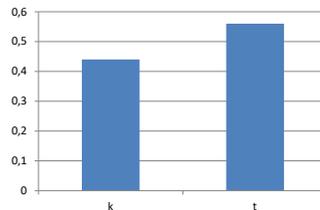
1560 7.4.1 Vecteur de phrases et mécanismes d’attention

1561 La valence, première composante du *vecteur de la phrase*, est définie pour chaque verbe (voir
1562 section 4). Les 4 autres composantes sont partiellement définies par les préfixes et suffixes et sont
1563 donc des distributions probabilistes. Comme le nombre total de préfixes et de suffixes est faible,
1564 la description de la fonction F est essentiellement une définition par cas.

1565 Les mécanismes d’attention [40, 60] présentés dans le chapitre 3.4 prévoient pour un mot donné
1566 w_i d’une phrase, la distribution des autres mots les plus corrélés. La corrélation d’un mot w_j avec
1567 w_i est approximativement la valeur $v(w_i).v(w_j)$. On peut alors calculer $\{v(w_i).v(w_j) : j \neq i\}$
1568 pour un w_i fixe et normaliser les valeurs pour obtenir une distribution. Le mécanisme d’attention
1569 est surtout utilisé pour la traduction automatique.

1570 Nous pouvons nous concentrer sur les préfixes et les suffixes et définir de la même façon
1571 l’attention des préfixes et des suffixes comme le produit scalaire des vecteurs associés aux préfixes
1572 ou aux suffixes. Dans l’exemple de la figure 7.9, considérons la phrase : ” *Mi-melaw k-u wawa*
1573 *t-u tilibi*,” avec les préfixes $mi-$, $k-$, $t-$. Calculons l’attention pour le préfixe $mi-$. Nous calculons
1574 $v_p(mi).v_p(k) = 1, 149$ et $v_p(mi).v_p(t) = 1, 446$, en utilisant les vecteurs donnés dans la matrice B
1575 de la section 7.2.2. Avec une normalisation, on obtient la distribution de la figure 7.9.

1576 Pour les préfixes importants tels que $mi-$ et $ma-$, la répartition de l’attention est importante
1577 et codifie certaines propriétés intentionnelles d’une phrase.



Attention of the prefix mi-

FIGURE 7.9: Distribution de l'attention pour *mi-* dans la phrase *Mi-melaw k-u wawa t-u tilibi*.

1578 7.5 Grammaires et statistiques

1579 Nous étudions maintenant l'impact des préfixes et suffixes d'une phrase sur les arbres de
 1580 dérivation possibles, définis par les analyseurs syntaxiques des grammaires classiques. Étant donné
 1581 une grammaire et une phrase, quels sont les arbres de dérivation les plus probables ? Il s'agit d'une
 1582 question centrale pour tout analyseur, y compris les analyseurs de dépendance [18] construits par
 1583 apprentissage à partir de données étiquetées. Dans ce cas, la distribution des exemples définit un
 1584 espace probabiliste.

1585 Un problème fondamental est d'associer à une séquence de symboles un arbre de dérivation
 1586 possible, ce problème existe dans plusieurs disciplines : en biologie le repliement de structure *ARN*
 1587 (Acide RiboNucléique)⁸, en informatique l'arbre de décomposition d'une phrase. Un repliement
 1588 de l'*ARN* est un arbre de dérivation possible avec une énergie potentielle et le repliement le
 1589 plus probable a un niveau d'énergie minimum. Dans notre contexte, nous disposons de diverses
 1590 statistiques sur les préfixes et les suffixes, et l'entropie de chaque distribution, obtenues à partir
 1591 des données linguistiques. Nous n'avons pas de définition simple de l'énergie associée à un arbre
 1592 et nous nous tournons vers l'interprétation linguistique, donnée en premier lieu dans la section
 1593 7.3.

1594 Une grammaire G pour les mots représentés par la structure *préfixe-racine-suffixe* peut être
 1595 représentée par des règles du type

Non – terminal \rightarrow regular expression

1596 Un exemple typique d'utilisation de Non-terminaux⁹ qui comprennent des mots, des préfixes,
 1597 des racines et des suffixes est :

$$1598 \quad S \rightarrow VP.KP + VP.KP^*$$

1599

8. Un repliement transforme une séquence linéaire de lettre en une structure avec des liens, un arbre qui est lui-même ensuite transformé en une structure tridimensionnelle.

9. KP stands for Case Phrase, DP stands for Determiner Phrase, K stands for Case, PossP stands for Possessive Phrase. K and Voice represent Case prefixes and Voice prefixes respectively.

1600 $VP \rightarrow Voice.V.KP^*$
 1601 $KP \rightarrow K.DP$
 1602 $DP \rightarrow D.N + D.N.PossP$
 1603 $PossP \rightarrow K.DP$

1604 Un autre ensemble de règles énumère les mots et les affixes possibles.

1605 $K \rightarrow t^- + \dots$
 1606 $V \rightarrow \textit{padang} + \dots$
 1607 $Voice \rightarrow \textit{mi}^- + \dots$
 1608 $N \rightarrow \textit{suwal} + \dots$
 1609 $D \rightarrow u + \dots$
 1610 $PossP \rightarrow \textit{n-ira} + \dots$

1611 Chaque symbole non terminal S, VP, KP, \dots génère une expression régulière simple pour la
 1612 première partie. Certains symboles non terminaux K, V, \dots énumèrent tous les mots possibles
 1613 pour la deuxième partie. On peut générer une phrase appliquant plusieurs règles à partir de la
 1614 racine S . L'arbre de dérivation est la construction inverse. Il existe plusieurs arbres de dérivation
 1615 possibles associés à une phrase, telle que

1616 $\textit{mi-padang t-u suwal n-ira tatakulaq,}$
 1617 '(he) supports the words of the frog'

1618 décrits dans la figure 7.10, et nous discutons de la meilleure représentation.

1619

1620 7.5.1 Les meilleurs arbres de dérivation

1621 Étant donné un *vecteur de phrase* s , nous pouvons alors décider que l'arbre de dérivation
 1622 (a) de la figure 7.10 est mieux adapté que l'arbre (b) pour la phrase présentée dans la section
 1623 précédente. Nous reprenons en partie l'explication des verbes de la section 7.3.

1624 Nous savons qu'un verbe dont le préfixe est *mi-* assigne un nominatif à l'Acteur 'voix' AV. La
 1625 deuxième composante du vecteur s est : $[1, 0, 0, 0]$. Par conséquent, la structure VP de l'arbre de
 1626 dérivation (a) encode tous les arguments alors qu'ils sont séparés dans l'arbre de dérivation (b).
 1627 Si on a un Acteur 'voix', il faut dériver le verbe avec son KP . Si on a le préfixe *mi-*, l'action doit
 1628 être attachée au verbe. Par conséquent, l'arbre de dérivation (a) est une meilleure représentation.

1629 7.5.2 Comparaison avec d'autres modèles stochastiques

1630 Dans une grammaire stochastique [41], les dérivations avec le même symbole non-terminal
 1631 ont une probabilité p telle que la somme des probabilités pour chaque symbole non-terminal est
 1632 de 1. L'espace probabiliste associé à l'ensemble (s, t) , où s est une phrase et t est un arbre de
 1633 dérivation, est le produit des probabilités des règles utilisées selon l'arbre, noté $p(s, t)$. Étant
 1634 donnée une phrase, une tâche classique consiste à prédire l'arbre de dérivation le plus probable,
 1635 et elle peut être réalisée en $O(n^3)$ pour une phrase de n mots .

1636 Dans notre contexte, l'espace probabiliste est entièrement différent. Les vecteurs structurés
 1637 des mots nous permettent de prédire le mot, préfixe ou suffixe le plus probable, compte tenu
 1638 du contexte des mots précédents. Pour une phrase donnée, les positions des préfixes et suffixes

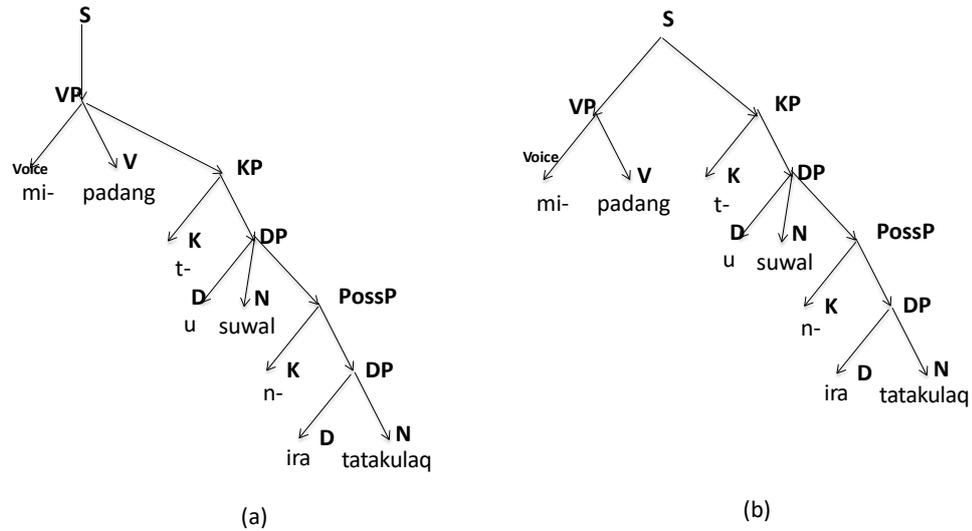


FIGURE 7.10: Arbre de dérivation de la phrase *mi-padang t-u suwal n-ira tatakulaq* pour une grammaire G .

1639 déterminent les distributions probabilistes du vecteur s introduite dans la section 7.2.4. Nous
 1640 examinons l'arbre de dérivation préférable, compte tenu de cette distribution des composantes
 1641 sémantiques.

1642 7.5.3 Conclusion

1643 Nous proposons un modèle statistique pour la morphologie, inspiré par la langue *Amis*. La
 1644 morphologie d'un mot w est la structure *préfixe-racine-suffixe*. Nous avons construit les distri-
 1645 butions classiques des préfixes, racines et suffixes les plus fréquents, et les distributions corres-
 1646 pondantes avec une racine, un préfixe ou un suffixe possible. On construit des vecteurs : $v_p(pre)$
 1647 pour le préfixe pre , $v_r(root)$ pour la racine $root$ et $v_s(suf)$ pour le suffixe suf . Le vecteur du mot
 1648 $v(w)$ est la concaténation de ces trois vecteurs. Nous avons défini un vecteur probabiliste associé
 1649 à une phrase. L'analyse des préfixes et suffixes, détermine la plupart des composantes du vecteur
 1650 de la phrase S , de manière probabiliste. Nous avons montré comment prédire la classe syntaxique
 1651 d'une phrase simple à partir du vecteur probabiliste, en utilisant les arbres de décision. Étant
 1652 donné une grammaire G et une phrase $S : w_1, w_2, \dots, w_n$, nous avons montré comment trouver
 1653 l'arbre de décomposition syntaxique préférable en utilisant le vecteur de la phrase. Toutes les
 1654 prédictions sont approximatives.

1655 8 Recherche et explications

1656 Dans ce chapitre, nous analysons plusieurs flux de textes ou Twitter et utilisons leurs corrélations
1657 de contenu, introduites dans la section 5.4.2.3. La distance entre deux flux est l'inverse de leur
1658 corrélation, de même pour la distance entre deux clusters. Ces distances vont nous permettre de
1659 définir une distance *dist* entre deux mots, dans le contexte de plusieurs flux.

1660 Nous introduisons alors le concept de *Recherche dans plusieurs flux* : c'est la généralisation
1661 d'un moteur de recherche, comme Google, dans le contexte de plusieurs flux dynamiques. Etant
1662 donnés des mots clés, $\sigma_1, \sigma_2, \dots, \sigma_l$, on va rechercher les mots σ qui minimisent $\sum_{i=1, \dots, l} dist(\sigma, \sigma_i)$.
1663 Chaque solution σ aura une trace qui permettra de suivre les composantes géantes et les flux,
1664 qui sera le témoin du minimum : ces sont les explications.

1665 A partir d'une matrice de distance entre les flux, nous pouvons utiliser les méthodes de
1666 phylogénie qui construisent un arbre de décomposition dont les feuilles sont les flux, avec des
1667 coûts sur les arêtes. Nous étendons cet arbre en un DAG (Directed Acyclic Graph) en plaçant un
1668 noeud pour chaque composante géante d'un flux et une arête entre le noeud de la composante et
1669 le noeud du flux. Nous plaçons ensuite un noeud pour chaque mot ou tag et une arête entre ce
1670 mot et chaque composante géante où ce mot se trouve. Les distances entre mots seront les plus
1671 courts chemins dans ce graphe.

1672 Nous stockons les composantes géantes de chaque fenêtre, pour chaque flux, plus exactement
1673 le 2-core de chaque composante, comme un ensemble de noeuds et un ensemble d'arêtes. Chaque
1674 cluster est le résultat d'un tirage uniforme ou pondéré réalisé par un Réservoir qui nous permet
1675 cependant d'approximer les distances.

