

## System Write-Up

### Part A: Descriptive Statistics

#### **Total words : 1800**

This number calculated by running `nlk.word_tokenize()` on the scrubbed file and counting the length of the returned list.

#### **Total sentences : 106**

This number calculated by running a sentence tokenizer, loaded from 'tokenizers/punkt/english.pickle', on the scrubbed file and counting the length of the returned list.

#### **Total Unique Words : 774**

#### **FREQUENCY TABLE (goes for a few pages) :**

Word	Occurrences	Frequency
-----	-----	-----
the	112	1
and	57	2
s	50	3
a	49	4
of	43	5
it	43	6
in	31	7
i	30	8
to	29	9
that	25	10
mr	25	11
return	23	12
lynch	20	13
is	19	14
on	18	15
like	15	16
was	14	17
this	14	18
twin	13	19
peaks	13	20

he	13	21
has	13	22
for	13	23
but	13	24
with	12	25
tv	11	26
one	11	27
we	10	28
its	10	29
t	9	30
not	9	31
his	9	32
as	9	33
who	8	34
frost	8	35
from	8	36
been	8	37
you	7	38
what	7	39
show	7	40
new	7	41
first	7	42
say	6	43
most	6	44
love	6	45
have	6	46
be	6	47
at	6	48
world	5	49
when	5	50
they	5	51
so	5	52
series	5	53
more	5	54

just	5	55
here	5	56
film	5	57
different	5	58
cooper	5	59
also	5	60
which	4	61
ve	4	62
two	4	63
think	4	64
there	4	65
seems	4	66
season	4	67
re	4	68
original	4	69
now	4	70
me	4	71
man	4	72
james	4	73
if	4	74
how	4	75
even	4	76
can	4	77
brando	4	78
before	4	79
an	4	80
again	4	81
would	3	82
wally	3	83
too	3	84
their	3	85
than	3	86
television	3	87
something	3	88

<b>some</b>	<b>3</b>	<b>89</b>
<b>since</b>	<b>3</b>	<b>90</b>
<b>see</b>	<b>3</b>	<b>91</b>
<b>palmer</b>	<b>3</b>	<b>92</b>
<b>over</b>	<b>3</b>	<b>93</b>
<b>out</b>	<b>3</b>	<b>94</b>
<b>our</b>	<b>3</b>	<b>95</b>
<b>old</b>	<b>3</b>	<b>96</b>
<b>off</b>	<b>3</b>	<b>97</b>
<b>no</b>	<b>3</b>	<b>98</b>
<b>movie</b>	<b>3</b>	<b>99</b>
<b>m</b>	<b>3</b>	<b>100</b>
<b>laura</b>	<b>3</b>	<b>101</b>
<b>into</b>	<b>3</b>	<b>102</b>
<b>had</b>	<b>3</b>	<b>103</b>
<b>going</b>	<b>3</b>	<b>104</b>
<b>familiar</b>	<b>3</b>	<b>105</b>
<b>episode</b>	<b>3</b>	<b>106</b>
<b>don</b>	<b>3</b>	<b>107</b>
<b>david</b>	<b>3</b>	<b>108</b>
<b>dale</b>	<b>3</b>	<b>109</b>
<b>course</b>	<b>3</b>	<b>110</b>
<b>by</b>	<b>3</b>	<b>111</b>
<b>box</b>	<b>3</b>	<b>112</b>
<b>because</b>	<b>3</b>	<b>113</b>
<b>are</b>	<b>3</b>	<b>114</b>
<b>always</b>	<b>3</b>	<b>115</b>
<b>all</b>	<b>3</b>	<b>116</b>
<b>years</b>	<b>2</b>	<b>117</b>
<b>year</b>	<b>2</b>	<b>118</b>
<b>work</b>	<b>2</b>	<b>119</b>
<b>word</b>	<b>2</b>	<b>120</b>
<b>wonder</b>	<b>2</b>	<b>121</b>
<b>women</b>	<b>2</b>	<b>122</b>

<b>while</b>	<b>2</b>	<b>123</b>
<b>were</b>	<b>2</b>	<b>124</b>
<b>way</b>	<b>2</b>	<b>125</b>
<b>watch</b>	<b>2</b>	<b>126</b>
<b>viewers</b>	<b>2</b>	<b>127</b>
<b>us</b>	<b>2</b>	<b>128</b>
<b>up</b>	<b>2</b>	<b>129</b>
<b>unspeakable</b>	<b>2</b>	<b>130</b>
<b>times</b>	<b>2</b>	<b>131</b>
<b>time</b>	<b>2</b>	<b>132</b>
<b>through</b>	<b>2</b>	<b>133</b>
<b>those</b>	<b>2</b>	<b>134</b>
<b>these</b>	<b>2</b>	<b>135</b>
<b>then</b>	<b>2</b>	<b>136</b>
<b>terrifying</b>	<b>2</b>	<b>137</b>
<b>take</b>	<b>2</b>	<b>138</b>
<b>sometimes</b>	<b>2</b>	<b>139</b>
<b>small</b>	<b>2</b>	<b>140</b>
<b>simple</b>	<b>2</b>	<b>141</b>
<b>shows</b>	<b>2</b>	<b>142</b>
<b>scene</b>	<b>2</b>	<b>143</b>
<b>room</b>	<b>2</b>	<b>144</b>
<b>read</b>	<b>2</b>	<b>145</b>
<b>quote</b>	<b>2</b>	<b>146</b>
<b>popping</b>	<b>2</b>	<b>147</b>
<b>poniewozik</b>	<b>2</b>	<b>148</b>
<b>person</b>	<b>2</b>	<b>149</b>
<b>part</b>	<b>2</b>	<b>150</b>
<b>other</b>	<b>2</b>	<b>151</b>
<b>or</b>	<b>2</b>	<b>152</b>
<b>open</b>	<b>2</b>	<b>153</b>
<b>often</b>	<b>2</b>	<b>154</b>
<b>needs</b>	<b>2</b>	<b>155</b>
<b>need</b>	<b>2</b>	<b>156</b>

