

Traitement du Langage Naturel et Linguistique

LE BELLEGO Victor, MARTZLOFF Alice, MOUSSOX Vincent

M2 IAAA 2021/2022

1 Certaines langues sont elles plus difficiles à analyser que d'autres ?

L'objectif de ce projet est d'analyser et de comprendre les raisons pour lesquelles des analyseurs syntaxiques partageant la même architecture et entraînés sur la même quantité de données obtiennent des performances très différentes sur différentes langues. On peut observer ce phénomène dans la Table 2 qui présente les performances calculées à l'aide des mesures LAS (Labeled Accuracy Score) et UAS (Unlabeled Accuracy Score) obtenues par un analyseur sur 36 langues différentes.

lang	las	uas	lang	las	uas	lang	las	uas
da	66.18	73.06	zh	46.76	55.60	pl	68.76	80.69
hr	58.88	69.67	lv	51.46	59.32	sv	64.52	71.11
id	69.99	75.63	he	64.82	70.06	cs	69.77	77.40
ar	60.88	69.78	ko	46.06	55.00	nl	54.92	63.69
eu	47.94	58.17	ja	71.54	82.71	hu	57.21	69.24
it	74.82	81.01	ca	64.79	71.13	bg	68.08	78.14
fa	47.70	56.17	en	66.57	71.22	vi	48.40	49.62
es	66.09	71.01	pt	70.84	74.55	sl	49.46	61.51
ro	57.23	68.61	et	59.04	73.09	nno	73.58	78.30
de	58.12	64.19	fr	68.40	73.46	nob	66.22	73.93
hi	64.03	72.54	el	69.96	76.58			

Table 1 – LAS/UAS calculées pour les différentes langues (sans les *configurational features*)

1.1 Variables explicatives

Nous faisons l'hypothèse qu'il est possible de prédire les capacités d'un analyseur entraîné sur une langue à partir de variables directement observables. La plupart sont observées sur le corpus d'apprentissage, mais nous utilisons aussi un score de complexité à partir du WALS.

1.1.1 Observations effectuées sur le corpus d'apprentissage

Ces variables observables sur le corpus sont les suivantes :

- Taille du vocabulaire
- Longueur moyenne d'une phrase
- Longueur moyenne d'un mot
- Taux de projectivité
- Longueur moyenne des dépendances
- Nombre moyen de dépendances
- Longueur moyenne de la plus longue dépendance

Taux de projectivité

La projectivité est une propriété de certains arbres de dépendance. Elle se traduit par le fait que lorsqu'on trace les dépendances au dessus des mots de la phrase, il n'y a jamais de croisement.

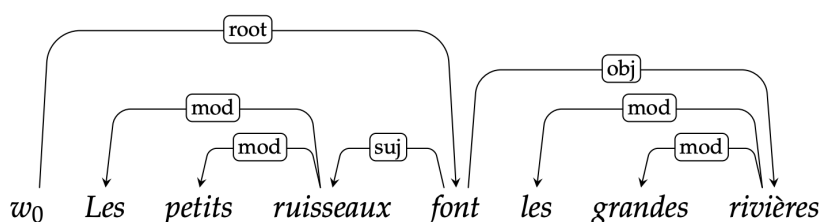


Figure 1 – Exemple d'une phrase française projective

La grande majorité des phrases françaises sont projectives, mais il existe certains cas de non projectivité.

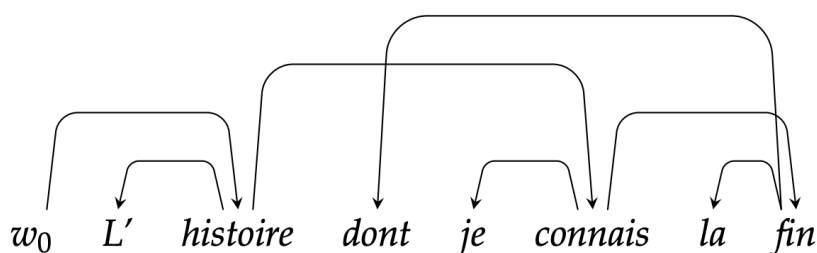


Figure 2 – Exemple d'une phrase française non projective

Le degré de non projectivité varie d'une langue à l'autre. Il est élevé pour le tchèque, le néerlandais ou le turc, par exemple.