1676 Dans la première section, nous rappelons la définition de la corrélation de flux. Dans la
1677 deuxième section, nous décrivons l'arbre phylogénique entre les flux. Dans la troisième section,
1678 nous étendons l'arbre en un DAG et définissons la distance entre deux mots, dans le contexte de
1679 flux. Dans la quatrième, nous présentons le moteur de recherche dans ce cadre. Dans la cinquième
1680 section, nous décrivons les explications disponibles pour les résultats de ce moteur.

1681 8.1 Corrélation de flux

Nous avons introduit cette notion dans la section 5.4.2.3, pour estimer la corrélation entre deux flux Twitter. La notion se généralise aux flux de textes, puisque nous suivons la même approche. Pour un flux i et un graphe G_i associé, soit $C_{i,j}$ le j -ème (γ, δ) -cluster de C_i et

$$C_i = \bigcup_j C_{i,j}$$

1682 l'ensemble des (γ, δ) -clusters du graphe G_i , pour $i = 1$ ou 2 .

1683 **Definition 8.** La corrélation de deux flux définis par G_1, G_2 est le coefficient $\rho = J(C_1, C_2)$.

1684 Il est alors possible d'approximer la corrélation en utilisant les composantes géantes des
1685 graphes. Il suffit de considérer le 2-core de chaque composante géante. En effet, comme le montre
1686 [47], le 2-core de la composante géante du Réservoir est une bonne approximation, au sens en-
1687 sembliste, de chaque $C_{i,j}$ et donc de l'union.

1688 Dans la suite, nous utilisons plutôt la distance entre deux flux.

1689 **Definition 9.** La distance $dist(C_1, C_2)$ entre deux flux définis par G_1, G_2 est le coefficient $1 - \rho$.

1690 Ces définitions se généralisent aux graphes G_1, G_2 définis par des fenêtres dans les flux.

1691 8.2 Phylogénie entre flux

1692 Supposons que nous observions plusieurs flux. Les distances et corrélations vont varier au
1693 cours du temps. Pour la fenêtre à l'instant t , soit A_t la matrice de distance entre les flux i et j :

$$A_t(i, j) = dist(i, j)$$

1694 La méthode Neighbor Joining [56], construit un arbre T avec des coûts sur les arêtes, de
1695 telle sorte que chaque flux apparaisse comme une feuille dans T et que la distance $d'(i, j)$ entre
1696 deux feuilles dans l'arbre soit approximativement la distance définie par la matrice A_t . Cette
1697 construction de l'arbre phylogénique suppose une propriété additive des distances, mais il existe
1698 toujours une solution approximative.

1699

1700 En quoi la phylogénie est-elle plus intéressante que la matrice de corrélation? Etant donné
1701 un flux i , i.e. une feuille de l'arbre, on peut trouver les flux proches plus efficacement en utilisant
1702 l'arbre phylogénique qu'en utilisant la matrice de distance. En effet on trouve un flux proche sans
1703 parcourir toute la matrice.

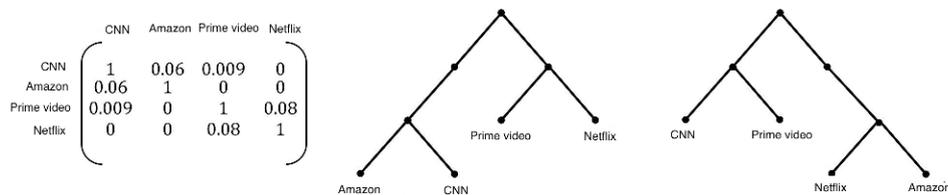


FIGURE 8.1: Arbre phylogénique T (centre) à partir de la matrice de distance, et un autre arbre proche T' (droite)

1704 Avec la matrice au temps t , on construit l'arbre T . Plus tard, nous avons une matrice différente
 1705 A' et un arbre différent T' comme dans la Figure 8.1. Étant donné un flux, les voisins d'une feuille
 1706 de T sont en général les feuilles les plus proches dans T , et nous utilisons cette propriété dans la
 1707 recherche par corrélation.
 1708

1709 8.3 Distance entre mots, relative à des flux

1710 La figure 8.2 indique comment compléter l'arbre phylogénétique introduit à la section précédente
 1711 8.2 pour définir un DAG dont les nœuds représentent les flux, les clusters et les mots. Chaque
 1712 flux est une feuille de l'arbre phylogénétique : on place tous les clusters d'un même flux comme des
 1713 nœuds connectés au flux par des arêtes de poids défini ci-dessous. Les mots d'un cluster sont des
 1714 nœuds connectés aux clusters dont ils font partie, à un niveau encore inférieur, connectés par des
 1715 arêtes de poids défini ci-dessous.
 1716

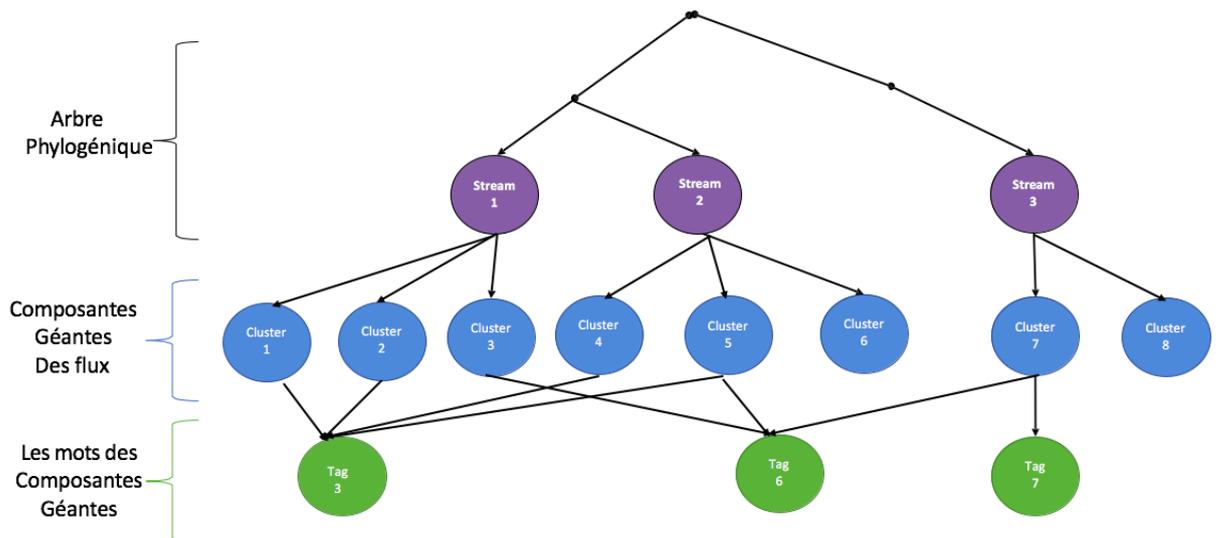


FIGURE 8.2: DAG de recherche qui regroupe les flux, les clusters et les mots.

On cherche des coûts sur les arêtes entre 0 et 1, qui privilégient les clusters le plus récents et qui tendent vers 1 quand l'âge du cluster tend vers l'infini. La fonction

$$f(t) = \frac{t}{t + 2}$$

1717 où $t = 0, 1, 2, \dots$ est l'âge de la fenêtre. La fenêtre courante correspond à $t = 0$ et $f(t)$ vaut zéro.
 1718 Pour $t = 1, 2, 3, \dots$, les fenêtres précédentes (décalée de $t * \lambda$), la fonction $f(t)$ vaut $1/3, 1/2, 3/5, \dots$
 1719 et tend vers 1 quand t croît.

1720 Les coûts des arêtes du graphe de la figure 8.2 combinent deux critères : privilégier les clusters
1721 récents et les noeuds de grand degré. Ils sont définis comme suit :

- 1722 — Le coût des arêtes de l'arbre phylogénique suit la construction de l'arbre phylogénique,
- 1723 — Le coût des arêtes entre les noeuds de flux et les noeuds de clusters est de $f(t)$ si t est l'âge
1724 du cluster,
- Le coût entre les mots et les clusters est inversement proportionnel aux degrés d des noeuds
dans le cluster et dépend de l'âge de cluster :

$$\frac{1}{d} + \frac{f(t)}{2}$$

1725 Le coût $\text{dist}_\mu(\sigma_1, \sigma_2)$ d'un chemin μ entre σ_1 et σ_2 est la somme des coûts des arêtes du
1726 chemin. Nous avons 3 cas possibles pour deux tags σ_1, σ_2 , qui sont illustrés par la figure 8.2 :

- 1727 — Cas 1 : σ_1, σ_2 sont dans le même flux et dans les mêmes clusters
- 1728 — Cas 2 : σ_1, σ_2 sont dans le même flux mais dans deux clusters différents
- 1729 — Cas 3 : σ_1, σ_2 sont dans deux flux différents et dans deux clusters différents

1730 Les tags σ_1 et σ_2 de degrés d_1 et d_2 appartiennent à des composantes géantes, C_{1,t_i} et C_{2,t_j}
1731 dans des flux f_1 et f_2 . Soit d_T la distance entre f_1 et f_2 dans l'arbre phylogénique. Si on fixe un
1732 chemin μ entre σ_1 et σ_2 la distance est :

$$\text{dist}_\mu(\sigma_1, \sigma_2) = \begin{cases} \frac{1}{d_1} + \frac{1}{d_2} + f(t) & \text{si } C_{1,t_i} = C_{1,t_j} \wedge f_1 = f_2 \\ \frac{1}{d_1} + \frac{1}{d_2} + 3/2 * f(t_i) + 3/2 * f(t_j) & \text{si } C_{1,t_i} \neq C_{1,t_j} \wedge f_1 = f_2 \\ \frac{1}{d_1} + \frac{1}{d_2} + 3/2 * f(t_i) + 3/2 * f(t_j) + d_T & \text{si } C_{1,t_i} \neq C_{1,t_j} \wedge f_1 \neq f_2 \end{cases}$$

1736

La distance entre deux tags σ_1 et σ_2 est :

$$\text{dist}(\sigma_1, \sigma_2) = \text{Min}_\mu \text{ dist}_\mu(\sigma_1, \sigma_2)$$

1737 Cette définition se généralise à deux noeuds arbitraires du DAG. Nous pouvons vérifier que la
1738 fonction dist est une fonction binaire symétrique qui satisfait l'inégalité triangulaire.

1739

1740 8.4 Moteur de recherche par similarité

1741 Nous présentons un algorithme de recherche pour trouver des tags proches d'une liste de tags
1742 $\sigma_1, \dots, \sigma_l$.

1743 **Recherche par similarité** ($\sigma_1, \dots, \sigma_l$) :

- Trouver les mots σ qui minimisent :

$$\sum_{i=1, \dots, l} \text{dist}(\sigma, \sigma_i)$$

1744 Etant donné un tag σ , on peut calculer par l'algorithme de Dijkstra les mots les plus proches
1745 et les plus courts chemins vers $\sigma_1, \sigma_2, \dots, \sigma_l$ et donc $\sum_{i=1, \dots, l} \text{dist}(\sigma, \sigma_i)$. Nous utilisons cette
1746 procédure comme une boîte noire.

1747 8.4.1 Algorithme de recherche

1748 **Algorithme de recherche**($\sigma_1, \dots, \sigma_l$) :

- 1749 — Pour chaque mot (σ_i), $i = 1, \dots, l$, soit C_{i,t_i} le cluster le plus récent qui contient σ_i ,
- 1750 — Soit $\sigma_{d_{max,i}}$ le nœud de plus grand degré (le centre) de la composante C_{i,t_i} ,
- 1751 — Pour chaque $\sigma_{d_{max,i}}$, on calcule $\sum_{j=1, \dots, l} dist(\sigma_{d_{max,i}}, \sigma_j)$ à l'aide de la boîte noire,
- 1752 — Les réponses sont les $\sigma_{d_{max,i}}$ qui minimisent la somme précédente.

1753 Notons que si σ est une solution optimale, ce n'est pas nécessairement le centre d'une com-
 1754 posante. En effet, σ est dans un cluster C mais son centre c n'est pas nécessairement la solution
 1755 minimum. La somme $\sum_{j=1, \dots, l} dist(c, \sigma_j)$ peut être supérieure à la somme $\sum_{j=1, \dots, l} dist(\sigma, \sigma_j)$,
 1756 selon les positions des arêtes du centre vers les clusters C_{i,t_i} .

1757 Nous allons montrer que les solutions $\sigma_{d_{max,i}}$ sont des solutions approchées dans le sens où
 1758 $\sum_{j=1, \dots, l} dist(\sigma_{d_{max,i}^*}, \sigma_j)$ est plus petit que l'optimum multiplié par une constante c , par exemple
 1759 $c = 2$.

