



# Projet Méthodologie

Création des features et évaluation de leur efficacité

GAO SHUAI

# Introduction

- But : Créer des « features » et évaluer leur efficacité dans la différenciation des textes de genres différents
- Corpus : Corpus Brown du Natural Language Toolkit
- Outils : Codes en Python fournis en cours et Weka (Waikato Environment for Knowledge Analysis)

# Définition du corpus

- Corpus Brown du NLTK
  - Constitution
    - « reviews », « editorial » → « news »
    - « science\_fiction », « fiction », « mystery » → « literature »
    - « learned » → « sciences »

# Choix des « features »

```
1 {  
2   "verbs_present_tense":  
3   [  
4     "be", "am", "are", "is", "have", "has",  
5     "do", "does", "say", "says", "go", "goes",  
6     "get", "gets", "make", "makes", "know", "knows",  
7     "think", "thinks", "take", "takes",  
8     "see", "sees", "come", "comes", "want", "wants",  
9     "look", "looks", "use", "uses", "find", "finds",  
10    "give", "gives", "tell", "tells", "work", "works",  
11    "call", "calls", "try", "tries", "ask", "asks",  
12    "need", "needs", "feel", "feels", "become", "becomes"  
13  ],  
14  
15  "verbs_past_tense":  
16  [  
17    "was", "were", "had",  
18    "did", "said", "went",  
19    "got", "made", "knew",  
20    "thought", "took",  
21    "saw", "came", "wanted",  
22    "looked", "used", "found",  
23    "gave", "told", "worked",  
24    "called", "tried", "asked",  
25    "needed", "felt", "became"  
26  ]  
27 }
```

# Tests de classification dans Weka

all\_features

	F-Mesure	Rappel	Précision	Instances correctement classifiés
J48	0,686	0,720	0,720	72%
NaiveBayes	0,600	0,620	0,613	62%
RandomForest	0,537	0,560	0,534	56%

verbs\_past\_tense

	F-Mesure	Rappel	Précision	Instances correctement classifiés
J48	0,686	0,720	0,720	72%
NaiveBayes	0,611	0,640	0,618	64%
RandomForest	0,548	0,560	0,545	56%

verbs\_present\_tense

	F-Mesure	Rappel	Précision	Instances correctement classifiés
J48	?	0,460	?	46%
NaiveBayes	0,450	0,480	0,481	48%
RandomForest	0,403	0,400	0,433	40%



# Discussion des résultats

- Résultats
  - « Features » peu efficaces
    - Confusion de formes des verbes au passé et des participes passés
    - Présence sous-estimés des verbes au présent
- Perspectives
  - Prise en considération des formes composées
  - Élimination des verbes au présent dans « features »



Merci