Nous supposons que le taux de non projectivité peut expliquer les performances d'un analyseur.

Le calcul du taux de projectivité est très simple. Il suffit de compter, pour une langue, le nombre de phrases contenues dans le fichier `lang.conllu` et le nombre de phrases contenues dans le fichier `lang_proj.conllu` (qui ne contient que les phrases projectives) et d'en faire le quotient (voir la fonction `taux_projectivite()` du fichier `featuresAnalyser.py` en annexe)

Longueur des dépendances

Nous supposons que la longueur des dépendances des phrases projectives ou non projectives peut aussi expliquer les performances de l'analyseur.

Pour calculer la longueur des dépendances, nous utilisons les fichiers `.mcf` et la fonction `long_depandances()` (voir `featuresAnalyser.py` en annexe). L'algorithme opère de la manière suivante :

1. découpage des phrases grâce à l'attribut EOS (*End Of Sentence*)
2. création liste de listes contenant l'attribut GOV (positions relatives des gouverneurs pour chaque mot)
3. conversion de ces positions relatives en positions absolues grâce à la fonction `rel2abs()`
4. création de l'arbre correspondant pour chaque liste grâce à la fonction `maketree()`
5. calcul de la profondeur de l'arbre

Nous avons extrait 3 variables explicatives pour chaque langue :

- Longueur moyenne des dépendances
- Nombre moyen de dépendances par phrase
- Longueur moyenne de la plus longue dépendance de chaque phrase

lang	Taille du vo- cabu- laire	Taille moyenne des phrases	Lon- gueur moyenne des mots	Taux de non projec- tivité	Nombre moyen de dépen- dances	Longueur moyenne des dépen- dances	Longueur de la plus longue dépendance
da	3545	15.71	7.62	0.2	10.77	3.15	4.56
hr	3974	19.2	7.65	0.08	13.06	3.64	5.59
id	3580	18.72	7.4	0.04	12.64	3.66	5.59
ar	5304	36.68	6.43	0.09	17.07	5.04	7.83
eu	2772	11.39	7.31	0.34	7.26	2.76	4.01
it	3354	20.23	8.06	0.05	13.34	3.49	5.11
nno	2161	14.9	7.42	0.08	10.87	3.11	4.59
fa	3904	24.02	5.42	0.07	14.64	3.78	6.11
es	3424	25.63	7.58	0.08	17.69	3.93	5.94
ro	3767	20.01	7.01	0.12	13.52	3.75	5.75
de	4121	16.38	9.1	0.11	11.5	3.18	4.52
hi	2453	19.43	5.33	0.14	12.14	3.44	5.0
zh	4289	21.62	1.99	0.0	15.34	3.45	5.25
lv	3359	11.68	7.34	0.08	8.39	2.98	4.15
he	6050	29.75	5.16	0.01	15.49	3.99	6.18
ko	10283	10.42	3.63	0.12	5.72	3.39	4.55
ja	4034	20.18	2.51	0.0	14.74	3.54	5.36
ca	2904	27.82	7.41	0.09	18.93	3.99	6.13
en	3375	14.61	6.57	0.08	10.32	3.07	4.4
pt	3980	23.93	7.5	0.13	15.18	3.75	5.55
et	3683	8.71	7.31	0.09	6.55	2.56	3.4
fr	3351	23.27	7.53	0.12	15.5	3.61	5.43
el	2734	21.28	7.86	0.11	13.76	3.82	5.63
pl	4762	8.51	7.79	0.01	6.09	2.85	3.9
sv	4451	14.83	7.66	0.11	9.79	3.11	4.51
cs	4875	15.52	7.35	0.13	10.22	3.45	5.02
nl	3552	12.52	7.62	0.27	7.73	2.82	3.75
hu	3677	19.69	8.05	0.23	12.5	3.37	5.1
nob	2155	13.1	7.41	0.08	9.42	3.05	4.42
bg	3736	11.68	7.8	0.03	8.4	3.19	4.42
vi	1798	12.36	3.42	0.02	8.81	3.1	4.68
sl	4219	13.38	7.3	0.14	9.34	2.97	4.2
moy	3863	18.04	6.73	0.10	11.77	3.41	5.02

Table 2 – Synthèse des variables explicatives pour les corpus *train* pour chaque langue

1.1.2 Complexité selon M. Parkvall

Dans son papier *The simplicity of creoles in a cross-linguistic perspective* (Parkvall 2008) sorti en 2008, Mikael Parkvall s'intéresse à quantifier la complexité des langues. Il part du postulat qu'une expression est d'autant plus complexe qu'elle implique de règles, c'est-à-dire qu'elle requiert une longue description. Ainsi, l'hypothèse de base de l'auteur est la suivante : une langue complexe est une langue avec des constructions plus complexes. Il

explore un aspect de complexité structurelle.

