



**ANALISIS BUKTI DIGITAL *CYBERBULLYING* PADA *JEJARING SOSIAL*
MENGUNAKAN *NAÏVE BAYES CLASSIFIER* (NBC)**

HARIANI

13917146

Tesis diajukan sebagai syarat untuk meraih gelar Magister Komputer

Konsentrasi Forensika Digital

Program Studi Magister Teknik Informatika

Program Pascasarjana Fakultas Teknologi Industri

Universitas Islam Indonesia

2017

Lembar Pengesahan Pembimbing

**ANALISIS BUKTI DIGITAL *CYBERBULLYING* PADA *JEJARING SOSIAL*
MENGUNAKAN *NAÏVE BAYES CLASSIFIER* (NBC)**



Pembimbing I

A handwritten signature in black ink, which appears to be 'Dr. Imam Riadi, M.Kom'. The signature is written in a cursive style and is positioned above the printed name.

Dr. Imam Riadi, M.Kom

Lembar Pengesahan Penguji

ANALISIS BUKTI DIGITAL CYBERBULLYING PADA JEJARING SOSIAL
MENGUNAKAN NAÏVE BAYES CLASSIFIER (NBC)

Nama: Hariani

NIM: 13917146

Yogyakarta, Maret, 2017

Tim Penguji,

Dr. Imam Riadi, M.Kom

Ketua

Yudi Prayudi, S.Si., M.Kom

Anggota I

Dr. Bambang Sugiantoro, MT

Anggota II



Mengetahui,

Ketua Program Pascasarjana Fakultas Teknologi Industri

Universitas Islam Indonesia



Dr. R. Deduh Dirgahayu, S.T., M.Sc.

ABSTRAK

Pesatnya perkembangan Jejaring Sosial sebagai salah satu kemajuan dalam teknologi tidak hanya membawa dampak positif tetapi juga memperkenalkan masalah-masalah baru saat tidak digunakan secara tepat atau menyalahi secara semestinya, hal ini disebut *cybercrime*. Twitter merupakan salah satu jejaring sosial yang sangat diminati saat ini karena kemudahan akses yang dapat dilakukan melalui smarthphone, laptop, tablet maupun berbagai layanan aplikasi lainnya. Pertumbuhan Jejaring Sosial Twitter tersebut membawa *trend* baru dalam masyarakat sebagai ajang untuk melakukan tindakan penindasan secara online atau yang lebih dikenal dengan *cyberbullying*.

Cyberbullying merupakan salah satu *cybercrime* yang sedang marak saat ini. Praktik *cyberbullying* tidak hanya terbatas pada anak-anak tetapi juga pada orang dewasa hal ini disebut *cyberstalking* atau *cyberharrasment*. Dampak yang ditimbulkan bagi pelaku bisa dijerat hukuman penjara sesuai dengan undang-undang yang berlaku, sementara dari sisi psikologi korban bisa mengalami depresi bahkan bunuh diri. Dengan fenomena tersebut, penelitian ini fokus untuk menganalisis perkembangan *cyberbullying* di Indonesia melalui Jejaring Sosial Twitter dengan pendekatan Data Mining menggunakan algoritma Naïve Bayes Classifier pada Machine Learning WEKA. Output penelitian ini berupa *knowledge* mengenai perkembangan *cyberbullying* pada Jejaring Sosial dan Jenis *cyberbullying* yang banyak digunakan sehingga pemerintah atau instansi terkait dapat melakukan *edukasi* dan memberi sanksi tegas bagi masyarakat dan pengguna Jejaring Sosial yang melakukan *bullying* untuk menghindari korban *bullying* lebih banyak.

Berdasarkan analisis hasil klasifikasi dapat diketahui bahwa yang positif mengandung konten *bullying* sebanyak 86.97%. Jenis *cyberbullying* yang banyak digunakan *related psychology* sebanyak 61.63%, *related animals* sebanyak 3.08%, *general bullying* sebanyak 19.57% dan *sexuality* sebanyak 3.08%. penelitian ini untuk periode November-Desember 2016 sehingga dari hasil klasifikasi tersebut dapat disimpulkan bahwa *cyberbullying* di Indonesia untuk periode tersebut cukup tinggi.

Kata kunci: *cybercrime, cyberbullying, bukti digital, jejaring sosial, data mining, Naïve Bayes*

ABSTRACT

Development of Sosial Media as one advances in technology not only bring positif impact but also introduce new problems when not used properly or violated undue, this case called Cybercrime. Twitter is one of the sosial networks in great demand nowadays because of ease to access that can be done through the smarthphone, laptop, tablet or a variety of other application services. The growth of sosial networking Twitter brings new trends in society as the arena to perform acts of oppression online, this called is cyberbullying.

Cyberbullying is one of cybercrime that are crowded today. The practice of cyberbullying is not just limited to children but also in adults and this is called cyberstalking or cyberharrasment. The impact to the offender can be sentenced to jail in accordance with the legislation in force, while from the side of psychology the victim's could get depressed even committed suicide. With the phenomenon, this research focus to analyze the development of cyberbullying in Indonesia through Sosial Media Twitter approach Data Mining with algorithm Naïve Bayes Classifier on the Machine Learning WEKA. Output of this research in form of knowledge about the development of cyberbullying on Sosial Media and Types of Cyberbullying widely used, so the Government or related institutions can undertake education and giving sanction expressly for communities and sosial media users who do the bullying to avoid more victims of bullying.

Based on the analysis results of classification can be a positif contains content of bullying as much 86.97%. A widely used type of cyberbullying is related psychology as much 61.63%, related animals as much 3.08%, general bullying as much 20% and sexuality as much 3.08%. This research for period November-December 2016 so that the classification of the results it can be concluded that cyberbullying in Indonesia for the period is quite high.

Keywords

Cybercrime, Cyberbullying, Digital Evidence, Sosial Media, Data mining, Naïve Bayes

Pernyataan keaslian tulisan

Dengan ini saya menyatakan bahwa tesis ini merupakan tulisan asli dari penulis, dan tidak berisi material yang telah diterbitkan sebelumnya atau tulisan dari penulis lain terkecuali referensi atas material tersebut telah disebutkan dalam tesis. Apabila ada kontribusi dari penulis lain dalam tesis ini, maka penulis lain tersebut secara eksplisit telah disebutkan dalam tesis ini.

Dengan ini saya juga menyatakan bahwa segala kontribusi dari pihak lain terhadap tesis ini, termasuk bantuan analisis statistik, desain survei, analisis data, prosedur teknis yang bersifat signifikan, dan segala bentuk aktivitas penelitian yang dipergunakan atau dilaporkan dalam tesis ini telah secara eksplisit disebutkan dalam tesis ini.

Segala bentuk hak cipta yang terdapat dalam material dokumen tesis ini berada dalam kepemilikan pemilik hak cipta masing-masing. Untuk material yang membutuhkan izin, saya juga telah mendapatkan izin dari pemilik hak cipta untuk menggunakan material tersebut dalam tesis ini.

Yogyakarta, Maret 2017



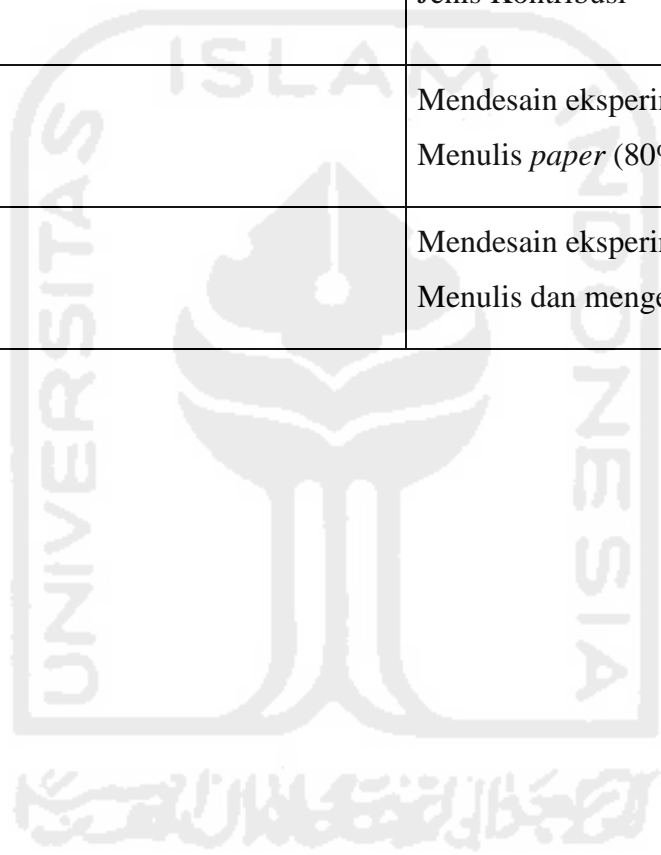
Hariani, S.Kom

Publikasi selama masa studi

Hariani., dan Riadi. I (2017). Detection of Cyberbullying on Sosial Media using Data Mining Techniques, IJCSIS vol. 15 No. 3

Publikasi yang menjadi bagian dari tesis

Kontributor	Jenis Kontribusi
Hariani	Mendesain eksperimen (80%) Menulis <i>paper</i> (80%)
Imam Riadi	Mendesain eksperimen (20%) Menulis dan mengedit <i>paper</i> (20%)



Kontribusi yang diberikan oleh pihak lain dalam tesis ini

Saran dan Masukan Oleh Pembimbing dan Penguji:

- Dr. Imam Riadi, M.Kom
- Yudi Prayudi, S.Si., M.Kom
- Dr. Bambang Sugiantoro, MT



Halaman Persembahan

BISMILLAHIRRAHMANIRRAHIM...

Segala Puji Syukur atas Nikmat Allah S.W.T atas semua kelancaran, kesehatan dan semua kebaikan yang diberikanNya sehingga saya bisa menyelesaikan Tesis ini, karya ini Special kupersembahkan kepada semua pihak yang telah memberi dukungan baik moril maupun materil selama pembuatan tesis ini.