mystery	2	157
much	2	158
movies	2	159
mind	2	160
men	2	161
medium	2	162
matter	2	163
many	2	164
made	2	165
mad	2	166
maclachlan	2	167
lynchian	2	168
logic	2	169
live	2	170
know	2	171
knew	2	172
itself	2	173
instead	2	174
inland	2	175
including	2	176
idea	2	177
hook	2	178
hell	2	179
head	2	180
go	2	181
felt	2	182
f.b.i	2	183
eyes	2	184
everyone	2	185
eternal	2	186
entertainment	2	187
empire	2	188
else	2	189
easy	2	190

<b>driven</b>	<b>2</b>	<b>191</b>
<b>down</b>	<b>2</b>	<b>192</b>
<b>dougie</b>	<b>2</b>	<b>193</b>
<b>doesn</b>	<b>2</b>	<b>194</b>
<b>directly</b>	<b>2</b>	<b>195</b>
<b>deep</b>	<b>2</b>	<b>196</b>
<b>critics</b>	<b>2</b>	<b>197</b>
<b>critic</b>	<b>2</b>	<b>198</b>
<b>crime</b>	<b>2</b>	<b>199</b>
<b>could</b>	<b>2</b>	<b>200</b>
<b>cool</b>	<b>2</b>	<b>201</b>
<b>close</b>	<b>2</b>	<b>202</b>
<b>characters</b>	<b>2</b>	<b>203</b>
<b>cartoon</b>	<b>2</b>	<b>204</b>
<b>call</b>	<b>2</b>	<b>205</b>
<b>built</b>	<b>2</b>	<b>206</b>
<b>black</b>	<b>2</b>	<b>207</b>
<b>big</b>	<b>2</b>	<b>208</b>
<b>bad</b>	<b>2</b>	<b>209</b>
<b>audience</b>	<b>2</b>	<b>210</b>
<b>assumed</b>	<b>2</b>	<b>211</b>
<b>another</b>	<b>2</b>	<b>212</b>
<b>am</b>	<b>2</b>	<b>213</b>
<b>agent</b>	<b>2</b>	<b>214</b>
<b>abc</b>	<b>2</b>	<b>215</b>
	<b>2</b>	<b>216</b>
<b>young</b>	<b>1</b>	<b>217</b>
<b>york</b>	<b>1</b>	<b>218</b>
<b>yet</b>	<b>1</b>	<b>219</b>
<b>writing</b>	<b>1</b>	<b>220</b>
<b>wrapped</b>	<b>1</b>	<b>221</b>
<b>wow</b>	<b>1</b>	<b>222</b>
<b>worlds</b>	<b>1</b>	<b>223</b>
<b>wore</b>	<b>1</b>	<b>224</b>

woodsman	1	225
woman	1	226
wizard	1	227
witness	1	228
wish	1	229
will	1	230
wild	1	231
wicked	1	232
why	1	233
whom	1	234
western	1	235
well	1	236
weekly	1	237
week	1	238
web	1	239
weakest	1	240
waylaid	1	241
watching	1	242
wasn	1	243
warren	1	244
wants	1	245
wall	1	246
walk	1	247
waking	1	248
vortexes	1	249
void	1	250
voice	1	251
visual	1	252
vision	1	253
viewership	1	254
very	1	255
velocity	1	256
<a href="http://vegas.it">vegas.it</a>	1	257
usual	1	258



using	1	259
used	1	260
use	1	261
until	1	262
unified	1	263
under	1	264
unconscious	1	265
turns	1	266
turning	1	267
turned	1	268
tropes	1	269
town	1	270
tortured	1	271
topped	1	272
took	1	273
together	1	274
timed	1	275
tidy	1	276
threat	1	277
though	1	278
things	1	279
them.that	1	280
theatrically	1	281
that/him/her	1	282
terry	1	283
takes	1	284
switched	1	285
sweep	1	286
swaggering	1	287
surrealists	1	288
surrealism	1	289
surprising	1	290
surprised	1	291
sure	1	292

