

helped helpful helping hero heroic heroically heroine heroize heros high-quality		
high-spirited hilarious holy homage honest honesty honor honorable honored honoring hooray hopeful		
hospitable hot hotcake hotcakes hottest hug humane humble humility humor humorous humorously		
humour humourous ideal idealize ideally idol idolize idolized idyllic illuminate illuminati		
illuminating illumine illustrious ilu imaculate imaginative immaculate immaculately immense impartial impartiality impartially		
<pre>impassioned impeccable impeccably important impress impressed impresses impressive impressively impressiveness improve</pre>		
<pre>improved improvement improvements improves improving incredible incredibly indebted individualized indulgence indulgent industrious</pre>		
industrious inestimable inestimably inexpensive infallibility infallible infallibly influential ingenious ingeniously ingenuity ingenuous		
ingenuous ingenuously innocuous innovation innovative inpressed insightful insightfully inspiration inspirational inspire inspiring		
<pre>instantly instructive instrumental integral integrated intelligence intelligent intelligible interesting interests intimacy</pre>		
<pre>intimate intricate intrigue intriguing intriguingly intuitive invaluable invaluablely inventive invigorate invigorating invincibility</pre>		
invincible inviolable inviolate invulnerable irreplaceable irreproachable irresistible irresistibly issue-free jaw-droping jaw-dropping		
jollify jolly jovial joy joyful joyfully joyous joyously jubilant jubilantly jubilate		
jubilation jubiliant judicious justly keen keenly keenness kid-friendly kindliness kindly		
knowledgeable kudos large-capacity laud laudable laudably lavish lavishly law-abiding lawful lawfully lead		
leading leads lean led legendary leverage levity liberate liberation liberty lifesaver		
light-hearted lighter likable like liked likes liking lionhearted lively logical long-lasting lovable		
lovable lovably love loved loveliness lovely lover loves loving low-cost low-price low-priced		
low-risk lower-priced loyal loyalty lucid lucidly luck luckier luckiest luckiness lucky		
lucrative luminous lush luster lustrous luxuriant luxuriate luxurious luxuriously luxury lyrical		
magic magical magnanimous magnanimously magnificence magnificent magnificently majestic majesty manageable maneuverable marvel		
marveled marvelled marvellous marvelous marvelously marvelousness marvels master masterful masterfully masterpiece		
masterpieces masters mastery matchless mature maturely maturity meaningful memorable merciful mercifully		
mercy merit meritorious merrily merriment merriness merry mesmerize mesmerized mesmerizes mesmerizing		
mesmerizingly meticulous meticulously mightily mighty mind-blowing miracle miracles miraculous miraculously miraculousness modern		
modern modest modesty momentous monumental monumentally morality motivated multi-purpose navigable neat neatest		
neatest neatly nice nicely nicer nicest nifty nimble noble nobly noiseless non-violence non-violent		
notably noteworthy nourish nourishing nourishment novelty nurturing oasis obsession obsessions		
obtainable openly openness optimal optimism optimistic opulent orderly originality outdo outdone outperform		
outperformed outperforming outperforms outshine outshone outsmart outstanding outstandingly outstrip outwit ovation overjoyed		
overjoyed overtake overtaken overtakes overtaking overtook overture pain-free painless painlessly palatial pamper		
pampered pamperedly pamperedness pampers panoramic paradise paramount pardon passion passionate passionately		
patience patient patiently patriot patriotic peace peaceable peaceful peacefully peacekeepers peach		
peerless pep pepped pepping peppy peps perfect perfection perfectly permissible perseverance		
persevere personages personalized phenomenal phenomenally picturesque piety pinnacle playful playfully pleasant pleasant		
pleased pleases pleasing pleasingly pleasurable pleasurably pleasure plentiful pluses plush plusses		
poetic poeticize poignant poise poised polished polite politeness popular portable posh		
positive positively positives powerful powerfully praise praiseworthy praising pre-eminent precious precise		
precisely preeminent prefer preferable preferably prefered preferes preferring prefers premier prestige prestigious		
prettily pretty priceless pride principled privilege privileged prize proactive problem-free problem-solver		
prodigious prodigiously prodigy productive productively proficient proficiently profound profoundly profuse profusion		
progress progressive prolific prominence prominent promise promised promises promising promoter prompt		
promptly proper properly propitious propitiously pros prosper prosperity prosperous prospros prospros prospros		
protection protective proud proven proves providence proving prowess prudence prudent prudently punctual		
pure purify purposeful quaint qualified qualify quicker quiet quieter radiance radiant rapid		
rapport rapt rapture raptureous raptureously rapturous rapturously rational razor-sharp reachable readable		
readily ready reaffirm reaffirmation realistic realizable reasonable reasonably reasoned reassurance reassure receptive		
receptive reclaim recomend recommend recommendation recommendations recommended reconcile reconciliation record-setting recover		
rectification rectify rectifying redeem redeeming redemption refine refined refinement reform reformed		
reforming reforms refresh refreshed refreshing refund refunded regal regally regard rejoice		
rejoicing rejoicingly rejuvenate rejuvenated rejuvenating relaxed relent reliable reliably relief relish remarkable		
remarkable remarkably remedy remission remunerate renaissance renewed renown renowned replaceable reputable reputation		
resilient resolute resound resounding resourceful resourcefulness respect respectable respectful respectfully respite resplendent		
resplendent responsibly responsive restful restored restructure restructured restructuring retractable revel revelation revere		
revere reverence reverent reverently revitalize revival revive revives revolutionary revolutionize revolutionized revolutionizes reward		
reward rewarding rewardingly rich richer richly richness right righten righteous righteously righteousness		
righteousness rightful rightfully rightly rightness risk-free robust rock-star rock-stars rockstar rockstars romantic		
romantically romanticize roomier roomy rosy safe safely sagacity sagely saint saintliness		
saintly salutary salute sane satisfactorily satisfactory satisfied satisfies satisfy satisfying satisified		
saver savings savior savvy scenic seamless seasoned secure securely selective self-determination		
self-respect self-satisfaction self-sufficiency self-sufficient sensation sensational sensationally sensations sensible sensibly sensitive serene		
serenity sexy sharp sharper sharpest shimmering shimmeringly shine shiny significant silent		
simpler simplest simplified simplifies simplify simplifying sincere sincerely sincerity skill		
skillful skillfully slammin sleek slick smart smarter smartest smartly smile smiles		
smiling smilingly smitten smooth smoother smoothes smoothest smoothly snappy snazzy sociable soft		
softer solace solicitous solicitously solid solidarity soothe soothingly sophisticated soulful soundly		
soundness spacious sparkle sparkling spectacular spectacularly speedily speedy spellbind spellbinding spellbindingly		
spellbound spirited spiritual splendid splendidly splendor spontaneous sporty spotless sprightly stability stabilize		
stabilize stable stainless standout state-of-the-art stately statuesque staunch staunchly staunchness steadfast steadfastly		
steadfastly steadfastness steadiest steadiness steady stellar stellarly stimulate stimulates stimulating stimulative		
stimulative stirringly straighten straightforward streamlined striking strikingly striving strong stronger strongest stunned		
stunned stunning stunningly stupendous stupendously sturdier sturdy stylish stylishly stylized suave suavely		
suavely sublime subsidize subsidized subsidizes subsidizing substantive succeed succeeded succeeding succeeds succeeds		
succes successes successes successful successfully suffice sufficed suffices sufficient sufficiently suitable sumptuous		
sumptuous sumptuously sumptuousness super superb superbly superior superiority supple support supported supporter		
supporter supporting supportive supports supremacy supreme supremely supurb supurbly surmount surpass surreal		
surreal survival survivor sustainability sustainable swank swankier swankiest swanky sweeping sweet sweeten		
sweeten sweetheart sweetly sweetness swift swiftness talent talented talents		
<pre>tantalize tantalizing tantalizingly tempt tempting</pre>		

