

In [40]:

```
1 import numpy as np
2 import pandas as pd
```

In [41]:

```
1 data = pd.read_csv("Kasus1.csv", delimiter=';', usecols=['text']) # read file csv nya dulu yang java heritage
2 data
```

Out[41]:

	text
0	Apa yg bisa diharapkan dari politisi gaya pela...
1	Kalau sy jadi @DennyJA_WORLD , Syahganda aka...
2	Di @PDemokrat banyak org2 yg punya akal sehat ...
3	Cypridophobia adalah fobia atau takut pada pel...
4	Yang RAMAI kerja waktu malam ni biasanya PELAC...
...	...
395	Letakan DiKepalak Sepatunya Klo Tidak Mau Koto...
396	Media pelacur lahir dari rahim yg tak jujur. h...
397	Jadi pelacur ko nanti dijawab nak. Jangan...
398	Kok jadi pelacur
399	/r/t/ no filter ini pake Samsung A50. Negara y...

400 rows × 1 columns

In [42]:

```
1 pip install Sastrawi
```

Requirement already satisfied: Sastrawi in c:\users\user\anaconda3\lib\site-packages (1.0.1)Note: you may need to restart the kernel to use updated packages.

In [43]:

```

1 import Sastrawi
2 import re
3 from Sastrawi.StopWordRemover.StopWordRemoverFactory import StopWordRemoverFactory
4 from Sastrawi.Stemmer.StemmerFactory import StemmerFactory
5
6 slangs={"@": "di", "abis": "habis", "ad": "ada", "adlh": "adalah", "afaik": "as far as i know",
7 "ahaha": "haha", "aj": "saja", "ajep-ajep": "dunia gemerlap", "ak": "saya", "akika": "aku",
8 "akkoh": "aku", "akuwh": "aku", "alay": "norak", "alow": "halo", "ambilin": "ambilkan",
9 "ancur": "hancur", "anjrit": "anjing", "anter": "antar", "ap2": "apa-apa", "apasih": "apa sih",
10 "apes": "sial", "aps": "apa", "aq": "saya", "aquwh": "aku", "asbun": "asal bunyi", "aseekk": "asyik",
11 "asekk": "asyik", "asem": "asam", "aspal": "asli tetapi palsu", "astul": "asal tulis", "ato": "atau",
12 "au ah": "tidak mau tahu", "awak": "saya", "ay": "sayang", "ayank": "sayang", "b4": "sebelum",
13 "bakalan": "akan", "bandes": "bantuan desa", "bangedh": "banget", "banpol": "bantuan polisi",
14 "banpur": "bantuan tempur", "basbang": "basi", "bcanda": "bercanda", "bdg": "bandung",
15 "begajulan": "nakal", "beliin": "belikan", "bencong": "banci", "bentar": "sebentar", "ber3": "bertiga",
16 "beresin": "membereskan", "bete": "bosan", "beud": "banget", "bg": "abang", "bgmn": "bagaimana",
17 "bgt": "banget", "bijimane": "bagaimana", "bintal": "bimbingan mental", "bkl": "akan",
18 "bknya": "bukannya", "blegug": "bodoh", "blh": "boleh", "bln": "bulan", "blum": "belum", "bnai": "benci",
19 "bnran": "yang benar", "bodor": "lucu", "bokap": "ayah", "boker": "buang air besar", "bokis": "bohong",
20 "boljug": "boleh juga", "bonek": "bocah nekat", "boye": "boleh", "br": "baru", "brg": "bareng",
21 "bro": "saudara laki-laki", "bru": "baru", "bs": "bisa", "bsen": "bosan", "bt": "buat", "btw": "ngomong-ngomong",
22 "buaya": "tidak setia", "bubbu": "tidur", "bubu": "tidur", "bumil": "ibu hamil", "bw": "bawa",
23 "bwt": "buat", "byk": "banyak", "byrin": "bayarkan", "cabal": "sabar", "cadas": "keren",
24 "calo": "makelar", "can": "belum", "capcus": "pergi", "caper": "cari perhatian", "ce": "cewek",
25 "cekal": "cegah tangkal", "cemen": "penakut", "cengengesan": "tertawa", "cepet": "cepat",
26 "cew": "cewek", "chuyunk": "sayang", "cimeng": "ganja", "cipika cipiki": "cium pipi kanan cium pipi kiri",
27 "ciyh": "sih", "ckepp": "cakep", "ckp": "cakep", "cmliw": "correct me if i'm wrong", "cmpur": "campur",
28 "cong": "banci", "conlok": "cinta lokasi", "cowwyy": "maaf", "cp": "siapa", "cpe": "capek", "cppe": "capek",
29 "cucok": "cocok", "cuex": "cuek", "cumi": "Cuma miscall", "cups": "culun",
30 "curanmor": "pencurian kendaraan bermotor", "curcol": "curahan hati colongan",
31 "cwek": "cewek", "cyin": "cinta", "d": "di", "dah": "deh", "dapet": "dapat", "de": "adik",
32 "dek": "adik", "demen": "suka", "deyh": "deh", "dgn": "dengan", "diancurin": "dihancurkan",
33 "dimaaftin": "dimaaftkan", "dimintak": "diminta", "disono": "di sana", "dket": "dekati",
34 "dkk": "dan kawan-kawan", "dll": "dan lain-lain", "dlu": "dulu", "dngn": "dengan", "dodol": "bodoh",
35 "doku": "uang", "dongs": "dong", "dpt": "dapat", "dri": "dari", "drmn": "darimana", "drted": "dari tadi",
36 "dst": "dan seterusnya", "dtg": "datang", "duh": "aduh", "duren": "durian", "ed": "edisi",
37 "egg": "emang gue pikirin", "eke": "aku", "elu": "kamu", "emangnya": "memangnya", "emng": "memang",
38 "endak": "tidak", "enggak": "tidak", "envy": "iri", "ex": "mantan", "fax": "facsimile",
39 "fiffo": "first in first out", "folbek": "follow back", "fyi": "sebagai informasi", "gaada": "tidak ada uang",
40 "gag": "tidak", "gaje": "tidak jelas", "gak papa": "tidak apa-apa", "gan": "juragan", "gaptek": "gagap teknologi",
41 "gatek": "gagap teknologi", "gawe": "kerja", "gbs": "tidak bisa", "gebetan": "orang yang disukai",
42 "geje": "tidak jelas", "gepeng": "gelandangan dan pengemis", "ghiy": "lagi", "gile": "gila",