<b>sunday</b>	<b>1</b>	<b>293</b>
<b>stuff</b>	<b>1</b>	<b>294</b>
<b>structure</b>	<b>1</b>	<b>295</b>
<b>streets</b>	<b>1</b>	<b>296</b>
<b>strange</b>	<b>1</b>	<b>297</b>
<b>straddled</b>	<b>1</b>	<b>298</b>
<b>story</b>	<b>1</b>	<b>299</b>
<b>storehouse</b>	<b>1</b>	<b>300</b>
<b>store</b>	<b>1</b>	<b>301</b>
<b>start</b>	<b>1</b>	<b>302</b>
<b>standards</b>	<b>1</b>	<b>303</b>
<b>stand</b>	<b>1</b>	<b>304</b>
<b>spooky</b>	<b>1</b>	<b>305</b>
<b>spirited</b>	<b>1</b>	<b>306</b>
<b>speaks</b>	<b>1</b>	<b>307</b>
<b>sort</b>	<b>1</b>	<b>308</b>
<b>sooty</b>	<b>1</b>	<b>309</b>
<b>somewhat</b>	<b>1</b>	<b>310</b>
<b>soaps</b>	<b>1</b>	<b>311</b>
<b>soap</b>	<b>1</b>	<b>312</b>
<b>snob</b>	<b>1</b>	<b>313</b>
<b>slippers</b>	<b>1</b>	<b>314</b>
<b>slice</b>	<b>1</b>	<b>315</b>
<b>skull</b>	<b>1</b>	<b>316</b>
<b>sized</b>	<b>1</b>	<b>317</b>
<b>sitcom</b>	<b>1</b>	<b>318</b>
<b>sit</b>	<b>1</b>	<b>319</b>
<b>siren</b>	<b>1</b>	<b>320</b>
<b>significant</b>	<b>1</b>	<b>321</b>
<b>shut</b>	<b>1</b>	<b>322</b>
<b>shrouded</b>	<b>1</b>	<b>323</b>
<b>shovels</b>	<b>1</b>	<b>324</b>
<b>shoveling</b>	<b>1</b>	<b>325</b>
<b>shouldn</b>	<b>1</b>	<b>326</b>

<b>should</b>	<b>1</b>	<b>327</b>
<b>shoes</b>	<b>1</b>	<b>328</b>
<b>ship</b>	<b>1</b>	<b>329</b>
<b>shifted</b>	<b>1</b>	<b>330</b>
<b>shaped</b>	<b>1</b>	<b>331</b>
<b>shape</b>	<b>1</b>	<b>332</b>
<b>shaking</b>	<b>1</b>	<b>333</b>
<b>sewn</b>	<b>1</b>	<b>334</b>
<b>setup</b>	<b>1</b>	<b>335</b>
<b>setting</b>	<b>1</b>	<b>336</b>
<b>set</b>	<b>1</b>	<b>337</b>
<b>serves</b>	<b>1</b>	<b>338</b>
<b>serve</b>	<b>1</b>	<b>339</b>
<b>sequel</b>	<b>1</b>	<b>340</b>
<b>sentimental</b>	<b>1</b>	<b>341</b>
<b>self</b>	<b>1</b>	<b>342</b>
<b>seen</b>	<b>1</b>	<b>343</b>
<b>secrecy</b>	<b>1</b>	<b>344</b>
<b>seasons</b>	<b>1</b>	<b>345</b>
<b>screwball</b>	<b>1</b>	<b>346</b>
<b>screen</b>	<b>1</b>	<b>347</b>
<b>sci</b>	<b>1</b>	<b>348</b>
<b>sarah</b>	<b>1</b>	<b>349</b>
<b>same</b>	<b>1</b>	<b>350</b>
<b>said</b>	<b>1</b>	<b>351</b>
<b>sacrificed</b>	<b>1</b>	<b>352</b>
<b>run</b>	<b>1</b>	<b>353</b>
<b>ruby</b>	<b>1</b>	<b>354</b>
<b>roles</b>	<b>1</b>	<b>355</b>
<b>rode</b>	<b>1</b>	<b>356</b>
<b>roads</b>	<b>1</b>	<b>357</b>
<b>roadhouse</b>	<b>1</b>	<b>358</b>
<b>road</b>	<b>1</b>	<b>359</b>
<b>right</b>	<b>1</b>	<b>360</b>