1760 8.4.2 Justification de l'algorithme de recherche

1761 Etant donné un tag σ , on peut calculer par l'algorithme de Dijkstra, aussi appelé A^* , les
 1762 plus courts chemins vers $\sigma_1, \sigma_2, \dots, \sigma_l$ et donc $\sum_{i=1, \dots, l} dist(\sigma, \sigma_i)$. Il existe trois types de chemin
 1763 μ entre deux tags :

1764

- 1765 • Chemins de type 1 : tag-cluster-tag-cluster-tag : $\sigma_1, C_1, \sigma_2, C_2, \sigma_3$, avec σ_2 un noeud in-
 1766 termédiaire. Pour un chemin de type 1, soit d'_1 le degré de σ_2 dans C_1 et d_2 le degré de σ_2 dans
 1767 C_2 , et soit d'_2 le degré de σ_3 dans C_2 . Le coût total est de :

$$\begin{aligned} & f(t_1)/2 + 1/d_1 + f(t_1)/2 + 1/d'_1 + f(t_2)/2 + 1/d_2 + f(t_2)/2 + 1/d'_2 \\ & = f(t_1) + 1/d_1 + 1/d'_1 + f(t_2) + 1/d_2 + 1/d'_2 \end{aligned}$$

En général pour un chemin qui a $i - 1$ noeuds intermédiaires, le coût sera de :

$$\sum_i f(t_i) + \sum_i (1/d_i + 1/d'_i)$$

- 1768 • Chemins de type 2 : tag-cluster-stream-cluster-tag $\sigma_1, C_1, s, C_2, \sigma_2$ pour des clusters du
 1769 même stream s . Le coût est de :

$$\begin{aligned} & f(t_1)/2 + 1/d_1 + f(t_1) + f(t_2) + 1/d_2 + f(t_2)/2 \\ & = 3/2.(f(t_1) + f(t_2)) + (1/d_1 + 1/d_2) \end{aligned}$$

- 1770 • Chemins de type 3 : tag-cluster-stream-phylogénie-stream-cluster-tag, $\sigma_1, C_1, s_1, T, s_2, C_2, \sigma_2$,
 1771 pour deux flux s_1, s_2 et T l'arbre phylogénique. Le coût est de :

$$\begin{aligned} & f(t_1)/2 + 1/d_1 + f(t_1) + f(t_2) + 1/d_2 + f(t_2)/2 + d_T(s_1, s_2) \\ & = 3/2.(f(t_1) + f(t_2)) + (1/d_1 + 1/d_2) + d_T(s_1, s_2) \end{aligned}$$

1772 Le coût générique entre deux tags σ_1 et σ_2 est le chemin de type 2 ou 3 où on prend les
1773 clusters les plus récents. Pour un chemin de type 1, la somme $\sum_i f(t_i)$ croît comme la longueur
1774 chemin et dépassera la coût d'un chemin de type 2 ou 3.

1775 **Lemme 5.** *Si un chemin μ de type 1 entre deux tags σ_1 et σ_2 est meilleur que le chemin générique,*
1776 *alors il est court.*

Démonstration. Un chemin de type 1 de longueur 6, avec deux noeuds intermédiaires a un coût de :

$$\sum_{i=1,2,3} f(t_i) + \sum_{i=1,2,3} (1/d_i + 1/d'_i)$$

1777 Pour un modèle où tous les $f(t_i)$ ont une même moyenne, le chemin générique de type 2 a un
1778 coût inférieur, $3/2.(f(t_1) + f(t_2)) + (1/d_1 + 1/d_2)$, d'un facteur environ 1/2. Pour un chemin de
1779 type 3, le coût additionnel $d_T(s_1, s_2)$ est constant. Il existe donc une longueur $2.p$ telle que le
1780 coût d'un chemin de type 1 de longueur $2.p$ sera supérieur au coût du chemin de type 3. \square

Soit σ l'optimum, qui minimise $\sum_{i=1,\dots,l} dist(\sigma, \sigma_i)$ et soit ρ^* la valeur minimale. Montrons que les solutions produites par l'algorithme ne sont pas très éloignées de ρ^* . L'algorithme produit les $\sigma_{d_{max,i}}$, les centres des composantes les plus récentes des σ_i . Soit i^* l'indice qui minimise

$$\sum_{j=1,\dots,l} dist(\sigma_{d_{max,i}}, \sigma_j)$$

1781 **Lemme 6.** *Il existe une constante c telle que $\sum_{j=1,\dots,l} dist(\sigma_{d_{max,i^*}}, \sigma_j) \leq c.\rho^*$.*

Démonstration. Montrons que la propriété est vraie pour chaque i et donc aussi pour i^* . Pour chaque i , $dist(\sigma_{d_{max,i}}, \sigma_i) \leq dist(\sigma, \sigma_i)$ car $d_{max,i}$ est le point le plus proche de σ_i : il appartient au même cluster et a un degré maximum. Pour les autres σ_j , où $j \neq i$, on applique l'inégalité triangulaire, $dist(\sigma_{d_{max,i}}, \sigma_j) \leq dist(\sigma_{d_{max,i}}, \sigma) + dist(\sigma, \sigma_j)$. D'après le lemme 5, il existe une constante c telle que :

$$dist(\sigma_{d_{max,i}}, \sigma_j) \leq c.dist(\sigma, \sigma_j)$$

1782 et on obtient bien le résultat annoncé.

1783

\square

1784 8.4.3 Structures de données

1785 **Ce qui est stocké dans le temps.** A certains moments discrets t_1, t_2, \dots , nous stockons les
1786 grandes composantes connectées des Réservoirs R_t . Il pourrait n'y en avoir aucune. Nous utilisons
1787 une base de données NoSQL, avec 4 tables (Clé, Valeurs) où la clé est toujours un tag (@x ou
1788 #y) et les valeurs stockent les noeuds des clusters. Notez qu'un flux est identifié par un tag (ou
1789 un ensemble de tags) et qu'un cluster est également identifié par un tag (son noeud le plus haut
1790 degré).

1791 — *Stream(tag, list(cluster, timestamp))* est la table qui fournit les clusters les plus récents d'un
1792 flux,

1793 — $Cluster(tag, list(stream, timestamp, list(high-degree\ nodes), list(nodes(nodes, degree))))$ est
1794 la table qui fournit la liste des noeuds de haut degré et la liste des noeuds avec leur degré,
1795 dans un cluster donné,

1796 — $Nodes(tag, list(stream, cluster, timestamp), tweet)$ est la table qui fournit pour chaque noeud
1797 la liste des flux, clusters et timestamps où le noeud apparaît,

1798 8.5 Explications

1799 L'algorithme de recherche prend le centre c_i des clusters de chacun des tags, puis regarde
1800 le meilleur des centres, celui qui a la distance la plus courte par-rapport à tous les autres tags.
1801 Chacun des centres va appliquer l'algorithme de *Dijkstra* pour calculer la somme de distance qui
1802 minimise la fonction de coût $\sum_{j=1, \dots, l} dist(\sigma_{d_{max, i^*}}, \sigma_j)$. On obtient via l'algorithme de Dijkstra
1803 les chemins qui vont être de type 1 ou de type 2 ou 3. Les chemins de type 1 ne peuvent pas
1804 être très longs. Par défaut l'algorithme nous donne le chemin de type 2. Nous observons trois
1805 chemins, comme le montre la figure 8.3. À partir du centre c nous obtenons deux chemins de type
1806 1 et un chemin de type 2. Nous avons trouvé les plus courts chemins de cette manière. Ce tag est
1807 une réponse intéressante aux questions de trois tags, parce qu'il a un indice pour cette fonction
1808 d'optimum qui n'est pas mauvais.

L'algorithme de recherche trouve le centre c_i du cluster le plus récent de chacun des tags, qui a la distance totale aux autres tags la plus courte. Chacun des centres applique l'algorithme de *Dijkstra* et estime la somme des distances

$$\sum_{j=1, \dots, l} dist(\sigma_{d_{max, i^*}}, \sigma_j).$$

1809 Prenons par exemple une recherche sur les trois tags *CNN*, *AMAZON*, *Prime video* qui sont
1810 appelés *tag 3*, *tag 6*, *tag 7* sur la figure 8.3. La solution obtenue est le centre du cluster 5 et les
1811 trois chemins trouvés par *Dijkstra* sont en rouge bleu et vert.

1812 L'algorithme de Dijkstra trouve des chemins de type 1 ou de type 2 ou 3. Les chemins de
1813 type 1 ne peuvent pas être très longs. La figure 8.3 montre trois plus courts chemins. À partir
1814 du centre c nous obtenons deux chemins de type 1 et un chemin de type 2. La solution n'est pas
1815 l'optimum garanti, mais est dans un rapport fixe à l'optimum, comme le montre le lemme 8.4.2.

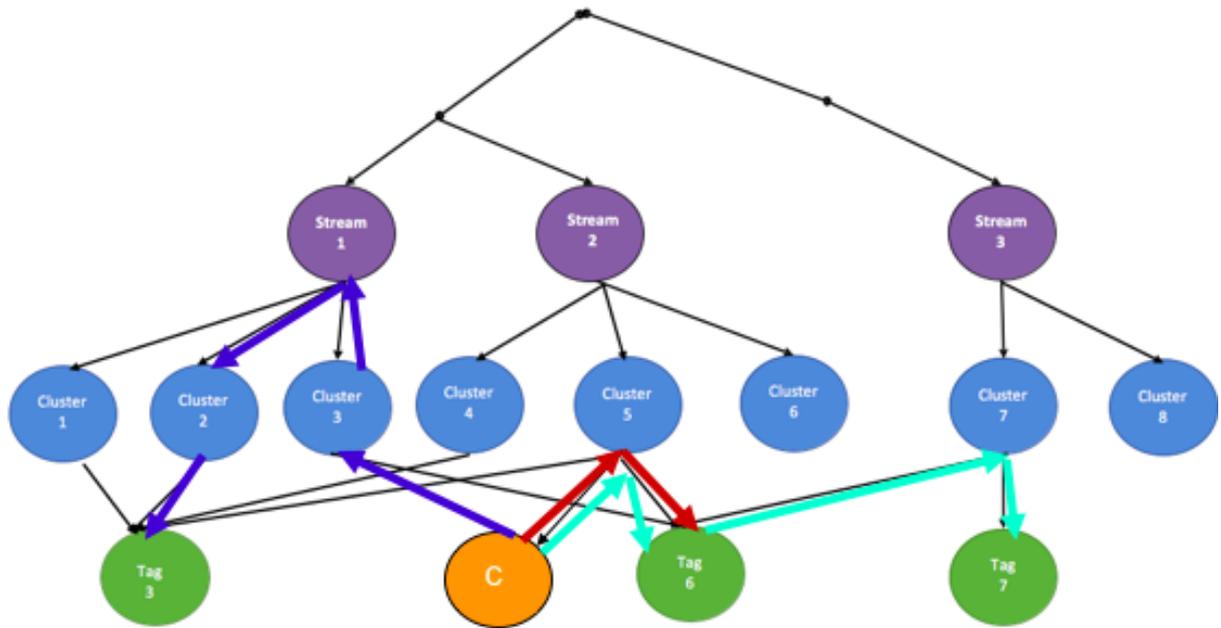


FIGURE 8.3: DAG de recherche.

1816 8.6 Conclusion

1817 La corrélation de contenus entre deux composantes définit une distance qui se généralise aux
 1818 flux. On peut donc associer un arbre phylogénique aux flux, puis l'étendre aux composantes
 1819 géantes et aux tags. Nous avons introduit la notion de recherche par corrélation basée sur ces
 1820 distances. Ces travaux conjuguant plusieurs domaines, ouvrent des champs nouveaux d'applica-
 1821 tions.

1822 9 Expériences

1823 Nous présentons dans ce chapitre les différentes expériences réalisées dans le cadre de cette
1824 thèse.

1825 Dans la première section, nous présentons les expériences dans le cadre d'un flux Twitter.
1826 Dans la deuxième section, nous montrons des expériences pour un flux de texte, en compa-
1827 rant l'échantillonnage uniforme avec l'échantillonnage pondéré. Dans la troisième section, nous
1828 présentons des expériences sur le chapitre recherche et explication.

1829 9.1 Flux Twitter

1830 Aujourd'hui, Twitter présente les tendances à un instant donné, en affichant les tags les plus
1831 fréquents. Nous avons étendu un programme présenté par [61], pour estimer les corrélations entre
1832 les tags et nous l'étendons pour les corrélations entre flux. Cela permet de voir comment les tags
1833 les plus fréquents sont reliés entre eux, de connaître la fréquence des paires de tags.