43 "gimana": "bagaimana", "gino": "gigi ngongol", "githu": "gitu", "gj": "tidak jelas", "gmana": "bagaimana",
44 "gn": "begini", "goblok": "bodoh", "golput": "golongan putih", "gowes": "mengayuh sepeda",
45 "gpony": "tidak punya", "gr": "gede rasa", "gretongan": "gratisan", "gtau": "tidak tahu",
46 "gua": "saya", "guoblok": "goblok", "gw": "saya", "ha": "tertawa", "haha": "tertawa",
47 "hallow": "halo", "hankam": "pertahanan dan keamanan", "hehe": "he", "helo": "halo", "hey": "hai",
48 "hlm": "halaman", "hny": "hanya", "hoax": "isu bohong", "hr": "hari", "hrus": "harus",
49 "hubdar": "perhubungan darat", "huff": "mengeluh", "hum": "rumah", "humz": "rumah",
50 "ilang": "hilang", "ilfil": "tidak suka", "imho": "in my humble opinion", "imoetz": "imut",
51 "item": "hitam", "itungan": "hitungan", "iye": "iya", "ja": "saja", "jadiin": "jadi",
52 "jain": "jaga image", "jayus": "tidak lucu", "jdi": "jadi", "jem": "jam", "jga": "juga",
53 "jgnkan": "jangankan", "jin": "anjing", "jln": "jalan", "jomblo": "tidak punya pacar",
54 "jubir": "juru bicara", "jutek": "galak", "k": "ke", "kab": "kabupaten", "kabor": "kabur",
55 "kacrut": "kacau", "kadiv": "kepala divisi", "kagak": "tidak", "kalo": "kalau", "kampret": "sialan",
56 "kamtibmas": "keamanan dan ketertiban masyarakat", "kamuwh": "kamu", "kanwil": "kantor wilayah",
57 "karna": "karena", "kasubbag": "kepala subbagian", "katrok": "kampung", "kayanya": "kayaknya",
58 "kbr": "kabar", "kdu": "harus", "kec": "kecamatan", "kejurnas": "kejuaraan nasional",
59 "kekeuh": "keras kepala", "kel": "kelurahan", "kemaren": "kemarin", "kepengen": "mau", "kepingin": "mau",
60 "kepsek": "kepala sekolah", "kesbang": "kesatuan bangsa", "kesra": "kesejahteraan rakyat",
61 "ketrima": "diterima", "kgiatan": "kegiatan", "kibul": "bohong", "kimpoi": "kawin", "kl": "kalau",
62 "klian": "kalian", "kloter": "kelompok terbang", "klw": "kalau", "km": "kamu", "kmps": "kampus",
63 "kmrn": "kemarin", "kna": "kenal", "knp": "kenapa", "kodya": "kota madya",
64 "komdis": "komisi disiplin", "komsov": "komunis soviet", "kongkow": "kumpul bareng teman-teman",
65 "kopdar": "kopi darat", "korup": "korupsi", "kpn": "kapan", "krenz": "keren", "krn": "kirim", "kt": "kita", "ktmu": "ketemu",
66 "kyk": "seperti", "la": "lah", "lam": "salam", "lamp": "lampiran", "lanud": "landasan udara",
67 "latgab": "latihan gabungan", "lebay": "berlebihan", "leh": "boleh", "lelet": "lambat",
68 "lemot": "lambat", "lgi": "lagi", "lgsg": "langsung", "liat": "lihat", "litbang": "penelitian dan pengembangan",
69 "lmyn": "lumayan", "lo": "kamu", "loe": "kamu", "lola": "lambat berfikir", "lough": "cinta", "low": "kalau",
70 "lp": "lupa", "luber": "langsung, umum, bebas, dan rahasia", "luchuw": "lucu", "lum": "belum", "luthu": "lucu",
71 "lwn": "lawan", "maacih": "terima kasih", "mabal": "bolos", "macem": "macam", "macih": "masih",
72 "maem": "makan", "magabut": "makan gaji buta", "maho": "homo", "mak jang": "kaget", "maksain": "memaksa",
73 "malem": "malam", "mam": "makan", "maneh": "kamu", "maniez": "manis", "mao": "mau", "masukin": "masukkan",
74 "melu": "ikut", "mepet": "dekati sekali", "mgu": "minggu", "migas": "minyak dan gas bumi",
75 "mikul": "minuman beralkohol", "miras": "minuman keras", "mlah": "malah", "mngkn": "mungkin",
76 "mo": "mau", "mokad": "mati", "moso": "masa", "mpe": "sampai", "msk": "masuk", "mslh": "masalah",
77 "mt": "makan teman", "mubes": "musyawarah besar", "mulu": "melulu", "mumpung": "selagi",
78 "munas": "musyawarah nasional", "muntaber": "muntah dan berak", "musti": "mesti", "muupz": "maaf",
79 "mw": "now watching", "n": "dan", "nanam": "menanam", "nanya": "bertanya", "napa": "kenapa",
80 "napi": "narapidana", "napza": "narkotika, alkohol, psikotropika, dan zat adiktif",
81 "narkoba": "narkotika, psikotropika, dan obat terlarang", "nasgor": "nasi goreng", "nda": "tidak",
82 "ndiri": "sendiri", "ne": "ini", "nekolin": "neokolonialisme", "nembak": "menyatakan cinta",
83 "ngabuburit": "menunggu berbuka puasa", "ngaku": "mengaku", "ngambil": "mengambil",
84 "nganggur": "tidak punya pekerjaan", "ngapah": "kenapa", "ngaret": "terlambat", "ngasih": "memberikan",
85 "ngebandel": "berbuat bandel", "ngegosip": "bergosip", "ngeklaim": "mengklaim", "ngeksis": "menjadi eksis",
86 "ngeles": "berkilah", "ngelidur": "menggigau", "ngerampok": "merampok", "ngga": "tidak", "ngibul": "berbohong",
87 "ngiler": "mau", "ngiri": "iri", "ngisiin": "mengisiskan", "ngmng": "bicara", "ngomong": "bicara",
88 "ngubek2": "mencari-cari", "ngurus": "mengurus", "nie": "ini", "nih": "ini", "niyh": "nih", "nmn": "nomor",
89 "nntn": "nonton", "nobar": "nonton bareng", "np": "now playing", "ntar": "nanti", "ntn": "nonton",
90 "numpuk": "bertumpuk", "nutupin": "menutupi", "nyari": "mencari", "nyekar": "menyekar", "nyicil": "mencicil",