<b>riffle</b>	<b>1</b>	<b>361</b>
<b>riff</b>	<b>1</b>	<b>362</b>
<b>riding</b>	<b>1</b>	<b>363</b>
<b>ride</b>	<b>1</b>	<b>364</b>
<b>revivals</b>	<b>1</b>	<b>365</b>
<b>revisit</b>	<b>1</b>	<b>366</b>
<b>reunion</b>	<b>1</b>	<b>367</b>
<b>retrospective</b>	<b>1</b>	<b>368</b>
<b>resist</b>	<b>1</b>	<b>369</b>
<b>reminds</b>	<b>1</b>	<b>370</b>
<b>remember</b>	<b>1</b>	<b>371</b>
<b>rematerialized</b>	<b>1</b>	<b>372</b>
<b>reintroduction s</b>	<b>1</b>	<b>373</b>
<b>red</b>	<b>1</b>	<b>374</b>
<b>recently</b>	<b>1</b>	<b>375</b>
<b>reboots</b>	<b>1</b>	<b>376</b>
<b>realized</b>	<b>1</b>	<b>377</b>
<b>real</b>	<b>1</b>	<b>378</b>
<b>rare</b>	<b>1</b>	<b>379</b>
<b>rainbow</b>	<b>1</b>	<b>380</b>
<b>radiator</b>	<b>1</b>	<b>381</b>
<b>rabbit</b>	<b>1</b>	<b>382</b>
<b>question</b>	<b>1</b>	<b>383</b>
<b>python</b>	<b>1</b>	<b>384</b>
<b>put</b>	<b>1</b>	<b>385</b>
<b>pushes</b>	<b>1</b>	<b>386</b>
<b>pull</b>	<b>1</b>	<b>387</b>
<b>preyed</b>	<b>1</b>	<b>388</b>
<b>premiere</b>	<b>1</b>	<b>389</b>
<b>predecessor</b>	<b>1</b>	<b>390</b>
<b>pouter</b>	<b>1</b>	<b>391</b>
<b>porous</b>	<b>1</b>	<b>392</b>
<b>poetry</b>	<b>1</b>	<b>393</b>

pleasure	1	394
pleasing	1	395
pleasantly	1	396
plays	1	397
playing	1	398
plated	1	399
plane	1	400
place	1	401
pilot	1	402
picture	1	403
photo	1	404
performances	1	405
perfectly	1	406
payoff	1	407
pay	1	408
past	1	409
parts	1	410
pandora	1	411
painter	1	412
oz	1	413
own	1	414
outside	1	415
ostensibly	1	416
opera	1	417
opens	1	418
openly	1	419
ooo	1	420
onto	1	421
only	1	422
once	1	423
odysseus	1	424
obviously	1	425
obligatory	1	426
o	1	427

nuclear	1	428
nowhere	1	429
novocain	1	430
novel	1	431
nostalgia	1	432
nor	1	433
nod	1	434
nightmare	1	435
night	1	436
nice	1	437
never	1	438
neither	1	439
neck	1	440
neat	1	441
naughty	1	442
narrative	1	443
namechecks	1	444
mysteries	1	445
museum	1	446
murdered	1	447
mulholland	1	448
muck	1	449
moves	1	450
motorcycle	1	451
motifs	1	452
monty	1	453
mirrors	1	454
miguel	1	455
michael	1	456
meta	1	457
met	1	458
mention	1	459
memorable	1	460
means	1	461

<b>maybe</b>	<b>1</b>	<b>462</b>
<b>May</b>	<b>1</b>	<b>463</b>
<b>masters</b>	<b>1</b>	<b>464</b>
<b>master</b>	<b>1</b>	<b>465</b>
<b>masks</b>	<b>1</b>	<b>466</b>
<b>masculine</b>	<b>1</b>	<b>467</b>
<b>mark</b>	<b>1</b>	<b>468</b>
<b>map</b>	<b>1</b>	<b>469</b>
<b>manohla</b>	<b>1</b>	<b>470</b>
<b>make</b>	<b>1</b>	<b>471</b>
<b>main</b>	<b>1</b>	<b>472</b>
<b>lure</b>	<b>1</b>	<b>473</b>
<b>lost</b>	<b>1</b>	<b>474</b>
<b>lose</b>	<b>1</b>	<b>475</b>
<b>look</b>	<b>1</b>	<b>476</b>
<b>long</b>	<b>1</b>	<b>477</b>
<b>lonely</b>	<b>1</b>	<b>478</b>
<b>lodge</b>	<b>1</b>	<b>479</b>
<b>locations</b>	<b>1</b>	<b>480</b>
<b>lived</b>	<b>1</b>	<b>481</b>
<b>literally</b>	<b>1</b>	<b>482</b>
<b>likes</b>	<b>1</b>	<b>483</b>
<b>life</b>	<b>1</b>	<b>484</b>
<b>let</b>	<b>1</b>	<b>485</b>
<b>legacy</b>	<b>1</b>	<b>486</b>
<b>leather</b>	<b>1</b>	<b>487</b>
<b>lays</b>	<b>1</b>	<b>488</b>
<b>laughed</b>	<b>1</b>	<b>489</b>
<b>late</b>	<b>1</b>	<b>490</b>
<b>last</b>	<b>1</b>	<b>491</b>
<b>las</b>	<b>1</b>	<b>492</b>
<b>larger</b>	<b>1</b>	<b>493</b>
<b>lane</b>	<b>1</b>	<b>494</b>
<b>landmark</b>	<b>1</b>	<b>495</b>