- ♥ Suami tercinta yang tak pernah menyerah dalam doa dan nyata dan menemaniku berjuang hingga akhir.
- ♥ Orang Tua dan Mertua tersayang yang memberiku motivasi, dukungan dan restunya.
- ♥ Adik2ku tersayang yang memberi motivasi agar selalu berusaha menjadi kakak yang baik dan bisa jadi panutan.
- ♥ Teman-teman Kinanti Angels atas kebersamaan, persaudaraan dan dukungannya selama tinggal bersama.
- ♥ Keluarga besar Forensika Digital UII, special FD angkatan VIII, seru mengenal kalian..
- ♥ Kak Medy dan Kak Rose Akatsuki yang dengan setia antar jemput kemana-mana, thanks a lot of kak.
- ♥ Kucing-kucing tersayang atas hiburannya dengan kelucuannya Babon, Mini, Miki, Mina, Mono, Mina dan Miko.

Dengan Setulus Hati ini,
Hariani

Kata Pengantar

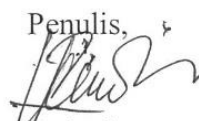
بِسْمِ اللَّهِ الرَّحْمَنِ الرَّحِيمِ

Alhamdulillah, segala puji bagi Allah *Subhanhu wa ta'ala* atas limpahan rahmat, hidayah, serta bimbingan-NYA. Shalawat dan salam tercurah kepada Nabi Muhammad *Shallallohu 'alaihi wa sallam*. Akhirnya penulis dapat menyelesaikan tesis ini dengan judul **“ANALISIS BUKTI DIGITAL CYBERBULLYING PADA JEJARING SOSIAL MENGGUNAKAN NAÏVE BAYES CLASSIFIER (NBC)”**, sebagai salah satu syarat untuk mendapatkan gelar Magister Komputer pada Program Pascasarjana (S-2) Fakultas Teknologi Industri, Universitas Islam Indonesia Yogyakarta. Pada kesempatan ini pula dengan segala kerendahan hati penulis mengucapkan terima kasih kepada:

1. Suamiku Ibnuhazm yang memberi dukungan sepenuhnya, terutama dukungan secara moril yang memberi semangat, motivasi dan kasih sayang tiada terkira.
2. Orang tuaku dan Mertuaku, atas segala doa, restunya. Semoga Allah merahmati kalian.
3. Bapak Dr. Imam Riadi, M.Kom selaku Dosen pembimbing, terimakasih pak atas kesabarannya selama membimbing saya.
4. Bapak Yudi Prayudi, S.Si, M.Kom selaku Dosen penguji
5. Bapak Dr. Bambang Sugiantoro, M.T selaku Dosen penguji
6. Ketua Program Pascasarjana FTI UII dan seluruh jajarannya.
7. Dosen Magister Teknik Informatika khususnya dosen Forensik Digital
8. Teman-teman seperjuangan Forensika Digital angkatan 8
9. Semua pihak yang memberi dukungan moril dan doa yang tidak bisa saya sebut satu persatu.

Semoga Allah memberikan balasan kebaikan atas segala bantuan yang diberikan kepada penulis. Tesis ini tentu masih jauh dari kata sempurna, untuk itu segala kritikan dan saran yang bersifat membangun guna menyempurnakan tesis ini sangat diharapkan. Semoga tesis ini dapat bermanfaat bagi kita semua. Aamiin.

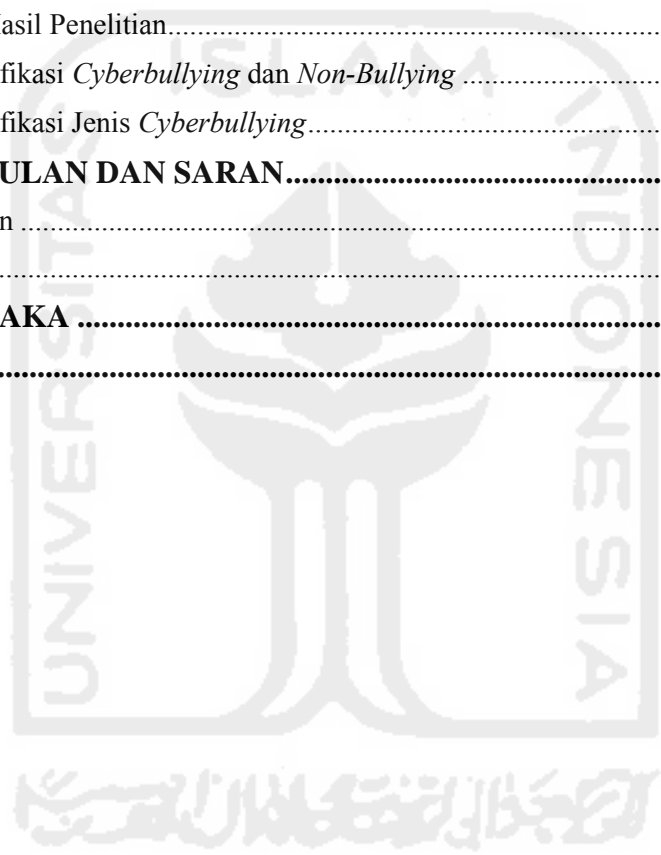
Yogyakarta, 27 Maret 2017

Penulis,

Hariani

DAFTAR ISI

LEMBAR PENGESAHAN PEMBIMBING	i
LEMBAR PENGESAHAN PENGUJI	ii
ABSTRAK	iii
ABSTRACT	iv
PERNYATAAN KEASLIAN TULISAN	v
PUBLIKASI SELAMA MASA STUDI	vi
KONTRIBUSI YANG DIBERIKAN OLEH PIHAK LAIN DALAM TESIS INI	vii
HALAMAN PERSEMBAHAN	viii
KATA PENGANTAR	ix
DAFTAR ISI	x
BAB 1 PENDAHULUAN	1
1.1 Latar Belakang	1
1.2 Rumusan Masalah	6
1.3 Batasan Masalah.....	6
1.4 Tujuan Penelitian	7
1.5 Manfaat Penelitian	7
1.6 Metodologi Penelitian	7
1.7 Sistematika Penulisan.....	8
BAB 2 LANDASAN TEORI	10
2.1 Kajian Penelitian Terdahulu.....	10
2.2 <i>Cybercrime</i>	16
2.3 Data Mining dan Digital Forensik.....	17
2.4 Jejaring Sosial	19
2.5 <i>Twitter</i>	20
2.6 <i>Cyberbullying</i>	21
2.7 Tinjauan Undang-Undang.....	22
2.8 Data Mining	23
2.9 <i>Text Mining</i>	25
2.10 <i>N-Gram</i>	27
2.11 Naïve Bayes Classifier	28
2.12 <i>Bag-of-words Model</i>	30
2.13 Teknik Pencarian Data	30
2.14 <i>API Twitter</i>	31
BAB 3 METODOLOGI PENELITIAN	33
3.1 Metodologi Penelitian	33
3.2 Perangkat Pendukung Penelitian.....	33

3.3	Pengumpulan Data <i>Log Tweet</i>	34
3.4	Preprocessing Data.....	40
3.4.1	Menghapus Special Karakter	41
3.4.2	Normalisasi Kalimat.....	42
3.4.3	Stemming.....	44
3.4.4	Penggunaan Lexicon	45
3.4.5	Data Training dan Data Testing	46
3.4.6	Pengolaan Data Menggunakan WEKA.....	48
3.5	Klasifikasi	51
BAB 4 HASIL DAN PEMBAHASAN.....		53
4.1	Deskripsi Penelitian	53
4.2	Evaluasi Hasil Penelitian.....	53
4.2.1	Klasifikasi <i>Cyberbullying</i> dan <i>Non-Bullying</i>	53
4.2.2	Klasifikasi Jenis <i>Cyberbullying</i>	55
BAB 5 KESIMPULAN DAN SARAN.....		61
5.1	Kesimpulan	61
5.2	Saran.....	62
DAFTAR PUSTAKA		63
LAMPIRAN		67



DAFTAR TABEL

Tabel 1. 1 Tabel Statistic Cyberbullying Beberapa Negara (Singhal 2013).....	3
Tabel 2. 1 Perbandingan Penelitian Terdahulu.....	13
Tabel 2. 2 Teknik Digital Forensik dan Data Mining.....	17
Tabel 3. 1 Tabel Find dan Replace.....	42
Tabel 3. 2 Contoh kata tidak baku menjadi baku.....	43
Tabel 3. 3 Pola Cyberbullying	47
Tabel 3. 4 Pola Negasi	47
Tabel 3. 5 Tabel Stopword Tala.....	49
Tabel 4. 1 Hasil Presentase Klasifikasi.....	59



DAFTAR GAMBAR

Gambar 1. 1 Grafik Analisis Cyberbullying (Bangsal 2013).....	3
Gambar 1. 2 Grafik Analisis Cyberbullying (2015)	4
Gambar 2. 1 Fase-fase Dalam Data Mining (Fayyad, 1996).....	24
Gambar 3. 1 Alur Penelitian	33
Gambar 3. 2 Teknik Pengumpulan Data.....	34
Gambar 3. 3 Access Token	35
Gambar 3. 4 Map Twitter.....	36
Gambar 3. 5 Hasil Ekstrak Data Setelah di Parsing.....	37
Gambar 3. 6 Script Program Pengumpulan Data.....	38
Gambar 3. 7 Hasil Validasi file Json Menggunakan JsonLint.....	40
Gambar 3. 8 Teknik Preprocessing Data	41
Gambar 3. 9 Script Stemming.....	44
Gambar 3. 10 Kamus Lexicon	45
Gambar 3. 11 Data Sebelum Penggunaan Lexicon.....	46
Gambar 3. 12 Data Setelah Penggunaan Lexicon.....	46
Gambar 3. 13 Interface WEKA.....	48
Gambar 3. 14 Flowchart Naive Bayes Classifier.....	51
Gambar 4. 1 Grafik Klasifikasi Data Training Menggunakan WEKA.....	53
Gambar 4. 2 Evaluasi Naive Bayes Pada Data Training Menggunakan WEKA.....	54
Gambar 4. 3 Prediksi Data Testing	55
Gambar 4. 4 Grafik Data Training Jenis Cyberbullying Menggunakan WEKA.....	56
Gambar 4. 5 Klasifikasi Data Training Jenis Cyberbullying	57
Gambar 4. 6 Prediksi Klasifikasi Jenis Cyberbullying	58
Gambar 4. 7 Grafik Prediksi Jenis Cyberbullying	58
Gambar 4. 8 Grafik Cyberbullying dan Jenis Cyberbullying	59

Bab 1 Pendahuluan

1.1 Latar Belakang

Perkembangan dan kemajuan teknologi selain membawa dampak positif juga memperkenalkan masalah-masalah baru saat digunakan secara tidak tepat atau menyalahi dari yang semestinya, hal ini sering disebut dengan kejahatan *cyber* atau *cybercrime*. Menurut (Hamzah 2012) *cybercrime* merupakan sebagai bentuk perbuatan melawan hukum yang dilakukan dengan menggunakan internet yang berbasis pada kecanggihan teknologi komputer dan telekomunikasi. Kementerian Komunikasi dan Informatika (Kemenkominfo) mengungkapkan pengguna Internet di Indonesia saat ini mencapai 63 juta orang. Dari angka tersebut 95 persennya menggunakan internet untuk mengakses Jejaring Sosial. Direktur Pelayanan Informasi Internasional Ditjen Informasi dan Komunikasi Publik (IKP), situs Jejaring Sosial yang paling banyak diakses adalah *Facebook* dan *Twitter*. Indonesia menempati peringkat 5 pengguna Twitter terbesar di dunia (di kutip dari www.kominfo.go.id, 2013).