thankful thinner thoughtful thoughtfully thoughtfulness

thrift

heartfelt heartily heartwarming



text_blob = """2-face 2-faces abnormal abolish abominable abominably abominate abomination abort aborted	d		
aborted aborts abrade abrasive abrupt abruptly abscond absence absent-minded absentee absurd absurdity			
absurdly absurdness abuse abused abuses abusive abysmal abysmally abyss accidental accost			
accursed accusation accusations accuse accuses accusing accusingly acerbate acerbic acerbically ache ached			
ached aches achey aching acrid acridly acridness acrimonious acrimony adamant			
adamantly addict addicted addicting addicts admonish admonisher admonishingly admonishment admonition adulterate			
adulterated adulteration adulterier adversarial adversary adverse adversity afflict affliction afflictive			
affront afraid aggravate aggravating aggravation aggression aggressive aggressiveness aggressor aggrieve aggrieved			
aggrivation aghast agonies agonize agonizing agonizingly agony aground ail ailing			
ailment aimless alarm alarmed alarming alarmingly alienate alienated alienation allegation allegations			
allege allergic allergies allergy aloof altercation ambiguity ambiguous ambivalence ambivalent			
ambush amiss amputate anarchism anarchist anarchistic anarchy anemic anger angrily angriness			
angry anguish animosity annihilate annihilation annoy annoyance annoyances annoyed annoying annoyingly			
annoys anomalous anomaly antagonism antagonist antagonistic antagonize anti- anti-american anti-israeli			
anti-occupation anti-proliferation anti-semites anti-social anti-us anti-white antipathy antiquated antithetical anxieties anxiety			
anxious anxiously anxiousness apathetic apathetically apathy apocalypse apocalyptic apologist apologists			
appal appall appalled appalling appallingly apprehension apprehensive apprehensively arbitrary arcane			
archaic arduous arduously argumentative arrogance arrogant arrogantly ashamed asinine asininely			
asinininity askance asperse aspersion aspersions assail assassin assassinate assault assult astray			
astray asunder atrocious atrocities atrocity atrophy attack attacks audacious audaciously audaciousness			
audacity audiciously austere authoritarian autocrat autocratic avalanche avarice avaricious avariciously			
avenge averse aversion aweful awful awfully awfulness awkward awkwardness ax babble			
back-logged back-wood back-woods backache backaches backaching backbite backbiting backward backward			
backwood backwoods bad badly baffle baffled bafflement baffling bait balk			
banal banalize bane banish banishment bankrupt barbarian barbaric barbarically barbarity barbarous			
barbarous barbarously barren baseless bash bashed bashful bashing bastard bastards battered battering			
battering batty bearish beastly bedlam bedlamite befoul beg beggar beggarly begging			
beguile belabor belated beleaguer belie belittle belittled belittling bellicose belligerence belligerent			
belligerent belligerently bemoan bemoaning bemused bent berate bereave bereavement bereft berserk			
berserk beseech beset besiege besmirch bestial betray betrayal betrayals betrayer betraying			
betraying betrays bewail beware bewilder bewildered bewildering bewilderingly bewilderment bewitch bias			
biased biases bicker bickering bid-rigging bigotries bigotry bitch bitchy biting			
bitingly bitter bitterly bitterness bizarre blab blabber blackmail blah blame			
blameworthy bland blandish blaspheme blasphemous blasphemy blasted blatant blatantly blather bleak			
bleakly bleakness bleed bleeding bleeds blemish blind blinding blindingly blindside blister			
blister blistering bloated blockage blockhead bloodshed bloodthirsty bloody blotchy blow blunder			
blundering blunders blunt blur bluring blurred blurring blurry blurs blurt			
boastful boggle bogus boil boiling boisterous bomb bombard bombardment bombastic bondage			
bonkers bore bored boredom bores boring botch bother bothered bothering			
bothers bothersome bowdlerize boycott braggart bragger brainless brainwash brash brashly brashness			
brat bravado brazen brazenly brazenness breach break break-up break-ups breakdown			
breaking breaks breakup breakups bribery brimstone bristle brittle broke broken broken-hearted			
brood browbeat bruise bruised bruises bruising brusque brutal brutalising brutalities			
brutality brutalize brutalizing brutally brute brutish bs buckle bug bugging buggy			
bugs bulkier bulkiness bulky bulkyness bull**** bull bullies bullshit bullshyt bully			
bullying bullyingly bum bump bumped bumping bumpping bumpy bumps bumpy			
bungler bungling bunk burden burdensome burdensomely burn burned burning burns bust			
busts busybody butcher butchery buzzing byzantine cackle calamities calamitous calamitously calamity			
callous calumniate calumniation calumnies calumnious calumniously calumny cancer cancerous cannibal cannibalize			
capitulate capricious capriciously capriciousness capsize careless carelessness caricature carnage carp			
cartoonish cash-strapped castigate castrated casualty cataclysm cataclysmic cataclysmically catastrophe catastrophes			
catastrophies catastrophic catastrophically catastrophies caustic caustically cautionary cave censure chafe chaff			
chagrin challenging chaos chaotic chasten chastise chastisement chatter chatter chatterbox cheap			
cheapen cheaply cheat cheated cheater cheating cheats checkered cheerless cheesy chide			
childish chill chilly chintzy choke choleric choppy chore chronic chunky			
clamor clamorous clash cliche cliched clique clog clogged clogs cloud clouding			
cloudy clueless clumsy clunky coarse cocky coerce coercion coercive cold coldly			
coldly collapse collude collusion combative combust comical commiserate commonplace commotion commotions			
complacent complain complained complaining complains complaint complaints complex complicated complication complicit			
compulsion compulsive concede conceded conceit conceited concens concern concern concerns			
concerns concession concessions condemn condemnable condemnation condemned condemns condescend condescending condescendingly			
condescendingly condescension confess confession confessions confined conflict conflicted conflicting conflicts confound			
confounded confounding confront confrontation confrontational confuse confused confuses confusing confusion confusions			
congested congestion cons conscons conservative conspicuous conspicuously conspiracies conspiracy conspirator conspirator			
conspiratorial conspire consternation contagious contaminate contaminated contaminates contaminating contamination contempt contemptible			
contrariness contravene contrive contrived controversial controversy convoluted corrode corrosion corrosions corrosive			
corrupt corrupted corrupting corruption corrupts corruptted costlier costly counter-productive counterproductive coupists			
coupists covetous coward cowardly crabby crack cracked cracks craftily craftly crafty			
cramp cramped cramping cranky crap crappy craps crash crashed crashes crashing			
crashing crass craven cravenly craze crazily craziness crazy creak creaking creaks			
creaks credulous creep creeping creeps creepy crept crime criminal cringe cringed cringes			
cringes cripple crippled cripples crippling crisis critic critical criticism criticisms criticize			
cruel crueler cruelest cruelly cruelness cruelties cruelty crumble crumbling crummy crumple			
crumple crumpled crumples crush crushed crushing cry culpable culprit cumbersome cunt			
cunt cunts cuplrit curse cursed curses curt cuss cust cussed cutthroat cynical			
damned damning damper danger dangerous dangerousness dark darken darkened darker darkness			
darkness dastard dastardly daunt daunting dauntingly dawdle daze dazed dead deadbeat			
deadlock deadly deadweight deaf dearth death debacle debase debasement debaser debatable			
debaucher debauchery debilitate debilitating debility debt debts decadence decadent decay			
deception deceptive deceptively declaim decline declines declining decrement decrepit			
decrepit decrepitude decry			