1834 Il s'agit d'un programme Python¹, exécutable comme un notebook sur *Google Colab*. La
1835 première version nécessite un compte de développeur Twitter² qui spécifie 4 clés. Le code python
1836 doit inclure les bonnes clés et lit les entrées suivantes qui spécifient les paramètres importants
1837 ainsi que ceux des fenêtres dynamique introduites dans la section 5.4.2.2 :

- 1838 1. La taille du Réservoir (k),
- 1839 2. Les mots-clés de Twitter,
- 1840 3. La longueur τ de la fenêtre,
- 1841 4. Le paramètre de pas λ ,
- 1842 5. La valeur seuil, la taille minimale des composantes géantes,
- 1843 6. La durée de l'observation.

1844 La sortie est un ensemble de fichiers, un pour chaque fenêtre. Chaque fichier contient des arêtes
1845 (u, v, texte) où u, v sont des nœuds et le texte est le contenu du tweet. Par exemple, Le "`@tho-`
1846 `masjacksonjr`" a envoyé un tweet qui est le texte ci-dessous, dans laquelle $v = \#PrimeVideo$ est
1847 un des 4 tags et l'arête :

1848 (`@thomasjacksonjr, #PrimeVideo, "Watched #ABCMurders from beginning to end and`
1849 `it was pure #mystery and delight. Excellent show and a must watch for #mysterylovers. Great`
1850 `show on #PrimeVideo."`).

1. <https://github.com/alassou/t/blob/master/topic.ipynb>

2. <https://developer.twitter.com>

1851 9.1.1 Expérience sur un flux

1852 Une exécution possible précise : $k = 400$, le mot clé : "CNN", $\tau = 12min$, $\lambda = 3min$, seuil
1853 $= 10$, la durée de l'observation $= 24min$. Ces paramètres définissent 5 fenêtres qui se chevauchent
1854 et la sortie contient 5 fichiers.

1855 Nous observons 3.10^3 arêtes, générés par 10^3 tweets par minute pendant $24min$. La figure
1856 9.1 décrit un Réservoir de 400 d'arêtes où tous les petites composantes (de taille inférieure au
1857 seuil=10) ont été enlevées. Pour une observation de $\tau = 12$, nous avons approximativement $m =$
1858 3.10^4 et $n = 8000$, comme $m = O(n. \log n)$. La taille k du Réservoir est $(\sqrt{n}. \log n) = 90.4 = 360$,
1859 soit moins de 400 dans l'expérience.

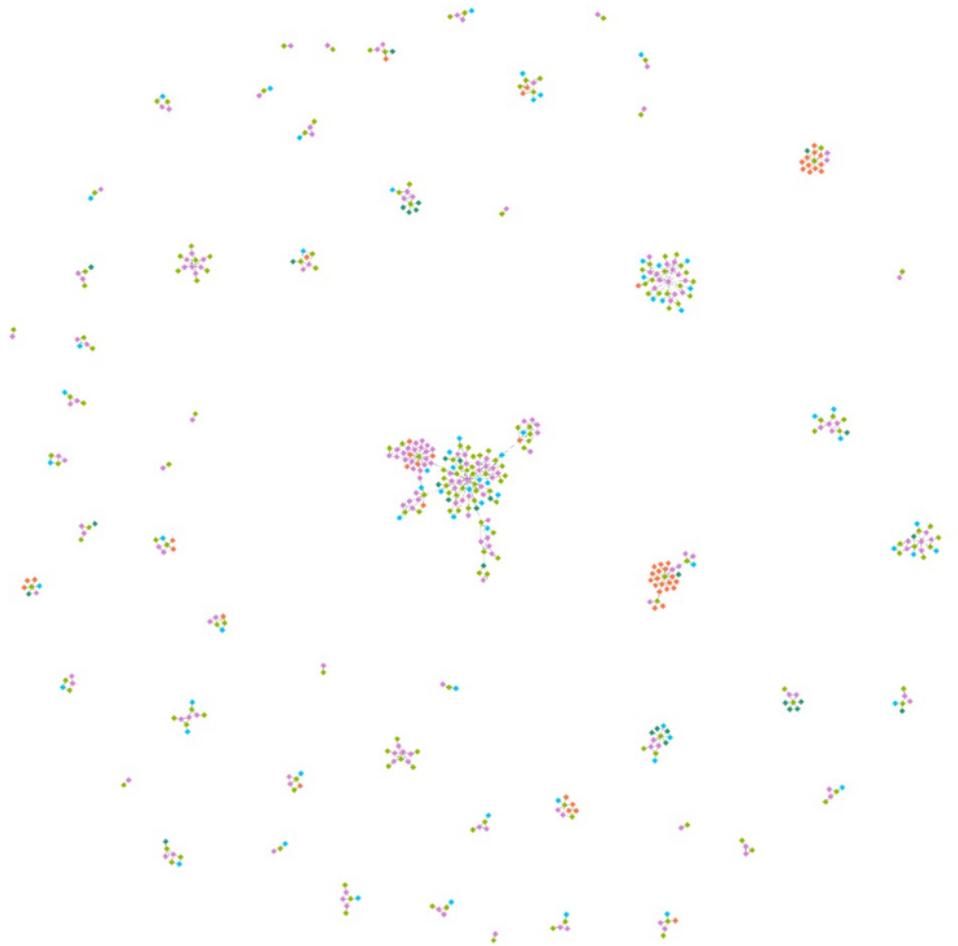


FIGURE 9.1: Les composantes géantes du Réservoir

1860 9.1.2 Expérience sur plusieurs flux

1861 Supposons une autre expérience avec 4 flux twitter sur les tags #CNN, #FoxNews, #Bitcoin,
 1862 et #Ripple pendant 24 heures avec un taille de fenêtre $\tau = 1h$ et un intervalle de temps de
 1863 $\lambda = 30mins$. La figure 9.3 indique le nombre d'arêtes dans une fenêtre, environ $m = 30 \cdot 10^3$ par
 1864 flux, pour chaque flux, donc 10^6 d'arêtes en 24 heures. Pour le #Bitcoin avec $k = 400$, la taille
 1865 de la composante géante est d'environ 100.

1866 Ces données influencent le développement d'un moteur de recherche par similarité, présenté
 1867 dans la section 8.2.

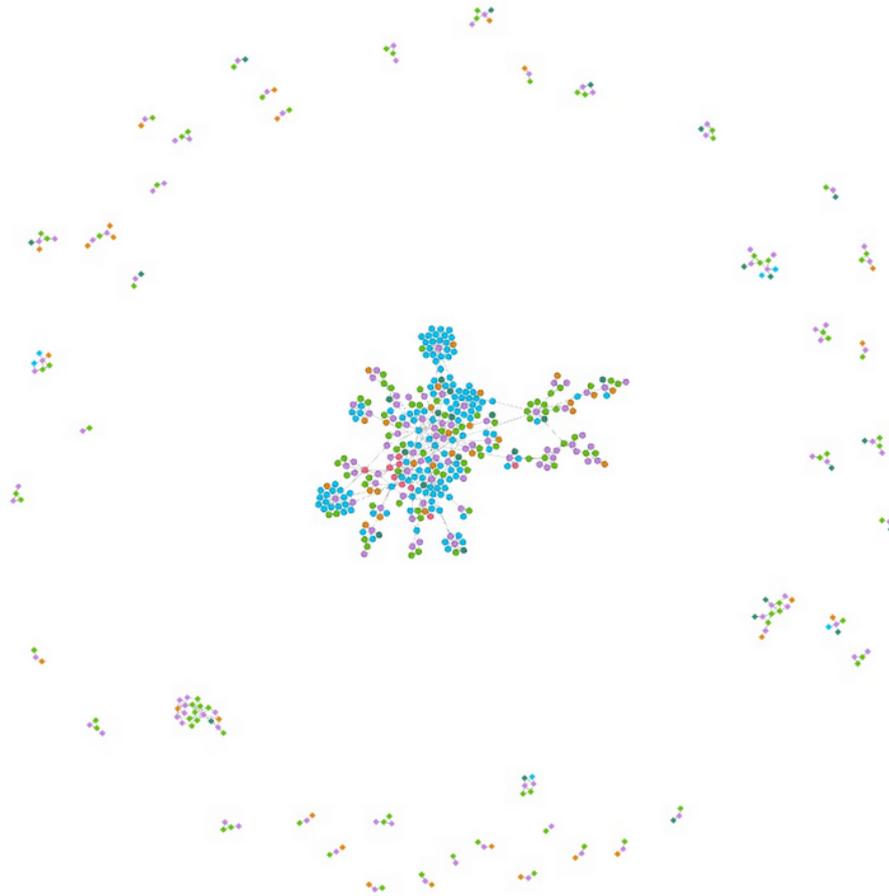


FIGURE 9.2: Les composantes géantes du Réservoir pour le mot clé *Bitcoin*.

1868 9.2 Analyse de texte

1869 Nous pouvons appliquer notre méthode à un corpus de textes standards. Pour chaque phrase,
 1870 nous appliquons d'abord *la lemmatisation, la reconnaissance des entités* [49] et éliminons les
 1871 *stop words* ³ et les mots les moins fréquents. Nous générons soit les *bigrammes* (paires de mots
 1872 contigus) [50] soit toutes les paires possibles d'une phrase donnée, comme des arêtes potentielles.
 1873 Pour l'échantillonnage uniforme, les poids des arêtes sont constants (par exemple égal à 1). Pour
 1874 l'échantillonnage pondéré, le poids d'une arête est la similarité entre deux mots de *Word2vec*
 1875 présentée dans la section 6.2. Dans les deux cas, nous traitons le texte comme un flux.

1876 Nous évaluons nos méthodes sur l'ensemble de données (*NIPS*)⁴, (*Googleplaystore_user_reviews*)

3. stop word est un mot commun qu'il est inutile d'indexer ou de l'utiliser dans une recherche. En français, par exemple « le », « la », « de », « du », « ce »...

4. <https://www.kaggle.com/benhamner/nips-paperspapers.csv>

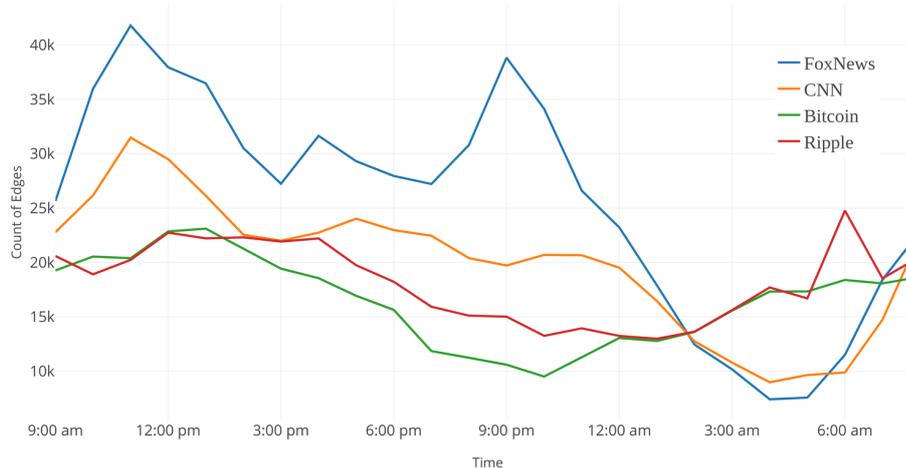


FIGURE 9.3: Nombre d'arêtes/par heure pour 4 flux pendant 24h

1877 *dataset*⁵ et le corpus de *Stanford Natural Language Inference (SNLI)* [13] qui consiste en une
 1878 collection de phrases en anglais. Nous utilisons uniquement le jeu de données *test*⁶, qui contient
 1879 3000 phrases. Nous avons obtenu 71269 de paires de mots.

1880 9.2.1 Échantillonnage uniforme ou pondéré

1881 Nous appliquons l'échantillonnage uniforme comme dans le cas de Twitter ou pondéré où le
 1882 poids est la valeur absolue de la similarité *Word2vec* entre deux mots. Nous introduisons une
 1883 notion de stabilité pour la construction du Réservoir introduite dans la section 6.3. Supposons
 1884 que nous effectuons deux passages indépendants sur le corpus D , avec deux Réservoirs de même
 1885 taille k . Si nous observons une composante géante dans chaque Réservoir, c'est-à-dire C_1 pour
 1886 la première passe et C_2 pour la deuxième passe, quelle est la taille de l'intersection $C_1 \cap C_2$ par
 1887 rapport à la taille de C_1 ou de C_2 ?

1888 Pour les figures [9.4,9.5] et le tableau 9.1, nous regardons les *bigrammes*, tandis que pour la
 1889 figure 9.6 nous regardons toutes les paires possibles. Pour l'échantillonnage pondéré, la figure 9.4
 1890 donne la taille de la composante géante en fonction de k .

1891 En regardant les *bigrammes* sur une phrase donnée, la figure 9.5 montre la comparaison entre
 1892 $\rho(D, k)$ pour le Réservoir uniforme et le Réservoir pondéré avec différentes tailles de Réservoir
 1893 $k = 100, 200, \dots, 500$. Nous avons observé que dans le cas du Réservoir uniforme, $\rho(D, k)$ augmente
 1894 avec la taille du Réservoir et dans le cas du Réservoir pondéré, $\rho(D, k)$ est presque stable comme
 1895 le montre la figure 9.5. Dans tous les cas, le Réservoir pondéré est toujours plus stable.

