

Journal Report 4

9/23/19-9/29/19

Abby Kim

Computer Systems Research Lab

Period 2, White

Daily Log

Monday September 9

Looked at how the tokens file tags people and how the book.html file lists different characters. trying to see how to standardize the names of different characters and sort them by gender.

Tuesday September 10

Did not feel well in the morning :(Felt better later though and ended up tracing the people tags through with the token numbers and references, so I kind of understand the way words are connected to people specifically.

Thursday September 12

Continuation of Tuesday. Traced Oliver Twist and Rose a couple times. Started coding the algorithm I (kind of) figured out.

Timeline

Date	Goal	Met
Two Weeks Ago	"Implement the algorithm I drew up. I'm kind of bad at coding so there will probably be a lot of bugs so this might take a lot of time."	done (toit)
Last week	Figure out the next part of the algorithm regarding filtering out characters by gender	yesh shir
This week	Implement that algorithm	N/A
Today plus 1 week	Graph the data points generated from the previous week and see what it looks like. (This seems simple but I'm also really bad at math so this will probably take a week)	N/A
Today plus 2 weeks	Figure out how to run other books with the pipeline	N/A

Reflection

Pretty cute week not going to lie. So the next part I want to code up this week has to do with extracting words that are associated with different people. So right now I have the grouped words with other words associated them by using their token id. Right now it's a giant dictionary with a root head as the key and a set of all words connected to it as the value. Example:

```
9, 196530, 196531, 196532, 196533, 196534, 196535, 196536, 196537]], 65436: set([65437, 65438]), 174747: set([174728, 174729, 174730, 174731, 174732, 174733, 174734, 174735, 174736, 174737, 174738, 174739, 174740, 174741, 174742, 174743, 174744, 174745, 174746, 174748, 174749, 174750, 174751, 174752, 174753]), 65446: set([65440, 65441, 65442, 65443, 65444, 65445, 65446, 65447, 65448, 65449, 65450, 65451, 65452, 65453, 65454, 65455, 65456, 65457, 65458, 65459, 65460, 65461, 65462, 65463, 65464, 65465, 65466, 65467, 65468, 65469, 65470, 65471, 65472, 65473, 65474, 65475, 65476, 65477, 65478, 65479, 65480, 65481, 65482, 65483, 65484, 65485, 65486, 65487, 65488, 65489, 65490, 65491, 65492, 65493, 65494, 65495, 65496, 65497, 65498, 65499, 65500, 65501, 65502, 65503, 65504, 65505, 65506, 65507, 65508, 65509, 65510, 65511, 65512, 65513, 65514, 65515, 65516, 65517, 65518, 65519, 65520, 65521, 65522, 65523, 65524, 65525, 65526, 65527, 65528, 65529, 65530, 65531, 65532, 65533, 65534, 65535, 65536, 65537, 65538, 65539, 65540, 65541, 65542, 65543, 65544, 65545, 65546, 65547, 65548, 65549, 65550, 65551, 65552, 65553, 65554, 65555, 65556, 65557, 65558, 65559, 65560, 65561, 65562, 65563, 65564, 65565, 65566, 65567, 65568, 65569, 65570, 65571, 65572, 65573, 65574, 65575, 65576, 65577, 65578, 65579, 65580, 65581, 65582, 65583, 65584, 65585, 65586, 65587, 65588, 65589, 65590, 65591, 65592, 65593, 65594, 65595, 65596, 65597, 65598, 65599, 65600, 65601, 65602, 65603, 65604, 65605, 65606, 65607, 65608, 65609, 65610, 65611, 65612, 65613, 65614, 65615, 65616, 65617, 65618, 65619, 65620, 65621, 65622, 65623, 65624, 65625, 65626, 65627, 65628, 65629, 65630, 65631, 65632, 65633, 65634, 65635, 65636, 65637, 65638, 65639, 65640, 65641, 65642, 65643, 65644, 65645, 65646, 65647, 65648, 65649, 65650, 65651, 65652, 65653, 65654, 65655, 65656, 65657, 65658, 65659, 65660, 65661, 65662, 65663, 65664, 65665, 65666, 65667, 65668, 65669, 65670, 65671, 65672, 65673, 65674, 65675, 65676, 65677, 65678, 65679, 65680, 65681, 65682, 65683, 65684, 65685, 65686, 65687, 65688, 65689, 65690, 65691, 65692, 65693, 65694, 65695, 65696, 65697, 65698, 65699, 65700, 65701, 65702, 65703, 65704, 65705, 65706, 65707, 65708, 65709, 65710, 65711, 65712, 65713, 65714, 65715, 65716, 65717, 65718, 65719, 65720, 65721, 65722, 65723, 65724, 65725, 65726, 65727, 65728, 65729, 65730, 65731, 65732, 65733, 65734, 65735, 65736, 65737, 65738, 65739, 65740, 65741, 65742, 65743, 65744, 65745, 65746, 65747, 65748, 65749, 65750, 65751, 65752, 65753, 65754, 65755, 65756, 65757, 65758, 65759, 65760, 65761, 65762, 65763, 65764, 65765, 65766, 65767, 65768, 65769, 65770, 65771, 65772, 65773, 65774, 65775, 65776, 65777, 65778, 65779, 65780, 65781, 65782, 65783, 65784, 65785, 65786, 65787, 65788, 65789, 65790, 65791, 65792, 65793, 65794, 65795, 65796, 65797, 65798, 65799, 65800, 65801, 65802, 65803, 65804, 65805, 65806, 65807, 65808, 65809, 65810, 65811, 65812, 65813, 65814, 65815, 65816, 65817, 65818, 65819, 65820, 65821, 65822, 65823, 65824, 65825, 65826, 65827, 65828, 65829, 65830, 65831, 65832, 65833, 65834, 65835, 65836, 65837, 65838, 65839, 65840, 65841, 65842, 65843, 65844, 65845, 65846, 65847, 65848, 65849, 65850, 65851, 65852, 65853, 65854, 65855, 65856, 65857, 65858, 65859, 65860, 65861, 65862, 65863, 65864, 65865, 65866, 65867, 65868, 65869, 65870, 65871, 65872, 65873, 65874, 65875, 65876, 65877, 65878, 65879, 65880, 65881, 65882, 65883, 65884, 65885, 65886, 65887, 65888, 65889, 65890, 65891, 65892, 65893, 65894, 65895, 65896, 65897, 65898, 65899, 65900, 65901, 65902, 65903, 65904, 65905, 65906, 65907, 65908, 65909, 65910, 65911, 65912, 65913, 65914, 65915, 65916, 65917, 65918, 65919, 65920, 65921, 65922, 65923, 65924, 65925, 65926, 65927, 65928, 65929, 65930, 65931, 65932, 65933, 65934, 65935, 65936, 65937, 65938, 65939, 65940, 65941, 65942, 65943, 65944, 65945, 65946, 65947, 65948, 65949, 65950, 65951, 65952, 65953, 65954, 65955, 65956, 65957, 65958, 65959, 65960, 65961, 65962, 65963, 65964, 65965, 65966, 65967, 65968, 65969, 65970, 65971, 65972, 65973, 65974, 65975, 65976, 65977, 65978, 65979, 65980, 6
```