defensive
defiance
defiant
defiantly
deficiencies
deficiency
deficient
defile

degeneration degradation degrade degrading degradingly dehumanization dehumanize deign deject dejected dejectedly dejection
delay delayed delaying delays delinquency delinquent delirious delirium delude delude deluded deluded deluge delusion
delusions demean demeaning demise demolish demolisher demon demonic demonize demonized demonizes
<pre>demonizing demoralize demoralizing demoralizingly denial denied denies denigrate denounce denounce dense dent dented</pre>
denunciation denunciations deny denying deplete deplorable deplorably deplore deploring deploring deplorings deplorings deplorings deplorably
depraved depravedly deprecate depress depressed depressing depressing depressing depressions depression depressions deprive deprived deride derision
derisive derisively derisiveness derogatory desecrate desert desertion desiccate desiccated desititute desolate desolate
despair despairing despairingly desperate desperately desperation despicable despicably despise despised despised
despoiler despondence despondency despondent despondently despot despotic despotism destabilisation destains destitute destitute destitute
destroyer destruction destructive desultory deter deteriorate deteriorating deterioration deterrent deterstable
<pre>detestably detested detesting detests detract detracted detracting detraction detracts detriment detrimental devastate</pre>
devastates devastating devastatingly devastation deviate deviation devil devilish devilishly devilishly devilment devilry devious
deviousness devoid diabolic diabolical diabolically diametrically diappointed diatribe diatribes dick dick dictator
dictatorial die die-hard died dies difficult difficulties difficulty diffidence dilapidated dilemma dilly-dally
<pre>dim dimmer din ding dings dinky dire direly direness dirt dirtbag dirtbags dirts</pre>
dirty disable disabled disaccord disadvantage disadvantaged disadvantageous disadvantages disaffect disaffected disafferm disagree
disagreeably disagreeing disagreement disagrees disallow disapointed disapointment disapointed disapointed disapointed disapointed disapointed disapointed disapointed
disappointing disappointment disappointments disappoints disapprobation disapproval disapprove disapproving disarm disarray disarter
disastrous disastrously disavow disavowal disbelief disbelieve disbeliever disclaim discombobulate discomfit discomfit
discompose disconcerted disconcerting disconcertingly disconsolate disconsolately disconsolately disconsolation discontent discontented discontented discontentedly discontinued
discontinuous discord discordance discordant discountenance discourage discourage discouragement discouraging discouraging discouragingly discourteous discourteous discourteous discourteous
discredit discrepant discriminate discrimination discriminatory disdain disdained disdainful disdainfully disfavor disgrace disgrace
disgracefully disgruntle disgruntled disgust disgusted disgustful disgustfully disgusting disgustingly dishearten
dishearteningly dishonest dishonestly dishonorstly dishonor dishonorable dishonorablely disillusion disillusioned disillusionment disillusions disinclination
disinclined disingenuous disingenuously disintegrate disintegrates disintegrates disintegration disinterest disinterest disinterested dislike dislikes dislikes dislikers
disliking dislocated disloyal disloyalty dismal dismally dismalness dismay dismayed dismaying dismaying dismayingly dismissive dismissively
disobedient disobey disoobedient disorder disordered disorderly disorganized disorient disorient disoriented disoriented disown disparage
disparagingly dispensable dispirit dispirited dispiritedly dispiriting displace displace displace displaced displease displease displeased displeasend displeasing displeasure
disproportionate disprove disputable dispute disputed disquiet disquieting disquietingly disquietude disregard disregardful disreputable
disrepute disrespect disrespectable disrespectablity disrespectful disrespectfully disrespectfulness disrespecting disrupt disrupt disruption disruptive diss dissapointed
dissappointing dissatisfaction dissatisfied dissatisfies dissatisfy dissatisfy dissatisfy dissatisfy dissatisfy dissatisfying dissed dissemble dissemble dissembler dissension
dissenter dissention disservice disses dissidence dissident dissidents dissidents dissioute dissolute dissolute dissolution dissonance
dissonant dissonantly dissuade dissuasive distains distaste distasteful distastefully distort distorted distortion distorts
distracting distraction distraught distraughtly distraughtness distress distressed distressing distressing distressingly distrust
distrusting disturb disturbance disturbed disturbing disturbingly disunity disvalue divergent divisive divisively divisiveness
dizzingly dizzy doddering dodgey dogged doggedly dogmatic doldrums domineer domineering donside doom
doomed doomsday dope doubt doubtful doubtfully doubts douchbag douchebags douchebags downbeat downcast
downfall downfallen downhearted downheartedly downhill downside downsides downturn downturns drab
draconic drag dragged dragging dragoon drags drain drained draining drains drastic drastically
drawback drawbacks dread dreadful dreadfully dreadfulness dreary dripped dripping drippy drips drones
droops drop-out drop-outs dropout dropouts dropouts drought drowning drunk drunkard drunken dubious
<pre>dubiously dubitable dud dull dullard dumb dumbfound dump dumped dumping dumpss dunce dungeon</pre>
<pre>dungeons dupe dust dusty dwindling dying earsplitting eccentric eccentricity effigy effrontery egocentric</pre>
egotism egotistical egotistically egregious egregiously election-rigger elimination emaciated emasculate embarrass embarrassing
<pre>embarrassment embattled embroil embroiled embroilment emergency emphatic emphatic emphatically emptiness encroach encroach encroachment endanger</pre>
enemies enemy enervate enfeeble enflame engulf enjoin enmity enrage enraged enraging enslave entangle
entangle entanglement entrap entrapment envious enviously enviousness epidemic equivocal erase erode erode erodes erosion err
evade evasion evasive evil evildoer evils eviscerate exacerbate exagerate exagerate exagerated exagerates exaggerate exaggerate exaggerate exaggerate
exaggeration exasperate exasperated exasperating exasperatingly exasperation excessive excessively exclusion excoriate excruciating excruciatingly excruciatingly excuse
excuses execrate exhaust exhausted exhaustion exhausts exhorbitant exhort exile exorbitant
expensive expire expired explode exploit exploitation explosive expropriate expropriate expropriation expulse expunge exterminate extermination
exterminate extermination extinguish extort extortion extraneous extravagance extravagant extravagant extravagantly extremism extremist extremist extremist extremists eyesore f**k
<pre>f**k fabricate fabrication facetious facetiously fail failed failing fails failure failures faint fainthearted</pre>
<pre>fainthearted faithless fake fall fallacies fallacious fallaciously fallaciousness fallacy fallen falling fallout falls</pre>
<pre>false falsehood falsely falsify falter faltered famine famished famatic fanatical</pre>
<pre>fanatically fanaticism fanatics fanciful far-fetched farce farcical farcical-yet-provocative farcically farfetched fascism</pre>
<pre>fanaticism fanatics fanciful far-fetched farce farcical farcical-yet-provocative farcically farfetched fascism fascism fascist fastidious fastidiously fastidiously fastidiously fattoat fat-cat fat-cat</pre>
fanaticism fanatics fanciful far-fetched farce farcical farcical-yet-provocative farcically farfetched fascism fascism fascist fastidiously fastidously fastuous fat fat-cat fat-cat fat-cats fatal fatalistic fatalistic fatalistic fatalistic fatalistic fatalistically fatally fatally fatally fatact fatcats fatect fatcats fatect fatcats fatect fat
famelicism famelicis famelicis famelicis famelicis far-fetched farce famelicis famelic
Amenics fancius fancius fancius farcical factor fascis fascis fascis fascis fascis fatcar fatca
Family Lise Family
Search Colombia Search Colombi
Case tribs Dar-feet Cook Dar-f
Grant of the Control
Section  Exercition  Exercitio
SAME AND
March   Marc
March   Marc
March   Marc