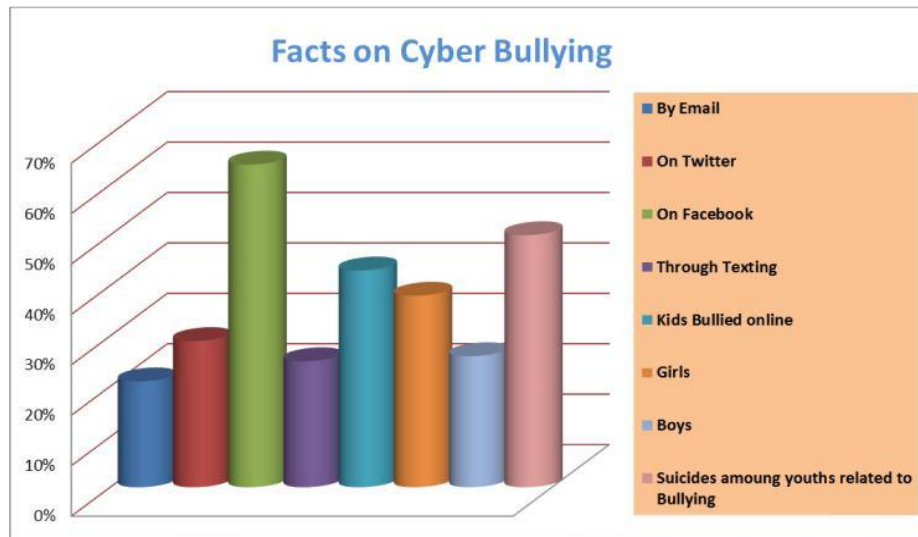
Pesatnya perkembangan jejaring sosial *Twitter* sebagai alat komunikasi yang mudah digunakan oleh siapa saja dan dapat diakses dimana saja membuat fenomena besar terhadap arus informasi, pertumbuhan jejaring sosial *Twitter* membawa *trend* baru dalam masyarakat sebagai ajang untuk melakukan tindakan penindasan secara online atau yang lebih dikenal dengan *Cyberbullying* (Yurnalita, 2016). *Cyberbullying* merupakan salah satu *cybercrime* yang sedang marak saat ini. *Cyberbullying* adalah tindakan seorang atau remaja secara sengaja mengintimidasi, mengancam, atau mempermalukan seseorang, atau sekelompok anak lain melalui teknologi informasi, seperti media sosial atau *mobile device* (Firman dan Ngazis, 2012). Sedangkan menurut *Urban Dictionary*, *Cyberbullying* melibatkan penggunaan teknologi informasi dan komunikasi seperti *e-mail*, telepon selular, *chat room*, dan jejaring sosial pribadi yang dilakukan secara sengaja dan berulang, dimaksudkan untuk menyakiti pihak lain (Yurnalita, 2016).

Banyaknya fenomena *Cyberbullying* di kalangan masyarakat yang mengakibatkan dampak negatif baik secara hukum maupun secara psikologi membuat peneliti tertarik untuk meneliti lebih dalam lagi mengenai *cyberbullying* di Indonesia. Secara hukum pelaku dapat dijerat sesuai undang-undang yang berlaku, sementara dari sisi psikologi terjadi pada korban

yang mengakibatkan depresi, sulit konsentrasi, merasa terisolasi, diperlakukan tidak manusiawi, penurunan kepercayaan diri, putusnya harapan, perasaan kesepian yang bisa mengakibatkan sampai melakukan bunuh diri (Rahayu 2012). Praktik *cyberbullying* tidak hanya terbatas pada anak-anak. Pada orang dewasa hal ini disebut *cyberstalking* atau *cyberharrasment*. Taktik yang umum digunakan oleh *cyberstalkers* adalah vandalisme mesin pencari atau ensklopedia, mengancam harta korban, pekerjaan, reputasi atau keselamatan. *Cyberstalking* adalah penggunaan komunikasi internet, e-mail atau elektronik lainnya dan umumnya mengacu pada pola perilaku mengancam atau berbahaya, sementara *Cyberharassment* biasanya berkaitan dengan mengancam atau melecehkan pesan *e-mail*, pesan instan, atau melalui *blog* atau *website* yang didedikasikan untuk menyiksa individu (Haryati 2014). Fenomena *bullying* di dunia maya, terjadi dibanyak kalangan, seperti dikalangan pelajar sekolah atau mahasiswa perguruan tinggi, hingga dengan kalangan selebritis, politikus, bahkan kalangan pejabat Negara.

Kasus *cyberbullying* yang pernah terjadi di Indonesia diantaranya adalah kasus Farhat Abbas S.H yang berkicau pada akun Twitternya mempermasalahkan penjualan plat mobil pribadi b 2 DKI yang dijual oleh polisi kepada umum. Berikut *tweets* Farhat Abbas “@farhatabbaslaw: Ahok sana sini plat pribadi B 2 DKI dijual polisi ke orang umum katanya ! Dasar Ahok plat aja diributin ! Apapun plat nya tetap cina”. Kicauan Farhat Abbas tersebut bersifat *Harrasment* dan berbau SARA karena mendiskriminasikan etnis dan ras tertentu (Yurnalita 2016). Florence Sihombing yang berdomisili di Yogyakarta memposting status yang menyinggung orang jogja melalui akun pribadinya di Path dengan mengatakan bahwa orang jogja “miskin, tolol, dan tak berbudaya”. Seseorang kemudian men *screenshot* statusnya dan menyebarkan ke media sosial. Hal ini menyebabkan Florence di *bully* di media sosial, dan diadukan ke pihak polisi oleh beberapa LSM yang ada di Yogyakarta. Kasus lainnya juga menimpa Walikota Bandung Ridwan Kamil, @Ridwankamil, men-*tweet* akan melaporkan pemilik akun @kemalsept, yang menghina Kota Bandung dengan sebutan kota yang penuh pelacur, selain penghinaan terhadap Kota Bandung @kemalsept juga menyebut Wali Kota Bandung Ridwan Kamil dengan kata “kunyuk” (Haryati 2014).

Kasus *Cyberbullying* lainnya yang terjadi pada Jejaring sosial dalam sebuah survey yang di lakukan oleh (Singhal, 2013) menunjukkan hasil survey berupa grafik analisis dari jejaring sosial diantaranya adalah *twitter*, *facebook* dan email, selain itu, hasil survey juga menganalisis pelaku dan korban *cyberbullying* baik laki-laki maupun perempuan serta penyebab dan akibat yang ditimbulkan dengan adanya kejahatan *bullying* tersebut. Adapun data statistik *cyberbullying* berdasarkan beberapa negara ditunjukkan dengan grafik pada **Gambar 1.1** dibawah ini:



Gambar 1. 1 Grafik Analisis Cyberbullying (Singhal, 2013)

Pada Gambar 1.1 menunjukkan survei presentase maksimum bullying dilakukan melalui media *Facebook* yaitu sekitar 70%, *Twitter* 27% dan minimum melalui *email* yaitu sekitar 25%. Sekitar 48% anak-anak mengalami gangguan bullying, 45% diantaranya adalah perempuan dan 30% anak-anak, sementara korban bunuh diri sekitar 55%. Selanjutnya, survei menunjukkan statistic *cyberbullying* dari berbagai negara yang mana Indonesia menempati urutan kedua setelah India. Persentase statistik dari beberapa Negara di tunjukkan pada Tabel 1.1:

Tabel 1. 1 Tabel Statistic Cyberbullying Beberapa Negara (Singhal 2013)

	My Child has experienced cyber bullying	A child in my community has experienced cyber bullying	Net parent awareness of cyber bullying in country
Total	12%	26%	38%
India	32%	45%	77%
Indonesia	14%	53%	47%
Sweden	14%	51%	65%
Canada	18%	31%	49%
Australia	13%	35%	48%
Brazil	20%	25%	45%
Saudia Arabia	14%	25%	41%
United States	15%	26%	41%
South Africa	10%	30%	41%
Turkey	5%	35%	40%
Mexico	8%	28%	36%
Argentina	9%	27%	36%
China	11%	25%	36%
Great Britain	11%	25%	36%
South Korea	8%	27%	35%
Poland	12%	20%	32%

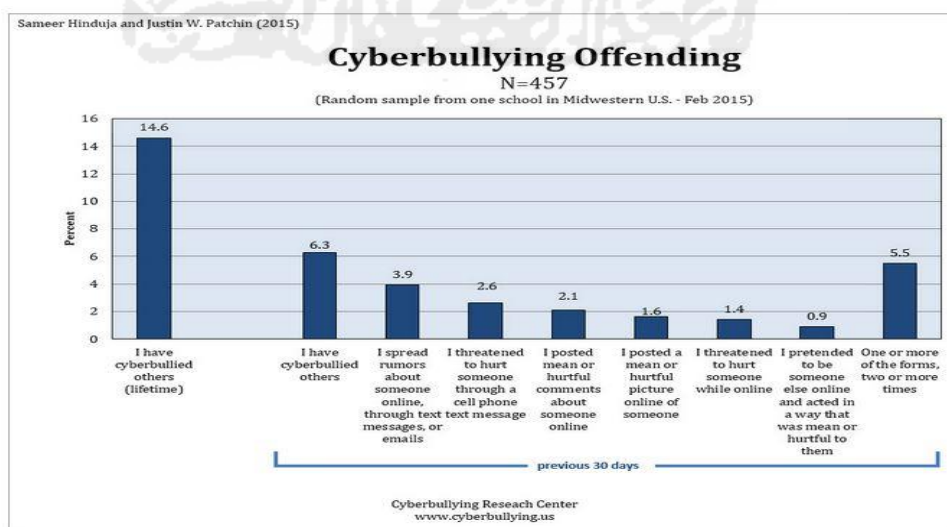
Tabel 1.1 Statistic Cyberbullying Beberapa Negara (Lanjutan)

	My Child has experienced cyber bullying	A child in my community has experienced cyber bullying	Net parent awareness of cyber bullying in country
Belgium	15%	13%	25%
Russia	5%	15%	20%
Germany	7%	12%	19%
Japan	7%	12%	19%
Hungary	7%	11%	18%
Italy	3%	15%	18%
Spain	5%	11%	16%
France	5%	10%	15%

Survey berikutnya dilakukan terhadap terhadap 126 pembaca majalah Kawanku dalam kutipan buku *Celebrate Your Wierdness* dalam (Nurjanah, 2014) bahwa:

- 69.84% pernah merasakan menjadi korban cyberbullying melalui twitter dan facebook
- 52% mengaku pernah menjadi pelaku cyberbullying di twitter dan facebook.
- Tindakan cyberbullying yang paling banyak diterima adalah cyberbullying berupa tulisan atau komentar, 53.97% mengaku pernah diejek dengan kata-kata kasar.
- Dampak yang paling nyata yang dialami korban adalah 38,10% mengaku merasa terasing dan merasa tidak punya teman.