1896 Nous lisons trois différents benchmarks et appliquons l'échantillonnage pondéré du Réservoir.

5. https://www.kaggle.com/lava18/google-play-store-appsgoogleplaystore_user_reviews.csv

6. <https://github.com/alassou/t/blob/master/snli.csv>

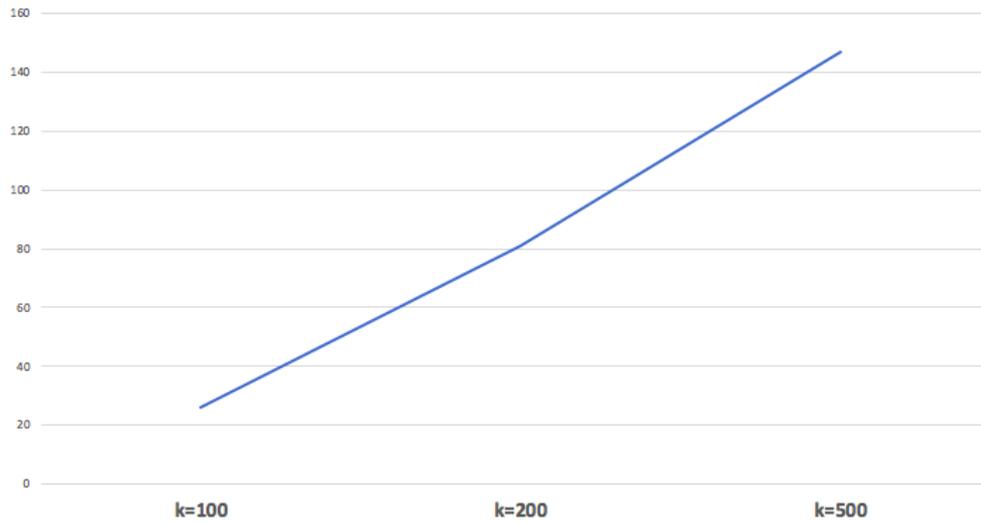


FIGURE 9.4: Taille de la composante géante pour le Réservoir pondéré

1897 Le tableau 9.1 montre la comparaison de $\rho(D, k)$ pour les trois différents benchmarks et pour des
1898 différentes valeurs de k .



FIGURE 9.5: Stabilité du Reservoir sampling et de l'échantillonnage pondéré.

Benchmarks	k=100	k=200	k=500
SNLI	0,42	0,56	0,77
NIPS PAPERS	0,58	0,65	0,72
Googleplaystore_user_reviews	0,15	0,52	0,61

TABLE 9.1: Stabilité du Réservoir pondéré

1899 La figure 9.6 montre la comparaison de stabilité entre l'échantillonnage uniforme et pondéré,
 1900 en regardant tous les paires possibles de la phrase avec différentes tailles de Réservoir $k =$
 1901 100, 200, ...500.

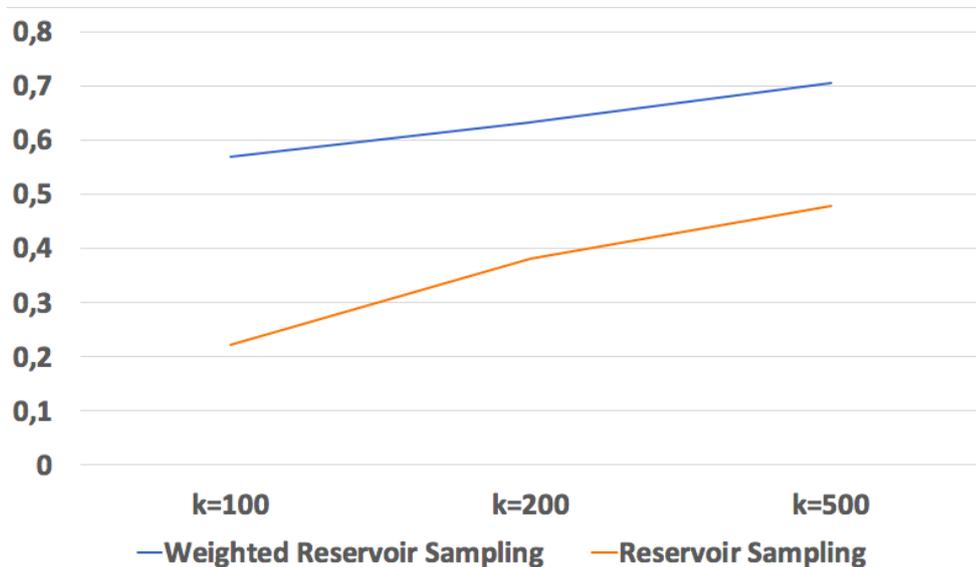


FIGURE 9.6: Stabilité du Reservoir sampling et de l'échantillonnage pondéré sur tous les paires.

1902 9.2.2 Classification

1903 Des mesures de qualité [64] ont été introduites pour comparer les différentes méthodes de clas-
 1904 sifications présentées dans la section 6.4. Dans notre cas, chaque expérience fournit un ensemble
 1905 de composantes géantes et chaque composante géante peut être interprétée comme un sujet. La
 1906 composante géante n'est pas une distribution mais un graphe avec sa propre structure.

1907 Nous pouvons comparer un support de distributions *LDA* présentées dans la section 3.3.2
 1908 avec l'ensemble des noeuds du graphe de la composante géante. Comme nous répétons différentes
 1909 expériences, l'espérance de la stabilité est très similaire à la stabilité du sujet introduit dans
 1910 [64]. Une différence importante est que *LDA* doit accéder à toutes les données, alors que nous
 1911 n'accédons qu'aux données en streaming en faisant 1 seule passe.

1912 Nous avons comparé notre méthode avec celle de *LDA* pour $k = 2$.

1913 Topic1 a été décrit par les mots significatifs suivants : ['man', 'woman', 'wear', 'boy', 'walk',
 1914 'girl', 'dog', 'stand', 'street', 'play', 'child', 'dress', 'hold', 'water'], tandis que Topic2 a été décrit
 1915 par les mots significatifs suivants ['man', 'sit', 'people', 'woman', 'look', 'stand', 'group', 'shirt',
 1916 'wear', 'table', 'hold', 'player', 'work', 'play', 'ball']. Les distributions sont données dans la figure
 1917 9.7.

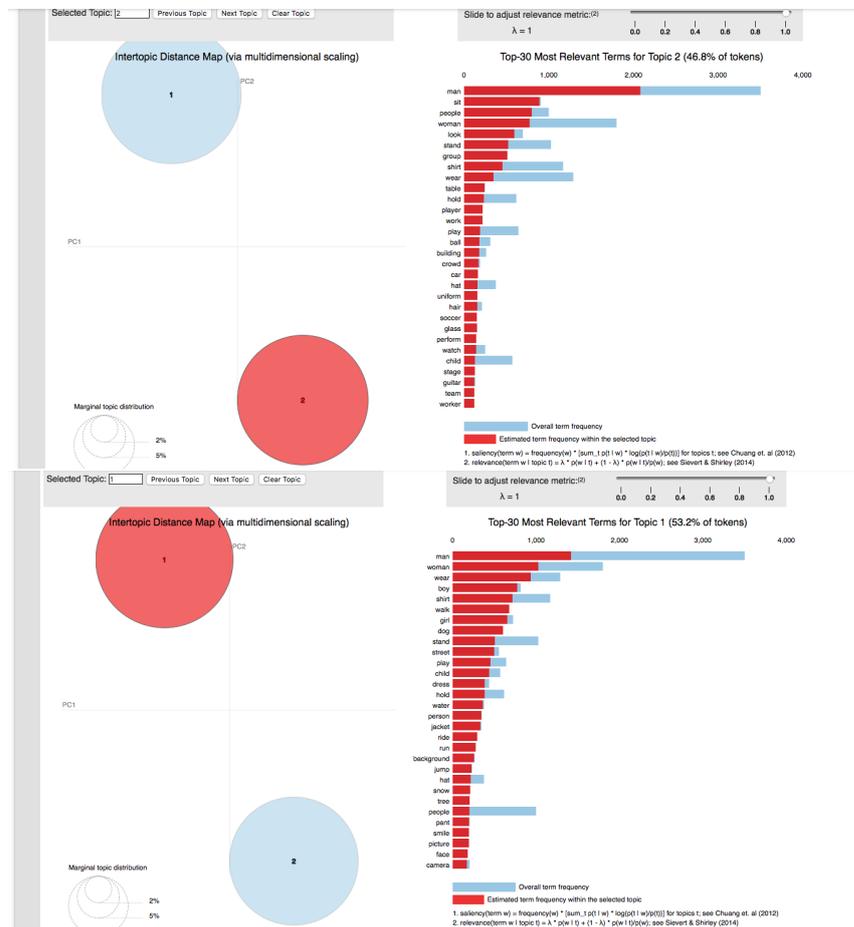


FIGURE 9.7: Distributions de 2 topics

1918 Si nous prenons le Réservoir pondéré stable, quelle partie de ces sujets couvrons-nous comme
 1919 noeuds? Cela est indiqué dans la figure 9.8. Plus k augmente, plus nous couvrons les mots du
 1920 Topic, pour atteindre 100% pour $k=500$.

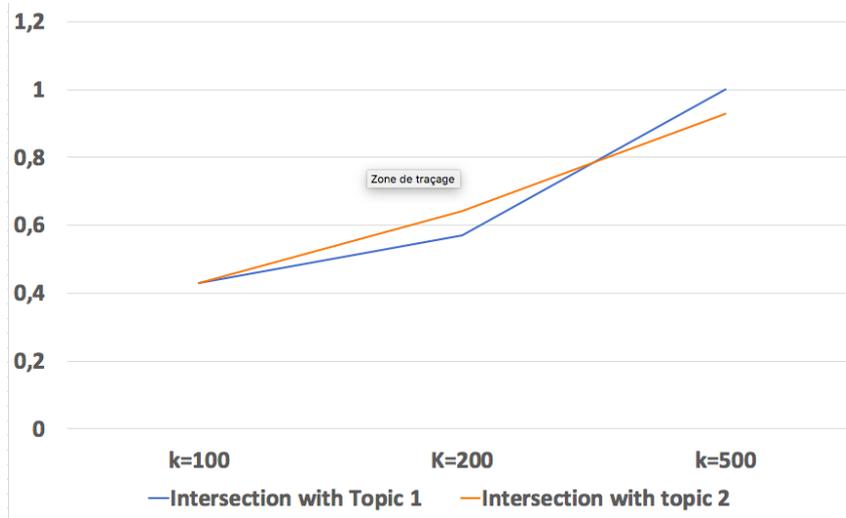


FIGURE 9.8: Fraction des thèmes de LDA comme des noeuds de la composante géante du réservoir pondéré

1921 Dans le cas général, si on prend le Réservoir pondéré dans des fenêtres, on va avoir un certain
 1922 nombre de composantes géantes pour lesquelles on a des distances. Si on applique l'algorithme
 1923 $k - means$, on aura $kclasses$. Finalement, on compare le représentant de chacune des classes
 1924 avec les topics de LDA . Dans le cas particulier représenté ci-dessus, on a une seule fenêtre et on
 1925 observe deux composantes géantes.

1926 9.3 Décompositions hiérarchiques des composantes géantes

1927 Chaque composante géante est le 2 - *Core* d'un graphe détaillé dans la section 2.3.4.1. Nous
 1928 pouvons explorer ce graphe en commençant par les noeuds de degré maximum et itérer jusqu'à ce
 1929 que nous explorions l'ensemble de la composante géante. Avec le graphe de la figure 9.10, nous
 1930 obtenons la décomposition en arbre de la figure 9.11. Nous commençons avec le noeud 1 de degré
 1931 maximum et nous explorons tous les noeuds à distance 1. À l'étape suivante, le noeud 6 est de
 1932 degré maximum et nous obtenons le noeud 3. La décomposition de l'arbre est de profondeur 2.

1933 Dans la figure 9.12, chaque noeud est un tag et le texte du tweet est associé à chaque arête.
 1934 Nous ne montrons que le texte de l'arête e_1 .

