

TIPE: AUTOCOMPLETION

MP 2016 - 2017



- I. Présentation du TIPE
- II. Réalisation
- III. Amélioration
- IV. Étude théorique
- V. Conclusion

SOMMAIRE



I - PRESENTATION DU TIPE



- ▶ **Texte** : suite de caractères (ici, éléments de l'alphabet latin ou de sa ponctuation) découpée en propositions (séparées par la ponctuation), elles-mêmes découpées en mots.
- ▶ **Autocomplétion** : prédiction d'un mot d'un texte en fonction de ce qui le précède

$$\mathcal{A} = \{a, b, \dots, z, \acute{e}, \grave{e}, \dots\}$$

$$\mathcal{D} = \{a, \grave{a}, \text{abaissa}, \dots, \text{zythums}\}$$

$$\mathcal{P} = \{, . ? ! : \dots " - \}$$

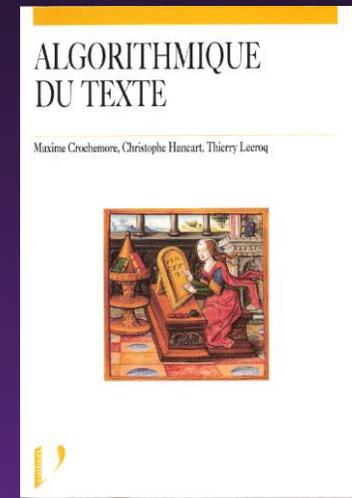
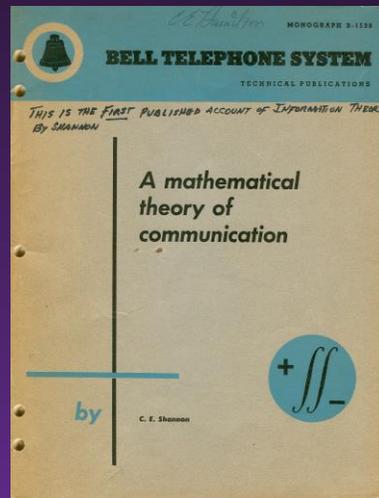
DEFINITIONS

Sous une forme écrite ou vocale, le texte est le seul véhicule fiable des concepts abstraits.

Algorithmique du texte, Crochemore, Hancart, Lecroq

- ▶ L'essor de technologies appelle celui des théories
- ▶ Le développement de la théorie de l'information est donc capital
- ▶ Pourtant l'autocomplétion est peu étudiée

CONTEXTE



- ▶ Écrire un programme Python d'autocomplétion en cours de saisie à partir d'une base de données
- ▶ Constituer automatiquement cette base de données
- ▶ Étudier et optimiser le programme
- ▶ Étudier les performances théoriques

OBJECTIFS



II - REALISATION



- ▶ Dans un premier temps, on compile des informations sur la succession des mots avec le programme « *compiler* »
- ▶ Ces informations seront stockées dans un fichier qui fait office de base de données
- ▶ Ensuite, on exploite cette base avec le programme principal, « *filler* »

ORGANISATION



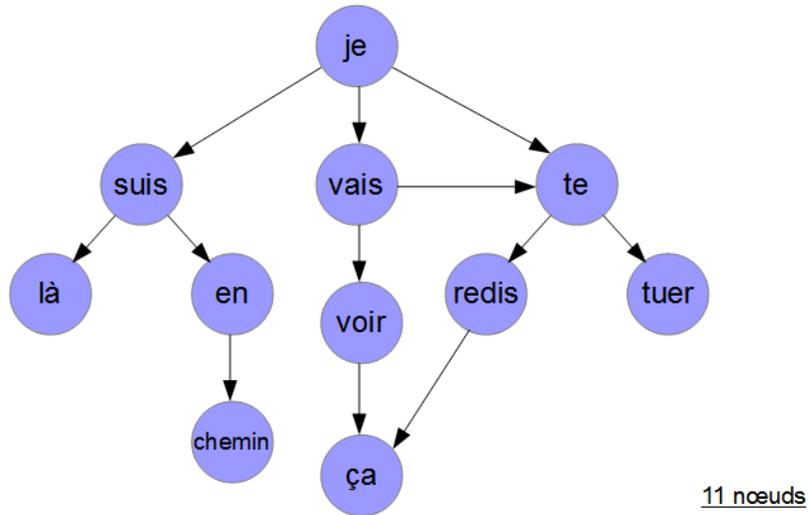


Figure 1 : Modélisation d'une phrase par une chaîne de Markov
(source sans mémoire)