```

91 "nyoblos": "mencoblos", "nyokap": "ibu", "ogah": "tidak mau", "ol": "online", "ongkir": "ongkos kirim",
 92 "oot": "out of topic", "org2": "orang-orang", "ortu": "orang tua", "otda": "otonomi daerah",
 93 "otw": "on the way, sedang di jalan", "pacal": "pacar", "pake": "pakai", "pala": "kepala",
 94 "pansus": "panitia khusus", "parpol": "partai politik", "pasutri": "pasangan suami istri", "pd": "pada",
 95 "pede": "percaya diri", "pelatnas": "pemusatan latihan nasional", "pemda": "pemerintah daerah",
 96 "pemkot": "pemerintah kota", "pemred": "pemimpin redaksi", "penjas": "pendidikan jasmani",
 97 "perda": "peraturan daerah", "perhatiin": "perhatikan", "pesenan": "pesanan", "pgang": "pegang", "pi": "tapi",
 98 "pilkada": "pemilihan kepala daerah", "pisan": "sangat", "pk": "penjahat kelamin", "plg": "paling",
 99 "pmrnth": "pemerintah", "polantas": "polisi lalu lintas", "ponpes": "pondok pesantren", "pp": "pulang pergi",
 100 "prg": "pergi", "prnh": "pernah", "psen": "pesan", "pst": "pasti", "pswt": "pesawat", "pw": "posisi nyaman",
 101 "qmu": "kamu", "rakor": "rapat koordinasi", "ranmor": "kendaraan bermotor", "re": "reply", "ref": "referensi",
 102 "rehab": "rehabilitasi", "rempong": "sulit", "repp": "balas", "restik": "reserse narkotika", "rhs": "rahasia",
 103 "rmh": "rumah", "ru": "baru", "ruko": "rumah toko", "rusunawa": "rumah susun sewa", "ruz": "terus",
 104 "saia": "saya", "salting": "salah tingkah", "sampe": "sampai", "samsek": "sama sekali", "sapose": "siapa",
 105 "satpam": "satuan pengamanan", "sbb": "sebagai berikut", "sbh": "sebuah", "sbnrny": "sebenarnya",
 106 "scr": "secara", "sdgkn": "sedangkan", "sdt": "sedikit", "se7": "setuju", "sebelas dua belas": "mirip",
 107 "sembako": "sembilan bahan pokok", "sempet": "sempat", "sendratari": "seni drama tari", "sgt": "sangat",
 108 "shg": "sehingga", "siech": "sih", "sikon": "situasi dan kondisi", "sinetron": "sinema elektronik",
 109 "siramin": "siramkan", "sj": "saja", "skalian": "sekalian", "sklh": "sekolah", "skt": "sakit",
 110 "slesai": "selesai", "sll": "selalu", "slna": "selama", "slsa": "selesai", "smpt": "sempat", "smw": "semua",
 111 "sndiri": "sendiri", "soljum": "sholat jumat", "songong": "sombong", "sory": "maaf", "sosek": "sosial-ekonomi",
 112 "sotoy": "sok tahu", "spa": "siapa", "sppa": "siapa", "spt": "seperti", "srtfkt": "sertifikat",
 113 "stiap": "setiap", "stlh": "setelah", "suk": "masuk", "sumpek": "sempit", "syg": "sayang", "t4": "tempat",
 114 "tajir": "kaya", "tau": "tahu", "taw": "tahu", "td": "tadi", "tdk": "tidak", "teh": "kakak perempuan",
 115 "telat": "terlambat", "telmi": "telat berpikir", "temen": "teman", "tengil": "menyebalkan", "tepar": "terkapan",
 116 "tgg": "tunggu", "tgu": "tunggu", "thankz": "terima kasih", "thn": "tahun", "tilang": "bukti pelanggaran",
 117 "tipiwan": "TvoOne", "tks": "terima kasih", "tlp": "telepon", "tls": "tulis", "tmbah": "tambah",
 118 "tmen2": "teman-teman", "tmpah": "tumpah", "tmpt": "tempat", "tngu": "tunggu", "tnyta": "nyata",
 119 "tokai": "tai", "toserba": "toko serba ada", "tpi": "tapi", "trdhulu": "terdahulu", "trima": "terima kasih",
 120 "trm": "terima", "trs": "terus", "trutama": "terutama", "ts": "penulis", "tst": "tahu sama tahu",
 121 "ttg": "tentang", "tuch": "tuh", "tuir": "tua", "tw": "tahu", "u": "kamu", "ud": "sudah", "udah": "sudah",
 122 "ujung": "ujung", "ul": "ulangan", "unyu": "lucu", "uplot": "unggah", "urang": "saya", "usah": "perlu",
 123 "utk": "untuk", "valas": "valuta asing", "w/": "dengan", "wadir": "wakil direktur", "wamil": "wajib militer",
 124 "warkop": "warung kopi", "warteg": "warung tegal", "wat": "buat", "wkt": "waktu", "wtf": "what the fuck",
 125 "xixixi": "tertawa", "ya": "iya", "yap": "iya", "yaudah": "ya sudah", "yawah": "ya sudah", "yg": "yang",
 126 "yl": "yang lain", "yo": "iya", "yowes": "ya sudah", "yup": "iya", "7an": "tujuan", "ababil": "abg labil",
 127 "acc": "accord", "adlah": "adalah", "adoh": "aduh", "aha": "tertawa", "aing": "saya", "aja": "saja",
 128 "ajj": "saja", "aka": "dikenal juga sebagai", "akko": "aku", "akku": "aku", "akyu": "aku",
 129 "aljasa": "asal jadi saja", "ama": "sama", "ambl": "ambil", "anjir": "anjing", "ank": "anak", "ap": "apa",
 130 "apaan": "apa", "ape": "apa", "aplot": "unggah", "apva": "apa", "aqu": "aku", "asap": "sesegera mungkin",
 131 "aseek": "asyik", "asek": "asyik", "aseknya": "asyiknya", "asoy": "asyik", "astrojim": "astagfirullahaladzim",
 132 "ath": "kalau begitu", "atu": "kalau begitu", "ava": "avatar", "aws": "awas", "ayang": "sayang",
 133 "ayok": "ayo", "bacot": "banyak bicara", "bales": "balas", "bangdes": "pembangunan desa", "bangkotan": "tua",
 134 "banpres": "bantuan presiden", "bansarkas": "bantuan sarana kesehatan",
 135 "bais": "badan amal, zakat, infak, dan sedekah", "bcot": "karena", "beb": "sayang", "bejibun": "banyak",
 136 "belom": "belum", "bener": "benar", "ber2": "berdua", "berdikari": "berdiri di atas kaki sendiri",
 137 "bet": "banget", "beti": "beda tipis", "beut": "banget", "bgd": "banget", "bgs": "bagus",