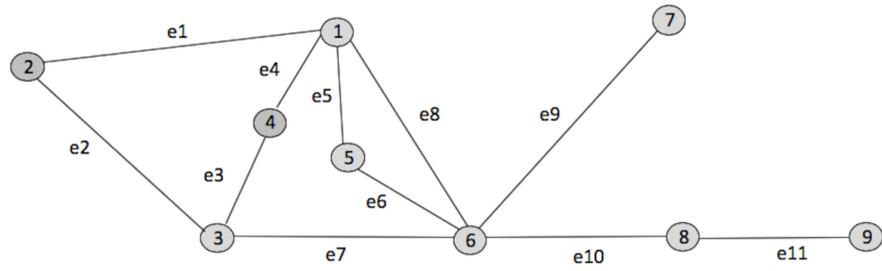


FIGURE 9.9: La grande composante connectée (C) dans le Réservoir

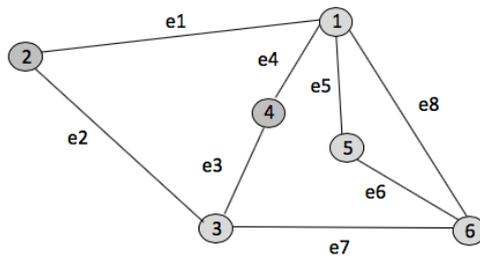


FIGURE 9.10: Le $2 - core(C)$

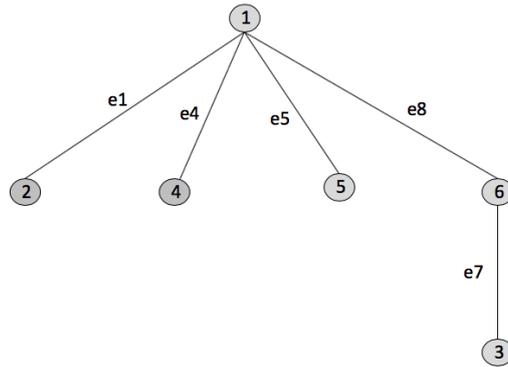


FIGURE 9.11: Arbre de décomposition

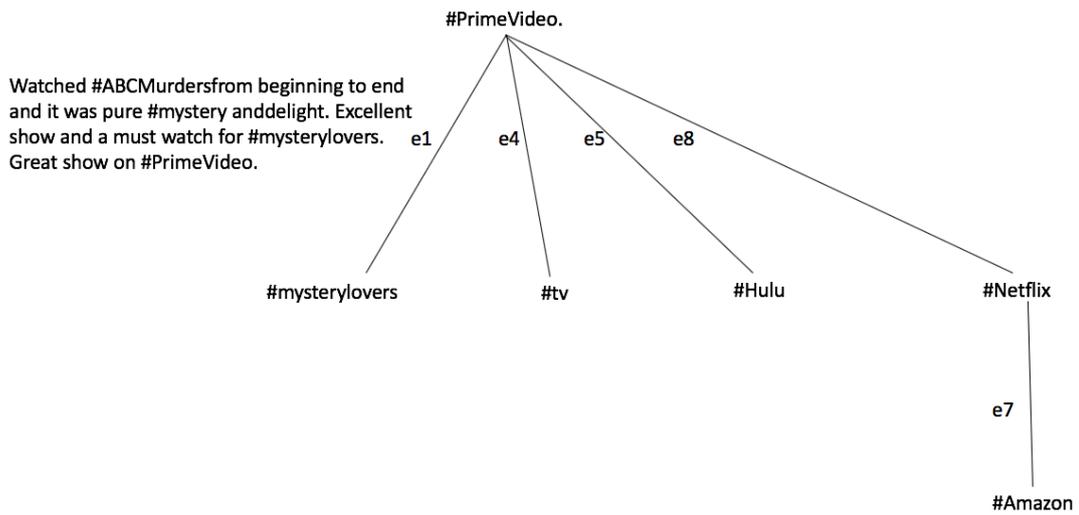


FIGURE 9.12: Arbre de décomposition où les noeuds sont des tags

1935 10 Conclusion

1936 L'objectif de cette thèse est de montrer comment approcher les propriétés de trois différents
1937 types de données (entrepôts de données au sein d'une base de données, graphes issus des réseaux
1938 sociaux, texte) à partir d'échantillons. Les résultats de cette thèse se décomposent en deux grandes
1939 parties. La première partie traite de l'échantillonnage de données afin d'approximer certaines
1940 propriétés et la deuxième partie approfondit un modèle morphologique de texte pour enrichir
1941 l'analyse du texte.

1942

1943 La première partie de la thèse couvre les sujets suivants :

1944 • Dans les bases de données, nous avons montré comment approximer des requêtes OLAP à l'aide
1945 d'un tirage non uniforme. Nous avons également introduit une notion de préférence comme
1946 une distribution sur les *Mesures*. Nous avons montré comment nous pouvons généraliser cette
1947 méthode d'approximation, pour toute préférence, c'est-à-dire toute distribution sur les *Mesures*.
1948

1949 • Pour un flux de Tweets, nous construisons un flux d'arêtes qui représente le graphe Twitter
1950 G . Nous avons montré que l'échantillonnage uniforme des arêtes reflète les grands clusters de
1951 G . L'échantillonnage uniforme est stable pour Twitter.
1952

1953 • Pour un flux de textes, nous construisons un graphe dont les arêtes sont les paires de mots
1954 de la même phrase. Nous avons montré que l'échantillonnage pondéré est stable pour analyser
1955 les composantes géantes de ce graphe. Chaque composante géante définit des mots centraux et
1956 des phrases centrales. Nous pouvons classifier les composantes avec l'algorithme *k-means*.
1957

1958 • On introduit la notion de moteur de recherche à partir des données d'échantillonnage pour
1959 donner des réponses à des requêtes définies par des mots clés. Les réponses sont obtenues par
1960 corrélation entre les composantes géantes et les échantillons. Pour chaque réponse, le moteur
1961 propose des explications.
1962

1963 La deuxième partie de la thèse introduit un modèle statistique pour la morphologie des
1964 langues naturelles, inspiré par *Amis*, une langue naturelle qui a une riche morphologie. Nous
1965 avons construit les distributions classiques des préfixes, racines et suffixes les plus fréquents, et
1966 les distributions correspondantes avec une racine, un préfixe ou un suffixe possible. A partir des
1967 seconds moments des distributions, nous avons construit des vecteurs pour les préfixes, les ra-
1968 cines et les suffixes qui capturent leurs corrélations. La morphologie d'un mot w est la structure

1969 *préfixe-racine-suffixe* . On construit des vecteurs : $v_p(pre)$ pour le préfixe *pre*, $v_r(root)$ pour la
1970 racine *root* et $v_s(suf)$ pour le suffixe *suf*. Le vecteur du mot $v(w)$ est la concaténation des trois
1971 vecteurs. Nous avons défini un vecteur probabiliste associé à une phrase. L'analyse des préfixes
1972 et suffixes, détermine la plupart des composantes du vecteur de la phrase S , de manière pro-
1973 babiliste. Nous avons montré comment prédire la classe syntaxique d'une phrase simple à l'aide
1974 d'un arbre de décision. Étant donné une grammaire G et une phrase $S : w_1, w_2, \dots, w_n$, nous avons
1975 montré comment trouver l'arbre de décomposition syntaxique préférable en utilisant le vecteur
1976 de la phrase. Toutes les prédictions sont approximatives.

1977

1978 Les résultats obtenus dans cette thèse ouvrent plusieurs pistes de recherches :

- 1979 ● Peut-on améliorer l'échantillonnage uniforme dans le cas des réseaux sociaux ? Nous n'utilisons
1980 pas d'historique sur les tags et utilisons la distribution uniforme par défaut. Nous pourrions
1981 construire la matrice de corrélation des tags, puis des vecteurs *Word2vec* comme dans le cas
1982 du langage naturel. Nous aurions alors un échantillonnage pondéré possible pour les paires de
1983 tags.
- 1984 ● Comment utiliser des vecteurs *Word2vec* structurés, comme ceux introduits pour l'analyse mor-
1985 phologique en ayant plusieurs distributions possibles sur lesquelles on peut faire de l'échantillonnage ?
1986 Les distributions des préfixes, racines et suffixes et leurs projections permettraient d'obtenir
1987 des échantillons plus variés.
- 1988 ● Comment classifier plusieurs flux à l'aide de l'algorithme *k-means*, appliqué aux composantes
1989 géantes de chaque fenêtre coulissante ? Plusieurs distances sont possibles entre les composantes,
1990 présentées au chapitre 6.4.1.
- 1991 ● Dans le domaine du *Question Answering*, comment utiliser les vecteurs de phrase basés sur
1992 la morphologie, introduits dans le chapitre 7 ? L'échantillonnage n'est pas adapté dans ce
1993 cadre, mais les vecteurs probabilistes associés aux phrases et les vecteurs *Word2vec* associés
1994 aux mots, donnent une information importante pour interpréter une question. Le mécanisme
1995 d'attention, basé sur les distributions de similarité entre un mot et les autres mots d'une phrase
1996 est à l'origine des techniques de *Transformers* utilisées pour la traduction d'une langue vers
1997 une autre. Comment traduire une question en langage naturel dans le contexte d'un schéma
1998 relationnel, en une requête SQL ?
1999 Dans le cas d'un ensemble de phrases, il peut être délicat d'isoler la question précise, en
2000 particulier en l'absence du symbole " ? ". Les vecteurs de phrase capturent les actes de langage
2001 de manière approchée et pourraient contribuer à mieux appréhender les questions. Un *indice*
2002 *de compréhension* pourrait alors mesurer expérimentalement la qualité des réponses sur des
2003 ensembles de tests.

2004 Bibliographie

- 2005 [1] Himmelmann N. (Ed.). Adelaar, K. (Ed.). *The Austronesian Languages of Asia and Mada-*
2006 *gascar*. London : Routledge,, 2005.
- 2007 [2] Charu C. Aggarwal. On biased reservoir sampling in the presence of stream evolution. In
2008 *Proceedings of the 32Nd International Conference on Very Large Data Bases, VLDB '06*,
2009 pages 607–618. VLDB Endowment, 2006.
- 2010 [3] Réka Albert and Albert-László Barabási. Statistical mechanics of complex networks. *Rev.*
2011 *Mod. Phys.*, 74 :47–97, Jan 2002.
- 2012 [4] Sanjeev Arora, Yuanzhi Li, Yingyu Liang, Tengyu Ma, and Andrej Risteski. A latent va-
2013 riable model approach to PMI-based word embeddings. *Transactions of the Association for*
2014 *Computational Linguistics*, 4 :385–399, 2016.
- 2015 [5] Bogdan Babych and Anthony Hartley. Improving machine translation quality with automatic
2016 named entity recognition. In *Proceedings of the 7th International EAMT workshop on MT*
2017 *and other language technology tools, Improving MT through other language technology tools,*
2018 *Resource and tools for building MT at EACL 2003*, 2003.
- 2019 [6] Bahman Bahmani, Ravi Kumar, and Sergei Vassilvitskii. Densest subgraph in streaming
2020 and mapreduce. *Proc. VLDB Endow.*, 5(5) :454–465, January 2012.
- 2021 [7] Y. Bengio, Réjean Ducharme, Pascal Vincent, and Christian Jauvin. *Neural Probabilistic*
2022 *Language Models*, volume 3, pages 137–186. 05 2006.
- 2023 [8] Sayan Bhattacharya, Monika Henzinger, Danupon Nanongkai, and Charalampos E. Tsou-
2024 rakakis. Space- and time-efficient algorithm for maintaining dense subgraphs on one-pass
2025 dynamic streams. *CoRR*, pages 173–182, 2015.
- 2026 [9] N. Biggs. *Discrete Mathematics*. Discrete Mathematics. OUP Oxford, 2002.
- 2027 [10] D. M. Blei, A. Y. Ng, and M. I. Jordan. Latent dirichlet allocation. *Journal of Machine*
2028 *Learning Research*, 3 :993-1022, 2003.
- 2029 [11] Piotr Bojanowski, Edouard Grave, Armand Joulin, and Tomas Mikolov. Enriching word
2030 vectors with subword information. *CoRR*, abs/1607.04606, 2016.
- 2031 [12] Bela Bollobas. *Graph theory : an introductory course / Bela Bollobas*. Springer Verlag,,
2032 1979.
- 2033 [13] Samuel Bowman, Gabor Angeli, Christopher Potts, and Christopher Manning. A large anno-
2034 tated corpus for learning natural language inference. *CoRR*, 08 2015.