Prenons par exemple la voix passive. Lorsqu'elle existe dans une langue, il faut pouvoir définir comment passer de la voie active à la voix passive, ce qui exige une explication de règle supplémentaire. Une langue qui possède une voix passive est donc, en ce qui concerne cette construction spécifique, plus complexe qu'une autre n'en possédant pas. Si on énumère donc un grand nombre de « constructions complexes », la langue la plus complexe sera celle qui en compte le plus grand nombre.

Pour extraire les « constructions complexes » qu'on peut trouver dans une langue, Parkvall utilise le set de données *World Atlas of Linguistic Structures* (WALS) publié en 2005 par Haspelmath et al. Il choisit 155 langues parmi plus de 2 500, et 47 caractéristiques parmi plus de 140.

Choix des caractéristiques

Parkvall exclut des caractéristiques selon un raisonnement défendu dans son papier et qui s'efforce de mettre la majorité des linguistes d'accord sur le fait qu'une caractéristique apporte ou non de la complexité. Il retient les caractéristiques suivantes :

Caractéristiques du WALS		
Size of consonant inventories	Distance contrast in demonstratives	Morphological imperative
Size of vowel quality inventories	Gender in pronouns	Morphological optative
Phonemic vowel nasalization	Politeness in pronouns	Grammaticalized evidentiality distinctions
Complexity of syllable structure	Person marking on adpositions	Both indirect and direct evidentials
Tone	Comitative ≠ instrumental	Non-neutral marking of full NPs
Overt marking of direct object	Ordinals exist as a separate class beyond 'first'	Non-neutral marking of pronouns
Double marking of direct object	Suppletive ordinals beyond 'first'	Subject marking as both free word and agreement
Possession by double marking	Obligatory numeral classifiers	Passive
Overt possession marking	Possessive classification	Antipassive
Reduplication	Conjunction 'and' ≠ adposition 'with'	Applicative
Gender	Difference between nominal and verbal conjunction	Obligatorily double negation
Number of genders	Grammaticalized perfective/imperfective	Asymmetric negation
Non-semantic gender assignment	Grammaticalized past/non-past	Equative copula ≠ Locative copula
Grammaticalized nominal plural	Remoteness distinctions of past	Obligatorily overt equative copula
Definite articles Indefinite articles	Morphological future	
Inclusivity (in either pronouns or verb morphology)	Grammaticalized perfect	

Table 3 – Liste des caractéristiques extraite directement du WALS

Il ajoute à ces caractéristiques, d'autres données « résiduelles » d'auteurs contributeurs au WALS. Ce sont les suivantes :

Caractéristiques d'auteurs du WALS		
Demonstratives marked for number	Demonstratives marked for gender	Demonstratives marked for case
Total amount of verbal suppletion	Alienability distinctions	

Table 4 – Liste de caractéristiques proposées par les auteurs du WALS

Enfin, il s'intéresse aussi à une donnée de Harley and Ritter (2002) à laquelle il a eu accès :

Caractéristique de Harley et Ritter
Number of pronominal numbers

Table 5 – Une caractéristique accessible, inspirée de Harley et al. (2002)

Les valeurs de ces caractéristiques ont toutes été traduites par l’auteur comme des valeurs comprises entre 0 et 1 :

- « Oui » ou « non » deviennent 0 ou 1 avec parfois l’introduction de valeurs intermédiaire 0,5 ;
- Des valeurs d’intensité comprises entre 1 et 4 sont comprimées en : 0 - 0,25 - 0,5 - 0,75 et 1 ;
- Des valeurs catégoriques comme la classification en « simple », « modérément complexe » et « complexe » sont traduites en 0 , 0,5 et 1.

