Insight into annotating process of UNIVERSITÉ Clermont brainstorming using machine learning Auvergne **Ecole doctorale Ducros Théo, Bouet Marinette and Toumani Farouk Sciences Pour** l'Ingénieur

LIMOS, UMR 6158, Clermont-Ferrand, France

ProFan

ProFan is a projet which takes part in the action « Innovation numérique pour l'excellence éducative ». The core of this project is an experiment conducted over 20 000 professional high-school students with the help of more than 1000 teachers. Those students are undertaking one of those high-school majors related to technical fields :

- « Accompagnement, Soin et Services à la Personne » (ASSP),
- « Métiers de l'ELectricité et de ses Environnements Connectés » (MELEC),
- « Commerce » (COM)

It aims to evaluate the effectiveness of collaborative learning through Jigsaw method.







Fig. 1 : Jigsaw Method

Along with scholar results, progression is also measured with activity like Text coping, Brainstorming...

Brainstorming data overview

Brainstorming activity produced chat conversation processed to determine the idea of each line. Based on that, three criteria evaluate creativity : flexibility, fluidity and originality.

- Large number of data available (> 50k annotated lines)
- Annotations provided by experts
- Features : Poor context
 - Unusual spelling (mistakes, abbreviations,...)
 - Reality representative

Elève	Contribution	ldée	Catégorie
4	un photophore	Photophore	Décoration
3	je connais pas		
Ą	tu sais c'est un trucs avec une bougie a l'interieur		
3	heummm nan je vois pas		
3	on pourrais faire des cannettes telephone	Téléphone	Electronique
C	les resicler	Recyclage	Autre

Fig. 2 : Annotated chat extract

Objectives

Computerize annotating process

The following sub-goals have been identified :

- Separate lines with idea from those without
- Find a way to give annotations to lines with idea
- Compute criteria according to expert's formula

A first approach using machine learning



Google tool allowing us to transform text data into vectors of floats.





Lines annotated with a category

une fleur métalique Flore

Results

In order to avoid, for now, noise from the classification mistake of the Filtering step, we extracted only the true positive data. They have been used to build a first classifier that aims to differentiate the categories. Among the 23 133 lines available, 4 627 composed the testing set.

The global accuracy reaches 58% which is low. However, the confusion matrix has an interesting feature : outside the first class, the mistake are very few as shown in Fig. 3.

After investigating, it was revealed that some of the first class predicted were in fact undetermined. They are easily identifiable so we can retrieve them which produces the Fig. 4.

Those results are encouraging because ~66% of the initial data are determined with high accuracy. Nonetheless, there are still ~34% to work on and other needs for the criteria that makes it unsatisfying.

912	2	1	6	0	2	1	0	2	0	1
207	261	3	3	0	2	1	0	1	0	1
160	0	149	2	0	8	1	0	1	0	0
228	2	0	151	0	1	0	0	1	0	1
46	0	0	0	39	2	0	0	0	0	1
361	3	2	6	1	659	0	1	3	0	0
134	0	0	0	1	1	246	0	0	0	1
70	1	0	0	0	1	0	122	0	0	0
183	2	0	2	0	0	1	0	180	0	0
75	0	1	1	0	0	1	0	0	160	0
104	1	0	0	1	1	0	0	3	1	98

667	2	1	6	0	2	1	0	2	0	1
12	261	3	3	0	2	1	0	1	0	1
6	0	149	2	0	8	1	0	1	0	0
9	2	0	151	0	1	0	0	1	0	1
1	0	0	0	39	2	0	0	0	0	1
11	3	2	6	1	659	0	1	3	0	0
3	0	0	0	1	1	246	0	0	0	1
1	1	0	0	0	1	0	122	0	0	0
4	2	0	2	0	0	1	0	180	0	0
2	0	1	1	0	0	1	0	0	160	0
5	1	0	0	1	1	0	0	3	1	98

Fig. 4 : Same matrix without Fig. 3 : Confusion matrix « undeterminded »