frustrated frustrates frustrating frustratingly frustration frustrations

fuck fucking fudge

defiler
deform
deformed
defrauding
defunct
defy
degenerate

fumble fume fumes fundamentalism funky funnily funny furious furiously furor fury			
fuss fussy fustigate fusty futile futilely futility fuzzy gabble gaff			
gaffe gainsay gainsayer gall galling gallingly galls gangster gape garbage garish			
gasp gauche gaudy gawk gawky geezer genocide get-rich ghastly ghetto ghosting			
gibber gibberish gibe giddy gimmick gimmicked gimmicking gimmicks gimmicky glare glaringly			
glaringly glib glibly glitch glitches gloatingly gloom gloomy glower glum glut			
gnawing goad goading god-awful goof goofy goon gossip graceless gracelessly graft			
grait grainy grapple grate grating gravely greasy greed greedy grief grievance			
grievances grieve grieving grievous grievously grim grimace grind gripe gripes			
grisly gritty gross grossly grotesque grouch grouchy groundless grouse growl			
grudge grudges grudging grudgingly gruesome gruesomely gruff grumble grumpier grumpiest			
grumpily grumpish grumpy guile guilt guiltily guilty gullible gutless gutter			
hack hacks haggard haggle hairloss halfhearted halfheartedly hallucinate hallucination hamper			
hampered handicapped hang hangs haphazard hapless harangue harass harassed harasses			
harassment harboring harbors hard hard-hit hard-line hard-liner hardball harden hardened			
hardheaded hardhearted hardliner hardliners hardship hardships harm harmed harmful harms			
harpy harridan harried harrow harsh harshly hasseling hassle hassled hassles			
haste hastily hasty hate hated hateful hatefully hatefulness hater haters			
hates hating hatred haughtily haughty haunt haunting havoc hawkish haywire			
hazard hazardous haze hazy head-aches headache headaches heartbreaker heartbreaking heartbreakingly			
heartless heathen heavy-handed heavyhearted heck heckle heckled heckles hectic hedge			
hedonistic heedless hefty hegemonism hegemonistic hegemony heinous hell hell-bent hellion			
hells helpless helplessly helplessness heresy heretic heretical hesitant hestitant hideous			
hideously hideousness high-priced hiliarious hinder hindrance hiss hissed hissing ho-hum			
hoard hoax hobble hogs hollow hoodium hoodwink hooligan hopeless hopelessly			
hopelessly hopelessness horde horrendous horrendously horrible horrid horrific horrified horrifies horrify			
horrifying horrifys hostage hostile hostilities hostility hotbeds hothead hotheaded hothouse			
hothouse hubris huckster hum humid humiliate humiliating humiliation humming hung hurt hurted			
hurted hurtful hurting hurts hustler hype hypocricy hypocrisy hypocrite hypocrites hypocritical			
hypocritical hypocritically hysteria hysteric hysterical hysterically hysterics idiocies idiocy idiot			
idiot idiotic idiotically idiots idle ignoble ignominious ignominiously ignominy ignorance ignorant			
ignore ill-advised ill-conceived ill-defined ill-designed ill-fated ill-favored ill-formed ill-mannered ill-natured			
<pre>ill-sorted ill-tempered ill-treated ill-treatment ill-usage ill-used illegal illegally illegitimate illicit</pre>			
<pre>illiterate illness illogic illogical illogically illusion illusions illusory imaginary imbalance</pre>			
imbecile imbroglio immaterial immature imminence imminently immobilized immoderate immoderately immodest			
immoral immorality immorally immovable impair impaired impasse impatience impatient impatiently			
impactencry impeach impedance impede impediment impending impenitent imperfect imperfect imperfection imperfections imperfectly			
imperiectly imperialist imperil imperious imperiously impermissible impersonal impertinent impetuous impetuous impetuously impiety			
implety impinge impious implacable implausible implausibly implicate implication implode impolite impolitely			
impolitic impolitic importunate importune impose imposers imposing imposition impossible impossiblity impossibly			
<pre>impotent impoverish impoverished impractical imprecate imprecise imprecisely imprecision imprison</pre>			
<pre>imprisonment improbability improbable improbably improper improperly impropriety imprudence imprudent impudence</pre>			
<pre>impudent impudently impugn impulsive impulsively impunity impure impurity inability inaccuracies inaccuracy</pre>			
inaccurate inaccurately inaction inactive inadequacy inadequate inadequately inadverent inadverently inadvisable			
<pre>inadvisably inane inanely inappropriate inappropriately inapt inaptitude inarticulate inattentive inaudible</pre>			
incapable incapably incautious incendiary incense incessant incessantly incite incitement incivility			
inclement incognizant incoherence incoherent incoherently incommensurate incomparable incomparably incompatability incompatibility			
incompatible incompetence incompetent incompetently incomplete incompliant incomprehensible incomprehension inconceivable inconceivably			
incongruous incongruously inconsequent inconsequential inconsequentially inconsequently inconsiderate inconsiderately inconsistence inconsistencies			
inconsistency inconsistent inconsolable inconsolably inconstant inconvenience inconveniently incorrect incorrectly incorrigible incorrigibly			
incredulous incredulously inculcate indecency indecent indecently indecision indecisive indecisively indecorum			
indefensible indelicate indeterminable indeterminate indifference indifferent indigent indignant indignantly			
indignation indignity indiscernible indiscreet indiscreetly indiscretion indiscriminate indiscriminate indiscriminately indiscriminating indistinguishable			
indoctrinate indoctrination indolent indulge ineffective ineffectively ineffectiveness ineffectual ineffectually ineffectualness			
<pre>inefficacious inefficacy inefficiency inefficient inefficiently inelegance inelegant ineligible ineloquent ineloquently</pre>			
<pre>inept ineptitude ineptly inequalities inequality inequitable inequitably inequities inescapable inescapably</pre>			
inessential inevitable inevitably inexcusable inexcusably inexorable inexorably inexperience inexperience inexpert inexpertly			
<pre>inexpertly inexpiable inexplainable inextricable inextricably infamous infamously infamy infected infection infections</pre>			
<pre>infections inferior inferiority infernal infest infested infidel infidels infiltrator infiltrators infiltrators</pre>			
<pre>inflame inflammation inflammatory inflammed inflated inflationary inflexible inflict infraction infringe</pre>			
<pre>infringement infringements infuriate infuriated infuriating infuriatingly inglorious ingrate ingratitude inhibit</pre>			
inhibition inhospitable inhospitality inhuman inhumane inhumanity inimical inimically iniquitous iniquity injudicious			
injudicious injure injurious injury injustice injustices innuendo inoperable inopportune inordinate inordinately			
inordinately insane insanely insanity insatiable insecure insecurity insensible insensitive insensitively insensitivity insidious			
instigator instigators insubordinate insubstantial insubstantially insufferable insufferably insufficiency insufficient insufficiently insufficiently			
insular insult insulted insulting insultingly insults insupportable insupportably insurmountable insurmountably insurrection			
insurrection intefere inteferes intense interfere interferece interferes intermittent interrupt interruption interruptions			
interruptions intimidate intimidating intimidatingly intimidation intolerable intolerablely intolerance intoxicate intractable intransigence			
intransigence intransigent intrude intrusion intrusive inundate inundated invader invalid invalidate invalidity invasive			
invasive invective inveigle invidious invidiously invidiousness invisible involuntarily involuntary			
<pre>irate irately ire irk irked irking irks irksome irksomey</pre>			
irksomeness irksomenesses ironic ironical ironically ironies irony irragularity irrational irrationalities irrationality			
<pre>itchy jabber jaded jagged jam jarring jaundiced jealous jealously jealousness</pre>			
<pre>jealousness jealousy jeer jeering jeeringly jeers jeopardize jeopardy jerk jerky</pre>			
judders jumpy junk junky junkyard jutter jutters kaput kill killed killer			
killer killing killjoy kills knave knife knock knotted kook kooky			
lack lackadaisical lacked lackey lackeys lacking lackluster lacks laconic lag			
lag lagged lagging laggy lags laid-off lambast lambaste lame lame-duck lament			
lament lamentable lamentably languid languish languor languorous languorously lanky lapse lapsed			
lapses lascivious last-ditch latency laughable laughably laughingstock			
lawbreaker lawbreaking lawless lawlessness			
lawbreaking lawless			