Choix des langues

L’auteur s’attèle à choisir des langues dont les annotations sont le moins lacunaires possible pour les caractéristiques décrites ci-dessus. Pour une langue i donnée :

$$\text{Score}_i = \frac{\sum_{k=1}^L \text{contribution}_k}{L} \quad (1)$$

où k représente une caractéristique. Chaque langue comptant un nombre différent L de caractéristiques (parmi celles choisies par l’auteur) effectivement annotées pour cette langue.

N.B. : en effet, de même que pour les caractéristiques que nous extrayons directement de nos données, Parkvall note que le set de données WALS n’est pas identiquement distribué : les langues ne sont pas identiquement annotées pour les caractéristiques proposées...

1.2 Cohérence des annotations

Dans leur article *Divergences entre annotations dans le projet Universal Dependencies* (Nivre et al. 2017) et leur impact sur l’évaluation de l’étiquetage morpho-syntaxique (Wisniewski et Yvon 2018), Guillaume Wisniewski et François Yvon montrent que la dégradation des performances observée lors de l’application d’un analyseur morpho-syntaxique à des données hors domaine comme ici, d’une langue à l’autre, résulte d’incohérences entre les annotations des ensembles de test et d’apprentissage. Ils montrent qu’appliquer le principe de variation des annotations de Dickinson & Meurers (Meurers 2003) permet d’identifier les erreurs d’annotation et donc les incohérences et évaluer leur impact. Nous souhaitons nous inspirer de ces méthodes mais n’avons pu en raison notamment du temps qui nous manque proposer une telle amélioration...

1.3 Mise en oeuvre

Nous réalisons une régression multiple à l'aide du modèle `LinearRegression` de `scikit-learn` pour prédire LAS après **normalisation** des variables explicatives.

Il est d'abord intéressant de regarder les potentielles corrélations entre nos variables explicatives.

