
Natural Language Processing (ENSAE 3A-MS)- Project

ENSAE 2021/2022

PROFESSEUR: BENJAMIN MULLER

MASTÈRE SPÉCIALISÉ - DATA SCIENCE

OSCAR FOSSEY JEFFREY VERDIÈRE

LAB SMALL GROUP PROFESSOR: GAËL GUIBON



Pôle Emploi Job offer classification

Contents

1	The problematic	2
2	Methodology	2
2.1	Data collection	2
2.2	Model litterature review for text classification	3
2.3	Models training and benchmarking	3
2.4	Model testing	3
2.5	Model picked	4
3	Product delivery	5
	References	6

1 The problematic

[1] Pôle emploi receives around one thousands of job offer per day on their plateforme and new to classify according to very rigid nomenclature.

B			Arts et Façonnage d'ouvrages d'art	
B	11		Arts plastiques	
B	11	01	Création en arts plastiques	
B	11	01	Aquarelliste	11101

Figure 1: Nomenclature of a job offer

As show on Figure 1, each job offer has to be classified following a ROME code with different levels. The first level is the letter, the second level is a job category which corresponds to a two digits number.

The aim of our project: is to find the best predictor for a job offer description which predicts the first level of the ROME Code. To be accurate, we want to predict for a chain of character the label between A to N (14 labels).

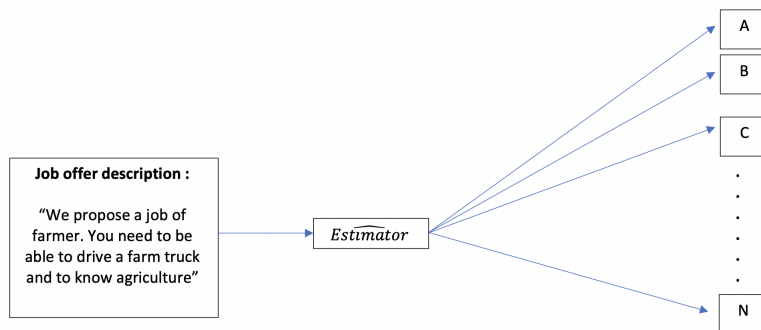


Figure 2: Data pipeline of our project

2 Methodology

2.1 Data collection

We collected the data on pôle emploi website API. It enables to collect 64 402 job description texts with the ROME code associated.

Table 1: Data Collection Organization

The total number of samples is	64 402 samples
The train dataset is	38640 samples
The validation dataset is	12 881 samples
The test dataset is	12 881 samples

Also our data set, has the following repartition for each labels.

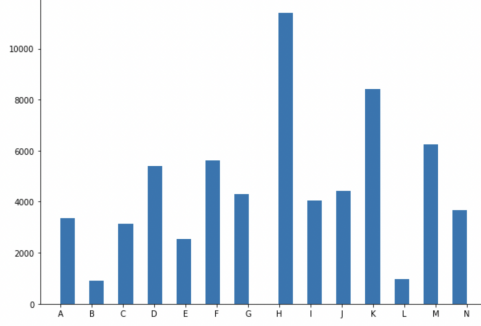


Figure 3: Repartition for the labels of our Data Set

2.2 Model litterature review for text classification

According to [3], the most common models for text classification are TF-IDF, LSTM and Transformers.

2.3 Models training and benchmarking

[2] explains that the best methods to solve the unbalanced data set problem released in our case with the label H are oversampling, undersampling or random weighted sampling.

For sampling techniques, we have tested four sampling techniques to build our data loaders and we displayed the results in the following tab.

Table 2: Model Accuracy on validation test for different sampling during training

Model	Undersampling	Oversampling	Weighted sampling	Normal Data Set
TD-IDF	Accuracy: 64%	Accuracy: 69%	Accuracy: -	Accuracy: 47%
Camembert	Accuracy: -	Accuracy:-	Accuracy:72%	Accuracy: 76%
LSTM	Accuracy: -	Accuracy:-	Accuracy:-	Accuracy: 69%

