

D3 - Authorship Attribution Rankings for AllTokens and Singleton

Lana Abdelmohsen, Robert Helck, Casey Lishko, and Alex Quezada

Below are the authorship attribution rankings for the “33913.txt” and “70535.txt” test files, respectively, according to AllTokens. They are based on the test and train directories that were generated for a particular run of the code; executing again will create new randomized lines of an author’s review that are separated into training and test files, but the resultant ranking differs only slightly from what is seen here.

The ranking is a list of tuples where the first value identifies the author by the name of the training file they wrote (without the .txt file extension), and the second is the system score (which is the geometric mean of unigram probabilities). The author of 33913 in the training set is the third most likely (more likely authors were calculated with a higher geometric mean) to have produced the 33913 test file. However, the author of the train set’s “70535.txt” ranks poorly in likelihood to have written the test set’s “70535.txt”:

AllTokens 33913:

```
[('562732', 0.0005339909155648803), ('3280905', 0.0005148154865674477), ('33913', 0.0005121692854479732), ('819382', 0.0005079819742607628), ('317399', 0.0004838080491027413), ('1173088', 0.00048076446022851725), ('2483625', 0.0004713200100564184), ('342623', 0.0004664884272208933), ('1609079', 0.0004660040130102563), ('453068', 0.00046321731956751306), ('391152', 0.0004618872489906928), ('1406078', 0.00045218215067499), ('9938570', 0.0004516563204245384), ('1355507', 0.0004485976640348497), ('2707735', 0.00044059018221805963), ('2248099', 0.00043576710368245015), ('2721657', 0.0004348010626299325), ('4532636', 0.000433026297118583), ('2020269', 0.00043097257631559855), ('1293485', 0.00043086204995644773), ('3109237', 0.00042712821152904903), ('7743887', 0.0004148527641808078), ('1391795', 0.00041164980810160905), ('2093818', 0.00041104649452126463), ('1416505', 0.0004102124084632225), ('1219578', 0.00040656721183498427), ('4445210', 0.00040572141275807593), ('3223254', 0.0004014704415779093), ('2467618', 0.0004009466144330946), ('102816', 0.0004001410763272033), ('2626332', 0.00039863194221984895), ('2911571', 0.00039480161601081865), ('1048771', 0.00038818621148026575), ('3717154', 0.00038809256174474847), ('1532177', 0.0003863932949904254), ('463200', 0.0003842836672947981), ('3079504', 0.0003814987548510017), ('2171244', 0.0003779153901491259), ('2567136', 0.00037672229793479767), ('2542703', 0.0003753295014942663), ('11228318', 0.0003734083955659491), ('4248714', 0.00036965281184789744), ('1111192', 0.00036600396061559537), ('1399158', 0.0003605127353683321), ('15896852', 0.0003437215718544325), ('4888011', 0.0003374115039352052), ('989035', 0.00032349693067250824), ('1162550', 0.000311141133504273), ('2488512', 0.00030928972179571643), ('1617546', 0.0003089830016127558), ('11423174', 0.0003047341646618741), ('663392', 0.00030207178585236056), ('306861', 0.0002952213801007826), ('449021', 0.0002929348160398321), ('8239592',
```

0.000287586423608339), ('583640', 0.0002835526076318813), ('453228',
0.00027894448887491596), ('865972', 0.00027421561185639537), ('1132073',
0.00026158287509734805), ('386241', 0.0002612839729980792), ('783721',
0.0002398459009342192), ('70535', 0.0002340666270100422)]

AllTokens 70535:

[('562732', 0.000518552280006532), ('819382', 0.0005127365111256938), ('317399',
0.00047962682650506456), ('1173088', 0.000477954743536225), ('453068',
0.000473912028225093), ('2707735', 0.0004622497739366479), ('3280905',
0.00046137785533725205), ('1355507', 0.00044931130078976027), ('9938570',
0.00044730889572186045), ('1609079', 0.0004464776991711101), ('2721657',
0.00043878909342899574), ('391152', 0.0004369544704420376), ('2248099',
0.00043508331365077485), ('2093818', 0.00042990110279326517), ('2020269',
0.00042917927403468336), ('3109237', 0.0004273220478737522), ('342623',
0.0004260743255749048), ('2483625', 0.0004236143460826501), ('4445210',
0.0004114701232930964), ('2911571', 0.000411764460902987), ('102816',
0.0004109111465599181), ('4532636', 0.0004106452899925199), ('1293485',
0.000407611839228008), ('3079504', 0.00040464900977160164), ('1391795',
0.0004039942067362228), ('3717154', 0.00040206832253084226), ('2626332',
0.0003984512836293251), ('1416505', 0.000398148757831135), ('2467618',
0.00039647655420892175), ('1219578', 0.0003949296492705918), ('7743887',
0.0003929266858533766), ('1406078', 0.00039062358659237997), ('1048771',
0.0003895441459811563), ('1111192', 0.00038691008206360826), ('1532177',
0.00038242545480822487), ('463200', 0.00037224892934137414), ('1399158',
0.00036415087267082144), ('3223254', 0.00036189849782383693), ('2567136',
0.000355958425833527), ('2171244', 0.00035143053951961005), ('11228318',
0.0003480728945379931), ('2542703', 0.0003426698268947622), ('33913',
0.000330245846343969), ('1162550', 0.0003238340219124098), ('989035',
0.0003219519666503292), ('4888011', 0.00031991273176206327), ('70535',
0.00031905006902833256), ('15896852', 0.0003176956728670301), ('4248714',
0.00031281531381471204), ('1617546', 0.0003126859039565837), ('11423174',
0.0003124483204252193), ('8239592', 0.00030769027717267625), ('306861',
0.00029662139520791924), ('449021', 0.0002934826970402373), ('663392',
0.0002910764621424473), ('453228', 0.0002892033362887093), ('583640',
0.00028593335329482), ('2488512', 0.000285920423798044), ('865972',
0.00027453621171517255), ('386241', 0.00027368485173031193), ('783721',
0.0002624817634399248), ('1132073', 0.0002334969608058662)]

This is a marked difference from the author attribution rankings of our Singleton system. When we consider only the tokens that appear once in the test file in our geometric mean calculations, Singleton correctly attributes the author of the training file “33913.txt” to that of the test file “33913.txt” (highest rank). The Singleton system also has much more belief in 70535 as the author of test file “70535.txt,” placing it second:

Singleton 33913:

[('33913', 0.0011152688395228841), ('449021', 0.0010912122942468528), ('15896852', 0.0010762234374651857), ('3223254', 0.0010659360669232159), ('989035', 0.0010648171875485498), ('8239592', 0.0010467701396546392), ('4532636', 0.001043870185605839), ('2483625', 0.0010411091840119412), ('2542703', 0.0010398239532272376), ('9938570', 0.001031863546269829), ('7743887', 0.001029982481164329), ('342623', 0.0010299143377291479), ('819382', 0.0010239923954030985), ('1391795', 0.0010172929913778884), ('4248714', 0.0010143104615314656), ('2567136', 0.0010131274572442637), ('3280905', 0.0010128726378285668), ('11228318', 0.0010087710971385564), ('2707735', 0.001005088819905897), ('2488512', 0.001003326107148229), ('317399', 0.0010032732456402127), ('2171244', 0.001002049658462618), ('1173088', 0.0009995380393001774), ('1048771', 0.0009951550883929021), ('2020269', 0.0009939572317921811), ('1406078', 0.0009922950850921163), ('4888011', 0.0009909296161622707), ('391152', 0.000986149481098138), ('102816', 0.000985044130843528), ('1162550', 0.0009840605750399774), ('1532177', 0.0009802795278763279), ('1609079', 0.0009768549555198634), ('4445210', 0.0009753485080889982), ('663392', 0.0009741945209856076), ('2721657', 0.0009740228258044097), ('3717154', 0.0009716008319212375), ('1399158', 0.0009704821201491339), ('463200', 0.0009701281164391568), ('1617546', 0.0009664931810862482), ('70535', 0.0009653989903800148), ('865972', 0.0009640107446472424), ('783721', 0.0009617935590750047), ('3109237', 0.000957256666309364), ('386241', 0.0009536975977709882), ('1219578', 0.0009468965637461766), ('562732', 0.0009417446055989072), ('2626332', 0.0009415109244506021), ('2467618', 0.0009374261000384796), ('3079504', 0.0009287261222149483), ('583640', 0.0009283421726709697), ('1355507', 0.0009236886047942063), ('453228', 0.0009234394516905311), ('453068', 0.00092086751363675), ('2911571', 0.0009187368193275926), ('306861', 0.0009121850923192054), ('1293485', 0.0009059591336627905), ('11423174', 0.0008915649885934761), ('2248099', 0.0008786511535687451), ('1111192', 0.0008524384057204831), ('2093818', 0.0008296925494560909), ('1132073', 0.0008295852162971716), ('1416505', 0.0008239925412901767)]