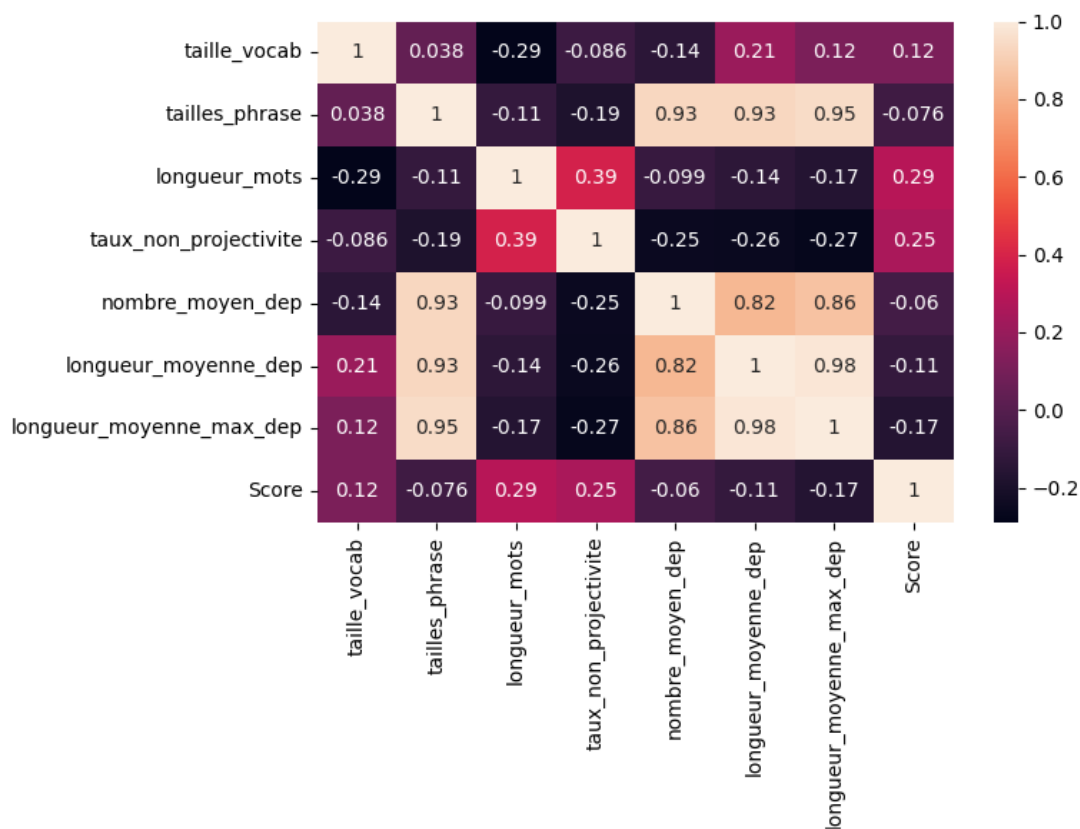


Figure 3 – Corrélations entre les variables explicatives

On remarque, en toute logique, que la taille des phrases est corélée au nombre et à la taille des dépendances. On constate aussi que la longueur moyenne des dépendances est corélée au nombre moyen des dépendances : plus il y a de dépendances, plus elles sont longues.

Nos variables normalisées sont notées comme suit :

- x_1 : Taille du vocabulaire
- x_2 : Longueur moyenne d'une phrase
- x_3 : Longueur moyenne d'un mot
- x_4 : Taux de projectivité
- x_5 : Longueur moyenne des dépendances
- x_6 : Nombre moyen de dépendances
- x_7 : Longueur moyenne de la plus longue dépendance
- x_8 : Score de M. Parkvall (WALS)

Les LAS sont notés y . On souhaite estimer les y tels que :

$$\exists [\alpha_1, \dots, \alpha_8, \beta] \in \mathbb{R}, \forall i \quad y_i = \beta + \sum_{j=1}^8 \alpha_j \cdot x_{i,j}$$

On utilise l'estimateur des moindres carrés :

$$(\hat{\alpha}, \hat{\beta}) = \arg \min_{\alpha, \beta} \sum_{i=1}^n (y_i - (\alpha x_i + \beta))^2$$

On obtient les résultats suivants :

$\beta = 61.66, \alpha_1 = -2.32, \alpha_2 = -0.75, \alpha_3 = 3.63, \alpha_4 = -3.56, \alpha_5 = 3.76, \alpha_6 = 17.40, \alpha_7 = -18.33, \alpha_8 = 0.02$

avec $R^2 = 0.476$, qui est moyen. L'erreur RMS s'élève à 6.198.

Du fait de la normalisation des variables d'entrée, il est directement possible de déduire l'impact de chaque variable. On remarque que le score de M. Parkvall (WALS) a un impact quasi nul sur la régression. En revanche, la longueur moyenne des dépendance et la longueur moyenne de la plus longue dépendance ont une importance majeure, bien qu'opposée.

1.4 Conclusion

1.4.1 Est-il possible de connaître a priori les performances de l'analyseur en fonction de certaines caractéristiques de la langue ?

Bien que les résultats obtenus ne soient pas à la hauteur des espérances, il semblerait certaines caractéristiques de la langue permettent de connaître a priori les performances de l'analyseur.

Notre étude porte sur un **corpus de faible taille qui comporte, par essence, de nombreux biais** ; il est donc difficile de généraliser. Une étude avec des corpus de plus grande taille serait nécessaire pour affiner la régression et ainsi conclure de manière plus certaine.

1.4.2 Est-il possible d'utiliser les conclusions de l'étude statistique pour améliorer les performances de l'analyseur ?

Il est a priori possible d'utiliser l'étude statistique pour améliorer les performances de l'analyseur. En effet, connaissant les caractéristiques de la langues, on peut choisir un analyseur moins sensible à la non projectivité ou aux longues dépendances, mais plus coûteux seulement pour les langues qui le requièrent. Les langues projectives (comme le japonais) peuvent être traités avec des analyseurs plus simples et moins coûteux.