2.4 Model testing

Table 3: Models Accuracy on validation set and confusion matrix

Model	Accuracy	Confusion matrix																																																																																																																																																																																																																																	
TD-IDF	Accuracy: 69%	<div><p>Confusion Matrix for tfidf model on test set</p><table><tr><th>Actual \ Predicted</th><th>A</th><th>B</th><th>C</th><th>D</th><th>E</th><th>F</th><th>G</th><th>H</th><th>I</th><th>J</th><th>K</th><th>L</th><th>M</th><th>N</th></tr><tr><th>A</th><td>457</td><td>3</td><td>19</td><td>12</td><td>6</td><td>27</td><td>7</td><td>47</td><td>22</td><td>14</td><td>22</td><td>2</td><td>22</td><td>20</td></tr><tr><th>B</th><td>0</td><td>124</td><td>1</td><td>7</td><td>11</td><td>7</td><td>0</td><td>29</td><td>3</td><td>0</td><td>1</td><td>0</td><td>1</td><td>1</td></tr><tr><th>C</th><td>0</td><td>0</td><td>483</td><td>23</td><td>3</td><td>5</td><td>7</td><td>7</td><td>7</td><td>2</td><td>6</td><td>1</td><td>74</td><td>1</td></tr><tr><th>D</th><td>8</td><td>8</td><td>112</td><td>732</td><td>10</td><td>8</td><td>42</td><td>19</td><td>14</td><td>11</td><td>7</td><td>1</td><td>57</td><td>17</td></tr><tr><th>E</th><td>1</td><td>4</td><td>18</td><td>11</td><td>263</td><td>7</td><td>9</td><td>92</td><td>7</td><td>2</td><td>9</td><td>6</td><td>45</td><td>15</td></tr><tr><th>F</th><td>4</td><td>6</td><td>26</td><td>14</td><td>12</td><td>773</td><td>4</td><td>78</td><td>51</td><td>2</td><td>8</td><td>1</td><td>62</td><td>13</td></tr><tr><th>G</th><td>2</td><td>2</td><td>22</td><td>30</td><td>7</td><td>1</td><td>661</td><td>6</td><td>3</td><td>9</td><td>52</td><td>1</td><td>51</td><td>16</td></tr><tr><th>H</th><td>24</td><td>104</td><td>35</td><td>47</td><td>42</td><td>131</td><td>13</td><td>1552</td><td>129</td><td>11</td><td>25</td><td>3</td><td>175</td><td>33</td></tr><tr><th>I</th><td>9</td><td>5</td><td>14</td><td>11</td><td>8</td><td>41</td><td>22</td><td>66</td><td>547</td><td>2</td><td>13</td><td>2</td><td>53</td><td>17</td></tr><tr><th>J</th><td>1</td><td>1</td><td>1</td><td>9</td><td>0</td><td>0</td><td>16</td><td>20</td><td>3</td><td>826</td><td>16</td><td>0</td><td>13</td><td>4</td></tr><tr><th>K</th><td>25</td><td>8</td><td>100</td><td>22</td><td>12</td><td>25</td><td>48</td><td>84</td><td>76</td><td>96</td><td>1006</td><td>7</td><td>170</td><td>39</td></tr><tr><th>L</th><td>2</td><td>2</td><td>3</td><td>7</td><td>17</td><td>3</td><td>6</td><td>10</td><td>12</td><td>2</td><td>4</td><td>111</td><td>17</td><td>2</td></tr><tr><th>M</th><td>1</td><td>1</td><td>156</td><td>46</td><td>35</td><td>13</td><td>19</td><td>25</td><td>17</td><td>28</td><td>45</td><td>6</td><td>651</td><td>33</td></tr><tr><th>N</th><td>4</td><td>0</td><td>11</td><td>15</td><td>3</td><td>23</td><td>14</td><td>46</td><td>10</td><td>3</td><td>14</td><td>1</td><td>75</td><td>509</td></tr></table></div>	Actual \ Predicted	A	B	C	D	E	F	G	H	I	J	K	L	M	N	A	457	3	19	12	6	27	7	47	22	14	22	2	22	20	B	0	124	1	7	11	7	0	29	3	0	1	0	1	1	C	0	0	483	23	3	5	7	7	7	2	6	1	74	1	D	8	8	112	732	10	8	42	19	14	11	7	1	57	17	E	1	4	18	11	263	7	9	92	7	2	9	6	45	15	F	4	6	26	14	12	773	4	78	51	2	8	1	62	13	G	2	2	22	30	7	1	661	6	3	9	52	1	51	16	H	24	104	35	47	42	131	13	1552	129	11	25	3	175	33	I	9	5	14	11	8	41	22	66	547	2	13	2	53	17	J	1	1	1	9	0	0	16	20	3	826	16	0	13	4	K	25	8	100	22	12	25	48	84	76	96	1006	7	170	39	L	2	2	3	7	17	3	6	10	12	2	4	111	17	2	M	1	1	156	46	35	13	19	25	17	28	45	6	651	33	N	4	0	11	15	3	23	14	46	10	3	14	1	75	509
Actual \ Predicted	A	B	C	D	E	F	G	H	I	J	K	L	M	N																																																																																																																																																																																																																					
A	457	3	19	12	6	27	7	47	22	14	22	2	22	20																																																																																																																																																																																																																					
B	0	124	1	7	11	7	0	29	3	0	1	0	1	1																																																																																																																																																																																																																					
C	0	0	483	23	3	5	7	7	7	2	6	1	74	1																																																																																																																																																																																																																					
D	8	8	112	732	10	8	42	19	14	11	7	1	57	17																																																																																																																																																																																																																					
E	1	4	18	11	263	7	9	92	7	2	9	6	45	15																																																																																																																																																																																																																					
F	4	6	26	14	12	773	4	78	51	2	8	1	62	13																																																																																																																																																																																																																					
