

HA#9 Report

- a.) While adding a classifier at the end of the pipeline after the sklearn has been trained the text. The cross validation gives us enough flexibility and to determine how many tress and how deep should we go. After I changed from 5 to 3 cross validation I started to realize the hyper parameters has been changed little bit. And also the prediction changed a little bit, but it increase the mean prediction when I turned on to the 3 cross validation.

```
Out[2]:
id text author ... num_noun_word num_verb_word
id26305 This process, however, afforded me no means of... EAP ... 28 8
id17569 It never once occurred to me that the fumbling... HPL ... 28 8
id11008 In his left hand was a gold snuff box, from wh... EAP ... 28 8
id27763 How lovely is spring As we looked from Windsor... MWS ... 28 8
id12958 Finding nothing else, not even gold, the Super... HPL ... 28 8

[5 rows x 11 columns]
Out[3]:
id processed length ... num_noun_word num_verb_word
id19417 this panorama is indeed glorious and i should ... 91 ... 28 8
id09522 there was a simple natural earnestness about h... 240 ... 28 8
id22732 who are you pray that i duc de lomelette princ... 387 ... 28 8
id10351 he had gone in the carriage to the nearest tow... 118 ... 28 8
id24580 there is no method in their proceedings beyond... 71 ... 28 8
```

- b.) In the previous HA#8, I added number of adjectives, nouns, and verbs.

Class important and least important features I found.

processed	length	words	words_not_stopword	avg_word_length
this process however afforded me no means of a...	224	41	21	6.380952
it never once occurred to me that the fumbling...	70	14	6	6.166667
in his left hand was a gold snuff box from whi...	195	36	19	5.947368
how lovely is				

Important features

commas	toritus
4	9A3
0	JQH
4	9A3

least important features

c.)

```
(1, 3582) 0.2289382705944011
(1, 5981) 0.29710947912130625
(1, 6628) 0.22357939073099628
(1, 7172) 0.23160082552507501
(1, 9317) 0.25681851547627504
(1, 10116) 0.2515703101082241
(1, 11058) 0.17834637520787042
(1, 12536) 0.21797103780515726
(1, 17252) 0.23449643109972057
(1, 19134) 0.22413626047542487
(1, 19189) 0.21036459238526695
:
(13115, 2214) 0.45585879012226876
(13115, 3752) 0.3901413850145597
(13115, 3829) 0.38313659891823865
(13115, 5089) 0.4771286675345845
(13115, 16794) 0.35133994372068006
(13115, 21516) -0.1444701537205625
(13115, 21517) -0.0915346596641655
(13115, 21518) -0.31912971679909746
(13115, 21519) 0.2076359596181049
(13115, 21520) -0.4510466952023476
(13116, 88) 0.3532883437037042
(13116, 1605) 0.3375391333269599
(13116, 3284) 0.24810135849814624
(13116, 5874) 0.33155156143046516
(13116, 8033) 0.29708663607294533
(13116, 13307) 0.3263648907347903
(13116, 14679) 0.3375391333269599
(13116, 15138) 0.28534534206566287
(13116, 16249) 0.3375391333269599
(13116, 16668) 0.29276637301797825
(13116, 21516) -0.3959318053066319
(13116, 21517) -0.4004391142282561
```

After training the data, some of the text have wrong prediction which if they have negative probabilities.