 138 "bhuhu": "tidur", "bimbuluh": "bimbingan dan penyuluhan", "bisi": "kalau-kalau", "bkn": "bukan", "bl": "beli",
 139 "blg": "bilang", "blm": "belum", "bls": "balas", "bnchi": "benci", "bngung": "bingung", "bnyk": "banyak",
 140 "bohay": "badan aduhai", "bokep": "porno", "bokin": "pacar", "bole": "boleh", "bolot": "bodoh",
 141 "bonyok": "ayah ibu", "bpk": "bapak", "brb": "segera kembali", "brngkt": "berangkat", "brp": "berapa",
 142 "brun": "saudara laki-laki", "bsa": "bisa", "bsk": "besok", "bu_bu": "tidur", "bubarin": "bubarkan",
 143 "buber": "buka bersama", "bujubune": "luar biasa", "buser": "buru sergap", "bwhn": "bawahan", "byar": "bayar",
 144 "byr": "bayar", "c8": "chat", "cabut": "pergi", "caem": "capek", "cama-cama": "sama-sama",
 145 "cangcut": "celana dalam", "cape": "capek", "caur": "jelek", "cekak": "tidak ada uang", "cekidot": "coba lihat",
 146 "cemplungin": "cemplungkan", "ceper": "pendek", "ceu": "kakak perempuan", "cewe": "cewek", "cibuk": "sibuk",
 147 "cin": "cinta", "ciye": "cie", "ckck": "ck", "clbk": "cinta lama bersemi kembali", "cmpr": "campur",
 148 "cnenk": "senang", "congor": "mulut", "cow": "cowok", "coz": "karena", "cpa": "siapa", "gokil": "gila",
 149 "gombal": "suka merayu", "gpl": "tidak pakai lama", "gpp": "tidak apa-apa", "gretong": "gratis", "gt": "begitu",
 150 "gtw": "tidak tahu", "gue": "saya", "guys": "teman-teman", "gws": "cepat sembuh", "haghaghag": "tertawa",
 151 "hakhak": "tertawa", "handak": "bahan peledak", "hansip": "pertahanan sipil", "hellow": "halo", "helow": "halo",
 152 "hi": "hai", "hlng": "hilang", "hnya": "hanya", "houm": "rumah", "hrs": "harus",
 153 "hubad": "hubungan angkatan darat", "hubla": "perhubungan laut", "huft": "mengeluh",
 154 "humas": "hubungan masyarakat", "idk": "saya tidak tahu", "ilfeel": "tidak suka", "imba": "jago sekali",
 155 "imoet": "imut", "info": "informasi", "itung": "hitung", "isengin": "bercanda", "iyala": "iya lah",
 156 "iyo": "iya", "jablay": "jarang dibelai", "jadul": "jaman dulu", "jancuk": "anjing", "jd": "jadi",
 157 "jdikan": "jadikan", "jg": "juga", "jgn": "jangan", "jijay": "jijik", "jkt": "jakarta", "jnj": "janji",
 158 "jth": "jatuh", "jurdil": "jujur adil", "jwb": "jawab", "ka": "kakak", "kabag": "kepala bagian",
 159 "kacian": "kasihan", "kadit": "kepala direktorat", "kaga": "tidak", "kaka": "kakak",
 160 "kamtib": "keamanan dan ketertiban", "kamuh": "kamu", "kamyu": "kamu", "kapt": "kapten",
 161 "kasat": "kepala satuan", "kasubbid": "kepala subbidang", "kau": "kamu", "kbar": "kabar",
 162 "kcian": "kasihan", "keburu": "terlanjur", "kedubes": "kedutaan besar", "kek": "seperti", "keknnya": "kayaknya",
 163 "keliatan": "kelihatan", "keneh": "masih", "kepikiran": "terpikirkan", "kepo": "mau tahu urusan orang",
 164 "kere": "tidak punya uang", "kesian": "kasihan", "ketauan": "ketahuan", "keukeuh": "keras kepala",
 165 "khan": "kan", "kibus": "kaki busuk", "kk": "kakak", "kliian": "kalian", "klo": "kalau", "kluarga": "keluarga",
 166 "klwrga": "keluarga", "kmari": "kemari", "kmpus": "kampus", "kn": "kan", "knl": "kenal", "knpa": "kenapa",
 167 "kog": "kok", "kompi": "komputer", "komtong": "komunis Tiongkok", "konjen": "konsulat jenderal", "koq": "kok",
 168 "kpd": "kepada", "kptsan": "keputusan", "krik": "garing", "krn": "karena", "ktauan": "ketahuan",
 169 "ktny": "katanya", "kudu": "harus", "kuq": "kok", "ky": "seperti", "kykny": "kayanya", "laka": "kecelakaan",
 170 "lambreta": "lambat", "lansia": "lanjut usia", "lapas": "lembaga pemasyarakatan", "lbur": "libur",
 171 "lekong": "laki-laki", "lg": "lagi", "lgkp": "lengkap", "lht": "lihat", "linmas": "perlindungan masyarakat",
 172 "lmyan": "lumayan", "lngkp": "lengkap", "loch": "loh", "lol": "tertawa", "lom": "belum", "loupz": "cinta",
 173 "lowh": "kamu", "lu": "kamu", "lucu": "lucu", "luff": "cinta", "luph": "cinta", "lw": "kamu", "lwt": "lewat",
 174 "maaciw": "terima kasih", "mabes": "markas besar", "macem-macem": "macam-macam", "madesu": "masa depan suram",
 175 "maen": "main", "mahatma": "maju sehat bersama", "mak": "ibu", "makasih": "terima kasih", "malah": "bahkan",
 176 "malu2in": "memalukan", "mamz": "makan", "manies": "manis", "mantep": "mantap", "markus": "makelar kasus",
 177 "mba": "mbak", "mending": "lebih baik", "mgkn": "mungkin", "mhn": "mohon", "miker": "minuman keras",
 178 "milis": "mailing list", "mksd": "maksud", "mls": "malas", "mnt": "minta", "moge": "motor gede",
 179 "mokat": "mati", "mosok": "masa", "msh": "masih", "mskpn": "meskipun", "msng2": "masing-masing",
 180 "muahal": "mahal", "muker": "musyawarah kerja", "mumet": "pusing", "muna": "munafik",
 181 "munaslub": "musyawarah nasional luar biasa", "musda": "musyawarah daerah", "muup": "maaf", "muuv": "maaf",
 182 "nal": "kenal", "nangis": "menangis", "naon": "apa", "napol": "narapidana politik", "naq": "anak",

```

183 "narsis": "bangga pada diri sendiri", "nax": "anak", "ndak": "tidak", "ndut": "gendut",
184 "nekolim": "neokolonialisme", "nelfon": "menelepon", "ngabis2in": "menghabiskan", "ngakak": "tertawa",
185 "ngambek": "marah", "ngampus": "pergi ke kampus", "ngantri": "mengantri", "ngapain": "sedang apa",