G	2	2	22	30	7	1	661	6	3	9	52	1	51	16																																																																																																																																																																																																																					
H	24	104	35	47	42	131	13	1552	129	11	25	3	175	33																																																																																																																																																																																																																					
I	9	5	14	11	8	41	22	66	547	2	13	2	53	17																																																																																																																																																																																																																					
J	1	1	1	9	0	0	16	20	3	826	16	0	13	4																																																																																																																																																																																																																					
K	25	8	100	22	12	25	48	84	76	96	1006	7	170	39																																																																																																																																																																																																																					
L	2	2	3	7	17	3	6	10	12	2	4	111	17	2																																																																																																																																																																																																																					
M	1	1	156	46	35	13	19	25	17	28	45	6	651	33																																																																																																																																																																																																																					
N	4	0	11	15	3	23	14	46	10	3	14	1	75	509																																																																																																																																																																																																																					
Camembert	Accuracy: 76%	<div><p>Confusion Matrix for the model trained on unbalanced data</p><table><tr><th>Actual \ Predicted</th><th>A</th><th>B</th><th>C</th><th>D</th><th>E</th><th>F</th><th>G</th><th>H</th><th>I</th><th>J</th><th>K</th><th>L</th><th>M</th><th>N</th></tr><tr><th>A</th><td>506</td><td>0</td><td>3</td><td>13</td><td>0</td><td>19</td><td>4</td><td>51</td><td>25</td><td>7</td><td>20</td><td>4</td><td>14</td><td>14</td></tr><tr><th>B</th><td>1</td><td>87</td><td>0</td><td>4</td><td>8</td><td>18</td><td>0</td><td>61</td><td>2</td><td>0</td><td>3</td><td>1</td><td>0</td><td>0</td></tr><tr><th>C</th><td>0</td><td>0</td><td>451</td><td>24</td><td>5</td><td>9</td><td>2</td><td>18</td><td>5</td><td>0</td><td>20</td><td>0</td><td>75</td><td>10</td></tr><tr><th>D</th><td>17</td><td>0</td><td>27</td><td>653</td><td>8</td><td>3</td><td>22</td><td>16</td><td>13</td><td>7</td><td>22</td><td>0</td><td>62</td><td>16</td></tr><tr><th>E</th><td>2</td><td>1</td><td>7</td><td>11</td><td>260</td><td>5</td><td>10</td><td>104</td><td>6</td><td>4</td><td>16</td><td>3</td><td>45</td><td>15</td></tr><tr><th>F</th><td>8</td><td>0</td><td>13</td><td>8</td><td>2</td><td>614</td><td>0</td><td>128</td><td>24</td><td>3</td><td>21</td><td>0</td><td>20</td><td>13</td></tr><tr><th>G</th><td>13</td><td>0</td><td>8</td><td>27</td><td>2</td><td>2</td><td>667</td><td>10</td><td>11</td><td>5</td><td>56</td><td>2</td><td>36</td><td>24</td></tr><tr><th>H</th><td>24</td><td>13</td><td>10</td><td>27</td><td>19</td><td>66</td><td>8</td><td>1903</td><td>83</td><td>8</td><td>45</td><td>0</td><td>68</td><td>50</td></tr><tr><th>I</th><td>3</td><td>0</td><td>1</td><td>7</td><td>4</td><td>31</td><td>15</td><td>104</td><td>568</td><td>3</td><td>25</td><td>2</td><td>28</td><td>19</td></tr><tr><th>J</th><td>7</td><td>2</td><td>0</td><td>7</td><td>0</td><td>2</td><td>4</td><td>22</td><td>2</td><td>794</td><td>55</td><td>0</td><td>4</td><td>11</td></tr><tr><th>K</th><td>23</td><td>2</td><td>31</td><td>27</td><td>11</td><td>11</td><td>19</td><td>78</td><td>41</td><td>26</td><td>1327</td><td>1</td><td>92</td><td>29</td></tr><tr><th>L</th><td>0</td><td>2</td><td>0</td><td>5</td><td>14</td><td>10</td><td>7</td><td>20</td><td>13</td><td>2</td><td>12</td><td>99</td><td>12</td><td>2</td></tr><tr><th>M</th><td>7</td><td>1</td><td>58</td><td>55</td><td>20</td><td>20</td><td>13</td><td>66</td><td>31</td><td>5</td><td>94</td><td>2</td><td>651</td><td>33</td></tr><tr><th>N</th><td>4</td><td>0</td><td>7</td><td>15</td><td>0</td><td>14</td><td>4</td><td>50</td><td>19</td><td>0</td><td>17</td><td>0</td><td>32</td><td>566</td></tr></table></div>	Actual \ Predicted	A	B	C	D	E	F	G	H	I	J	K	L	M	N	A	506	0	3	13	0	19	4	51	25	7	20	4	14	14	B	1	87	0	4	8	18	0	61	2	0	3	1	0	0	C	0	0	451	24	5	9	2	18	5	0	20	0	75	10	D	17	0	27	653	8	3	22	16	13	7	22	0	62	16	E	2	1	7	11	260	5	10	104	6	4	16	3	45	15	F	8	0	13	8	2	614	0	128	24	3	21	0	20	13	G	13	0	8	27	2	2	667	10	11	5	56	2	36	24	H	24	13	10	27	19	66	8	1903	83	8	45	0	68	50	I	3	0	1	7	4	31	15	104	568	3	25	2	28	19	J	7	2	0	7	0	2	4	22	2	794	55	0	4	11	K	23	2	31	27	11	11	19	78	41	26	1327	1	92	29	L	0	2	0	5	14	10	7	20	13	2	12	99	12	2	M	7	1	58	55	20	20	13	66	31	5	94	2	651	33	N	4	0	7	15	0	14	4	50	19	0	17	0	32	566