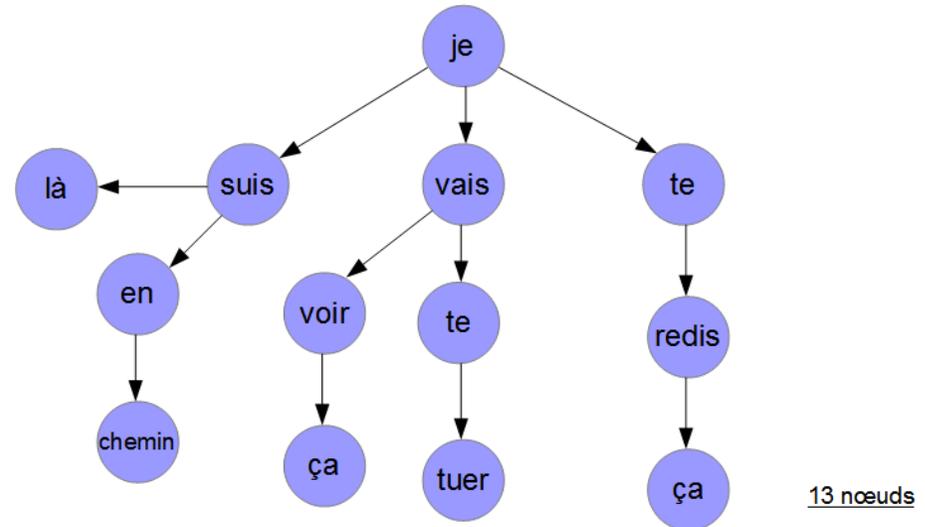


Figure 2 : Modélisation d'une phrase par une chaîne quelconque
(source avec mémoire)

STRUCTURE DE LA LANGUE

FILLER

```
# -*- coding: utf-8 -*-

#Ouverture du fichier et acquisition des données-----
with open("Successeurs.txt", encoding="utf-8", mode="r") as fichier:
    arbre = fichier.read()
    noeuds = arbre.splitlines()
    successeurs = []
    # "successeurs" est une liste de listes, contenant chacune un mot et ses successeurs
    for i in range(0, len(noeuds)):
        successeurs.append(noeuds[i].split())

#Fonction de succession-----
#renvoie la liste des mots suivant probablement w, par ordre de probabilité
def succ(w):
    if type(w) != str:
        print("ERREUR : vous devez entrer une chaîne de caractères !")
    else:
        #si début de phrase, file = première liste
        if w in "[\n\u0022.?!;,:*-(—) ]":
            file = successeurs[0][-3:]
        else:
            k=1
            #file sera renvoyé tel quel si la boucle se termine sans l'avoir modifiée
            file = [ 'Aucun successeur trouvé' ]
            while k in range(1, len(successeurs)) and file == [ 'Aucun successeur trouvé' ]:
                #évitte de modifier successeurs
                noeud = successeurs[k].copy()
                #quand on trouve le mot recherché
                if noeud[0] == w:
                    #pour ne garder que les successeurs et pas w
                    noeud.remove(w)
                    file = noeud
                k=k+1
            return(file)

#Saisie intuitive-----

saisie = [ ' ' ]
#condition permettant d'arrêter la saisie si besoin
running = True
while running:
    #affichage des suggestions
    S = succ(saisie[-1].lower())
    print(S)
    i = input("Tapez 1, 2 ou 3 (0 si mot absent)\n")
    if i == 'FinDuTest':
        running = False
    #si le mot voulu est proposé
    elif int(i) != 0:
        saisie.append(S[int(i)-1])
    #sinon on laisse l'utilisateur le taper
    else:
        saisie.append(input())
```