186 "ngaruh": "berpengaruh", "ngawur": "berbicara sembarangan", "ngeceng": "kumpul bareng-bareng",
187 "ngeh": "sadar", "ngekos": "tinggal di kos", "ngelamar": "melamar", "ngeliat": "melihat",
188 "ngemeng": "bicara terus-terusan", "ngerti": "mengerti", "nggak": "tidak", "ngikut": "ikut",
189 "nginep": "menginap", "ngisi": "mengisi", "ngmg": "bicara", "ngocol": "lucu", "ngomongin": "membicarakan",
190 "ngumpul": "berkumpul", "ni": "ini", "nyasar": "tersesat", "nyariin": "mencari", "nyiapin": "mempersiapkan",
191 "nyiram": "menyiram", "nyok": "ayo", "o/": "oleh", "ok": "ok", "priksa": "periksa", "pro": "profesional",
192 "psn": "pesan", "psti": "pasti", "puanas": "panas", "qmo": "kamu", "qt": "kita", "rame": "ramai",
193 "raskin": "rakyat miskin", "red": "redaksi", "reg": "register", "rejekin": "rezeki", "renstra": "rencana strategis",
194 "reskrim": "reserse kriminal", "sni": "sini", "somse": "sombong sekali", "sorry": "maaf", "sosbud": "sosial-budaya",
195 "sospol": "sosial-politik", "sowry": "maaf", "spd": "sepeda", "spri": "seperti", "spy": "supaya",
196 "stelah": "setelah", "subbag": "subbagian", "sumbangin": "sumbangkan", "sy": "saya", "syp": "siapa",
197 "tabanas": "tabungan pembangunan nasional", "tar": "nanti", "taun": "tahun", "tawh": "tahu", "tdi": "tadi",
198 "te2p": "tetap", "tekor": "rugi", "telkom": "telekomunikasi", "telp": "telepon", "temen2": "teman-teman",
199 "tengok": "menjenguk", "terbitin": "terbitkan", "tgl": "tanggal", "thanks": "terima kasih",
200 "thd": "terhadap", "thx": "terima kasih", "tipi": "TV", "tkg": "tukang", "tll": "terlalu", "tlpn": "telepon",
201 "tman": "teman", "tmbh": "tambah", "tmn2": "teman-teman", "tmph": "tumpah", "tnda": "tanda", "tnh": "tanah",
202 "togel": "toto gelap", "tp": "tapi", "tq": "terima kasih", "trgntg": "tergantung", "trims": "terima kasih",
203 "cb": "coba", "y": "ya", "munafik": "munafik", "reklamuk": "reklamasi", "sma": "sama", "tren": "trend",
204 "ngehe": "kesal", "mz": "mas", "analisis": "analisis", "sadaar": "sadar", "sept": "september",
205 "nmenarik": "menarik", "zonk": "bodoh", "rights": "benar", "simiskin": "miskin", "ngumpet": "sembunyi",
206 "hardcore": "keras", "akhirx": "akhirnya", "solve": "solusi", "watuk": "batuk", "ngebully": "intimidasi",
207 "masy": "masyarakat", "still": "masih", "tauk": "tahu", "mbual": "bual", "tioghoa": "tioghoa",
208 "ngentotin": "senggama", "kentot": "senggama", "faktakta": "fakta", "sohib": "teman", "rubahnn": "rubah",
209 "trlalu": "terlalu", "nyela": "cela", "heters": "pembenci", "nyembah": "sembah", "most": "paling",
210 "ikon": "lambang", "light": "terang", "pndukung": "pendukung", "setting": "aturan", "seting": "akting",
211 "next": "lanjut", "waspadalah": "waspada", "gantengsaya": "ganteng", "parte": "partai", "nyerang": "serang",
212 "nipu": "tipu", "ktipu": "tipu", "jentelmen": "berani", "buangbuang": "uang", "tsangka": "tersangka",
213 "kurng": "kurang", "ista": "nista", "less": "kurang", "koan": "teriak", "paranoid": "takut",
214 "problem": "masalah", "tahi": "kotoran", "tiran": "tiran", "tilep": "tilap", "happy": "bahagia",
215 "tak": "tidak", "penertiban": "tertib", "uasai": "kuasa", "mnolak": "tolak", "trending": "trend",
216 "taik": "tahi", "wkwkw": "tertawa", "ahokncc": "ahok", "istaa": "nista", "benarjujur": "jujur",
217 "mgkin": "mungkin"}
218
219 processed_comments = []
220
221 for sentence in data['text']:
222     # Remove all the special characters
223     processed_comment = re.sub(r'\W| ', ' ', str(sentence))
224
225     # Converting to Lowercase
226     processed_comment = processed_comment.lower()
227
228     #Remove number
229     processed_comment = re.sub(r'\d+', ' ', processed_comment)
230
231     # remove all single characters
232     processed_comment = re.sub(r'\s+[a-zA-Z]\s+', ' ', processed_comment)
233
234     #remove duplicate character
235     pattern=re.compile(r"(\.)\1{1,}",re.DOTALL)
236     processed_comment=pattern.sub(r"\1",processed_comment)
237
238     #Corrected Slang words
239     words = processed_comment.split()
240     rfrm=[slangs[word] if word in slangs else word for word in words]
241     processed_comment= " ".join(rfrm)
242
243     #remove stopwords
244     factory = StopWordRemoverFactory()
245     f = open("stopwords-tala.txt", "r")
246     more_stopword = [] #menambahkan stopwords
247     for line in f:
248         stripped_line = line.strip()
249         line_list = stripped_line.split()
250         more_stopword.append(line_list[0])
251     f.close()
252     stopwords = factory.get_stop_words() + more_stopword
253     temp = [t for t in re.findall(r'\b[a-z]+-?[a-z]+\b',processed_comment) if t not in stopwords]
254     processed_comment = ' '.join(temp)
255
256     #stemming
257     stemmer = StemmerFactory().create_stemmer()
258     processed_comment = stemmer.stem(processed_comment)
259
260     #Subtitusing multiple spaces with single space
261     processed_comment = re.sub(r'\s+', ' ', processed_comment, flags=re.I)
262
263     processed_comments.append(processed_comment)