leech leer

lemon lengthy

leery
left-leaning

fugitive
full-blown
fulminate

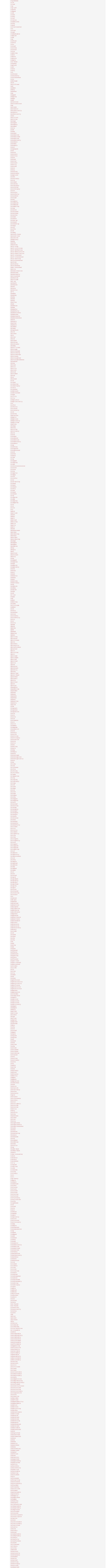
lethargy lewd lewdly lewdness liability liable liar liars licentious licentiously licentiousness		
lied lier lies life-threatening lifeless limit limitation limitations limited limits limp listless		
litigious little-known livid lividly loath loathe loathing loathly loathsome loathsomely lone loneliness		
lonely loner lonesome long-time long-winded longing longingly loophole loopholes loose loot		
lorn lose loser losers loses losing loss losses lost loud louder		
lousy loveless lovelorn low-rated lowly ludicrous ludicrously lugubrious lukewarm lull lumpy lunatic		
lunaticism lurch lure lurid lurk lurking lying macabre mad madden maddening		
maddeningly madder madly madman madness maladjusted maladjustment malady malaise malcontent malcontented		
maledict malevolence malevolent malevolently malice malicious maliciously maliciousness malign malignant malodorous maltreatment		
mangle mangled mangling mania maniac maniacal manic manipulate manipulation manipulative		
manipulative manipulators mar marginal marginally martyrdom martyrdom-seeking mashed massacre massacres matte mawkish		
mawkishly mawkishness meager meaningless meanness measly meddle meddlesome mediocre mediocrity melancholy		
melodramatic melodramatically meltdown menace menacing menacingly mendacious mendacity menial merciless mercilessly		
mess messed messes messing messy midget miff militancy mindless mindlessly mirage mire		
misalign misaligned misaligns misapprehend misbecome misbecoming misbegotten misbehave misbehavior miscalculate miscalculation		
miscellaneous mischief mischievous mischievously misconception misconceptions miscreant miscreants misdirection miser miserable		
miserableness miserably miseries miserly misery misfit misfortune misgiving misgivings misguidance misguide		
misguided mishandle mishap misinform misinformed misinterpret misjudge misjudgment mislead misleading misleadingly mislike		
mismanage mispronounce mispronounced mispronounces misread misreading misrepresent misrepresentation miss missed misses		
misstatement mist mistake mistaken mistakenly mistakes mistified mistress mistrust mistrustful mistrustfully		
mists misunderstand misunderstanding misunderstandings misunderstood misuse moan mobster mock mocked mocked mockeries mockery		
mocking mockingly mocks molest molestation monotonous monotony monster monstrosities monstrosity monstrous		
monstrously moody moot mope morbid morbidly mordant mordantly moribund moron moronic		
morons mortification mortified mortify mortifying motionless motley mourn mourner mournful mournfully		
muddle muddy mudslinger mudslinging mulish multi-polarization mundane murder murderer murderous murderously murky		
murky muscle-flexing mushy musty mysterious mysteriously mystery mystify myth nag nagging naive		
naively narrower nastily nastiness nasty naughty nauseate nauseates nauseating nauseatingly naïve nebulous		
nebulous nebulously needless needlessly needy nefarious nefariously negate negation negative negatives negativity		
negativity neglect neglected negligence negligent nemesis nepotism nervous nervously nervousness nettle nettlesome		
nettlesome neurotic neurotically niggle niggles nightmare nightmarish nightmarishly nitpick nitpicking noise noises		
noisier noisy non-confidence nonexistent nonresponsive nonsense nosey notoriety notorious notoriously noxious		
nuisance numb obese object objection objectionable objections oblique obliterate obliterated oblivious		
obnoxious obnoxiously obscene obscenely obscenity obscure obscured obscures obscurity obsess obsessive		
obsessively obsessiveness obsolete obstacle obstinate obstinately obstruct obstructed obstructing obstruction obstructs obtructs		
obtuse occlude occluded occludes occluding odd odder oddest oddities oddity oddly		
odor offence offend offender offending offenses offensive offensively offensiveness officious ominous		
ominously omission omit one-sided onerous onerously onslaught opinionated opponent opportunistic oppose		
opposition oppositions oppress oppression oppressive oppressively oppressiveness oppressors ordeal orphan ostracize outbreak		
outburst outcast outcry outlaw outmoded outrage outraged outrageous outrageously outrageousness		
outrages outsider over-acted over-awe over-balanced over-hyped over-priced over-valuation overact overacted overawe		
overbalance overbalanced overbearing overbearingly overblown overdo overdone overdue overemphasize overheat overkill overloaded		
overloaded overloaded overpaid overpayed overplay overpower overpriced overrated overreach overrun overshadow oversight		
oversight oversights oversimplification oversimplified oversimplify oversize overstate overstated overstatement overstatements overstates overstates		
overthrow overthrows overturn overweight overwhelm overwhelmed overwhelming overwhelmingly overwhelms overzealous overzealously		
overzelous pain painful painfull painfully pains pale pales paltry pan pandemonium pander		
pandering panders panic panick panicked panicking panicky paradoxical paradoxically paralize paralyzed		
paralyzed paranoia paranoid parasite pariah parody partiality partisan partisans passe passive passive		
pathetic pathetically patronize paucity pauper paupers payback peculiar peculiarly pedantic peeled		
peeve peeved peevish peevishly penalize penalty perfidious perfidity perfunctory peril perilous perilous perilous		
periously perish pernicious perplex perplexed perplexing perplexity persecute persecution pertinacious pertinaciously pertinacity		
perturb perturbed pervasive perverse perversely perversion perversity pervert pervert perverted perverts pessimism		
pessimism pessimistic pessimistically pest pestilent petrified petrify pettifog petty phobia phobic phony		
picket picketed picketing pickets picky pig pigs pillage pillory pimple pinch pique		
pique pitiable pitiful pitifully pitiless pitilessly pittance pity plagiarize plague plasticky plaything		
plea pleas plebeian plight plot plotters ploy plunder plunderer pointless poison		
poison poisonous poisonously pokey poky polarisation polemize pollute polluter polluters polution pompous		
poor poorer poorest poorly posturing pout poverty powerless prate pratfall prattle precarious		
precarious precariously precipitate precipitous predatory predicament prejudge prejudice prejudices prejudicial premeditated preoccupy		
preoccupy preposterous preposterously presumptuous presumptuously pretence pretend pretense pretentious pretentious prevaricate pricey pricier		
prick prickle prickles prideful prik primitive prison prisoner problem problematic		
problems procrastinate procrastinates procrastination profane profanity prohibit prohibitive prohibitively propaganda propagandize		
proprietary prosecute protest protested protesting protests protracted provocation provocative provoke pry pugnacious		
pugnaciously pugnacity punch punish punishable punitive punk puny puppet puppets puzzled		
puzzlement puzzling quack qualm qualms quandary quarrel quarrellous quarrellously quarrels quarrelsome		
quash queer questionable quibble quibbles quitter rabid racism racist racists racy radical		
radical radicalization radically radicals rage raged raging rail raked rampage rampant ramshackle		
ramshackle rancor randomly rankle rant ranted ranting rantingly rants rape raped raping rascal		
raping rascal rascals rash rattle rattled rattles ravage raving reactionary rebellious rebuff		
rebuff rebuke recalcitrant recant recession recessionary reckless recklessly recklessness recoil recourses redundancy redundant		
redundant refusal refuse refused refuses refusing refutation refute refuted refutes refuting refuting regress		
regress regression regressive regret regreted regretful regretfully regrets regrettable regrettably regretted reject		
reject rejected rejecting rejection rejects relapse relentless relentlessly relentlessness reluctance reluctant reluctantly		
reluctantly remorse remorseful remorsefully remorseless remorselessly remorselessness renounce renunciation repel repetitive reprehensible		
reprehensible reprehensibly reprehension reprehensive repress repression repressive reprimand reproach		
reproachful reprove reprovingly repudiate		
reprove reprovingly repudiate repudiation repugn repugnance repugnant repugnantly repulse repulsed repulsing repulsive repulsively		
reprove reprovingly repudiate repudiation repugn repugnance repugnant repugnantly repulse repulsed repulsing repulsive repulsively repulsiveness resent resentful resentment resignation resigned resistance restless restrict restricted restriction		
reprove reprovingly repudiate repudiation repugn repugnance repugnantly repulse repulsed repulsing repulsive repulsively repulsiveness resent resentful resentment resignation resigned resistance restless restrict restricted		
reprove reprovingly repudiate repudiation repugn repugnance repugnant repugnantly repulse repulsed repulsing repulsive repulsively repulsiveness resent resentful resentment resignation resigned resistance restless restless restrict restricted restriction restrictive resurgent retaliate retaliatory retard retarded retarded retract retract retreat		