- 2035 [14] Isabelle Bril. Roots and stems : Lexical and functional flexibility in Amis and Nêlêmwa.
2036 *Studies in Language*, 41(2) :358–407, 2017.
- 2037 [15] Isabelle Bril. Roots and stems : Lexical and functional flexibility in amis and nêlêmwa, 2017.
- 2038 [16] Isabelle Bril, Achraf Lassoued, and Michel de Rougemont. A statistical model for morphology
2039 inspired by the amis language. In *Proceedings of Language, Ontology, Terminology and*
2040 *Knowledge Structures Workshop (LOTKS 2017)*, Montpellier, France, 2017. Association for
2041 Computational Linguistics.
- 2042 [17] Kris Cao and Marek Rei. A joint model for word embedding and word morphology. *CoRR*,
2043 abs/1606.02601, 2016.
- 2044 [18] Danqi Chen and Christopher Manning. A fast and accurate dependency parser using neural
2045 networks. In *Proceedings of the 2014 Conference on Empirical Methods in Natural Language*
2046 *Processing (EMNLP)*, pages 740–750. Association for Computational Linguistics, October
2047 2014.
- 2048 [19] Teresa Chen. Verbal constructions and verbal classifications in nataoran-amis. In *Series C.*
2049 *Canberra : Pacific linguistics*, 1987.
- 2050 [20] Kenneth Ward Church and Patrick Hanks. Word association norms, mutual information,
2051 and lexicography. *Comput. Linguist.*, 16(1) :22–29, March 1990.
- 2052 [21] Ronan Collobert and Jason Weston. A unified architecture for natural language processing :
2053 deep neural networks with multitask learning. In *ICML '08*, 2008.
- 2054 [22] M. de Rougemont and G. Vimont. The content correlation of streaming edges. In *IEEE*
2055 *International Conference on Big Data*, pages 1101–1106, 2018.
- 2056 [23] Scott Deerwester, Susan T. Dumais, George W. Furnas, Thomas K. Landauer, and Richard
2057 Harshman. Indexing by latent semantic analysis. 41(6) :391–407, 1990.
- 2058 [24] Steve Cassidy, Diego Mollá, Menno van Zaanen. Named entity recognition in question
2059 answering of speech data. *Proceedings of the Australasian Language Technology Workshop*,
2060 pages 57–65, 2007.
- 2061 [25] Alessandro Epasto, Silvio Lattanzi, and Mauro Sozio. Efficient densest subgraph compu-
2062 tation in evolving graphs. In *Proceedings of the 24th International Conference on World*
2063 *Wide Web, WWW '15*, pages 300–310, Republic and Canton of Geneva, Switzerland, 2015.
2064 International World Wide Web Conferences Steering Committee.
- 2065 [26] Alessandro Epasto, Silvio Lattanzi, and Mauro Sozio. Efficient densest subgraph computa-
2066 tion in evolving graphs. In *Proceedings of the 24th International Conference on World Wide*
2067 *Web, WWW '15*, pages 300–310, 2015.
- 2068 [27] Paul Erdős and Alfréd Rényi. On random graphs i. *Publicationes Mathematicae (Debrecen)*,
2069 6 :290–297, 1959 1959.
- 2070 [28] P. Erdős and A Rényi. On the evolution of random graphs. pages 17–61, 1960.
- 2071 [29] Hossein Esfandiari, MohammadTaghi Hajiaghayi, and David P. Woodruff. Brief announce-
2072 ment : Applications of uniform sampling : Densest subgraph and beyond. In *Proceedings*
2073 *of the 28th ACM Symposium on Parallelism in Algorithms and Architectures, SPAA 2016,*
2074 *Asilomar State Beach/Pacific Grove, CA, USA, July 11-13, 2016*, pages 397–399, 2016.

- 2075 [30] Anastasia Giachanou and Fabio Crestani. Like it or not : A survey of twitter sentiment
2076 analysis methods. *ACM Computing Surveys (CSUR)*, 49(2) :1–41, 2016.
- 2077 [31] J. Hastad. Clique is hard to approximate within $n^{1-\epsilon}$. In *Proceedings of the 37th Annual*
2078 *Symposium on Foundations of Computer Science, FOCS '96*, pages 627–. IEEE Computer
2079 Society, 1996.
- 2080 [32] Sepp Hochreiter and Jürgen Schmidhuber. Long short-term memory. *Neural computation*,
2081 9 :1735–80, 12 1997.
- 2082 [33] W. Hoeffding. Probability inequalities for sums of bounded random variables. *Journal of*
2083 *the American Statistical Association*, 58(2) :13–30, 1963.
- 2084 [34] Pili Hu and Wing Cheong Lau. A survey and taxonomy of graph sampling. *CoRR*,
2085 abs/1308.5865, 2013.
- 2086 [35] William H. Inmon. *Building the Data Warehouse*. John Wiley , Sons, Inc., USA, 4rth edition,
2087 2005.
- 2088 [36] Rajeev Motwani John E. Hopcroft and Jeffrey D. Ullman. *Introduction to Automata Theory,*
2089 *Languages, and Computation*. Springer Verlag,, 2001.
- 2090 [37] Sara Keretna, Chee Lim, and Doug Creighton. A hybrid model for named entity recognition
2091 using unstructured medical text, 06 2014.
- 2092 [38] Efthymios Kouloumpis, Theresa Wilson, and Johanna Moore. Twitter sentiment analysis :
2093 The good the bad and the omg! In *Fifth International AAI conference on weblogs and*
2094 *social media*, 2011.
- 2095 [39] Omer Levy and Yoav Goldberg. Linguistic regularities in sparse and explicit word represen-
2096 tations. In *Proceedings of the Eighteenth Conference on Computational Natural Language*
2097 *Learning*, pages 171–180, June 2014.
- 2098 [40] Zhouhan Lin, Minwei Feng, Cícero Nogueira dos Santos, Mo Yu, Bing Xiang, Bowen Zhou,
2099 and Yoshua Bengio. A structured self-attentive sentence embedding. *CoRR*, abs/1703.03130,
2100 2017.
- 2101 [41] Christopher D. Manning and Hinrich Schütze. *Foundations of Statistical Natural Language*
2102 *Processing*. MIT Press, Cambridge, MA, USA, 1999.
- 2103 [42] Claire Mathieu and Michel de Rougemont. Large very dense subgraphs in a stream of edges.
2104 In *ACM-IMS Foundations of Data Science*, pages –, 2020.
- 2105 [43] Andrew McGregor, David Tench, Sofya Vorotnikova, and Hoa T. Vu. Densest subgraph in
2106 dynamic graph streams. *CoRR*, abs/1506.04417, 2015.
- 2107 [44] Tomas Mikolov, Kai Chen, Greg Corrado, and Jeffrey Dean. Efficient estimation of word
2108 representations in vector space. *CoRR*, abs/1301.3781, 2013.
- 2109 [45] Tomas Mikolov, Ilya Sutskever, Kai Chen, Greg Corrado, and Jeffrey Dean. Distributed
2110 representations of words and phrases and their compositionality. *CoRR*, abs/1310.4546,
2111 2013.
- 2112 [46] Tomas Mikolov, Scott Wen-tau Yih, and Geoffrey Zweig. Linguistic regularities in continuous
2113 space word representations. In *Proceedings of the 2013 Conference of the North American*
2114 *Chapter of the Association for Computational Linguistics : Human Language Technologies*
2115 *(NAACL-HLT-2013)*. Association for Computational Linguistics, May 2013.

- 2116 [47] M. Molloy and B. Reed. A critical point for random graphs with a given degree sequence.
2117 *Random Struct. Algorithms*, 6 :161–180, 1995.
- 2118 [48] Michael Molloy and Bruce Reed. The size of the giant component of a random graph with
2119 a given degree sequence. *Comb. Probab. Comput.*, 7(3) :295–305, 1998.
- 2120 [49] David Nadeau and Satoshi Sekine. A survey of named entity recognition and classification.
2121 *Linguisticae Investigationes*, 30(1) :3–26, January 2007.
- 2122 [50] Rao Muhammad Adeel Nawab, Mark Stevenson, and Paul Clough. Comparing Medline
2123 citations using modified N-grams. *Journal of the American Medical Informatics Association*,
2124 21(1) :105–110, 05 2013.
- 2125 [51] Jeffrey Pennington, Richard Socher, and Christopher D. Manning. Glove : Global vectors
2126 for word representation. In *In EMNLP*, 2014.
- 2127 [52] P. Ratinaud and S Déjean. Iramuteq : Implémentation de la méthode alceste d’analyse de
2128 texte dans un logiciel libre. *HAL*, 2009.
- 2129 [53] Douglas Rohde, Laura Gonnerman, and David Plaut. An improved model of semantic
2130 similarity based on lexical co-occurrence. *Cognitive Science - COGSCI*, 8, 01 2006.
- 2131 [54] D.E. Rumelhart, G.E. Hinton, and R.J. Williams. Learning representations by back propa-
2132 gating errors. *cogn. Nature*, 5, 01 1986.
- 2133 [55] D.E. Rumelhart and James McClelland. On learning the past tenses of english verbs. parallel
2134 distributed processing : explorations in the microstructure of cognition. *Psychological and*
2135 *Biological Models*, 2 :216–271, 01 1986.
- 2136 [56] N Saitou and M Nei. The neighbor-joining method : a new method for reconstructing
2137 phylogenetic trees. *Molecular Biology and Evolution*, 4(4) :406–425, 1987.
- 2138 [57] Sebastian Schuster and Christopher D. Manning. Enhanced english universal dependencies :
2139 An improved representation for natural language understanding tasks. In *LREC*, 2016.
- 2140 [58] Richard Socher, John Bauer, Christopher D. Manning, and Andrew Y. Ng. Parsing with
2141 compositional vector grammars. In *In Proceedings of the ACL conference*, 2013.
- 2142 [59] Richard Socher, Christopher D. Manning, and Andrew Y. Ng. Learning continuous phrase
2143 representations and syntactic parsing with recursive neural networks. In *In Proceedings of*
2144 *the NIPS-2010 Deep Learning and Unsupervised Feature Learning Workshop*, 2010.
- 2145 [60] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez,
2146 Lukasz Kaiser, and Illia Polosukhin. Attention is all you need. *CoRR*, abs/1706.03762, 2017.
- 2147 [61] Guillaume Vimont. Approximation dynamique de clusters dans un graphe social : Méthodes
2148 et applications. *Université Panthéon - Assas*, pages 120–122, 2019.
- 2149 [62] Jeffrey S. Vitter. Random sampling with a reservoir. *ACM Trans. Math. Softw.*, 11(1) :37–57,
2150 March 1985.
- 2151 [63] Yingce Xia, Tianyu He, Xu Tan, Fei Tian, Di He, and Tao Qin. Tied transformers : Neural
2152 machine translation with shared encoder and decoder. In *AAAI*, 2019.
- 2153 [64] Linzi Xing, Michael J. Paul, and Giuseppe Carenini. Evaluating topic quality with posterior
2154 variability. In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language*
2155 *Processing (EMNLP-IJCNLP)*, pages 3471–3477. Association for Computational Linguistics,
2156 2019.

- 2157 [65] Tom Young, Devamanyu Hazarika, Soujanya Poria, and Erik Cambria. Recent trends in
2158 deep learning based natural language processing. *CoRR*, abs/1708.02709, 2017.
- 2159 [66] Zhengyan Zhang, Xu Han, Zhiyuan Liu, Xin Jiang, Maosong Sun, and Qun Liu. Ernie :
2160 Enhanced language representation with informative entities, 2019.

2161 Table des figures

2162	2.1	Espace probabiliste	19
2163	2.2	Espace probabiliste cas général	21
2164	2.3	Exemple du modèle de configuration	27
2165	2.4	Exemple de Composante géante	28
2166	2.5	2-Core d'une composante géante.	29
2167	3.1	Phrase lemmatisée.	32
2168	3.2	Outil NER de Stanford avec la librairie SpaCy.	34
2169	3.3	Un simple exemple d'une Grammaires hors-contexte, l'ensemble des non-terminaux	
2170		N contient un ensemble de base de catégories syntaxiques : S = Sentence, VP =	
2171		Verb Phrase, NP = Noun Phrase, PP = Prepositional Phrase, DT = Determiner,	
2172		Vi = Intransitive Verb, Vt = Transitive Verb, NN = Noun, IN = Preposition. Σ	
2173		est l'ensemble des mots possibles dans la langue.	35
2174	3.4	Treebank	36
2175	3.5	Vecteur one hot	37
2176	3.6	Matrice 10 000 lignes et 300 colonnes	38
2177	3.7	Matrice 10 000 lignes et 300 colonnes	38
2178	3.8	Vecteur de similarité	39
2179	3.9	L'attention pour un mot spécifique.	42
2180	4.1	Schéma relationnel	44
2181	4.2	Schéma OLAP.	45
2182	4.3	Une requête OLAP de dimension 1 et une requête de dimension deux.	46
2183	4.4	Entrepôt de données.	49
2184	4.5	Échantillons provenant d'un entrepôt de données pour les trois <i>Mesures</i>	51
2185	5.1	Exemple de Twitter JSON	56
2186	5.2	Exemple d'arborescence Twitter	57
2187	5.3	D'un tweet à un graphe	58
2188	5.4	Exemple de graphe non symétrique d'utilisateurs sur Twitter	58
2189	5.5	Graphe Twitter	59
2190	5.6	Distributions de degrés	60
2191	5.7	Exemple de la distribution de degrés dans un graphe Twitter	61
2192	5.8	Step reservoir sampling.	65