<b>lady</b>	<b>1</b>	<b>496</b>
<b>label</b>	<b>1</b>	<b>497</b>
<b>kyle</b>	<b>1</b>	<b>498</b>
<b>kind</b>	<b>1</b>	<b>499</b>
<b>killed</b>	<b>1</b>	<b>500</b>
<b>keys</b>	<b>1</b>	<b>501</b>
<b>kept</b>	<b>1</b>	<b>502</b>
<b>keeping</b>	<b>1</b>	<b>503</b>
<b>kafka</b>	<b>1</b>	<b>504</b>
<b>judges</b>	<b>1</b>	<b>505</b>
<b>jones</b>	<b>1</b>	<b>506</b>
<b>jokes</b>	<b>1</b>	<b>507</b>
<b>johnny</b>	<b>1</b>	<b>508</b>
<b>jim</b>	<b>1</b>	<b>509</b>
<b>jacket</b>	<b>1</b>	<b>510</b>
<b>isn</b>	<b>1</b>	<b>511</b>
<b>irrelevant</b>	<b>1</b>	<b>512</b>
<b>invokes</b>	<b>1</b>	<b>513</b>
<b>invoke</b>	<b>1</b>	<b>514</b>
<b>intoning</b>	<b>1</b>	<b>515</b>
<b>interpretive</b>	<b>1</b>	<b>516</b>
<b>interested</b>	<b>1</b>	<b>517</b>
<b>instructive</b>	<b>1</b>	<b>518</b>
<b>instance</b>	<b>1</b>	<b>519</b>
<b>inside</b>	<b>1</b>	<b>520</b>
<b>inflected</b>	<b>1</b>	<b>521</b>
<b>ineffable</b>	<b>1</b>	<b>522</b>
<b>impossibly</b>	<b>1</b>	<b>523</b>
<b>images</b>	<b>1</b>	<b>524</b>
<b>imagery</b>	<b>1</b>	<b>525</b>
<b>ideas</b>	<b>1</b>	<b>526</b>
<b>hurley</b>	<b>1</b>	<b>527</b>
<b>humorous</b>	<b>1</b>	<b>528</b>
<b>hours</b>	<b>1</b>	<b>529</b>



<b>hour</b>	<b>1</b>	<b>530</b>
<b>horrors</b>	<b>1</b>	<b>531</b>
<b>horror</b>	<b>1</b>	<b>532</b>
<b>horrific</b>	<b>1</b>	<b>533</b>
<b>horne</b>	<b>1</b>	<b>534</b>
<b>home</b>	<b>1</b>	<b>535</b>
<b>hits</b>	<b>1</b>	<b>536</b>
<b>history</b>	<b>1</b>	<b>537</b>
<b>hired</b>	<b>1</b>	<b>538</b>
<b>hinged</b>	<b>1</b>	<b>539</b>
<b>himself</b>	<b>1</b>	<b>540</b>
<b>him</b>	<b>1</b>	<b>541</b>
<b>highway</b>	<b>1</b>	<b>542</b>
<b>helped</b>	<b>1</b>	<b>543</b>
<b>heels</b>	<b>1</b>	<b>544</b>
<b>heaven/everyt hing</b>	<b>1</b>	<b>545</b>
<b>heaven</b>	<b>1</b>	<b>546</b>
<b>heads</b>	<b>1</b>	<b>547</b>
<b>headlights</b>	<b>1</b>	<b>548</b>
<b>harley</b>	<b>1</b>	<b>549</b>
<b>hard</b>	<b>1</b>	<b>550</b>
<b>happy</b>	<b>1</b>	<b>551</b>
<b>happily</b>	<b>1</b>	<b>552</b>
<b>happens</b>	<b>1</b>	<b>553</b>
<b>great</b>	<b>1</b>	<b>554</b>
<b>gratification</b>	<b>1</b>	<b>555</b>
<b>grateful</b>	<b>1</b>	<b>556</b>
<b>gordon</b>	<b>1</b>	<b>557</b>
<b>good</b>	<b>1</b>	<b>558</b>
<b>golly</b>	<b>1</b>	<b>559</b>
<b>gold</b>	<b>1</b>	<b>560</b>
<b>glorious</b>	<b>1</b>	<b>561</b>
<b>glad</b>	<b>1</b>	<b>562</b>

<b>girls</b>	<b>1</b>	<b>563</b>
<b>gilliam</b>	<b>1</b>	<b>564</b>
<b>ghastly</b>	<b>1</b>	<b>565</b>
<b>get</b>	<b>1</b>	<b>566</b>
<b>genres</b>	<b>1</b>	<b>567</b>
<b>genre</b>	<b>1</b>	<b>568</b>
<b>funny</b>	<b>1</b>	<b>569</b>
<b>friend</b>	<b>1</b>	<b>570</b>
<b>freud</b>	<b>1</b>	<b>571</b>
<b>free</b>	<b>1</b>	<b>572</b>
<b>franchise</b>	<b>1</b>	<b>573</b>
<b>found</b>	<b>1</b>	<b>574</b>
<b>formulaic</b>	<b>1</b>	<b>575</b>
<b>forms</b>	<b>1</b>	<b>576</b>
<b>form</b>	<b>1</b>	<b>577</b>
<b>forced</b>	<b>1</b>	<b>578</b>
<b>folk</b>	<b>1</b>	<b>579</b>
<b>flying</b>	<b>1</b>	<b>580</b>
<b>flavor</b>	<b>1</b>	<b>581</b>
<b>flashed</b>	<b>1</b>	<b>582</b>
<b>fixations</b>	<b>1</b>	<b>583</b>
<b>fixated</b>	<b>1</b>	<b>584</b>
<b>fire</b>	<b>1</b>	<b>585</b>
<b>finishing</b>	<b>1</b>	<b>586</b>
<b>fine/in</b>	<b>1</b>	<b>587</b>
<b>finally</b>	<b>1</b>	<b>588</b>
<b>finale</b>	<b>1</b>	<b>589</b>
<b>fi</b>	<b>1</b>	<b>590</b>
<b>ferrer</b>	<b>1</b>	<b>591</b>
<b>feel</b>	<b>1</b>	<b>592</b>
<b>fascinating</b>	<b>1</b>	<b>593</b>
<b>far</b>	<b>1</b>	<b>594</b>
<b>face</b>	<b>1</b>	<b>595</b>
<b>express</b>	<b>1</b>	<b>596</b>