FILLER

```
#!/usr/bin/env python3
# -*- coding: utf-8 -*-

#Ouverture du fichier et acquisition des données-----
with open("Successeurs.txt", encoding="utf-8", mode="r") as fichier:
    arbre = fichier.read()
    noeuds = arbre.splitlines()
    successeurs = []
    # "successeurs" est une liste de listes, contenant chacune un mot et ses successeurs
    for i in range(0, len(noeuds)):
        successeurs.append(noeuds[i].split())

#Fonction de succession-----
#renvoie la liste des mots suivant probablement w, par ordre de probabilité
def succ(w):
    if type(w) != str:
        print("ERREUR : vous devez entrer une chaîne de caractères !")
    else:
        #si début de phrase, file = première liste
        if w in "[\n\u0022.?!,;,:*~(---)]":
            file = successeurs[0][-3:]
        else:
            k=1
            #file sera renvoyé tel quel si la boucle se termine sans l'avoir modifiée
            file = [ 'Aucun successeur trouvé' ]
            while k in range(1, len(successeurs)) and file == [ 'Aucun successeur trouvé' ]:
                #évite de modifier successeurs
                noeud = successeurs[k].copy()
                #quand on trouve le mot recherché
                if noeud[0] == w:
                    #pour ne garder que les successeurs et pas w
                    noeud.remove(w)
                    file = noeud
                k=k+1
            return(file)

#Saisie intuitive-----
saisie = [ ' ' ]
#condition permettant d'arrêter la saisie si besoin
running = True
while running:
    #affichage des suggestions
    S = succ(saisie[-1].lower())
    print(S)
    i = input("Tapez 1, 2 ou 3 (0 si mot absent)\n")
    if i == "FinDuTest":
        running = False
    #si le mot voulu est proposé
    elif int(i) != 0:
        saisie.append(S[int(i)-1])
    #sinon on laisse l'utilisateur le taper
    else:
        saisie.append(input())
```

LECTURE

FONCTION DE
SUCCESSION

INTERFACE

- ▶ En testant *filler*, le programme propose le mot attendu dans 17,3% des cas

RESULTATS

Je ne pense pas changer l'histoire de ce domaine car il y a eu beaucoup de grands hommes ici avant moi, et je ne travaille pas dessus depuis assez longtemps.

III - AMELIORATION



- ▶ Afin d'accélérer le parcours de la base de données, on la trie dans l'ordre alphabétique
- ▶ On doit donc également modifier *filler* en implémentant une recherche par dichotomie
- ▶ Un tri par fréquence d'apparition serait possible, mais moins intéressant au final

STRUCTURE DE LA BASE DE DONNEES

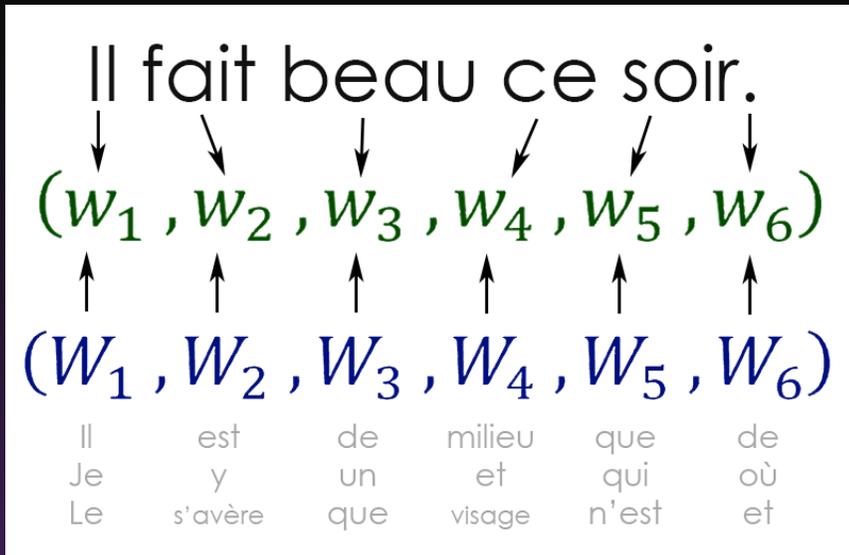


IV - ETUDE THEORIQUE



- ▶ Avec n le nombre de mots enregistrés, la complexité de *filler* est en $O(n)$ pour le fichier non trié, et en $O(\log(n))$ une fois l'amélioration effectuée.