2193	6.1	L'analyse de la phrase : <i>A boy is jumping on skateboard in the middle of a red</i>	
2194		<i>bridge.</i>	71
2195	6.2	Deux ensembles de vecteurs dans 2 dimensions	72
2196	6.3	Deux Réservoirs pour un texte T	73
2197	6.4	La composante géante du Réservoir R_α	74
2198	7.1	Préfixes et suffixes les plus fréquents.	78
2199	7.2	Les plus fréquents (préfixes ;suffixes) de la racine <i>banag.</i> (' know').	79
2200	7.3	Les vecteurs des 7 préfixes les plus fréquents <i>k-,ka-,n-,ni-,mi-,pa-,t-</i> en deux di-	
2201		mensions.	81
2202	7.4	Vecteur d'une phrase	83
2203	7.5	Les phrases et leurs encodages.	84
2204	7.6	Fréquences des différentes classes syntaxiques existantes dans le corpus.	85
2205	7.7	Arbre de décision pour prédire la classe syntaxique.	85
2206	7.8	Statistiques des prédictions.	86
2207	7.9	Distribution de l'attention pour <i>mi-</i> dans la phrase <i>Mi-melaw k-u wawa t-u tilibi.</i>	87
2208	7.10	Arbre de dérivation de la phrase <i>mi-padang t-u suwal n-ira tatakulaq</i> pour une	
2209		grammaire G	89
2210	8.1	Arbre phylogénique T (centre) à partir de la matrice de distance, et un autre arbre	
2211		proche T' (droite)	92
2212	8.2	DAG de recherche qui regroupe les flux, les clusters et les mots.	93
2213	8.3	DAG de recherche.	98
2214	9.1	Les composantes géantes du Réservoir	101
2215	9.2	Les composantes géantes du Réservoir pour le mot clé <i>Bitcoin.</i>	102
2216	9.3	Nombre d'arêtes/par heure pour 4 flux pendant 24h	103
2217	9.4	Taille de la composante géante pour le Réservoir pondéré	104
2218	9.5	Stabilité du Reservoir sampling et de l'échantillonnage pondéré.	105
2219	9.6	Stabilité du Reservoir sampling et de l'échantillonnage pondéré sur tous les paires.	106
2220	9.7	Distributions de 2 topics	107
2221	9.8	Fraction des thèmes de LDA comme des noeuds de la composante géante du	
2222		réservoir pondéré	108
2223	9.9	La grande composante connectée (C) dans le Réservoir	109
2224	9.10	Le $2 - core(C)$	109
2225	9.11	Arbre de décomposition	110
2226	9.12	Arbre de décomposition où les noeuds sont des tags	110

2227 Liste des tableaux

2228	3.1	Outils pour la reconnaissance des entités	33
2229	3.2	Notation utilisée dans le cadre de LDA	40
2230	5.1	Paramètres pour le point d'accès "POST statuses/filter"	55
2231	9.1	Stabilité du Réservoir pondéré	105

2232 A Phrases de la langue Amis

2233 Comme dans le texte, les affixes sont marqués par “-”, les infixes par “<...>”, les clitiques
2234 par “=”.

2235

2236 1. Itiya satu a remiad, mi-kasui i lakelal ci Sadaban hananay a tamdaw.

2237 2. Sulinay, adihay=tu k-u ni-pi-kasui.

2238 3. Ta-lumaq=tu cira.

2239 4. Si-pa-qeses cira i liyal.

2240 5. A tubu=tu k-ita.

2241 6. Na sa k-u balucuq n-ira.

2242 7. Pa-tangsul han=tu i liyal.

2243 8. Sulinay, ma-patay k-u balucuq n-ira.

2244 9. Pa-keda han=tu n-ira k-u tireng.

2245 10. U nuka n-u Kawas, si-pa-cahcah cira i niyaruq n-u Balaisan.

2246 11. I cuwa k-ita anini hakira ?

2247 12. Saasaan, ira=tu k-u Balaisan.

2248 13. Lepel han=tu cira.

2249 14. U maan a tamdaw-an k-u ma-tini-ay ?

2250 15. U pulut han-an-ay iri, u panga k-u han.

2251 16. Pa-habay-ay n-umita.

2252 17. Tuduh-en=ita.

2253 18. Tuwa mi-sangaq t-u pa-pulul-an n-i Sadaban.

2254 19. Sa-kapah-en k-u rangat.

2255 20. A pa-habay-an cira n-u Balaisan a niyaruq.

2256 21. Tuwa si-balucuq k-iya pa-habay-ay i ci Sadaban-an.

2257 22. Araw si-balucuq-ay cira.

2258 23. Tayra-an n-ira t-u labi.

2259 24. Sasaan ma-puyapuy cira.

2260 25. Tuwa ka-tangasa-an=tu k-u ka-subuc-an.

- 2261 26. Sulinay ma-subuc n-ira.
2262 27. U niyaruq hantu n-uhni, ma-banaq t-u ni-ka-subuc n-ira.
2263 28. ka-ta-tengteng han amin k-u panga n-iya wawa.
2264 29. Ma-susu=tu.
2265 30. Ma-dengay=tu.
2266 31. A patay-en k-ina babuy numita.
2267 32. Uruma k-u talaw n-ira t-u a=patay-en.
2268 33. A maan-en=ita ?
2269 34. Tuwa, tumitumian pa-ka-labi sa=tu cira-an.
2270 35. Melaw han n-i Sadaban, u puut.
2271 36. Tanu ulah sa=tu cira.
2272 37. Ira=tu k-iya Balaisan.
2273 38. Mi-kilim i cira-an.
2274 39. Tanu dungudungus sa=tu k-iya Balaisan.
2275 40. U mi-maan-ay k-isu ?
2276 41. Kiya tireng itini ?
2277 42. Tuwa meduk sa=tu tayra i ni-ka-sadak-an n-ira.
2278 43. Pa-sa-tip sa cira.
2279 44. Ira=tu k-iya Balaisan.
2280 45. A ma-lepel=tu k-ita.
2281 46. Tangasa sa cira i ungcung iri.
2282 47. Ci Sadaban hantu iri ma-tini k-u kakawaw numaku.
2283 48. U mi-laliw-ay k-aku t-u Balaisan.
2284 49. U mai-dudu-ay itakuwan k-ira Balaisan.
2285 50. Na ma-lepel k-aku.
2286 51. Urasisa tayni k-aku.
2287 52. Urasisa tangic k-aku.
2288 53. A maan-en=ita.
2289 54. Pa-ta-lumaq-aw-aku k-isu.
2290 55. Kalat-i k-u tangila numaku haw.
2291 56. Han=tu n-ira ci Sadaban.
2292 57. Hay-i sa=tu.
2293 58. Tangasa sa=tu k-uhni.
2294 59. Palita-an n-i Sadaban k-iya isu.
2295 60. Cima k-isu haw aru ?

- 2296 61. Han=tu n-ira ci Sadaban.
2297 62. Ma-pa-tangasa n-umisu k-aku iri.
2298 63. A caay=tu k-aku ka-pawan t-u ni-pa-urip numisu.
2299 64. ka-si-ayam t-u buhcal-ay.
2300 65. Han=tu n-ira.
2301 66. Sulinay dutuc han=tu nu niyaruq a paysin k-u na-u-suwal n-ira.
2302 67. Ala-en k-u besi numaku.
2303 68. Han=tu n-ira.
2304 69. Sulinay ma-patay sa=tu ci Sadaban.
2305 70. Sa cira.
2306 71. Sa=tu cira.
2307 72. Sa k-aku.
2308 73. Sa cira.
2309 74. Sa=tu k-uheni.
2310 75. Ma-wacay-ay k-uhni.
2311 76. Awaay k-u buduy
2312 77. Ma-'pud=tu k-uhni namaka lutuk haw.
2313 78. Ma-ruqruq itini i 'enar.
2314 79. Mi-awit tu kidudung.
2315 80. Adihay k-u ni-urung t-u kidudung saan iri.
2316 81. Mi-kapet tu sa-sait.
2317 82. Mi-kapet=tu t-u kidudung sa.
2318 83. Mi-sangaq tu talip.
2319 84. Araw ma-banaq k-iya niyaruq itira saan.
2320 85. Ma-banaq tu ngudu.
2321 86. Na caay hen ka-banaq tu ngudu.
2322 87. Dihku=tu k-uhni t-u ka-siqnaw-an iri.
2323 88. Ma-hemek k-uhni.
2324 89. Pa-beli-en k-aku t-u belac iri.
2325 90. Pa-beli-en-aku k-amu t-unian u kidudung.
2326 91. Pa-beli han 'amin n-uhni t-iya taw-an pa-cakay-ay t-u ciyapu.
2327 92. Manay si-buduy sa ku Pangcah.
2328 93. Pa-se-banaq n-iya pa-cakay-ay tu ciyapo sananay.
2329 94. Ma-rimurak ku niyaroq.
2330 95. Mi-cakay t-iya kidudung-an.

2331 Phrases négatives.

- 2332 1. Caay ka-demec k-u nanum.
2333 2. Caay ka-tala-mana k-u hadui n-ira.
2334 3. Caay=tu ka-wacay iri.

2335 **Résumé :**

2336 Nous étudions des données de nature diverse sous forme de flux, en particulier :

- 2337 • Base de données,
- 2338 • Réseaux sociaux,
- 2339 • Données de texte.

2340 Pour une base de données qui suit un schéma relationnel, un schéma d'analyse *OLAP* (Online
2341 Analytical Processing) définit une des tables de la base de données comme une table d'analyse.
2342 Nous supposons que les tuples de la table d'analyse arrivent sous forme d'un flux. Nous étudions
2343 l'approximation des requêtes *OLAP*, en échantillonnant de manière non uniforme les tuples du
2344 flux sans stocker les données d'analyse et donnons un modèle de préférence dans ce cadre.

2345 Dans le cas du réseau social *Twitter*, nous observons un flux de tweets qui contiennent un
2346 tag donné et le transformons en un flux d'arêtes d'un graphe. Nous souhaitons étudier l'existence
2347 des grands clusters dans le graphe ainsi obtenu. Nous proposons une méthode d'échantillonnage
2348 uniforme qui va associer au graphe un sous-graphe aléatoire et étudions les composantes géantes
2349 de ce sous-graphe aléatoire comme témoin des grands clusters du graphe d'origine.

2350 Pour un flux de texte, nous considérons les paires de mots dans une phrase lemmatisée comme
2351 des arêtes d'un graphe où les nœuds sont les mots. Nous transformons le flux de texte en flux
2352 d'arêtes. Nous échantillonnons les arêtes proportionnellement à la similarité *Word2vec* des mots.
2353 Nous analysons ensuite les composantes géantes.

2354 Nous étendons les vecteurs *Word2vec* en prenant en compte la morphologie d'une langue,
2355 en particulier la structure des préfixes et des suffixes d'un mot. Les nouveaux vecteurs d'un mot
2356 ont 3 composantes : un vecteur pour le préfixe, un vecteur pour la racine et un vecteur pour le
2357 suffixe. Nous définissons un vecteur probabiliste pour les phrases.

2358 Sur les trois types de données, nous avons échantillonné selon des distributions précises. Les
2359 résultats principaux de cette thèse montrent comment approcher les propriétés de ces données
2360 avec ces échantillons.

2361 *Descripteurs : Algorithmes de streaming, Approximation, Requêtes OLAP, Préférences, Cluste-*
2362 *ring, Graphes dynamiques, TAL, Morphologie*

2363 **Abstract :**

2364 We study different types of data as a stream :

- 2365 • Databases,
- 2366 • Social Networks,
- 2367 • Text data.

2368 For a database that follows a relational schema, an *OLAP* (Online Analytical Processing)
2369 analysis schema defines one of the database tables as an analysis table. We suppose that the
2370 tuples of the analysis table arrive as a stream. We study the approximation of *OLAP* queries,
2371 by sampling with weights the tuples of the stream without storing the analysis data. We give a
2372 *preference* model in this context.

2373 For the social network *Twitter*, we observe a stream of tweets that contain any given tag and
2374 transform it into a stream of edges of a graph. We study the existence of large clusters in the
2375 generated graph. We propose a uniform sampling method that will associate a random subgraph
2376 of the graph and study the giant components of this random subgraph as a witness of the large
2377 clusters of the original graph.

2378 For a stream of text, we consider the pairs of words in a lemmatized sentence as edges of a
2379 graph where the nodes are the words. We transform the stream of text into a stream of edges.
2380 We sample the edges proportionally to the *Word2vec* similarity of the words. We then analyze
2381 the giant components.

2382 We extend the classical *Word2vec* vectors, in order to take into account the morphology of the
2383 words, i.e. the prefixes, roots and suffixes. A new vector has 3 components, one for the prefix, one
2384 for the root and one for the suffix. We define a probabilistic vector associated with the sentences.

2385 With the three types of data, we sampled according to specific distributions. The main results
2386 of this thesis show how the properties of these data can be approximated with these samples.

2387 *Keywords : Streaming algorithms, Approximation, OLAP queries, Preferences, Clustering, Dy-*
2388 *namic graphs, NLP, Morphology*

2389 Nota : cette page, dernière de couverture, sera retournée avant reliure.