I need to start by importing the character tags from the html file that looks like this:

858 Oliver (778) Oliver Twist (48) OLIVER (32)
 357 Mr. Bumble (346) Bumble (11)
 341 Sikes (256) Mr. Sikes (79) Mr. William Sikes (4) MR. WILLIAM SIKES (2)
 301 Fagin (271) Mr. Fagin (26) MR. FAGIN (2) FAGIN (2)
 173 Brown (168) Mr. BROWNLOW (4) Brownlow (1)
 162 Noah (113) Mr. Claypole (32) Noah Claypole (15) Mr. Noah Claypole (2)
 143 Rose (130) Miss Rose (5)
 123 Charley (47) Bates (40) Charley Bates (36)
 119 Mr. Giles (84) Giles (35)
 116 Nancy (104) Miss Nancy (8) NANCY (2) MISS NANCY (2)
 95 Bill (87) Bill Sikes (7) Sikes (11)
 73 Toby (39) Crackit (26) Mr. Crackit (8) Crackit (6) Mr. Toby Crackit (2)
 65 Mr. Grimwig (60) Grimwig (3) MR. GRIMWIG (2)
 63 Mrs. Corney (62) Corney (1)
 57 Mrs. Mann (57)
 56 Harry (39) Harry Maylie (15) HARRY MAYLIE (2)
 55 Mrs. Maylie (55)
 47 Mr. Sowerberry (33) Sowerberry (14)
 46 Mr. Loscombe (45) Losberny (1)
 36 Mrs. Sowerberry (36)
 33 Barney (33)
 31 Mr. Fang (24) Fang (9)
 32 Mr. Bolter (26) Morris Bolter (3) Mr. Morris Bolter (2) Bolter (1)
 29 Mr. Chitting (26) Chitting (3)
 28 Mrs. Bedwin (24) Bedwin (4)
 25 Mrs. Gamfield (19) Gamfield (6)
 23 Mrs. Bumble (23)
 19 Duff (19)
 18 Tom (16) Tom White (2)

I'm going to try two ways of filtering the people out. The first is just looking at the root head and determining if it's lemma is in the list of names imported. If it is keep it and add it to a new bucket of words for that character, if not ignore it.

The second way involves looking at every token, not just the root heads. The html tags whether a character is the agent of patient, so I would count how many times they are each of those things then bucket words accordingly.