ridicules ridiculous ridiculously

rife rift rifts rigid rigidity

less-developed
lesser-known

letch
lethal
lethargic



	unspeakablely unspecified unstable unsteadily unsteadiness unsteady unsuccessful
	<pre>unsuccessfully unsupported unsupportive unsure unsuspecting unsustainable untenable untested</pre>
	unthinkable unthinkably untimely untouched untrue untrustworthy untruthful unusable
	unuseable unuseably unusual unusually unviewable unwanted unwarranted
	<pre>unwatchable unwelcome unwell unwieldy unwilling unwillingsy unwillingness unwise</pre>
	unwisely unworkable unworthy unyielding upbraid upheaval uprising uproar
	uproarious uproariously uproarous uproarously uproot upset upset
	<pre>upsets upsetting upsettingly urgent useless usurp usurper utterly</pre>
	<pre>vagrant vague vagueness vain vainly vanity vehement vehemently</pre>
	<pre>vengeance vengeful vengefully vengefulness venom venomous venomous venomously</pre>
	vestiges vex vexation vexing vexingly vibrate vibrated vibrates
	<pre>vibrating vibration vice vicious viciously viciousness victimize</pre>
	<pre>vile vileness vilify villainous villainously villains villian</pre>
	villianously villify vindictive vindictively vindictiveness violate violation
	<pre>violator violators violent violently viper virulence virulent</pre>
	<pre>virulently virus vociferous vociferously volatile volatility vomit vomited</pre>
	<pre>vomiting vomits vulgar vulnerable wack wail wallow wane</pre>
	<pre>waning wanton war-like warily wariness warlike warned</pre>
	<pre>warning warp warped wary washed-out waste wasted wasteful</pre>
	wastefulness wasting water-down watered-down wayward weak weaken
	<pre>weakening weaker weakness weaknesses weariness wearisome weary wedge</pre>
	weed weep weird weirdly wheedle whimper whine
	whining whiny whips whore whores wicked wickedly wickedness
	<pre>wild wildly wiles wilt wily wimpy wimpy</pre>
	wobbled wobbles wobegone woeful woefully
	womanizing worn worried worriedly worrier worries worrisome
	worry worrying worryingly worse worsen worsening worst worthless
	worthlessly worthlessness wound wounds wrangle wrath wreak
	wreaked wreaks wreck wrest wrestle wretch wretched
	<pre>wretchedly wretchedness wrinkle wrinkled wrinkles wrip wripped wripping</pre>
	writhe wrong wrongful wrongly wrought yawn zap
	<pre>zapped zaps zealot zealous zealously zombie """</pre>
	<pre>text_list = text_blob.split('\n') negative_words = ', '.join(text_list)  print(negative_words)  2-faced 2-faces abnormal abolish abominable abominably abominate abomination abort aborts aborts aborts aborts aborts aborts aborts aborts aborts aborts.</pre>
	2-faced, 2-faces, abnormal, abolish, abominable, abominably, abominate, abomination, abort, aborted, aborts, ab rade, abrasive, abrupt, abruptly, abscond, absence, absent-minded, absentee, absurd, absurdity, absurdly, absurdness, abuse, abused, abuses, abusive, abysmal, abysmally, abyss, accidental, accost, accursed, accusation, accusations, accuse, accuses, accusing, accusingly, acerbate, acerbic, acerbically, ache, ached, aches, achey, ach ing, acrid, acridly, acridness, acrimonious, acrimoniously, acrimony, adamant, adamantly, addict, addicted, addicting, addicts, admonish, admonisher, admonishingly, admonishment, admonition, adulterate, adulterated, adulteration, adulterier, adversarial, adversary, adverse, adversity, afflict, affliction, afflictive, affront, afraid, aggravate, aggravating, aggravation, aggression, aggressive, aggressiveness, aggressor, aggrieve, aggrieved,
	kier, spockiest, spocky, spochy, spoch, spoch, spoth, spoch, spothy, sporadic, spotty, spurious, spurn, sputter, squabbling, squander, stagnation, stasid, staid, stai
	nleash, unlicensed, unlikely, unlucky, unmoved, unnatural, unnaturally, unnecessary, unneeded, unnerve, unnerve
	violently, viper, virulence, virulent, virulently, virus, vociferous, vociferously, volatile, volatility, vomit, vomited, vomiting, vomits, vulgar, vulnerable, wack, wail, wallow, wane, waning, wanton, war-like, warily, wariness, warlike, warned, warning, warp, warped, wary, washed-out, waste, wasted, wasteful, wastefulness, wasting, water-down, watered-down, wayward, weak, weaken, weakening, weaker, weakness, weaknesses, weariness, wearis ome, weary, wedge, weed, weep, weird, weirdly, wheedle, whimper, whine, whining, whiny, whips, whore, whores, wicked, wickedly, wickedness, wild, wildly, wiles, wilt, wily, wimpy, wince, wobble, wobbled, wobbles, woe, woeb egone, woeful, woefully, womanizer, womanizing, worn, worried, worriedly, worrier, worries, worrisome, worry, worrying, worryingly, worse, worsen, worsening, worst, worthless, worthlessly, worthlessness, wound, wounds, wranted.