Singleton 70535:

[('9938570', 0.0005522214465047972), ('70535', 0.0005501557189218748), ('8239592', 0.0005491879033396891), ('819382', 0.0005483263645717414), ('2707735', 0.0005480141336345917), ('449021', 0.0005425840919459293), ('317399', 0.0005305111214578003), ('2020269', 0.0005293765012683157), ('2483625', 0.0005262517716984451), ('4532636', 0.0005256704761485837), ('989035', 0.000522156812974952), ('1391795', 0.0005220173502045796), ('1048771', 0.0005196384188213195), ('102816', 0.0005154526732634948), ('1173088', 0.0005134477083401226), ('15896852', 0.000511412938332664), ('1532177', 0.000508366448721294), ('3223254', 0.0005075161135669329), ('3717154', 0.0005074080032469161), ('1399158', 0.0005049758590305695), ('2721657',

0.0005032241029653155), ('4445210', 0.000501755697768795), ('562732',
0.0005014936308391629), ('3280905', 0.0004993352417101815), ('3109237',
0.0004967826100243339), ('2567136', 0.0004962327255176554), ('7743887',
0.0004961509906644764), ('391152', 0.0004946569792043787), ('342623',
0.0004934021619680744), ('2542703', 0.0004931726120872555), ('783721',
0.0004912445884635099), ('453068', 0.0004911610251478994), ('2171244',
0.0004909645244528675), ('1162550', 0.0004905266615409942), ('1609079',
0.0004896716195921394), ('1355507', 0.0004890386060162852), ('865972',
0.0004876636598158628), ('2911571', 0.00048753796078929563), ('11228318',
0.00048721077864150806), ('3079504', 0.00048672638996439513), ('1219578',
0.0004831620224075807), ('33913', 0.0004819644911621755), ('2467618',
0.00048084094587546314), ('2626332', 0.00047991663229289896), ('1617546',
0.0004765844020130472), ('386241', 0.0004750729334704502), ('463200',
0.000469987250618484), ('2248099', 0.0004660222986406398), ('4888011',
0.00046571307425175273), ('453228', 0.00046480327823914203), ('1406078',
0.0004634723174227236), ('2488512', 0.00046298833960207623), ('663392',
0.000457939436354479), ('4248714', 0.00045657250452509675), ('306861',
0.000455412241900875), ('1293485', 0.000453259433425138), ('583640',
0.00045038665033381544), ('1111192', 0.00044995816421382237), ('11423174',
0.00044659087115621186), ('2093818', 0.00042446521019174607), ('1416505',
0.00041067383143243697), ('1132073', 0.00037773379384001005)]

The difference makes sense intuitively. Perhaps a system that matches a test file to the author of a training file by the use of distinct words better captures the uniqueness of a writing style. But we also noticed that 70535 contains longer strings of abbreviating ellipses, which other files share and reduce the AllTokens system's confidence in 70535's authorship. If the training set's "70535.txt" contained several ellipses but the test set's did not, then the system is more likely to attribute the latter to an author of a training file with fewer ellipses. On the other hand, Singleton is not confused by repeated occurrences of the same token.