COMPLEXITE



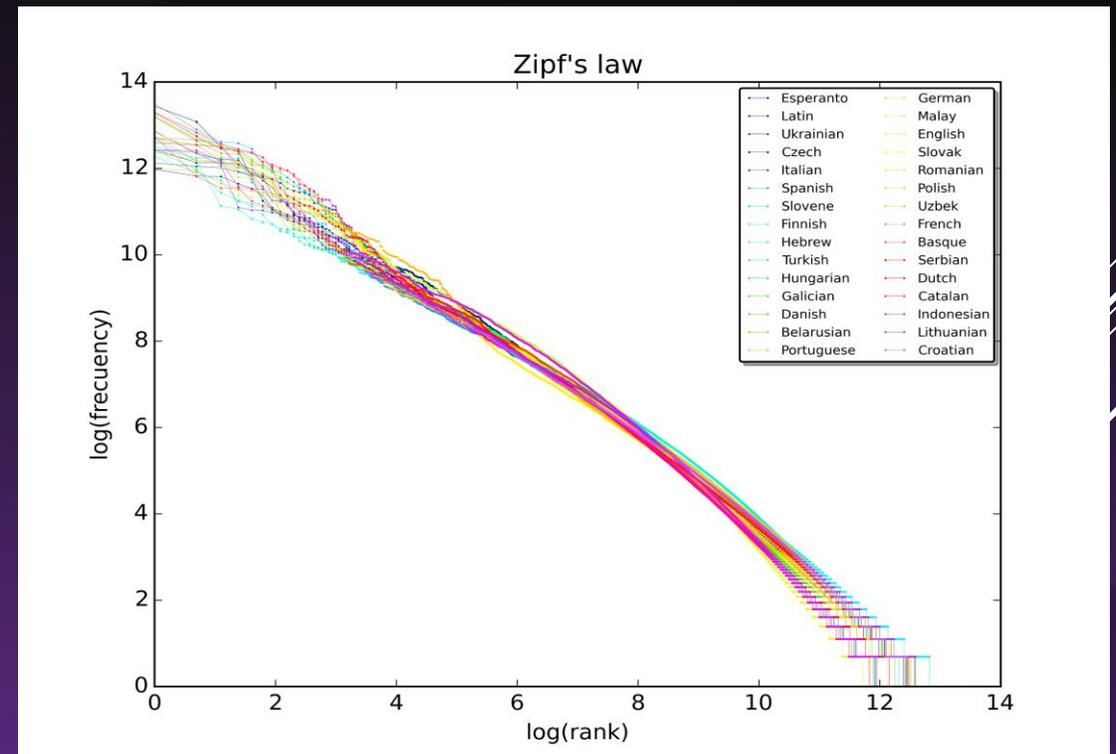
- ▶ Soit $(w_k)_{k \in \llbracket 1;n \rrbracket}$ une phrase et $(W_k)_{k \in \llbracket 1;n \rrbracket}$ la famille de variables aléatoires correspondante.
- ▶ On cherche la probabilité $P(W_k = w_k | W_{k-1} = w_{k-1})$

MODELISATION DU LANGAGE

- ▶ La répartition des mots suit la loi de Zipf : en notant m_j le j^{e} mot le plus courant, on a :