```

```
1 processed_comments
```

'harap politis gaya lacur pramagtis warna politik prinsip dlm politik milik integritas politis model suka rakyat',
'syahganda gugat syahganda mg mg lacur intelektual syahganda buka propaganda hasil survei denyawaslu https politik rmlol
co read syahganda buka propaganda hasil survei denyawaslu',
'pdemokrat org akal sehat maju bangsa orang begundal serta lacur politik grasak grusuk salah satu heran https twitter com
ferdinandhaean status',
'cypridophobia fobia takut lacur tular sakit kelamin',
'ramai kerja malam lacur siapa idea suh kerja malam tu takfahamtakapa',
'lacur intelektual kpcuranakyatmelawan kpungejarquickcounthttps twitter com suwandaben status',
'lembaga lacur lsi denyja germo kejar target atas propinsi dlm opini publik simultan dg operasi ampo suara hina hidup me
rekahttps twitter com status',
'bangun gada orang sukses kasur lacur selfreminder',
'abai lacur agama hidup nkri hidup presiden joko widodo',
'sat bersih jenis kotor bangsa masive awal laknat lacur agama biadab ancam selamat nkri',
'sukses capai kasur lacur',
'harap politis gaya lacur pramagtis warna politik prinsip dlm politik milik integritas politis model suka rakyat',
'syahganda gugat syahganda mg mg lacur intelektual syahganda buka propaganda hasil survei denyawaslu https politik rmlol
co read syahganda buka propaganda hasil survei denyawaslu',
'ndemokrat org akal sehat maju bangsa orang begundal serta lacur politik grasak grusuk salah satu heran https twitter com

```
1 from sklearn.feature_extraction.text import CountVectorizer
2 tf_vectorizer = CountVectorizer(max_df=1.0, min_df=1)
3 tf = tf_vectorizer.fit_transform(processed_comments)
4 #hasil representasi
5 tf_terms = tf_vectorizer.get_feature_names()
6 print(tf_vectorizer.get_feature_names())
7 matrix = tf.toarray()
8 print(matrix)
```

```
1 # panggil class LDA
2 from sklearn.decomposition import LatentDirichletAllocation as LDA
3 n_topics = 10 #untuk mendapatkan jumlah topik terbaik perlu trial
4 lda = LDA(n_components=n_topics, learning_method='batch', random_state=0).fit(tf)
5 lda
```

```
LatentDirichletAllocation(random_state=0)
```

In [47]:

```

1 #training LDA
2 vsm_topics = lda.transform(tf)
3 #tampilkan hasil
4 print(vsm_topics.shape)
5 vsm_topics

```

(400, 10)

Out[47]:

```

array([[0.00588256, 0.00588279, 0.00588328, ..., 0.00588236, 0.0058827 ,
        0.00588253],
       [0.0037039 , 0.00370397, 0.00370391, ..., 0.00370374, 0.96666464,
        0.00370386],
       [0.00454598, 0.00454641, 0.00454595, ..., 0.00454562, 0.00454585,
        0.00454574],
       ...,
       [0.02000242, 0.02000153, 0.02000101, ..., 0.02000006, 0.02000044,
        0.02000217],
       [0.05001691, 0.05001395, 0.0500139 , ..., 0.0500008 , 0.05000612,
        0.54987628],
       [0.00526458, 0.00526548, 0.00526366, ..., 0.00526317, 0.95262346,
        0.00526355]])

```

In [48]:

```

1 #Tampilkan nilai-nilai setiap fitur
2 print(lda.components_)

```

```

[[2.1      0.1      0.1      ... 0.1      0.1      0.1      ]
 [0.1      0.1      1.09999143 ... 0.1      0.1      0.1      ]
 [0.1      0.1      1.09999925 ... 0.1      0.1      0.1      ]
 ...
 [0.1      0.1      0.1      ... 0.1      0.1      0.1      ]
 [0.1      0.1      0.1      ... 0.1      0.1      0.1      ]
 [0.1      0.1      1.09996271 ... 0.1      0.1      0.10000178]]

```

In [49]:

```

1 #hasil label topic model untuk setiap dokumen
2 topics = np.argmax(vsm_topics, axis=1)
3 topics

```

Out[49]:

```

array([4, 8, 3, 1, 0, 3, 1, 9, 0, 6, 9, 4, 8, 3, 0, 3, 1, 0, 6, 1, 1, 7,
        6, 8, 1, 9, 5, 6, 9, 9, 4, 5, 6, 0, 0, 8, 5, 9, 6, 1, 4, 9, 3, 4,
        5, 5, 4, 9, 5, 4, 5, 6, 9, 4, 4, 0, 9, 9, 5, 9, 0, 0, 0, 2, 9, 2,
        4, 9, 5, 6, 4, 6, 9, 5, 5, 6, 8, 2, 4, 9, 6, 0, 4, 6, 8, 5, 9, 0,
        1, 9, 9, 5, 6, 9, 1, 6, 5, 6, 4, 5, 5, 4, 2, 6, 8, 4, 2, 9, 6, 4,
        5, 9, 6, 6, 2, 8, 8, 6, 1, 5, 0, 6, 6, 2, 5, 0, 6, 4, 5, 1, 2, 1,
        0, 0, 0, 9, 0, 4, 1, 1, 6, 4, 0, 5, 1, 2, 1, 0, 0, 0, 4, 2, 6, 9,
        0, 3, 1, 8, 6, 4, 2, 2, 6, 4, 0, 9, 0, 9, 1, 5, 9, 4, 4, 0, 2, 9,
        1, 9, 5, 2, 1, 4, 5, 4, 9, 1, 4, 0, 6, 3, 5, 6, 3, 9, 6, 3, 8, 4,
        4, 6, 5, 4, 0, 4, 4, 0, 6, 9, 0, 4, 6, 9, 1, 9, 8, 2, 2, 6, 9, 1,
        9, 9, 4, 2, 5, 5, 3, 5, 6, 5, 6, 1, 6, 5, 0, 5, 3, 4, 6, 4, 9, 2,
        9, 5, 9, 1, 6, 6, 2, 5, 4, 6, 6, 4, 8, 0, 5, 6, 4, 9, 9, 7, 1, 4,
        1, 4, 0, 2, 9, 3, 9, 1, 9, 8, 6, 6, 1, 6, 9, 6, 1, 9, 6, 3, 9, 2,
        1, 6, 9, 7, 9, 6, 9, 7, 8, 9, 1, 9, 6, 8, 6, 8, 4, 4, 2, 8, 2, 4,
        2, 4, 0, 6, 7, 4, 3, 6, 1, 2, 1, 6, 6, 9, 3, 8, 8, 9, 4, 2, 0, 1,
        3, 2, 2, 9, 8, 9, 0, 0, 2, 6, 1, 1, 6, 4, 4, 0, 4, 2, 4, 8, 5, 2,
        6, 1, 1, 8, 3, 4, 5, 8, 1, 1, 1, 1, 3, 0, 9, 1, 9, 1, 1, 5, 2, 5,
        9, 7, 1, 5, 8, 5, 1, 4, 0, 6, 9, 1, 9, 6, 9, 9, 9, 6, 6, 6, 6,
        1, 4, 9, 8], dtype=int64)