Actual \ Predicted	A	B	C	D	E	F	G	H	I	J	K	L	M	N																																																																																																																																																																																																																					
A	506	0	3	13	0	19	4	51	25	7	20	4	14	14																																																																																																																																																																																																																					
B	1	87	0	4	8	18	0	61	2	0	3	1	0	0																																																																																																																																																																																																																					
C	0	0	451	24	5	9	2	18	5	0	20	0	75	10																																																																																																																																																																																																																					
D	17	0	27	653	8	3	22	16	13	7	22	0	62	16																																																																																																																																																																																																																					
E	2	1	7	11	260	5	10	104	6	4	16	3	45	15																																																																																																																																																																																																																					
F	8	0	13	8	2	614	0	128	24	3	21	0	20	13																																																																																																																																																																																																																					
G	13	0	8	27	2	2	667	10	11	5	56	2	36	24																																																																																																																																																																																																																					
H	24	13	10	27	19	66	8	1903	83	8	45	0	68	50																																																																																																																																																																																																																					
I	3	0	1	7	4	31	15	104	568	3	25	2	28	19																																																																																																																																																																																																																					
J	7	2	0	7	0	2	4	22	2	794	55	0	4	11																																																																																																																																																																																																																					
K	23	2	31	27	11	11	19	78	41	26	1327	1	92	29																																																																																																																																																																																																																					
L	0	2	0	5	14	10	7	20	13	2	12	99	12	2																																																																																																																																																																																																																					
M	7	1	58	55	20	20	13	66	31	5	94	2	651	33																																																																																																																																																																																																																					
N	4	0	7	15	0	14	4	50	19	0	17	0	32	566																																																																																																																																																																																																																					
LSTM	Accuracy: 70%	<div><p>Confusion Matrix for the model trained on unbalanced data</p><table><tr><th>Actual \ Predicted</th><th>A</th><th>B</th><th>C</th><th>D</th><th>E</th><th>F</th><th>G</th><th>H</th><th>I</th><th>J</th><th>K</th><th>L</th><th>M</th><th>N</th></tr><tr><th>A</th><td>506</td><td>0</td><td>3</td><td>13</td><td>0</td><td>19</td><td>4</td><td>51</td><td>25</td><td>7</td><td>20</td><td>4</td><td>14</td><td>14</td></tr><tr><th>B</th><td>1</td><td>87</td><td>0</td><td>4</td><td>8</td><td>18</td><td>0</td><td>61</td><td>2</td><td>0</td><td>3</td><td>1</td><td>0</td><td>0</td></tr><tr><th>C</th><td>0</td><td>0</td><td>451</td><td>24</td><td>5</td><td>9</td><td>2</td><td>18</td><td>5</td><td>0</td><td>20</td><td>0</td><td>75</td><td>10</td></tr><tr><th>D</th><td>17</td><td>0</td><td>27</td><td>653</td><td>8</td><td>3</td><td>22</td><td>16</td><td>13</td><td>7</td><td>22</td><td>0</td><td>62</td><td>16</td></tr><tr><th>E</th><td>2</td><td>1</td><td>7</td><td>11</td><td>260</td><td>5</td><td>10</td><td>104</td><td>6</td><td>4</td><td>16</td><td>3</td><td>45</td><td>15</td></tr><tr><th>F</th><td>8</td><td>0</td><td>13</td><td>8</td><td>2</td><td>614</td><td>0</td><td>128</td><td>24</td><td>3</td><td>21</td><td>0</td><td>20</td><td>13</td></tr><tr><th>G</th><td>13</td><td>0</td><td>8</td><td>27</td><td>2</td><td>2</td><td>667</td><td>10</td><td>11</td><td>5</td><td>56</td><td>2</td><td>36</td><td>24</td></tr><tr><th>H</th><td>24</td><td>13</td><td>10</td><td>27</td><td>19</td><td>66</td><td>8</td><td>1903</td><td>83</td><td>8</td><td>45</td><td>0</td><td>68</td><td>50</td></tr><tr><th>I</th><td>3</td><td>0</td><td>1</td><td>7</