<b>explosion</b>	<b>1</b>	<b>597</b>
<b>exploited</b>	<b>1</b>	<b>598</b>
<b>explanatory</b>	<b>1</b>	<b>599</b>
<b>explanations</b>	<b>1</b>	<b>600</b>
<b>expanse</b>	<b>1</b>	<b>601</b>
<b>existential</b>	<b>1</b>	<b>602</b>
<b>existence</b>	<b>1</b>	<b>603</b>
<b>evolving</b>	<b>1</b>	<b>604</b>
<b>evil</b>	<b>1</b>	<b>605</b>
<b>everywhere</b>	<b>1</b>	<b>606</b>
<b>every</b>	<b>1</b>	<b>607</b>
<b>ever</b>	<b>1</b>	<b>608</b>
<b>events</b>	<b>1</b>	<b>609</b>
<b>especially</b>	<b>1</b>	<b>610</b>
<b>eraserhead</b>	<b>1</b>	<b>611</b>
<b>episodic</b>	<b>1</b>	<b>612</b>
<b>episodes</b>	<b>1</b>	<b>613</b>
<b>entirely</b>	<b>1</b>	<b>614</b>
<b>enjoyed</b>	<b>1</b>	<b>615</b>
<b>enigmatic</b>	<b>1</b>	<b>616</b>
<b>engaging</b>	<b>1</b>	<b>617</b>
<b>engagement</b>	<b>1</b>	<b>618</b>
<b>ends</b>	<b>1</b>	<b>619</b>
<b>eluded</b>	<b>1</b>	<b>620</b>
<b>eerie</b>	<b>1</b>	<b>621</b>
<b>earth</b>	<b>1</b>	<b>622</b>
<b>eagerly</b>	<b>1</b>	<b>623</b>
<b>each</b>	<b>1</b>	<b>624</b>
<b>e.</b>	<b>1</b>	<b>625</b>
<b>dusty</b>	<b>1</b>	<b>626</b>
<b>dust</b>	<b>1</b>	<b>627</b>
<b>drve</b>	<b>1</b>	<b>628</b>
<b>drowning</b>	<b>1</b>	<b>629</b>
<b>drive</b>	<b>1</b>	<b>630</b>

<b>dream</b>	<b>1</b>	<b>631</b>
<b>draw</b>	<b>1</b>	<b>632</b>
<b>drama</b>	<b>1</b>	<b>633</b>
<b>dr</b>	<b>1</b>	<b>634</b>
<b>douglas</b>	<b>1</b>	<b>635</b>
<b>doppelgänger</b>	<b>1</b>	<b>636</b>
<b>does</b>	<b>1</b>	<b>637</b>
<b>divide</b>	<b>1</b>	<b>638</b>
<b>distributed</b>	<b>1</b>	<b>639</b>
<b>distinctly</b>	<b>1</b>	<b>640</b>
<b>disjointed</b>	<b>1</b>	<b>641</b>
<b>discussing</b>	<b>1</b>	<b>642</b>
<b>discuss</b>	<b>1</b>	<b>643</b>
<b>director</b>	<b>1</b>	<b>644</b>
<b>directing</b>	<b>1</b>	<b>645</b>
<b>dim</b>	<b>1</b>	<b>646</b>
<b>digestion</b>	<b>1</b>	<b>647</b>
<b>did</b>	<b>1</b>	<b>648</b>
<b>dice</b>	<b>1</b>	<b>649</b>
<b>detours</b>	<b>1</b>	<b>650</b>
<b>detonated</b>	<b>1</b>	<b>651</b>
<b>dern</b>	<b>1</b>	<b>652</b>
<b>demons</b>	<b>1</b>	<b>653</b>
<b>defies</b>	<b>1</b>	<b>654</b>
<b>declared</b>	<b>1</b>	<b>655</b>
<b>decency</b>	<b>1</b>	<b>656</b>
<b>decaf</b>	<b>1</b>	<b>657</b>
<b>deaths</b>	<b>1</b>	<b>658</b>
<b>daytime</b>	<b>1</b>	<b>659</b>
<b>dark</b>	<b>1</b>	<b>660</b>
<b>dargis</b>	<b>1</b>	<b>661</b>
<b>dance</b>	<b>1</b>	<b>662</b>
<b>dammit</b>	<b>1</b>	<b>663</b>