```

In [50]:

```

1 #mencetak word fitur dengan nilai tertinggi pada setiap topik
2 n_top_words = 10 # jumlah fitur tertinggi yang kita tentukan
3 topic_words = {}
4 for topic, comp in enumerate(lda.components_):
5     word_idx = np.argsort(comp)[-1:][n_top_words]
6     # store the words most relevant to the topic
7     topic_words[topic] = [tf_vectorizer.get_feature_names()[i]+' '+str(comp[i]) for i in word_idx]

```

In [51]:

```
1 for topic, words in topic_words.items():
2     print('Topic: %d' % topic)
3     print(' %s' % ', '.join(words))
```

Topic: 0

lacur 36.858302877433104, ahok 14.100063547015367, com 10.33124212517124, twiter 9.3115701625257, negeri 9.10005257926472
3, kerja 7.100019285676232, agama 6.100051317894472, presiden 6.100040681345306, tu 6.100037092189047, laki 6.1000116111161
855

Topic: 1

lacur 34.17682018061304, status 20.95321637287863, twitter 20.72725394922268, com 20.633781250824796, https 16.6273719501033
8, leceh 9.002413895605685, tutup 8.100028846486598, tusuk 8.100019253365828, ras 8.099903939505213, verbal 7.9899674768977
045

Topic: 2

lacur 23.100425806660525, ras 11.788535040417313, orang 9.83517213284242, politik 6.100118994339808, iya 6.10002453701038,
suka 5.944528899123657, dasar 5.100115092764825, jaya 5.0999999995441625, gitu 5.09999999923475, maf 4.10008184079854

Topic: 3

lacur 12.784085463794048, com 7.100041782956527, twiter 6.100040680309208, orang 5.100091145426324, status 5.1000283347639
49, ga 4.100110571668511, ras 4.099959526783606, gak 3.957519079901731, politik 3.1000609901631715, abdu 3.1000392953393807
193071923

Topic: 4

lacur 43.823716325696914, ras 16.100238531340743, politik 16.055280738731454, gak 7.98385730774608, anjing 6.1000224106218
83, politis 5.100032068752101, khianat 4.100031667465745, kayak 3.9825075811719395, rakyat 3.755852239722175, yah 3.1000448
193071923

Topic: 5

lacur 41.35840756780545, politik 10.100082105401173, com 10.10001285224174, twiter 9.100020097764915, https 9.1000192583666
63, leceh 8.100048691293653, verbal 7.100032011533481, status 7.0999990467894145, nama 6.479139947659218, anjing 6.10002202
4478947

Topic: 6

lacur 56.97048436460543, agama 10.10004876004643, ras 9.411608944711196, mesum 8.099994396168306, twiter 8.09995653113680
4, com 8.099930792909472, media 7.100021220351084, orang 6.364847741523899, jual 6.100053142647907, demokrasi 6.10002990096
1755

Topic: 7

mesum 2.1001448956667024, org 2.100064838133289, skrg 2.100011354295569, duduk 2.100005646965668, iseng 2.099999999803583
4, negara 2.0999999997168217, hasil 1.100040421259895, nya 1.1000272434582734, cina 1.100025188701681, parah 1.100023937424
0526

Topic: 8

lacur 10.488518025497777, ras 8.100149221777906, syahganda 8.099999999695108, buka 5.100057243928322, orang 5.100054285997
409, emang 5.1000509453611125, kpop 5.09999999964672, gak 4.9446389126115395, https 4.447129214571693, hasil 4.1001277366085
69

Topic: 9

lacur 56.33942820903878, nya 10.8824002150346, orang 9.053454869108535, mesum 8.519297213112779, gak 7.100095266646303, ot
ak 7.040117200387565, iya 6.889353722015582, hina 5.81791365053057, beda 5.100014331414758, juang 5.100014231682571

LSA

In [52]:

```
1 from sklearn.feature_extraction.text import TfidfVectorizer
2 vektor = TfidfVectorizer(max_features = 1000)
```


In [53]:

```

1 #menghitung tf-idf dengan TfidfTransformer
2 vektor_dt = vektor.fit_transform(data['text'])
3 print(vektor_dt)
4 print(vektor_dt.shape)