Ainsi, l'étude statistique permet d'améliorer les performances de l'analyseur en orientant l'investissement en ressources pour les langues les plus « complexes »

Références

- Meurers, Detmar (jan. 2003). “Detecting Errors in Part-of-Speech Annotation”. In : p. 107-114.
- Nivre, Joakim et al. (2017). *Universal Dependencies 2.1*. LINDAT/CLARIAH-CZ digital library at the Institute of Formal and Applied Linguistics (ÚFAL), Faculty of Mathematics and Physics, Charles University. url : <http://hdl.handle.net/11234/1-2515>.
- Parkvall, Mikael (jan. 2008). “The simplicity of creoles in a cross-linguistic perspective”. In : p. 265-285. doi : 10.1075/slcs.94.17par.
- Wisniewski, Guillaume et François Yvon (mai 2018). “Divergences entre annotations dans le projet Universal Dependencies et leur impact sur l’évaluation de l’étiquetage morpho-syntaxique”. In : *Conférence sur le Traitement Automatique des Langues Naturelles*. Rennes, France. url : <https://hal.archives-ouvertes.fr/hal-01793784>.

Annexes

```

../resultsTable.py
1  import os
2  from os import listdir
3  from subprocess import run
4  import pandas as pd
5  from tabulate import tabulate
6
7  base_dir = os.getcwd()
8  os.chdir(base_dir)
9
10 list_lang = [k.replace('.conllu','').replace('train_','') for k in
    ↳ listdir('data') if '.conllu' in k and 'train' in k]
11 print('Training '+str(len(list_lang))+ ' languages')
12
13 os.chdir('expe/')
14 i=1
15 nbre = len(list_lang)
16 for lang in list_lang:
17     try :
18         print('Testing '+lang+' '+str(i)+'/'+str(nbre))
19         command = 'make lang=' + lang
20         run(command.split(), capture_output=True)
21         print('Success !')
22         i += 1
23     except :
24         print("Erreur " + lang)
25         pass
26
27 print('Reading results...')
28 os.chdir('out/')
29 df = pd.DataFrame(columns=["lang", "las", "uas"])
30 for lang in list_lang:
31     try:
32         with open(lang+'.res','r') as f:
33             lines = f.readlines()
34             to_append = lines[-1].split()
35             df_length = len(df)
36             df.loc[df_length] = to_append
37     except :
38         print("Erreur " + lang)
39         pass
40
41 df.set_index('lang', inplace=True)
42 df.to_csv('out.csv')
43 print(tabulate(df, headers='keys', tablefmt='psql'))

../featuresAnalyser.py

```

```

1  import os
2  from os import listdir
3  import pandas as pd
4  import numpy as np
5
6
7  base_dir = os.getcwd()
8
9  def extract_features(lang, filetype, proj=False):
10     if proj == True:
11         lang += "_proj"
12     os.chdir(base_dir + '/data/')
13     taille_phrases = []
14     liste_mots = []
15     with open(filetype + "_" + lang + ".conllu", 'r', encoding='utf8') as f:
16         lines = f.readlines()
17         empty_lines = [-1]
18         for i, line in enumerate(lines):
19             if line == "\n":
20                 empty_lines.append(i)
21         for indice in range(len(empty_lines)-1):
22             phrase = lines[empty_lines[indice]+1:empty_lines[indice+1]-1]
23             taille_phrase = 0
24             for word in phrase:
25                 if "PUNCT" not in word: #On enlève la ponctuation pour compter
26                     taille_phrase += 1
27                 if "PROPN" not in word and "NUM" not in word: #On enlève
28                     ↪ les nombres et les noms propres pour la diversité du
29                     ↪ vocab
30                 word = word.split()
31                 if word[2] != "_":
32                     liste_mots.append(word[2])
33                 else :
34                     liste_mots.append(word[1])
35             taille_phrases.append(taille_phrase)
36         liste_mots = list(set(liste_mots))
37         return {"liste_mots" : liste_mots, "taille_phrases":taille_phrases}
38
39  def var_exp_df(filetype):
40     os.chdir(base_dir)
41     list_lang = [k.replace('.conllu', '').replace('train_', '') for k in
42     ↪ listdir('data') if
43     ↪ '.conllu' in k and 'train' in k]
44     tailles_vocab, tailles_phrases, longueur_mots, tx_non_projectivite,
45     ↪ nombre_moyen_dep, longueur_moyenne_dep, longueur_moyenne_max_dep =
46     ↪ [], [], [], [], [], [], []
47     for lang in list_lang:
48         print('Analysing lang : ', lang)
49         my_dict = extract_features(lang, filetype)
50         tailles_vocab.append(len(my_dict["liste_mots"]))
51         tailles_phrases.append(np.mean(my_dict["taille_phrases"]))