<b>customizing</b>	<b>1</b>	<b>664</b>
<b>crushing</b>	<b>1</b>	<b>665</b>
<b>crowd</b>	<b>1</b>	<b>666</b>
<b>creators</b>	<b>1</b>	<b>667</b>
<b>cranking</b>	<b>1</b>	<b>668</b>
<b>craft</b>	<b>1</b>	<b>669</b>
<b>cozily</b>	<b>1</b>	<b>670</b>
<b>cox</b>	<b>1</b>	<b>671</b>
<b>count</b>	<b>1</b>	<b>672</b>
<b>coulson</b>	<b>1</b>	<b>673</b>
<b>core</b>	<b>1</b>	<b>674</b>
<b>copying</b>	<b>1</b>	<b>675</b>
<b>cops</b>	<b>1</b>	<b>676</b>
<b>conventions</b>	<b>1</b>	<b>677</b>
<b>convenience</b>	<b>1</b>	<b>678</b>
<b>contemporary</b>	<b>1</b>	<b>679</b>
<b>consistently</b>	<b>1</b>	<b>680</b>
<b>confrontational</b>	<b>1</b>	<b>681</b>
<b>compliment</b>	<b>1</b>	<b>682</b>
<b>complicated</b>	<b>1</b>	<b>683</b>
<b>common</b>	<b>1</b>	<b>684</b>
<b>comes</b>	<b>1</b>	<b>685</b>
<b>comedy</b>	<b>1</b>	<b>686</b>
<b>come</b>	<b>1</b>	<b>687</b>
<b>collection</b>	<b>1</b>	<b>688</b>
<b>coffee</b>	<b>1</b>	<b>689</b>
<b>co</b>	<b>1</b>	<b>690</b>
<b>click</b>	<b>1</b>	<b>691</b>
<b>clearly</b>	<b>1</b>	<b>692</b>
<b>cited</b>	<b>1</b>	<b>693</b>
<b>churning</b>	<b>1</b>	<b>694</b>
<b>children</b>	<b>1</b>	<b>695</b>
<b>character</b>	<b>1</b>	<b>696</b>

<b>channel</b>	<b>1</b>	<b>697</b>
<b>certain</b>	<b>1</b>	<b>698</b>
<b>cera</b>	<b>1</b>	<b>699</b>
<b>catherine</b>	<b>1</b>	<b>700</b>
<b>cash</b>	<b>1</b>	<b>701</b>
<b>career</b>	<b>1</b>	<b>702</b>
<b>car</b>	<b>1</b>	<b>703</b>
<b>cannes</b>	<b>1</b>	<b>704</b>
<b>called</b>	<b>1</b>	<b>705</b>
<b>cable</b>	<b>1</b>	<b>706</b>
<b>busted</b>	<b>1</b>	<b>707</b>
<b>broken</b>	<b>1</b>	<b>708</b>
<b>brilliant</b>	<b>1</b>	<b>709</b>
<b>brilliance</b>	<b>1</b>	<b>710</b>
<b>brand</b>	<b>1</b>	<b>711</b>
<b>boys</b>	<b>1</b>	<b>712</b>
<b>boxes</b>	<b>1</b>	<b>713</b>
<b>bowie</b>	<b>1</b>	<b>714</b>
<b>boundary</b>	<b>1</b>	<b>715</b>
<b>both</b>	<b>1</b>	<b>716</b>
<b>born</b>	<b>1</b>	<b>717</b>
<b>booking</b>	<b>1</b>	<b>718</b>
<b>bomb</b>	<b>1</b>	<b>719</b>
<b>bittersweetly</b>	<b>1</b>	<b>720</b>
<b>between</b>	<b>1</b>	<b>721</b>
<b>best</b>	<b>1</b>	<b>722</b>
<b>benefits</b>	<b>1</b>	<b>723</b>
<b>beloved</b>	<b>1</b>	<b>724</b>
<b>behind</b>	<b>1</b>	<b>725</b>
<b>begins</b>	<b>1</b>	<b>726</b>
<b>began</b>	<b>1</b>	<b>727</b>
<b>beautiful</b>	<b>1</b>	<b>728</b>
<b>awesome</b>	<b>1</b>	<b>729</b>
<b>away</b>	<b>1</b>	<b>730</b>