In [71]:	<pre>orrying, worryingly, worse, worsen, worsening, worst, worthless, worthlessly, worthlessness, wound, wounds, wrangle, wrath, wreak, wreaked, wreaks, wreck, wrest, wrestle, wretch, wretched, wretchedly, wretchedness, wrinkle, wrinkled, wrinkles, wrip, wripped, wripping, writhe, wrong, wrongful, wrongly, wrought, yawn, zap, zapped, zaps, zealot, zealous, zealously, zombie,  # Function to clean text and remove stopwords def clean_text(text):     stop_words = set(stopwords.words("english"))</pre>
In [72]:	<pre>tokens = nltk.word_tokenize(text)   clean_tokens = [word for word in tokens if word.lower() not in stop_words]   clean_text = " ".join(clean_tokens)   return clean_text  # Function to create a dictionary of positive and negative words def create_sentiment_dict(text, positive_words, negative_words):</pre>
	<pre>positive_dict = {} negative_dict = {} negative_dict = {} tokens = nltk.word_tokenize(text) for word in tokens:     if word in positive_words:         positive_dict[word] = positive_dict.get(word, 0) + 1     elif word in negative_words:         negative_dict[word] = negative_dict.get(word, 0) + 1</pre>
In [73]:	<pre>negative_dict[word] = negative_dict.get(word, 0) + 1 return positive_dict, negative_dict</pre>
In [74]: In [75]:	<pre># Function to calculate the negative score def calculate_negative_score (negative_dict):     negative_score = -sum(negative_dict.values())     return negative_score  # Function to calculate the polarity score</pre>
In [75]: In [76]:	<pre>def calculate_polarity_score (positive_score, negative_score):     polarity_score = (positive_score - negative_score) / (positive_score + negative_score + 0.000001)     return polarity_score  # Function to calculate the subjectivity score def calculate_subjectivity_score (positive_score, negative_score, total_words):</pre>
In [76]: In [77]:	<pre>def calculate_subjectivity_score (positive_score, negative_score, total_words):     subjectivity_score = (positive_score + negative_score) / (total_words + 0.000001)     return subjectivity_score  # Function to calculate average sentence length def calculate_avg_sentence_length(text):     sentences = nltk.sent_tokenize(text)     words = nltk.word_tokenize(text)</pre>
In [78]:	<pre>sentences = nltk.sent_tokenize(text) words = nltk.word_tokenize(text) avg_sentence_length = len(words) / len(sentences) return avg_sentence_length  # Function to calculate the percentage of complex words def calculate_percentage_complex_words(text):     words = nltk.word_tokenize(text)</pre>
In [79]:	<pre>words = nltk.word_tokenize(text) complex_word_count = sum(1 for word in words if len(word) &gt; 2) percentage_complex_words = (complex_word_count / len(words)) * 100 return percentage_complex_words</pre>
In [80]:	<pre>fog_index = 0.4 * (avg_sentence_length + percentage_complex_words)     return fog_index  # Function to calculate the average number of words per sentence def calculate_avg_words_per_sentence(text):     sentences = nltk.sent_tokenize(text)     words = nltk.word_tokenize(text)     avg_words_per_sentence = len(words) / len(sentences)</pre>
In [81]:	<pre>avg_words_per_sentence = len(words) / len(sentences) return avg_words_per_sentence</pre>
In [82]:	<pre>return complex_word_count  # Function to count total words  def count_total_words(text):     words = nltk.word_tokenize(text)     total_words = len(words)     return total_words</pre>
In [83]:	<pre># Function to calculate syllables per word  def calculate_syllables_per_word(text):     d = nltk.corpus.cmudict.dict()      def count_syllables(word):         if word.lower() in d:             return max([len(list(y for y in x if y[-1].isdigit()))) for x in d[word.lower()]])         else:</pre>
In [84]:	<pre>else:     return 0  words = nltk.word_tokenize(text) syllables_per_word = [count_syllables(word) for word in words] return syllables_per_word</pre>
In [84]: In [85]:	<pre>def count_personal_pronouns(text):     pronoun_pattern = r'\b(I we my ours us)\b'     personal_pronouns_count = len(re.findall(pronoun_pattern, text, flags=re.IGNORECASE))     return personal_pronouns_count</pre>
In [85]:	<pre>def calculate_avg_word_length(text):     words = nltk.word_tokenize(text)     total_characters = sum(len(word) for word in words)     avg_word_length = total_characters / len(words)     return avg_word_length</pre>
٠1:	<pre>def scrape_text_from_url(url):     try:     response = requests.get(url)  if response.status_code == 200:</pre>
	<pre>soup = BeautifulSoup(response.text, 'html.parser')  body_text = ' '.join([p.text for p in soup.find_all('p')])  return body_text else:</pre>
In [87]:	<pre>else:         print(f"Failed to retrieve content from {url}. Status code: {response.status_code}") except Exception as e:         print(f"An error occurred: {str(e)}") return None  # Function to scrape text from a column of URLs in a DataFrame def scrape_text_from_column(dataframe, url_column_name):</pre>
	<pre>def scrape_text_from_column(dataframe, url_column_name):     dataframe['Scraped_Text'] = dataframe[url_column_name].apply(scrape_text_from_url)</pre>