Selanjutnya, menurut data terbaru yang dipelopori Hinduja dkk melalui statistic *Cyberbullying 2015* (<http://cyberbullying.org/statistics/>) yang melakukan penelitian terhadap sampel acak dari 457 siswa antara usia 11 dan 15 tahun dari sekolah menengah di Midwestern Amerika Serikat. dikumpulkan pada bulan Februari tahun 2015 seperti terlihat pada **Gambar 1.2** dibawah ini:



Gambar 1. 2 Grafik Analisis Cyberbullying (2015)

Menurut penelitian ini, lima dari enam remaja mengakui menggunakan perangkat selular mereka untuk mengakses jejaring sosial, seperti Instagram dan Facebook, kemudian twitter dan Ask.fm. Remaja tersebut berkisar 11 dan 15 tahun. Menurut penelitian tersebut, Sekitar 34% siswa mengalami cyberbullying dalam hidup mereka. Ketika ditanya tentang jenis cyberbullying yang dialami, data menunjukkan dalam waktu 30 hari, komentar yang menyakitkan (12,8%), pencemaran nama baik (19,4%). Sementara menurut jenis kelamin, remaja perempuan lebih cenderung mengalami *cyberbullying* dari pada laki-laki dengan perbandingan 40,6% dan 28,2%.

Di Indonesia, terdapat beberapa aturan yang sudah berlaku dan secara substansi dapat digunakan untuk menuntut pelaku *Cyberbullying*. Aturan tersebut adalah KUHP, Undang-Undang perlindungan anak, Undang-Undang Pornografi dan Undang-Undang Informasi dan Teknologi. Namun diantara keempat undang-undang tersebut, belum ada secara khusus membahas secara spesifik tindakan *cyberbullying* dan hukuman yang tepat untuk dijatuhkan jika pelaku diproses melalui jalur hukum. Beberapa pasal yang relevan dan memenuhi semua unsur yang lazimnya *cyberbullying* adalah pasal 27 ayat (1) dan (3) Undang-Undang Informasi dan Transaksi Elektronik. Pasal tersebut dapat dikenakan karena mengatur persebaran informasi dan dokumen elektronik yang bermuatan penghinaan dan kesusilaan, dua hal ini yang seringkali dijadikan senjata utama oleh pelaku *Cyberbullying* untuk menyerang korbannya (Taibah 2013).

Penelitian mengenai cyberbullying di Indonesia banyak ditinjau dari sisi sosiologi, ilmu komunikasi, hukum, psikologi dan informatika dengan metode yang hampir sama yaitu melakukan penelitian di berbagai sekolah SMP maupun SMA/SMK di Indonesia menggunakan angket dan interview langsung. Untuk penelitian di bidang teknologi Informasi masih sedikit, khususnya mengenai cyberbullying. Berdasarkan hal tersebut maka Penelitian ini menggunakan sudut pandang yang berbeda, karena *cyberbullying* banyak terjadi di jejaring sosial seperti uraian diatas maka peneliti tertarik untuk melakukan analisis *cyberbullying* pada jejaring sosial Twitter dengan melakukan pengambilan data *Log Tweet* langsung pada database Twitter. Pengumpulan *Log Tweet* dilakukan pada periode November sampai dengan Desember 2016. Metode yang digunakan adalah dengan pendekatan teknik *Data mining*. Data Mining digunakan untuk melakukan analisis pada *Log Tweet* yang telah dikumpulkan. secara sederhana data mining dapat diartikan sebagai ekstraksi informasi atau pola yang penting atau menarik dari data yang ada di *datatabase* yang besar. Dalam jurnal Ilmiah, *Data Mining* juga dikenal dengan nama *Knowledge Discovery in Databases* (KDDI) (Suchyo 2013).

Log Tweet yang telah di ekstrak dari database Twitter kemudian di proses menggunakan *Text Mining*, *Text Mining* merupakan penerapan konsep dan teknik *Data Mining* untuk

mencari pola dalam teks, proses penganalisisan teks guna menyarikan informasi yang bermanfaat untuk tujuan tertentu (Susanto 2010) sehingga teknik ini cocok diterapkan untuk menganalisis *Log Tweet* untuk mendapatkan *Knowledge* yang dibutuhkan. Ada Beberapa metode yang digunakan untuk proses *Teks Mining* diantaranya adalah *Naïve Bayes Classifier (NBC)*, *SVM*, *C45*, *K-Nearest Neighbor*, *K-Means* dan algoritma genetika. Semua algoritma ini mempunyai kelebihan dan kekurangan masing-masing tergantung kasus yang akan diselesaikan. Namun pada penelitian ini peneliti memilih menggunakan *Naive Bayes Classifier (NBC)* untuk mengklasifikasikan *Cyberbullying* pada jejaring sosial karena NBC sederhana tetapi memiliki akurasi yang tinggi. Berdasarkan hasil eksperimen, NBC terbukti dapat digunakan secara efektif untuk mengklasifikasikan berita secara otomatis dengan akurasi mencapai 90.23%. Algoritma NBC yang sederhana dan kecepatannya yang tinggi dalam proses pelatihan dan klasifikasi membuat algoritma ini menarik untuk digunakan (Saraswati 2011). Kasus lainnya, NBC dapat mengklasifikasikan tweet yang berisi informasi kemacetan lalu lintas dengan akurasi 93.58% (Rodiyansyah 2012), penelitian selanjutnya, NBC dapat diterapkan untuk menilai kelayakan kredit pada BCA Finance dengan akurasi yang baik sehingga termasuk kategori *Excellent*, kelayakan kredit pada BCA menghasilkan akurasi akhir 92,54% (Ciptohartono 2014). NBC juga efektif diterapkan untuk mengidentifikasi *email spam*, dengan menggunakan beberapa syarat dan kondisi, seperti klasifikasi yang dilakukan secara online dan offline tingkat error lebih kecil, sementara *error* besar jika terdapat selisih pada jumlah keyword yang ada pada data training, namun secara keseluruhan penerapan NBC untuk identifikasi spam email ini cukup akurat (Anugroho 2016).

1.2 Rumusan Masalah

Berdasarkan uraian pada latar belakang maka dapat dirumuskan suatu permasalahan yaitu:

1. Bagaimana teknik dalam mengumpulkan bukti digital pada Jejaring Sosial *Twitter*?
2. Bagaimana proses klasifikasi bukti digital *cyberbullying* pada Jejaring Sosial *Twitter* menggunakan *Naïve Bayes Classifier*?

1.3 Batasan Masalah

Untuk menjaga fokus dalam penelitian maka beberapa batasan yang diberikan dalam penelitian adalah:

1. Penelitian ini menggunakan Jejaring sosial *Twitter* dengan *library API Developer* untuk mengekstrak *Log Twitter*.
2. Penelitian hanya berfokus pada komentar berbahasa Indonesia

3. Penelitian ini hanya berfokus pada komentar atau text yang mengandung *bullying* bukan pada gambar yang mengandung *bullying*.
4. Metode yang digunakan untuk mengklasifikasikan adalah *Naïve Bayes Classifier* dari *machine learning* WEKA.

1.4 Tujuan Penelitian

Tujuan dari penelitian ini adalah:

1. Mengembangkan teknik dalam mengumpulkan bukti digital *Log Tweet* pada Jejaring Sosial Twitter.
2. Mengklasifikasikan bukti digital *Log Tweet* untuk mendapatkan sebuah *knowledge* berupa informasi perkembangan *Cyberbullying* di Indonesia berdasarkan jenisnya menggunakan *Naïve Bayes Classifier*.

1.5 Manfaat Penelitian

Adapun manfaat dari penelitian ini adalah:

1. Memberi informasi kepada pihak-pihak yang membutuhkan mengenai perkembangan *Cyberbullying* di Indonesia untuk dijadikan referensi khususnya bagi pendidik, psikolog, pemerintah dan orang tua. Dengan Informasi tersebut pihak yang dimaksud dapat memberikan edukasi yang benar bagi pelaku maupun korban *cyberbullying* disekitarnya.
2. Menambah pengetahuan dalam menganalisis bukti digital *Log Tweet* dengan teknik klasifikasi menggunakan *Naïve Bayes Classifier*.

1.6 Metodologi Penelitian

Penelitian ini dilakukan dengan tahap-tahap sebagai berikut:

1. Studi kepustakaan
Pengumpulan bahan referensi, seperti jurnal penelitian, prosiding, tesis, buku-buku teori dan sumber-sumber lain termasuk informasi yang diperoleh melalui internet.
2. Pengumpulan Data
Pengumpulan data dilakukan dengan mengekstrak *Log* data dari jejaring sosial Twitter kemudian melakukan proses *preprocessing* guna memenuhi kebutuhan data yang diinginkan.
3. Pemeriksaan Data
Pada tahap ini peneliti melakukan pemeriksaan data dari *log Tweet* yang telah dikumpulkan apakah data tersebut sudah relevan dengan yang diinginkan.