<b>aware</b>	<b>1</b>	<b>731</b>
<b>aw</b>	<b>1</b>	<b>732</b>
<b>auteuristic</b>	<b>1</b>	<b>733</b>
<b>auteur</b>	<b>1</b>	<b>734</b>
<b>audrey</b>	<b>1</b>	<b>735</b>
<b>attenuated</b>	<b>1</b>	<b>736</b>
<b>atom</b>	<b>1</b>	<b>737</b>
<b>astonishing</b>	<b>1</b>	<b>738</b>
<b>assume</b>	<b>1</b>	<b>739</b>
<b>asides</b>	<b>1</b>	<b>740</b>
<b>artists</b>	<b>1</b>	<b>741</b>
<b>artist</b>	<b>1</b>	<b>742</b>
<b>art</b>	<b>1</b>	<b>743</b>
<b>around</b>	<b>1</b>	<b>744</b>
<b>applied</b>	<b>1</b>	<b>745</b>
<b>apparent</b>	<b>1</b>	<b>746</b>
<b>anything</b>	<b>1</b>	<b>747</b>
<b>answered</b>	<b>1</b>	<b>748</b>
<b>amusing</b>	<b>1</b>	<b>749</b>
<b>amp</b>	<b>1</b>	<b>750</b>
<b>america</b>	<b>1</b>	<b>751</b>
<b>ambitious</b>	<b>1</b>	<b>752</b>
<b>alternative</b>	<b>1</b>	<b>753</b>
<b>along</b>	<b>1</b>	<b>754</b>
<b>almost</b>	<b>1</b>	<b>755</b>
<b>allusions</b>	<b>1</b>	<b>756</b>
<b>alienation</b>	<b>1</b>	<b>757</b>
<b>airs</b>	<b>1</b>	<b>758</b>
<b>airing</b>	<b>1</b>	<b>759</b>
<b>ain</b>	<b>1</b>	<b>760</b>
<b>against</b>	<b>1</b>	<b>761</b>
<b>again.i</b>	<b>1</b>	<b>762</b>
<b>after</b>	<b>1</b>	<b>763</b>
<b>aesthetic</b>	<b>1</b>	<b>764</b>

admired	1	765
admire	1	766
abruptly	1	767
above	1	768
25	1	767

### Part B : Data Scrubbing Log

TAKE 1: Data found to not immediately work without user-end scrubbing.

- Spent hours fighting to convert text into a string decodable in Python
- Several 'ascii codec cannot decode' errors encountered
- Looked into the guts of the python 'read' function
- Consider emailing Blake that there's a problem

TAKE 2: Code-internal regex-based data scrubbing begins.

- Working manually in the file:
  - "Curly" quotes, apostrophes, and dashes replaced with simpler variants
  - Fixed the 'ascii codec' error
- Working in code:
  - Attempted to segment into sentences based on punctuation
  - Program fails due to Mr.'s and Dr.'s

TAKE 3 : Added more regex-based scrubbing to:

- Place all text in lowercase
- Replaces all final punctuation with periods
- Removes all other punctuation
  - Removal of apostrophes levels spelling contrast between "its" and "it's"
  - Replacement of dashes with spaces creates two distinct words from a single hyphenated word, possibly creating non-words in the process

### Part C : Table of Perplexities

This table shows the perplexities of unigrams, bigrams, and trigrams in each of the three test sets: testing on 10%, 20%, and 30%, respectively.

The perplexities did conform to our expectations: the perplexity equation is (small decimal)  $^{-1 / \text{large number}}$ , which approaches the value 1 the larger the denominator of the exponent gets. Therefore obtaining values of about 1 for our perplexity was expected.

Building this model, we ran through several iterations of this training function: originally, we had no functions and were planning to build an inordinate amount of individual dictionaries to calculate the perplexity values for each. Over time, we streamlined our approach using a set of generalized functions and the Ngram-making tools from the nltk toolkit to create a more concise and efficient parser.



The different train-test splits all showed an odd pattern: the more data we trained on, the higher our perplexity got. Since we did not shuffle the data and test on different chunks, perhaps this trend be explained by the relatively high complexity level of the final portion of the test set: if the last 10th were inordinately complex, then the models which trained on 90% would do worse, on average, than the models which trained on 70% and tested on a relatively easy 20% in addition to the same complex final 10%. These statistical biases could be avoided in the future by shuffling the data prior to training and testing to ensure that no particular piece of the data skews the results in this way.

Test%	Unigram Perplexity	Bigram Perplexity	Trigram Perplexity
Test10	1.02852	1.05033	1.05431
Test20	1.01478	1.02679	1.02682
Test30	1.0154	1.02435	1.02438

#### Part D: Test Sentence MLE Tables

Once the perplexity values were calculated, choosing the best model to evaluate the input sentences was straightforward. We simply reused the evaluation function made for calculating the perplexity values in Part C to calculate the MLE values for each sentence with each model.

As expected, the unigram model has the highest probability values, followed by the bigram model, and then by the trigram model. We anticipated this trend because unigrams are the least unique n-gram, allowing the same n-gram to appear more frequently than in the bigram and trigram models.

Sentence Number	Best Unigram	Best Bigram	Best Trigram
1	5.83E-25	4.64E-54	1.00E-60
2	1.56E-26	7.18E-52	1.00E-54
3	2.60E-32	9.87E-53	3.72E-64
4	4.23E-58	3.71E-100	5.18E-103
5	1.49E-20	1.03E-41	2.88E-45
6	4.15E-24	5.16E-37	7.19E-40
7	9.21E-112	7.65E-178	2.07E-186
8	1.05E-66	1.00E-96	1.00E-96
9	3.42E-29	1.00E-36	1.00E-36
10	2.60E-62	1.00E-90	1.00E-90