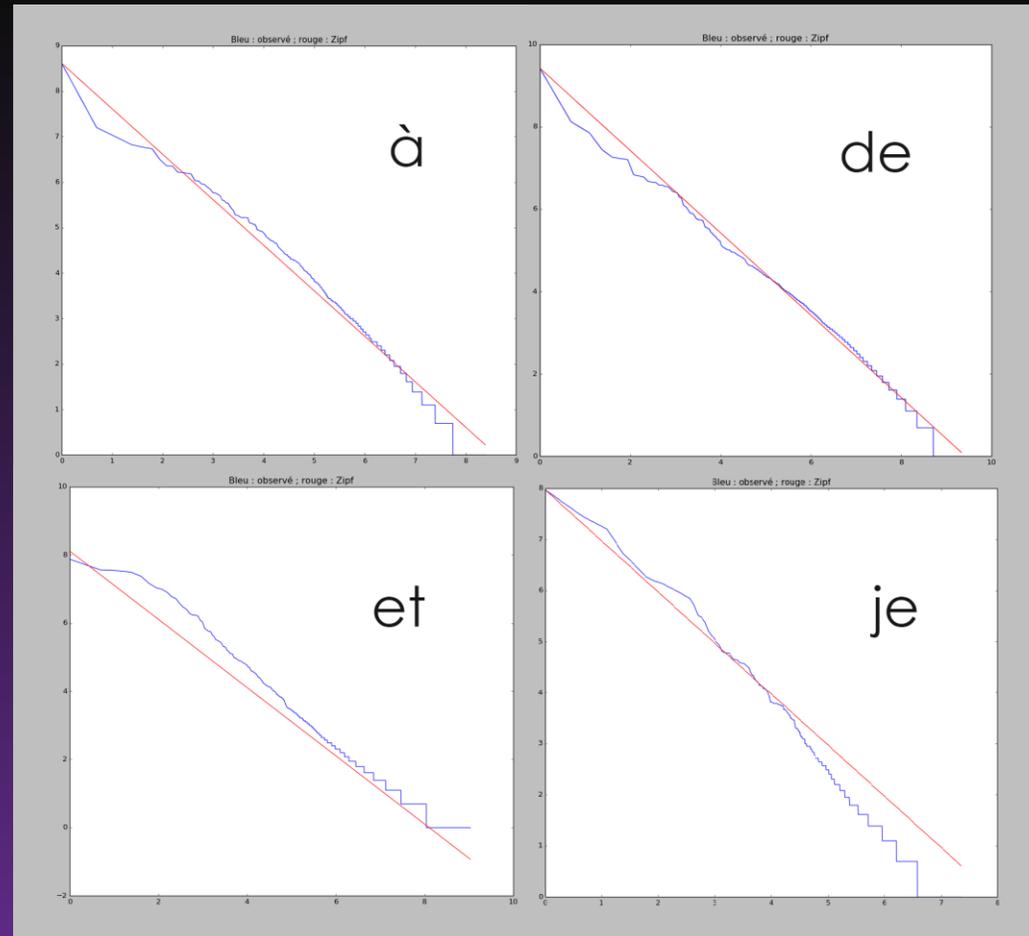
$$P(W_k = m_j) = \frac{P(W=m_1)}{j} = \frac{1}{j\mathcal{S}_N} \quad \text{où } \mathcal{S}_N = \sum_{k=1}^N \frac{1}{k} \quad \text{avec } N = \text{card}(W_k(\Omega))$$

LOI DE ZIPF



- ▶ Expérimentalement, les successeurs d'un mot donné suivent aussi une loi de Zipf

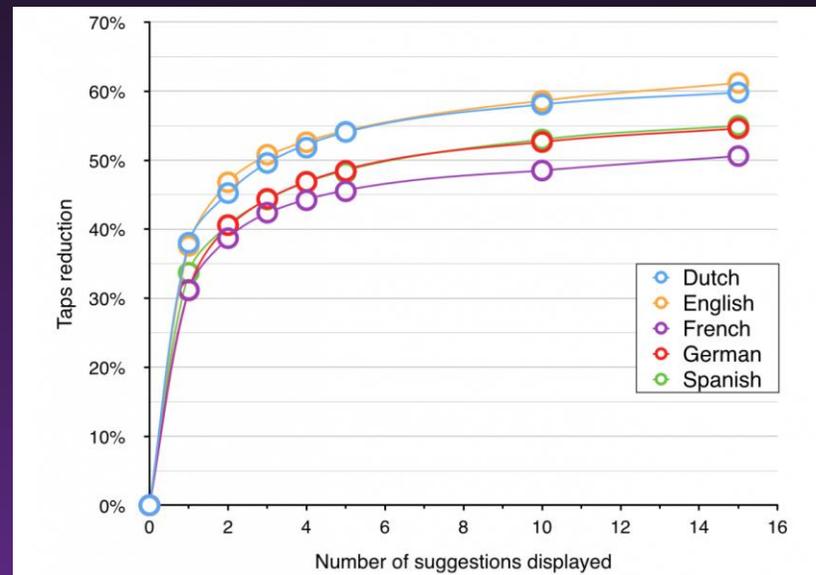
LOI DE ZIPF



- ▶ D'où $P(W_k = m_j | W_{k-1} = w_{k-1}) = \frac{1}{j\mathcal{S}_N}$, où N est le nombre de successeurs possibles
- ▶ Notant $S = \{s_1; s_2; s_3\}$ le triplet de suggestions affichées, on trouve en moyenne :

$$P(W_k \in S) = \left(1 + \frac{1}{2} + \frac{1}{3}\right) \frac{1}{\mathcal{S}_N} \approx 49,7\%$$

PROBABILITES



V - CONCLUSION



- ▶ Nous avons bien créé un programme d'autocomplétion fonctionnel
- ▶ Toutefois ses capacités ne sont pas encore optimales

BILAN



- ▶ Élargir la base de données peut permettre de meilleures prédictions
- ▶ Utiliser une variante du modèle de Zipf
- ▶ Prendre un modèle de Markov d'ordre 2

PISTES D'AMÉLIORATION