4. Analisis Data

Pada tahap ini dilakukan analisis data dengan melakukan analisa kebutuhan-kebutuhan untuk proses klasifikasi. Adapun proses klasifikasi dilakukan dengan langkah-langkah sebagai berikut:

- a. Menghapus duplikat *ID Tweet*
- b. Menghapus *Special character* (tanda baca: .,?!;:'[}]”&%)
- c. Menghapus *url*
- d. Menghapus Retweet (RT), namun teks RT tidak dihapus.
- e. Menghapus gambar
- f. Menghapus *Hastag*
- g. Melakukan Normalisasi pada teks, bahasa slang dijadikan baku contoh kata “Setaaan” kemudian di ubah menjadi “setan”
- h. Mengganti sinonim (yg=yang, mn=mana).
- i. Melakukan proses *stemming*
- j. Menggunakan *Lexicon*
- k. Menghilangkan *stopword*.
- l. *Casefolding*
- m. Menerapkan *N-Grams*

5. Klasifikasi

Mengimplementasikan algoritma *Naïve Bayes Classifier* dengan menggunakan *machine learning* WEKA pada data yang sudah di bersihkan pada proses *Preprocessing*.

1.7 Sistematika Penulisan

Dalam penyusunan penelitian ini, sistematika penulisan terbagi dalam beberapa bab yaitu:

BAB I PENDAHULUAN

Pendahuluan, merupakan pengantar terhadap permasalahan yang akan dibahas. Di dalamnya menguraikan tentang gambaran suatu penelitian yang terdiri dari latar belakang, rumusan masalah, batasan masalah, tujuan penelitian, manfaat penelitian, metodologi penelitian, serta sistematika penulisan.

BAB II LANDASAN TEORI

Pada Bab ini menjelaskan teori-teori yang digunakan untuk mendukung dalam memecahkan masalah pada penelitian ini.

BAB III METODOLOGI PENELITIAN

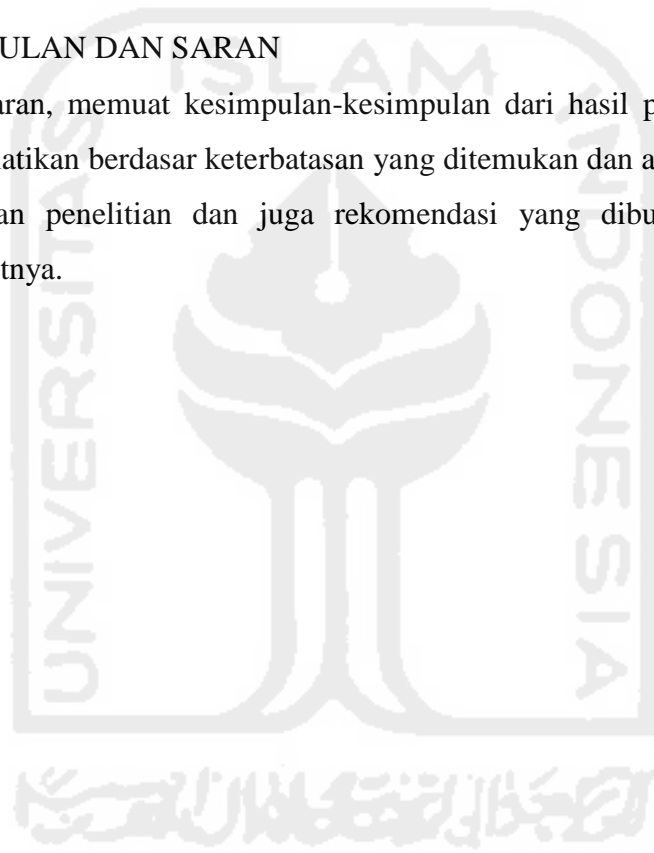
Bab ini membahas tentang langkah-langkah penelitian, kebutuhan perangkat keras dan perangkat lunak yang akan digunakan dalam proses klasifikasi untuk mencapai tujuan yang diinginkan.

BAB IV HASIL DAN PEMBAHASAN

Hasil dan Pembahasan, berisi tentang analisis dari algoritma yang digunakan untuk pengujian probabilitas klasifikasi data komentar *cyberbullying* pada Twitter.

BAB V KESIMPULAN DAN SARAN

Simpulan dan Saran, memuat kesimpulan-kesimpulan dari hasil penelitian dan saran-saran yang perlu diperhatikan berdasar keterbatasan yang ditemukan dan asumsi-asumsi yang dibuat selama melakukan penelitian dan juga rekomendasi yang dibuat untuk pengembangan penelitian selanjutnya.



Bab 2 Landasan Teori

2.1 Kajian Penelitian Terdahulu

Orebaugh dan Allnut (2009) melakukan penelitian dalam analisis forensik untuk mengidentifikasi pengguna *Instant Messenger* (IM) berdasarkan tingkah laku pengguna IM. Karakteristik yang digunakan yaitu *stylometric features* dengan kategori identifikasi kata menggunakan kata singkatan, struktur kalimat, karakter khusus, dan emoticon. Data yang digunakan dalam penelitian ini terdiri dari *Log* percakapan IM dengan 4 kategori *class* user yaitu user1, user2, user3 dan user4. Data tersebut diuraikan untuk menghitung struktur kalimat dan frekuensi yang telah ditetapkan pada kategori karakter khusus, emoticon, dan singkatan. Selanjutnya untuk klasifikasi digunakan metode *J48 decision tree, IBk nearest neighbor, and Naïve Bayes classifiers*. Hasilnya, akurasi pada Naif Bayes dengan atribut singkatan (97.85%), J48 (97.86%) dan IBk classifier (97.14%), hasil ketika semua atribut dikombinasikan menggunakan Naif Bayes memberikan akurasi terbaik yaitu 99,29%.

Penelitian Dinakar (2012) untuk mendeteksi *cyberbullying* pada jejaring sosial *Youtube* dan *Formspring.me* yang merupakan situs jejaring sosial yang populer di kalangan remaja. Penelitian ini menggunakan pendekatan state-of-the-art NLP dan sense knowledge base yang mana topiknya sangat luas pada kehidupan sehari-hari. Subyek yang diteliti seperti penampilan, *knowledge, racial* dan *ethnic shrul*. Penelitian ini menghasilkan aplikasi yang disebut *Sosial Network Dashboard* untuk menganalisis dan mendeteksi *Cyberbullying*. Contoh penggunaannya pada *Youtubube* yang mana komentar pada *youtube* dikategorikan dalam tiga kelompok pada 1500 komentar yang telah terunduh, diantaranya yang mengandung *cyberbullying* adalah kelompok *sexuality* 627 komentar, *Race and culture* 841 komentar dan *Intelligence* 809 komentar.

Penelitian Nalini dan Sheela (2014), mengusulkan sebuah framework berbasis *contextual* dan *word level feature* untuk mendeteksi *content* yang menyinggung dan mengidentifikasi pelaku *user* didalam *IRC logs*. Metode ini dibagi menjadi dua fase yaitu pertama, untuk mendeteksi level kata yang menyinggung, kedua memperoleh level user sebagai pelaku. Pada fase pertama, peneliti menerapkan teknik NLP seperti *word level feature* dan *contextual level feature*. Untuk fase kedua, peneliti menerapkan *user level feature* dengan

menggunakan *style*, *sturcture* dan fitur *cyberbullying*. Hasil dari metode tersebut akan di uji dengan menggunakan algoritma *classification* yang disediakan oleh WEKA, diantaranya *Random Forest*, *J48* dan *Sequential Minimal Optimization*.

Setty et al (2014) melakukan penelitian pada *facebook* untuk mengklasifikasikan berita ke beberapa kategori seperti *friends posts* dan *liked pages posts*. Kemudian *Friends posts* dikategorikan lagi menjadi *life events posts* dan *entertainment posts* pada halaman facebook. Penerapannya bertujuan untuk mengevaluasi klasifikasi *news feed* dan *sentiment analys* agar halaman berita facebook lebih terorganaisir dan terlihat lebih menarik. Penelitian ini Menggunakan aplikasi WEKA dan *Learning Model* untuk membandingkan beberapa algoritma klasifikasi pada 2000 data yang diekstrak dari facebook. Hasilnya, akurasi SVM dan Bayes Net lebih baik dari pada algoritma yang lain.

Kansara (2015), penelitian ini membangun sebuah *framework* dengan menggabungkan analisa teks dan gambar untuk mendeteksi ancaman *Cyberbullying*. Metode yang digunakan untuk analisa gambar adalah kombinasi antara *bag of visual word* (BoVW), *Local binary Pattern* (LBP) dan *SVM Classifier*. Sedangkan untuk analisa text menggunakan kombinasi *bag of word* (BoW) dan *Naïve Bayes Classifier*. Hasil analisa dari gambar dan text ini di proses lagi menggunakan *Boolean System* untuk menentukan apakah conten tersebut mrengandung *cyberbullying* atau tidak.

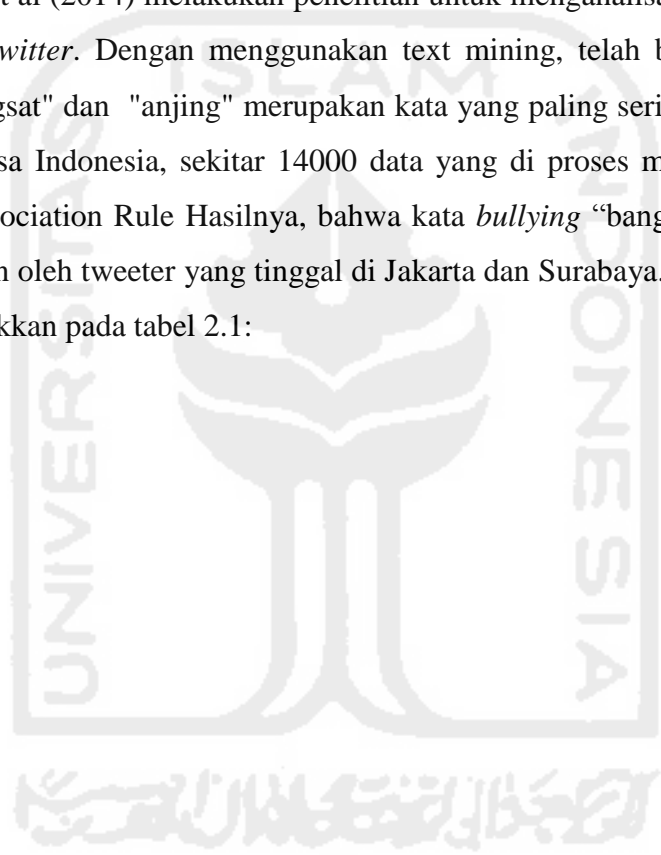
Penelitian selanjutnya, Hafilizara dan Adisantoso (2014) melakukan klasifikasi email spam dengan teknik Naïve Bayes, Pada teknik Naïve Bayes tersebut terdapat metode *smoothing* yang telah umum digunakan yaitu *Add-One smoothing* atau *Laplace smoothing*, namun masih ada metode lainnya yaitu *Jelinek-Mercer* (JM) *smoothing*, *Dirichlet* (Dir) *smoothing*, *Absolute Discounting* (AD) *smoothing*, dan *Two-Stage* (TS) *smoothing* yang dapat digunakan sebagai fungsi klasifikasi dan diduga mampu meningkatkan akurasi melebihi *Laplace smoothing*. Penelitian ini mengimplementasikan fungsi klasifikasi metode *smoothing* tersebut. Pada hasil percobaan terlihat akurasi yang dihasilkan metode *Laplace smoothing* lebih rendah dari metode *smoothing* lainnya. Dari hasil evaluasi terhadap nilai akurasi, *false rate*, dan *miss rate* terlihat metode *smoothing* Dirichlet memiliki nilai *miss rate* yang paling rendah sebesar 1.4%, nilai *false rate* 13.4%, dan akurasi 94.82%.