```

```

(0, 762)      0.15272779359201613
(0, 242)      0.19192313773966022
(0, 36)       0.14177393727307744
(0, 110)      0.18255764560035231
(0, 607)      0.19192313773966022
(0, 601)      0.19192313773966022
(0, 817)      0.14679231387154673
(0, 190)      0.1752932055260729
(0, 359)      0.19192313773966022
(0, 578)      0.19192313773966022
(0, 925)      0.29358462774309346
(0, 802)      0.1496249699362145
(0, 805)      0.12857401747832858
(0, 134)      0.18255764560035231
(0, 245)      0.1643393492071342
(0, 742)      0.19192313773966022
(0, 756)      0.13545318669050188
(0, 733)      0.19192313773966022
(0, 979)      0.18255764560035231
(0, 389)      0.15999223366629556
(0, 941)      0.24078494856409796
(0, 737)      0.19192313773966022
(0, 695)      0.04390090115685355
(0, 295)      0.19192313773966022
(0, 734)      0.3286786984142684
:             :
(397, 626)    0.5441070786149599
(397, 382)    0.3955651449392686
(397, 475)    0.4969607890072219
(397, 623)    0.4427114345470066
(397, 376)    0.2983654185284531
(397, 695)    0.1244601946296941
(398, 477)    0.8143305631184415
(398, 376)    0.5356650091115894
(398, 695)    0.22344738079620147
(399, 792)    0.35686995865224147
(399, 155)    0.2859150865896594
(399, 712)    0.3339049400532245
(399, 629)    0.2716809628530919
(399, 674)    0.2657134797002912
(399, 726)    0.2657134797002912
(399, 358)    0.16915875560193464
(399, 851)    0.23242200425197623
(399, 552)    0.26031523455388494
(399, 992)    0.1830994138898145
(399, 187)    0.1642305438990429
(399, 960)    0.16713572966240264
(399, 923)    0.20739716938742805
(399, 96)     0.24665612798854372
(399, 217)    0.3403960541036873
(399, 695)    0.07637811647778488
(400, 1000)

```

In [54]:

```

1 idf = vektor.idf_
2 dd = dict(zip(vektor.get_feature_names(), idf))
3 l = sorted(dd, key = (dd).get)
4 print(l[0], l[-1])
5 print(dd['pelacur'])
6 print(dd['wilde'])

```

```

pelacur wilde
1.3485145296633316
6.300814246746624

```

In [55]:

```

1 from sklearn.decomposition import TruncatedSVD
2 lsa_model = TruncatedSVD(n_components=10, algorithm='randomized', n_iter=10, random_state=42)
3 lsa_top = lsa_model.fit_transform(vektor_dt)

```

In [56]:

```

1 # jumlah dokumen * jumlah topik
2 print(lsa_top.shape)

```

```

(400, 10)

```

In [57]:

```
1 # top lsa
2 print(lsa_top)
```

```
[[ 1.48659449e-01 -9.07423341e-02 -9.32955306e-02 ... -1.12854271e-01
 5.71319799e-02 -1.41289657e-01]
 [ 1.17667791e-01 -6.98253277e-03 -9.08301957e-03 ... 1.49051278e-01
 4.12538022e-02 3.94359639e-03]
 [ 2.63552696e-01 -1.22096489e-01 1.58809497e-01 ... -1.89116206e-02
 8.87302284e-02 -8.13867092e-02]
 ...
 [ 1.43796548e-01 -6.33206813e-03 -9.86829494e-02 ... 4.92812579e-02
 8.67508362e-02 1.10058605e-01]
 [ 2.17107931e-01 -1.82523631e-04 -1.11976866e-01 ... 8.87682837e-02
 -2.13490196e-02 7.28173096e-02]
 [ 2.65659717e-01 -1.95614557e-01 3.84057286e-02 ... 1.87547711e-01
 -9.76421356e-02 2.07392673e-02]]
```

In [58]:

```
1 # memunculkan nilai lsa setiap topik
2 l = lsa_top[0]
3 print('Topik - Topik :')
4 for i, topic in enumerate(l):
5     print('Topic ', i, ' : ', topic*100)
```

```
Topik - Topik :
Topic 0 : 14.865944900587117
Topic 1 : -9.074233411497108
Topic 2 : -9.329553056265286
Topic 3 : -1.5971514002380256
Topic 4 : -0.6461870061277563
Topic 5 : -5.2107829719236065
Topic 6 : -18.1891684898388
Topic 7 : -11.28542705524751
Topic 8 : 5.713197985199807
Topic 9 : -14.128965680771529
```

In [59]:

```
1 #memunculkan jumlah kata-kata dalam setiap topik
2 print(lsa_model.components_.shape)
3 print(lsa_model.components_)
```

```
(10, 1000)
[[ 0.00725228 0.00192549 0.00399414 ... 0.00162078 0.0292714
 0.00869898]
 [ 0.00018748 -0.00018537 -0.004088 ... -0.00342359 0.03029904
 -0.00146251]
 [-0.00159705 -0.00032499 -0.00182139 ... -0.0034003 0.00997559
 -0.00775767]
 ...
 [-0.01669732 0.0085105 -0.00188842 ... -0.01050738 -0.02939351
 0.00095881]
 [-0.01419099 0.00259324 0.006313 ... 0.00109504 -0.01150681
 0.0032488 ]
 [ 0.00805032 0.00020352 -0.0013523 ... -0.00040257 0.02123483
 0.02292166]]
```

In [60]:

```
1 # word / kata paling penting dalam setiap topik
2 vocab = vektor.get_feature_names()
3 for i, comp in enumerate(lsa_model.components_):
4     vocab_comp = zip(vocab, comp)
5     sorted_words = sorted(vocab_comp, key = lambda x:x[1], reverse=True)[:10]
6     print('Topic '+str(i)+' : ')
7     for t in sorted_words:
8         print(t[0], end=', ')
9     print('\n')
```

Topic 0:

pelacur, politik, com, twitter, yg, ini, status, di, dan, https,

Topic 1:

politik, pelacur, negri, la, yusrilnorakdankasarmainnya, tutup01tusuk02, jabatan, cinta, para, yim,

Topic 2:

twitter, com, status, https, politik, pic, tutup01tusuk02, kpucurangrakyatmelawan, kpungejarquickcounthttps, suwandaben,

Topic 3:

politik, rasis, mesum, aja, ini, gak, lu, yg, ya, otak,

Topic 4:

mesum, aku, otak, ini, juga, babi, izinkan, cabul, ulama, gak,

Topic 5:

aku, izinkan, menjadi, tuhan, politik, ini, dan, binatang, jalang, filosofi,

Topic 6:

rasis, pelacur, aja, aku, mereka, kamu, tu, mesum, gak, kau,

Topic 7:

di, ahok, mesum, kecuali, kasur, jadi, negeri, tutup01tusuk02, otak, pic,

Topic 8:

kau, tu, politik, kata, orang, yang, la, kecuali, hina, tidak,

Topic 9:

kau, ini, anjing, muka, eh, tu, jadi, la, nak, pki,

Kesimpulan

- Terdapat beberapa perbedaan hasil Topic Modelling menggunakan LDA dan LSA
- Topik pertama dari kedua metode LDA maupun LSA sama dan menghasilkan output "pelacur" sebagai kata teratas
- Sedangkan di metode LDA, kata teratas yang dominan dari 10 topik menghasilkan "pelacur"