```

```

47     longueur_mots.append(np.mean([len(x) for x in my_dict["liste_mots"]]))
48     tx_non_projectivite.append(1-taux_projectivite(filetype,lang))
49     my_dict_dep = long_dependances(filetype,lang)
50     nombre_moyen_dep.append(my_dict_dep['nombre_moyen_dep'])
51     longueur_moyenne_dep.append(my_dict_dep['longueur_moyenne_dep'])
52
53     ↪ longueur_moyenne_max_dep.append(my_dict_dep['longueur_moyenne_max_dep'])
54 df = pd.DataFrame({'lang':list_lang,
55                   'taille_vocab' : tailles_vocab,
56                   'tailles_phrase' : np.round(tailles_phrases,2),
57                   'longueur_mots' : np.round(longueur_mots,2),
58                   'taux_non_projectivite':
59                   ↪ np.round(tx_non_projectivite,2),
60                   'nombre_moyen_dep': np.round(nombre_moyen_dep,2),
61                   'longueur_moyenne_dep':
62                   ↪ np.round(longueur_moyenne_dep,2),
63                   'longueur_moyenne_max_dep':
64                   ↪ np.round(longueur_moyenne_max_dep,2),
65                   })
66
67 return df
68
69 def make_df(filetype):
70     df = var_exp_df(filetype)
71     df.set_index('lang', inplace=True)
72     df2 = pd.read_csv(base_dir + '/out.csv', index_col='lang')
73     df = df.join(df2)
74     df2 = pd.read_csv(base_dir + '/complexityScore.csv', index_col='Code')
75     df2.drop(['Language', 'Extrapolation'], axis = 1, inplace=True)
76     df = df.join(df2)
77     os.chdir(base_dir)
78     df.to_csv('results_'+filetype+'.csv')
79     return df
80
81 def taux_projectivite(filetype,lang):
82     os.chdir(base_dir + '/data/')
83     with open(filetype + "_" + lang + ".conllu", 'r', encoding='utf8') as f:
84         lines = f.readlines()
85         empty_lines = [-1]
86         for i, line in enumerate(lines):
87             if line == "\n":
88                 empty_lines.append(i)
89             nbre_phrases = len(empty_lines)
90     os.chdir(base_dir + '/expe/out/')
91     with open(filetype + "_" + lang + "_proj.conllu", 'r', encoding='utf8') as
92     ↪ f:
93         lines = f.readlines()
94         empty_lines = [-1]
95         for i, line in enumerate(lines):
96             if line == "\n":
97                 empty_lines.append(i)
98             nbre_phrases_proj = len(empty_lines)
99     return nbre_phrases_proj/nbre_phrases

```

```

94
95 def long_dependances(filetype,lang):
96     os.chdir(base_dir + '/expe/out/')
97     long, nbre, longmax = [],[],[]
98     with open(filetype + "_" + lang + "_pgle.mcf", 'r', encoding='utf8') as f:
99         ## construction des indices de phrases ##
100         lines = f.readlines()
101         sentences = [0]
102         for i,line in enumerate(lines):
103             line = line.split("\t")
104             if int(line[4]) == 1:
105                 sentences.append(i+1)
106         liste = []
107         for i in range(len(sentences) - 1):
108             liste2 = []
109             for line in lines[sentences[i]:sentences[i + 1]]:
110                 word = line.split("\t")
111                 if word[3] == 'root':
112                     liste2.append(0)
113                 else :
114                     liste2.append(int(word[2]))
115             liste.append(liste2)
116         liste_abs = rel2abs(liste)
117         for liste3 in liste_abs:
118             arbre = maketree(liste3)
119             nbre.append(len(arbre))
120             longueurs = [len(dep) for dep in arbre]
121             long.append(np.mean(longueurs))
122             longmax.append(np.max(longueurs))
123         return {"nombre_moyen_dep" : np.mean(nbre),
124             ↪ "longueur_moyenne_dep":np.mean(long), "longueur_moyenne_max_dep" :
125             ↪ np.mean(longmax)}
124
125 def rel2abs(liste):
126     rep = []
127     for liste1 in liste:
128         rep.append([i+j if j != 0 else -1 for i,j in enumerate(liste1)])
129     return rep
130
131 def maketree(liste):
132     arbre = [[liste.index(-1)]]
133     profondeur = 1
134     while True:
135         nouvel_arbre = []
136         liste_recherche = [j[-1] for j in arbre if len(j)==profondeur]
137         if liste_recherche == []:
138             return arbre
139         for l in arbre:
140             if l[-1] in liste_recherche and l[-1] in liste:
141                 indices = [i for i, x in enumerate(liste) if x == l[-1]]
142                 for k in indices:
143                     nouvel_arbre.append(l+[k])