Penelitian Hosseinmardi et al (2015) pada sosial media Instagram menggunakan *Snowbal Sampling Method* telah mengidentifikasi 41K Ids user. Sekitar 61% user mempunyai *profile* yang *public* jadi sekitar 25K user. Setiap *user public* dikumpulkan dengan postingan gambar beserta 150 komentar, *id user* yang meng-*follow*, *id user* yang di-*follow*, *id user* yang mengomentari atau menyukai media yang di *share* oleh user tersebut. Hasilnya, 48% dari media tersebut bukan *cyberbullying*, melainkan hanya *cyberagression*, namun menggunakan

Linear SVM classifier dapat meningkatkan akurasi hingga 87% dengan menggabungkan fitur teks, gambar dan meta data dari media tersebut.

Akaichi (2013) melakukan penelitian dengan fokus pada penggunaan *teks mining* untuk *sentiment classification*. Ilustrasi dilakukan pada pengguna *Facebook* di Tunisia. metode yang diusulkan didasarkan pada dukungan Vektor mesin (SVM) dan Naif Bayes. peneliti juga membangun sebuah lexicon sentimen, berdasarkan emoticon, kata bahasa dan akronim, dari status data yang telah diekstrak. Selanjutnya, membandingkan dua *machine learning* yaitu algoritma SVM dan Naif Bayes melalui model pelatihan untuk klasifikasi sentiment. Hasilnya, data set dari data training 60% dan sisanya yaitu data test 40%.

Margono et al (2014) melakukan penelitian untuk menganalisa kata bullying berbahasa Indonesia pada *twitter*. Dengan menggunakan text mining, telah berhasil mengidentifikasi bahwa kata "bangsat" dan "anjing" merupakan kata yang paling sering digunakan untuk pola bullying berbahasa Indonesia, sekitar 14000 data yang di proses menggunakan metode Fp-Growth dan Association Rule Hasilnya, bahwa kata *bullying* "bangsat" dan "anjing" paling banyak digunakan oleh tweeter yang tinggal di Jakarta dan Surabaya. Perbandingan penelitian terdahulu ditunjukkan pada tabel 2.1:



Tabel 2. 1 Perbandingan Penelitian Terdahulu

No	Nama Peneliti	Objek	Judul Penelitian	Metode Penelitian	Hasil Penelitian
1.	Dinakar et al (2012)	Youtube dan Formspring.me	<i>Common Sense Reasoning for Detection, Prevention, and Mitigation of Cyberbullying</i>	NLP dan Common Sense Knowledgege Base	Menghasilkan aplikasi yang diberi nama <i>Sosial Network Dashboard</i> dan salah satu aplikasinya digunakan untuk Mendeteksi <i>Cyberbullying</i> dari komentar <i>Youtube</i> . Dari 15000 ribu komentar di kategorikan ke dalam tiga kelompok <i>cyberbullying</i> yaitu <i>Sexuality</i> sebanyak 627 komentar, <i>Race & Culture</i> 841 komentar dan <i>Inteliigence</i> 809 komentar.
2.	Orebaugh dan Allnutt (2009)	Instant Messanging	<i>Classification of Instant Messaging Communications for Forensics Analysis</i>	<i>J48</i> , <i>IBK classifier</i> dan <i>Naïve Bayes Classifier</i>	Data yang digunakan dalam penelitian ini terdiri dari <i>Log</i> percakapan IM dengan 4 kategori <i>class</i> user yaitu user1, user2, user3 dan user4. Data tersebut diuraikan untuk menghitung struktur kalimat dan frekuensi yang telah ditetapkan pada kategori karakter khusus, emoticon, dan singkatan. Selanjutnya untuk klasifikasi digunakan metode <i>J48 decision tree</i> , <i>IBk nearest neighbor</i> , and <i>Naïve Bayes classifiers</i> . Hasilnya, akurasi pada Naif Bayes dengan atribut singkatan (97.85%) , <i>J48</i> (97.86%) dan <i>IBk classifier</i> (97.14%), hasil ketika semua atribut dikombinasikan menggunakan Naif Bayes memberikan akurasi terbaik yaitu 99,29%.
3.	Hendro Margono. et al (2014)	Twitter	<i>Mining Indonesian Cyber Bullying Patterns in Sosial Networks</i>	FP-Growth, Assosiation Rule,	Data yang digunakan sebanyak 14000 <i>tweet</i> . Dengan fokus pada <i>tweet</i> bullying yang menggunakan bahasa Indonesia. Hasilnya, bahwa kata <i>bullying</i> “bangsat” dan “anjing” paling banyak digunakan oleh tweeter yang tinggal di Jakarta dan Surabaya.
3.	Shankar Setty. et al (2014)	Facebook	<i>Classification of Facebook News Feeds and Sentiment Analysis</i>	Menggunakan Algoritma BayesNet, <i>J48</i> , <i>Naïve Bayes</i> , <i>SVM</i> , <i>Logistic Regression</i> dari WEKA.	Menggunakan aplikasi WEKA dan <i>Learning Model</i> untuk membandingkan beberapa algoritma klasifikasi pada 2000 data yang diekstrak dari facebook. Hasilnya, akurasi <i>SVM</i> dan <i>Bayes Net</i> lebih baik dari pada algoritma yang lain.

Tabel 2.1 Perbandingan Penelitian Terdahulu (Lanjutan)

No	Nama Peneliti	Objek	Judul Penelitian	Metode Penelitian	Hasil Penelitian
4.	Kansara et al (2015)	Jejaring Sosial (Twitter, Facebook, Instagram dll)	<i>A Framework for Cyberbullying Detection in Sosial Network</i>	Untuk deteksi Gambar menggunakan metode <i>Bag of visual word (BoVW)</i> , <i>Local Binnary Pattern (LBP)</i> , dan <i>SVM Classification</i> . Untuk text metode yang digunakan yaitu <i>bag of word (BoW)</i> dan <i>Naïve Bayes Classifier</i> .	Penelitian ini membangun sebuah <i>Framework</i> untuk mendeteksi <i>Cyberbullying</i> menggunakan metode <i>Boolean System</i> dengan cara mengevaluasi hasil gambar dan hasil analisa teks.
5.	Mutia hafilizara dan Julio adisantoso (2014)	Email	Metode Smoothing dalam Naïve Bayes untuk Klasifikasi Email Spam.	Naïve Bayes, <i>Jelinek-Mercer (JM) smoothing</i> , <i>Dirichlet (Dir) smoothing</i> , <i>Absolute Discounting (AD) smoothing</i> , dan <i>Two-Stage (TS)</i>	Data email yang digunakan adalah korpus <i>email public</i> dengan 1396 <i>spam</i> , 250 <i>hard ham</i> , dan 1400 <i>ham</i> . Dari hasil evaluasi terhadap nilai akurasi, <i>false rate</i> , dan <i>miss rate</i> terlihat metode <i>smoothing Dirichlet</i> memiliki nilai <i>miss rate</i> yang paling rendah sebesar 1.4% , nilai <i>false rate</i> 13.4%, dan akurasi 94.82%.

Tabel 2.1 Perbandingan Penelitian Terdahulu (Lanjutan)

No	Nama Peneliti	Objek	Judul Penelitian	Metode Penelitian	Hasil Penelitian
6.	Hosseinmardi et al (2015)	Instagram	<i>Detection of Cyberbullying Incidents on The Instagram Sosial Network</i>	Metode yang digunakan untuk identifikasi user pada instagram menggunakan <i>Snowball sampling method</i> . Untuk klasifikasi menggunakan <i>Naive Bayes</i> dan <i>Linear SVM</i> .	Penelitian ini menggabungkan label komentar dan gambar. Hasil yang diperoleh, sekitar 48% dianggap bukan <i>cyberbullying</i> , sebagian menunjukkan hanya <i>cyberaggression</i> dan bukan <i>cyberbullying</i> , tetapi ketika menggunakan <i>Linear SVM</i> akurasi meningkat menjadi 87% dengan memasukkan multi-modal dari text, gambar, dan meta data dari media tersebut.
7.	Jalel Akaichi (2013)	Facebook	<i>Sosial Networks' Facebook' Statutes Updates Mining for Sentiment Classification</i>	Algoritma yang digunakan adalah <i>Naive Bayes multinomial event driven model</i> , <i>SVM</i> , <i>SMO (Sequential Minimal optimization)</i> .	Data yang digunakan sebanyak 260 status, dibagi menjadi tiga <i>lexicon</i> yaitu <i>Interjections</i> , <i>acronym</i> , dan <i>emoticons</i> . Hasilnya, jika menggunakan bigram dalam <i>Naive bayes</i> akurasinya tinggi, dan jika menggunakan unigram dalam <i>SVM</i> akurasinya lebih tinggi daripada algoritma <i>Naive Bayes</i> .
8.	Nalini et al (2014)	<i>Sosial Networking sites dan web site</i>	<i>A Survey on Datamining in Cyberbullying</i>	<i>Naive Bayes, SVM, Classifier term frequency</i> dan <i>tf-idf</i> .	Penelitian ini mengusulkan metode secara <i>contextual</i> dan <i>Level word</i> untuk mendeteksi <i>content</i> yang menyinggung dan <i>user</i> yang menyinggung. Sebanyak 325 transkrip yang digunakan sebagai <i>dataset</i> . Pengujiannya dilakukan dengan membandingkan beberapa algoritma klasifikasi yang di sediakan oleh WEKA. Metode ini dianggap efektif untuk mendeteksi <i>cyberbullying</i> pada <i>Sosial Network</i> dan <i>Web site</i> .
9	Hariani (2015)	Twitter	<i>Analisis Bukti Digital Cyberbullying Pada Jejaring Sosial Menggunakan Naive Bayes Classifier (NBC)</i>	<i>Tf-Idf, 10 fold cross validation, Naive Bayes Classifier</i>	Metode <i>Tf-Idf</i> diterapkan pada data yang telah dibersihkan, kemudian menggunakan pengujian <i>10 fold cross validation</i> sebagai validasi data, selanjutnya dilakukan klasifikasi menggunakan <i>Naive Bayes</i> untuk mendapatkan informasi perkembangan <i>Cyberbullying</i> dan jenis <i>cyberbullying</i> yang biasa digunakan di Indonesia.