unorthodox unorthodoxy unpleasant unpleasantries unpopular unpredictable unprepared unproductive unprofitable unprove

unprove unproved unproven unproving unqualified unravel

unraveled
unreachable
unreadable
unrealistic
unreasonable
unrelenting
unrelentingly
unreliability
unreliable
unresolved
unresponsive

unrest unruly unsafe

unsate
unsatisfactory
unsavory
unscrupulous
unscrupulously
unsecure

unsophisticated

unseemly
unsettle
unsettled
unsettling
unsettlingly
unskilled

unsound unspeakable

1 1 1	2 http: 3 http: 4 http: 09 http: 10 http: 11 https: 12 https:	s://insights s://insights s://insights s://insights ://insights. ://insights. /insights.b	s.blackcoffer s.blackcoffer s.blackcoffer s.blackcoffer blackcoffer blackcoffer.co	com/rise-of- com/rise-of- com/rise-of- com/corona com/corona com/what-ar	e-hea Au telem Au telem Au avirus-i Au te-the Au ing-dri Au	itomate the Citomate the Citoma	Data Manageme	ent Process Rea ent Process Rea ent Process Rea ent Process Rea ent Process Rea ent Process Rea ent Process Rea	Itime Itime Itime Itime Itime Itime Itime			
2]: 7	clean posit posit negat polar total subje avg_s perce fog_i avg_w compl total sylla perso	late_al ed_text ive_dic ive_sco ive_sco ity_sco _words ctivity entence ntage_c ndex = ords_pe ex_word _word_c bles_pe nal_pro	<pre>l_variab   = clean t, negat re = cal re = cal re = cal = len(cl _score = _length omplex_w calculat r_senten _count = ount = c r_word = noun_cou</pre>	cles(text, text)  ive_dict  culate_po  culate_po  eaned_text  calculate  rords = calculate  ce = calculate  ce = calculate  count_co  count_tota  int = count  count = count  count = calculate  count_tota  count = count  count = cou	xt) = create ositive_s egative_s olarity_s xt.split( te_subjec ate_avg_s alculate_ dex(avg_s culate_av omplex_wo al_words( te_syllab	re_words,  e_sentimer core (posi- core (po	negative_went_dict(cleative_dict) ative_dict) ative_dict) itive_score core(positive_score) length(text) length, percept per_sentence word(text) uns(text)	aned_text, , negative ve_score, ) words(text centage_co	_score) negative_so	words, negaticore, total_west		
7 5]: #	#applying new_df[[' ' lambd  #examinin new_df=ne new_df.he	perce total  the fu Positiv Avg_Sen Avg_Wor Syllabl a text:  g the n  w_df.dr ad()	ntage_co _word_co  nction o e_Score' tence_Le ds_Per_S es_Per_W pd.Seri  ew datas op(['Scr	on the new 'Negation 'Nega	rds, fog_lables_pe  w dataset ive_Score Percentag , 'Comple ersonal_P late_all_  t'],axis=	index, aver_word, preserved, pres	rity_Score' x_Words', ': ount', 'Toto ount', 'Avg s(text, pos.	r_sentence onoun_coun  , 'Subject Fog_Index' al_Word_Co _Word_Leng itive_word	<pre>, complex_v t, avg_word ivity_Score , unt', th']] = new s, negative</pre>	e', w_df['Scraped	_Text'].a	apply(
S	https://inshipshttps://insh	sights.black sights.black sights.black only the sight of	kcoffer.com/ of-e-l kcoffer.com/ of-tel kcoffer.com/ of-tel	/rise- hea /rise- hea /rise- lem /rise- lem 2]	162 177 264 264	URL	-68 2.24 -69 1.70	34146 37706 37692	0.218472 0.115957 0.182927 0.182927	30.758 19.767 24.968 24.968	3621 7123 3750	
	50921. 110 51382. 111 51844. 112 52306. 113 52768.	0 http 0 http 0 http 0 http 0 http 6 https 4 https://	s://insights.k s://insights.k s://insights.k s://insights.k s://insights.k :://insights.b ://insights.b	blackcoffer.co blackcoffer.co blackcoffer.co blackcoffer.co blackcoffer.co blackcoffer.co	om/rise-of-te om/rise-of-e- om/rise-of-e- om/rise-of-te om/coronavir om/coronavir om/what-are- om/marketing	-heahea elem elem rus-i the						
n n	url_ID  123.0  321.0	the Da = pd.m : head() https://ins https://ins	taFrames erge (sub	set_df, r	SRL Positive se- m se- ea se- ea	n='URL',	how='inner egative_Score -93 -39 -68 -69		4 0. 6 0. 6 0.	y_Score Avg_Sen .202653 .218472 .115957	24.77381 30.75862 19.76712 24.96875	21
7	#dropping merged_df  URL_II  0 123.  1 321.  2 2345.  3 4321.	na val dropna  http  http	ues ()  ttps://insight  ps://insights.	blackcoffer.c blackcoffer.c	n	F-    	-69  Score Negative  368  162  177  264	1.70769. ve_Score Pol -93 -39 -68 -69		.182927 Subjectivity_Score 0.202653 0.218472 0.115957 0.182927		
	4 432 50921. 107 50921. 108 51382. 109 51844. 110 52306. 111 52768.	https:// https:// https://	/insights.bladinsights.bladinsi/ s://insights.bladinsights.bladinsights.bl	ckcoffer.com ckcoffer.com blackcoffer.co lackcoffer.col	r.com/rise-of telem.  n/coronavirus i. n/coronavirus i. om/what-are the m/marketing dri. m/continued dem.	 :- :- :- :- :- :-	264  92 254 385 275	-69  -28 -154 -87 -88	1.707692  1.875000 4.080000 1.583893 1.941176 2.718310	0.182927 0.120983 0.068634 0.215786 0.161765 0.104412		24.9687 29.6538 39.1666 28.7702 28.8095 31.9354
7	12 rows × 1	ing th	e DataFr	rame to so								