td><td>4</td><td>31</td><td>15</td><td>104</td><td>568</td><td>3</td><td>25</td><td>2</td><td>28</td><td>19</td></tr><tr><th>J</th><td>7</td><td>2</td><td>0</td><td>7</td><td>0</td><td>2</td><td>4</td><td>22</td><td>2</td><td>794</td><td>55</td><td>0</td><td>4</td><td>11</td></tr><tr><th>K</th><td>23</td><td>2</td><td>31</td><td>27</td><td>11</td><td>11</td><td>19</td><td>78</td><td>41</td><td>26</td><td>1327</td><td>1</td><td>92</td><td>29</td></tr><tr><th>L</th><td>0</td><td>2</td><td>0</td><td>5</td><td>14</td><td>10</td><td>7</td><td>20</td><td>13</td><td>2</td><td>12</td><td>99</td><td>12</td><td>2</td></tr><tr><th>M</th><td>7</td><td>1</td><td>58</td><td>55</td><td>20</td><td>20</td><td>13</td><td>66</td><td>31</td><td>5</td><td>94</td><td>2</td><td>651</td><td>33</td></tr><tr><th>N</th><td>4</td><td>0</td><td>7</td><td>15</td><td>0</td><td>14</td><td>4</td><td>50</td><td>19</td><td>0</td><td>17</td><td>0</td><td>32</td><td>566</td></tr></table></div>	Actual \ Predicted	A	B	C	D	E	F	G	H	I	J	K	L	M	N	A	506	0	3	13	0	19	4	51	25	7	20	4	14	14	B	1	87	0	4	8	18	0	61	2	0	3	1	0	0	C	0	0	451	24	5	9	2	18	5	0	20	0	75	10	D	17	0	27	653	8	3	22	16	13	7	22	0	62	16	E	2	1	7	11	260	5	10	104	6	4	16	3	45	15	F	8	0	13	8	2	614	0	128	24	3	21	0	20	13	G	13	0	8	27	2	2	667	10	11	5	56	2	36	24	H	24	13	10	27	19	66	8	1903	83	8	45	0	68	50	I	3	0	1	7	4	31	15	104	568	3	25	2	28	19	J	7	2	0	7	0	2	4	22	2	794	55	0	4	11	K	23	2	31	27	11	11	19	78	41	26	1327	1	92	29	L	0	2	0	5	14	10	7	20	13	2	12	99	12	2	M	7	1	58	55	20	20	13	66	31	5	94	2	651	33	N	4	0	7	15	0	14	4	50	19	0	17	0	32	566
Actual \ Predicted	A	B	C	D	E	F	G	H	I	J	K	L	M	N																																																																																																																																																																																																																					
A	506	0	3	13	0	19	4	51	25	7	20	4	14	14																																																																																																																																																																																																																					
B	1	87	0	4	8	18	0	61	2	0	3	1	0	0																																																																																																																																																																																																																					
C	0	0	451	24	5	9	2	18	5	0	20	0	75	10																																																																																																																																																																																																																					
D	17	0	27	653	8	3	22	16	13	7	22	0	62	16																																																																																																																																																																																																																					
E	2	1	7	11	260	5	10	104	6	4	16	3	45	15																																																																																																																																																																																																																					
F	8	0	13	8	2	614	0	128	24	3	21	0	20	13																																																																																																																																																																																																																					
G	13	0	8	27	2	2	667	10	11	5	56	2	36	24																																																																																																																																																																																																																					
H	24	13	10	27	19	66	8	1903	83	8	45	0	68	50																																																																																																																																																																																																																					
I	3	0	1	7	4	31	15	104	568	3	25	2	28	19																																																																																																																																																																																																																					
J	7	2	0	7	0	2	4	22	2	794	55	0	4	11																																																																																																																																																																																																																					
K	23	2	31	27	11	11	19	78	41	26	1327	1	92	29																																																																																																																																																																																																																					
L	0	2	0	5	14	10	7	20	13	2	12	99	12	2																																																																																																																																																																																																																					
M	7	1	58	55	20	20	13	66	31	5	94	2	651	33																																																																																																																																																																																																																					
N	4	0	7	15	0	14	4	50	19	0	17	0	32	566																																																																																																																																																																																																																					