2.2 Cybercrime

Cybercrime merupakan bentuk-bentuk kejahatan yang timbul karena pemanfaatan teknologi internet. Menurut Hamzah (2012) dalam tulisannya mengartikan kejahatan komputer sebagai: "Kejahatan di bidang komputer secara umum dapat diartikan sebagai penggunaan komputer secara illegal". Dari pengertian di atas, secara ringkas dapat dikatakan bahwa *cybercrime* dapat didefinisikan sebagai perbuatan melawan hukum yang dilakukan dengan menggunakan internet yang berbasis pada kecanggihan teknologi komputer dan telekomunikasi. Selama ini dalam kejahatan konvensional, dikenal adanya dua jenis kejahatan sebagai berikut:

a. Kejahatan kerah biru (*blue collar crime*)

Kejahatan ini merupakan jenis kejahatan atau tindak kriminal yang dilakukan secara konvensional seperti misalnya perampokan, pencurian, pembunuhan dan lain-lain.

b. Kejahatan kerah putih (*white collar crime*)

Kejahatan jenis ini terbagi dalam empat kelompok kejahatan, yakni kejahatan korporasi, kejahatan birokrat, malpraktek, dan kejahatan individu.

Perkembangan penggunaan Internet telah menggeser banyak aktivitas masyarakat ke dunia maya. Kemajuan teknologi ini memberikan dampak positif dari segi kemudahan masyarakat untuk memperoleh pengetahuan dan informasi. Namun perkembangan internet tidak selamanya menghasilkan hal-hal positif. Internet dalam pengaruh negatifnya muncul dalam kasus-kasus semacam *cyberbullying*, *hoax*, *phishing* sampai pada *cybercrime*. Beberapa kasus dampak negatif internet antara lain (Haryati, 2014):

1. Pencurian dan penggunaan account internet milik orang lain. Salah satu kesulitan dari sebuah TSP (*Internet Service Provider*) pelanggan mereka yang "dicuri" dan digunakan secara tidak sah adalah adanya *account*.
2. Membajak situs web, Salah satu kegiatan yang sering dilakukan oleh *cracker* adalah mengubah halaman *web*, yang dikenal dengan istilah *deface*.
3. Probing dan port scanning, Salah satu langkah yang dilakukan *cracker* sebelum masuk ke *server* yang ditargetkan adalah melakukan pengintaian. Cara yang dilakukan adalah dengan melakukan "*port scanning*" atau "*probing*" untuk melihat servis-servis apa saja yang tersedia di *server target*.
4. *Virus*, seperti halnya di tempat lain, virus computer pun menyebar di Indonesia. Penyebaran umumnya dilakukan dengan menggunakan email.
5. Kejahatan yang berhubungan dengan nama *domain*. Nama domain (*domain name*) digunakan untuk mengidentifikasi perusahaan dan merek dagang.

6. *Spam*, yaitu email sampah yang berisi penawaran produk/jasa, penipuan berkedok kerja sama bisnis, penyebaran virus, dll.
7. *Hoax*, Informasi palsu yang menyesatkan (*Hoax*).
8. *Illegal Content*. Konten ilegal, seperti perjudian dan pornografi.
9. *Phising*, Pencurian identitas untuk melakukan pencurian, penipuan, dll.
10. *Cyberbullying*. Kekerasan menggunakan media internet, seperti pelecehan, ancaman, dan tuduhan.
11. *Copyright Plagiarism*. hak cipta. *infringement*, Pelanggaran.

2.3 Data Mining dan Digital Forensik

Data mining merupakan aplikasi yang terdiri dari beberapa algoritma untuk mengekstrak pola data pengetahuan yang berguna dalam pengambilan keputusan. Data mining memiliki beberapa aplikasi digital forensic diantaranya adalah korelasi dalam forensic data (*association*), menemukan dan menyortir data ke dalam kelompok tertentu (*classification*), mengidentifikasi kelompok tertentu (*clustering*) dan menemukan pola data untuk memprediksikan sesuatu (*forecasting*). Data Mining tidak terbatas dalam digital forensic, *tools* dan model yang dikembangkan dapat membantu penyidik untuk menemukan bukti digital lebih efisien dan lebih cepat (Kayarkar 2014). Klasifikasi Tools digital forensic dalam data mining dapat dilihat pada table 2.2 dibawah ini:

Tabel 2. 2 Teknik Digital Forensik dan Data Mining

Digital Forensic Techniques	Data Mining Techniques	Tools
Data Recovery, Data generation and Pre-Processing	Statistical test Analysis Bartlett's test of sphericity Kaiser-Meyer-Olkin (KMO)	Recuva FTK Encase Sleuth kit/Autopsy ProDiscover
Data Analysis	Clustering –K-means, EM, Hierarchical Clustering	Weka
	Classification Supervised learning - Decision Tree, Neural Networks, SVM, Naïve Bayesian	Weka
	Unsupervised learning PCA, Karnohuen Map	-

Tabel 2.2 Teknik Digital Forensics dan Data Mining (Lanjutan)

Digital Forensic Techniques	Data Mining Techniques	Tools
Data Analysis	Frequent Pattern Mining/Association rule Mining - Apriori, Eclat	Weka
	Named Entity recognition	Lingepipe
	Visualization	Cyber Forensics Time Lab
	Statistical Analysis and Anomaly Detection	EMT/MET
	Recursive data Mining	-
	Phishing	Invisible Witness
	Regression	-

Digital Forensik merupakan aplikasi bidang ilmu pengetahuan dan teknologi komputer untuk kepentingan pembuktian hukum (*pro justice*), yang dalam hal ini adalah untuk membuktikan kejahatan berteknologi tinggi atau *computer crime* secara ilmiah *scientific* hingga bisa mendapatkan bukti-bukti digital untuk menjerat pelaku kejahatan tersebut. Pencarian bukti-bukti digital untuk menjerat pelaku sering kali merupakan pekerjaan yang sangat kompleks dimana seorang *digital forensics analyst/investigator* harus mengikuti prosedur-prosedur yang diakui secara hukum baik nasional maupun Internasional, termasuk juga mereka harus memahami secara teoritis hal-hal yang berkaitan dengan bukti digital yang ditemukan, disamping juga memahami bagaimana penggunaan *software-software* forensic untuk mencari bukti-bukti digital tersebut dengan benar (Nuh, 2012).

Bukti digital diekstrak dari barang bukti elektronik. Beberapa contoh bukti digital yang dapat dianalisis adalah:

1. *Logical File*, yaitu file-file yang masih ada dan tercatat di file system di suatu partisi. Contoh: file aplikasi, library, office logs, multimedia dan lain-lain.
2. *Deleted File*, disebut juga *unallocated cluster* yang merujuk pada tempat penyimpanan yang telah terhapus.
3. *Lost File*, yaitu file yang sudah tidak tercatat lagi di file system yang sedang berjalan di suatu partisi.
4. *File slack*, yaitu sector penyimpanan yang berada diantara *end of file* (EoF) dengan *End of Cluster* (EoC).

5. *Log File*, yaitu file-file yang merekam aktivitas dari suatu keadaan tertentu misalnya log dari system operasi, internet browser, aplikasi, internet traffic, jejaring sosial dan lain-lain.
6. *Encrypted File*, yaitu file yang isinya sudah dilakukan enkripsi dengan menggunakan algoritma kriptografi sehingga tidak bias dilihat atau dibaca secara normal.

Dalam melakukan tindakan penegak hukum mempunyai acuan yang benar dan prosedural di dalam melakukan investigasi *computer crime* dan *computer-related crime*, serta memeriksa dan menganalisis barang bukti mempunyai *guidelines* yang telah diakui secara Internasional diantaranya adalah (Nuh, 2012):

1. *Good Practice Guide For Computer-Based Electronic Evidence* yang dikeluarkan *Association Of Chief Police Officers (ACPO)* yang merupakan asosiasi para pimpinan kepolisian di Inggris yang bekerja sama dengan 7safe.
2. *Forensics Examination of Digital Evidence: Guide for Law Enforcement*, yang dikeluarkan oleh *National Institute Of Justice* yang berada di bawah *U.S. Department of Justice*.
3. *Electronic Crime Scene Investigation: A Guide for First Responders*, yang juga dikeluarkan oleh *National Institute of Justice* yang berada di bawah *U.S. Departement of Justice*.

2.4 Jejaring Sosial

Seperti halnya di dunia nyata, menjalin hubungan persahabatan bisa juga dilakukan di dunia maya (internet). Bedanya, melalui internet ini tidak bisa bertatap muka secara langsung untuk berjabat tangan dan menanyakan siapa namanya. Di internet, seseorang bisa berkenalan dengan siapa pun, kapan pun, latar belakang, suku bangsa yang berbeda, bahkan antar negara dengan bahasa yang berbeda pula. Tentu saja, seperti halnya di dunia nyata, etika ketika mengajak berkenalan harus tetap dijaga. Meskipun yang diajak adalah teman lama, namun jangan sampai memaksa orang lain untuk menjadi teman kita. Proses untuk menjalin hubungan di dunia maya (internet) seperti itu, sering disebut *sosial networking* (jejaring sosial) Sudarma dalam (Akbar 2013).

Menurut modul CHFI (2011), jejaring sosial adalah semua orang berkomunikasi dengan berbagai perangkat dan media komunikasi yang digunakan. Seseorang dapat melakukan hubungan dengan orang yang lain, organisasi dan organisasi yang lain hanya bermodalkan sebuah perangkat. Berkomunikasi dengan kesamaan hobi, ide, teman dan sebagainya.