```

```

144         else:
145             nouvel_arbre.append(1)
146             arbre = nouvel_arbre
147             profondeur += 1

```

../RegressionMultiple.py

```

1  import numpy as np
2  import pandas as pd
3  from featuresAnalyser import make_df
4  import matplotlib.pyplot as plt
5  import seaborn as sns
6
7  ### Plot ###
8  def my_plot(df,x_value,y):
9      y = y["las"].values.tolist()
10     x = df[x_value].values.tolist()
11
12     plt.figure(figsize=(7, 4))
13     sns.set_style("whitegrid")
14     ax = sns.lineplot(x=x, y=y, color='red')
15     sns.despine(left=True)
16     ax.axhline(0, color='grey')
17     # ax.set(xlabel='Année', ylabel="Taux (%)")
18     # ax.legend(['Arabie Saoudite', 'Moyen-Orient et Afrique du Nord'])
19     plt.show()
20
21     #### Data ####
22
23     #df_train = make_df("train")
24     #df_train.sort_values(by=['uas'], inplace=True)
25     #df_test = make_df("dev")
26     #df_test.sort_values(by=['uas'], inplace=True)
27
28     def mean_norm(df_input):
29         return df_input.apply(lambda x: (x-x.mean())/ x.std(), axis=0)
30
31     df_train =
32     ↪ pd.read_csv('results_train.csv', index_col='lang').sort_values(by=['uas'])
33     df_test =
34     ↪ pd.read_csv('results_dev.csv', index_col='lang').sort_values(by=['uas'])
35
36     X_train = df_train[["taille_vocab" , "tailles_phrase" , "longueur_mots" ,
37     ↪ "taux_non_projectivite" , "nombre_moyen_dep" , "longueur_moyenne_dep" ,
38     ↪ "longueur_moyenne_max_dep", "Score"]]
39     X_train.fillna(X_train.mean(), inplace=True)
40     X_train = mean_norm(X_train)
41     y_train = df_train[["las"]]
42     X_test = df_test[["taille_vocab" , "tailles_phrase" , "longueur_mots" ,
43     ↪ "taux_non_projectivite" , "nombre_moyen_dep" , "longueur_moyenne_dep" ,
44     ↪ "longueur_moyenne_max_dep", "Score"]]
45     X_test.fillna(X_test.mean(), inplace=True)
46     X_test = mean_norm(X_test)

```

```

41 y_test = df_test[["las"]]
42
43 ### Corrélations ###
44
45 corr_df = X_train.corr(method='pearson')
46
47 plt.figure(figsize=(8, 6))
48 sns.heatmap(corr_df, annot=True)
49 plt.tight_layout()
50 plt.show()
51
52 ### Regression ###
53
54 # importing module
55 from sklearn.linear_model import LinearRegression
56 # creating an object of LinearRegression class
57 LR = LinearRegression()
58 # fitting the training data
59 LR.fit(X_train,y_train)
60
61 print(LR.intercept_)
62 print(LR.coef_)
63
64 y_prediction = LR.predict(X_test)
65 y_prediction
66
67 # importing r2_score module
68 from sklearn.metrics import r2_score
69 from sklearn.metrics import mean_squared_error
70 # predicting the accuracy score
71 score=r2_score(y_test,y_prediction)
72
73 print('r2 score is ', score)
74 print('MSE is ',mean_squared_error(y_test,y_prediction))
75 print('RMS error of is ',np.sqrt(mean_squared_error(y_test,y_prediction)))
76
77 for var in ["taille_vocab" , "tailles_phrase" , "longueur_mots" ,
↵ "taux_non_projectivite" , "nombre_moyen_dep" , "longueur_moyenne_dep" ,
↵ "longueur_moyenne_max_dep", "Score"]:
78     my_plot(X_train,x_value=var,y=y_train)

```