2.5 Model picked

Despite the large amount of parameters for the camembert models, it is the one which performs the best and can draw predictions for on text in less than a second which is fast enough. So we keep this model for our product delivery.

3 Product delivery

The product delivery is github link: <https://github.com/oscarfossey/NLP-Job-classifier-based-on-description>.

The data and the models were registered on huggingface cloud platform for convinient reasons. Link to data: <https://huggingface.co/datasets/oscarfossey/NLPpoleemploi/tree/main>

Link to models: https://huggingface.co/oscarfossey/job_classification

On this repository, you can find all the used notebooks. Particularly, we divided our works in three main differents files: scrapping, trainings of the models and pipelines.

The "main" notebook should be run using a google collab GPU.

References

- [1] Romarik Le Dourneuf. “Le nombre d’offres explose sur Pôle Emploi, voici les secteurs qui recrutent le plus”. In: (1982).
- [2] Gao. “Data Augmentation in Solving Data Imbalance Problems”. In: (2020).
- [3] Cambria Nanyang Nikzad Chenaghlu Gao Minaee Kalchbrenner. “Deep Learning Based Text Classification: A Comprehensive Review”. In: (2020).

List of Tables

1	Data Collection Organization	2
2	Model Accuracy on validation test for different sampling during training	3
3	Models Accuracy on validation set and confusion matrix	4

List of Figures

1	Nomenclature of a job offer	2
2	Data pipeline of our project	2
3	Repartition for the labels of our Data Set	3