Pendapat lainnya, Jejaring sosial atau terjemahan Inggrisnya adalah *social networking* yang merupakan hubungan antar individu atau organisasi yang di bentuk karena adanya kesamaan, misalnya kesamaan visi, misi, pertemanan, keturunan, suku, dan sebagainya Utomo dalam (Akbar 2013).

Dengan demikian sebuah situs jejaring sosial adalah situs yang digunakan sebagai media untuk berinteraksi sosial. Selain itu situs jejaring sosial berfungsi sebagai media komunikasi antar anggotanya. Saat ini banyak situs jejaring sosial yang digunakan dan dikenal oleh para pengguna internet di Indonesia, antara lain *Friendster, Facebook, Twitter, Myspace*. Selain itu, masih banyak situs jejaring sosial yang lainnya.

Sampai saat ini, banyak sekali situs yang menyediakan khusus untuk menjalin hubungan di dunia maya. Meskipun inti tujuannya sama, masing-masing situs memiliki fitur yang berbeda. Ada yang khusus untuk menjaring pertemanan saja, menjaring pertemanan dengan lebih interaktif dan menguak memori dengan teman lama, atau lebih menonjolkan komunikasi dan interaksi dengan teman lewat *blog*.

2.5 Twitter

Twitter sebagai salah satu *micro-blogging system* merupakan media jejaring sosial yang memberikan fasilitas bagi para penggunanya untuk mem-*posting* segala hal terkait aktivitas, opini, dan segala hal yang terjadi kepada public atau suatu grup melalui pesan yang dikenal dengan sebutan *tweet*. Karakteristik yang dimiliki oleh *tweet* diantaranya Go dalam (Hidayatullah 2014):

1. Length

Hal yang membedakan Twitter dengan jejaring sosial lainnya adalah adanya batasan dalam mem-*posting tweet* yaitu maksimum sebanyak 140 karakter sehingga pengguna dituntut untuk dapat mengekspresikan pandangan mereka menggunakan satu atau dua kata kunci.

2. Data Availability

Keberadaan Twitter API memberikan kemudahan dalam mengumpulkan jutaan *tweet* atau lebih untuk *training*.

3. *Languange Model*

Model bahasa yang disampaikan dalam *tweet* berasal dari berbagai macam media yang berbeda sehingga dimungkinkan banyak terdapat ejaan yang salah ataupun penggunaan bahasa '*slang*' di dalamnya.

4. *Domain*

Pengguna twitter mengirimkan pesan singkat tentang bermacam-macam topik. Berbeda dengan situs *review* yang memberika opini yang difokus kan pada topik tertentu.

2.6 **Cyberbullying**

Cyberbullying merupakan istilah yang merujuk kepada penggunaan teknologi Informasi untuk menggertak orang dengan mengirim atau posting teks yang bersifat mengintimidasi atau mengancam. Istilah ini telah ditambahkan ke dalam kamus OED (Oxford English Dictionary). OED menunjukkan penggunaan istilah ini pertama kali di Canberra pada tahun 1998, namun sebelumnya istilah telah di gunakan pada artikel New York Times 1995 dimana banyak sarjana dan penulis Besley seorang Kanada yang meluncurkan website *cyberbullying* tahun 2013 dengan istilah coining (Machsun 2016).

Cyberbullying adalah tindakan seorang atau remaja secara sengaja mengintimidasi, mengancam, atau mempermalukan seseorang, atau sekelompok anak lain melalui teknologi informasi, seperti media sosial atau *mobile device*. Hasil penelitian Ipsos Global menunjukkan sebanyak 60 persen responden mengatakan *cyberbullying* terjadi di sejumlah laman media sosial terkemuka seperti *facebook* (Firman dan Ngazis, 2012). *Cyberbullying* tidak hanya berlaku pada remaja dan anak-anak namun juga hal yang sama berlaku pada orang dewasa. Perbedaan dalam kelompok umur ini disebut *cyberstalking* atau *cyberharassment*. *Cyberstalking* adalah penggunaan komunikasi internet, *e-mail*, atau elektronik lainnya dan umumnya mengacu pada pola perilaku mengancam atau berbahaya. Sementara *Cyberharassment* berbeda dari *Cyberstalking*, *Cyberharassment* biasanya berkaitan dengan mengancam atau melecehkan pesan *e-mail*, pesan instan, atau *blog* dan *website* yang didekasikan sepenuhnya untuk menyiksa individu (Haryati 2014).

Professor Dan Olweus pada tahun 1993 dalam (Rudi 2010) telah mendefinisikan bullying yang mengandung tiga unsur mendasar perilaku bullying, yaitu:

1. Bersifat menyerang (agresif) dan negatif.
2. Dilakukan secara berulang kali.
3. Adanya ketidakseimbangan kekuatan antara pihak yang terlibat.

Olweus kemudian mengidentifikasi dua subtype bullying, yaitu perilaku secara langsung (*Direct bullying*), misalnya penyerangan secara fisik dan perilaku secara tidak langsung (*Indirect bullying*), misalnya pengucilan secara sosial. Underwood, Galen, dan Paquette di tahun 2001, mengusulkan istilah “Sosial Aggression“ untuk perilaku menyakiti secara tidak langsung. Riset menunjukkan bahwa bentuk bullying tidak langsung, seperti pengucilan atau penolakan secara sosial, lebih sering digunakan oleh perempuan daripada laki-laki. Sementara anak laki-laki menggunakan atau menjadi korban tipe bullying secara langsung, misalnya penyerangan secara fisik, dalam Nansel et al. 2001; Olweus 1997 (Rudi 2010).

Menurut survey global yang diadakan oleh *Latitude News*, Indonesia merupakan Negara dengan kasus bullying tertinggi kedua dunia setelah Jepang, dan mengalahkan kasus bullying di Amerika Serikat yang menempati urutan ketiga. Bullying paling banyak dilakukan melalui jejaring sosial. Sebagai Negara dengan jumlah populasi terbanyak keempat di dunia, Indonesia memiliki jumlah pengguna *Facebook* terbesar ketiga di dunia dan penyumbang 15 persen *tweet* setiap hari untuk *Twitter*. Berdasarkan penelitian 91% responden asal Indonesia mengaku telah melihat kasus *Cyberbullying*, dan paling sering terjadi melalui media sosial. Di Indonesia, 74% responden menunjuk *Facebook* sebagai biangnya *Cyberbullying*, dan 44% menyebut media website yang lain (Satalina 2014).

2.7 Tinjauan Undang-Undang

Cyberbullying di Indonesia diatur dalam Pasal 27 UU ITE, namun terdapat beberapa tindakan yang termasuk *cyberbullying* yakni *Flaming*, *Harassment* (gangguan), *Impersonation* (peniruan), *Outing* (menyebarkan rahasia orang lain), *Trickery* (tipu daya), *Exclusion*(pengeluaran), *Cyberstalking*. UU ITE hanya memuat unsur penghinaan dan pengancaman, padahal tindakan *cyberbullying* lainnya juga kerap kali terjadi dan menjadi langkah awal tindak pidana lain. Dengan berkembangnya situs jejaring sosial maka hal tersebut akan memudahkan pelaku *cyberbullying* melakukan tindakannya (satyawati 2014).

Didalam Pasal 27 ayat (3) UU ITE yang menyatakan bahwa Setiap Orang dengan sengaja dan tanpa hak mendistribusikan dan/atau mentransmisikan dan/atau membuat dapat diaksesnya Informasi Elektronik dan/atau Dokumen Elektronik yang memiliki muatan penghinaan dan/atau pencemaran nama baik. Kemudian dalam Pasal 27 ayat (4) UU ITE yang menyatakan bahwa Setiap Orang dengan sengaja dan tanpa hak mendistribusikan dan/atau mentransmisikan dan/atau membuat dapat diaksesnya Informasi Elektronik dan/atau Dokumen Elektronik *Cyber bullying* dalam UU ITE tidak terdapat unsur yang jelas. Hanya

terdapat unsur penghinaan, pencemaran nama baik, pengancaman dan pemerasan. Sedangkan jenis *cyberbullying* tidak hanya mengandung unsur penghinaan, pencemaran nama baik, pengancaman dan pemerasan saja. Pasal 27 ayat (3) dan (4) UU ITE belum menyangkut unsur dari *Flaming*, *Harassment* (gangguan), *Impersonation* (peniruan), *Outing* (menyebarkan rahasia orang lain), *Trickery* (tipu daya), *Exclusion* (pengeluaran), *Cyberstalking*. yang memiliki muatan pemerasan dan pengancaman (satyawati 2014).

2.8 Data Mining

Istilah data mining memiliki beberapa padanan, seperti *knowledge discovery* ataupun *pattern recognition*. Kedua istilah tersebut sebenarnya memiliki ketepatannya masing-masing. Istilah *knowledge discovery* atau penemuan pengetahuan digunakan karena tujuan utama dari data mining memang untuk mendapatkan pengetahuan yang masih tersembunyi di dalam bongkahan data. Istilah *pattern recognition* atau pengenalan pola pun tepat untuk digunakan karena pengetahuan yang hendak digali memang berbentuk pola-pola yang mungkin juga masih perlu digali dari dalam bongkahan data yang tengah dihadapi. Jadi apakah sebenarnya data mining itu? Banyak definisi bagi istilah ini dan belum ada yang dibakukan atau disepakati semua pihak. Namun demikian, istilah ini memiliki hakikat (notion) sebagai disiplin ilmu yang tujuan utamanya adalah untuk menemukan, menggali, atau menambang pengetahuan dari data atau informasi yang kita miliki (Susanto dan Suryadi, 2010).

Data mining bukan hanya pelengkap saja dalam suatu database, melainkan mempunyai fungsi yang penting untuk membantu penggunanya mendapatkan informasi yang berguna serta meningkatkan pengetahuan bagi sang pengguna itu sendiri dan dapat nantinya berguna untuk orang banyak. Pada dasarnya, data mining mempunyai empat fungsi dasar yaitu (Berson 2000):

1. Fungsi Klasifikasi (*classification*)

Data mining dapat digunakan untuk mengelompokkan data-data yang jumlahnya besar menjadi data-data yang lebih kecil

2. Fungsi Segmentasi (*Segmentation*)

Disini data mining juga digunakan untuk melakukan segmentasi (pembagian) terhadap data berdasarkan karakteristik tertentu.

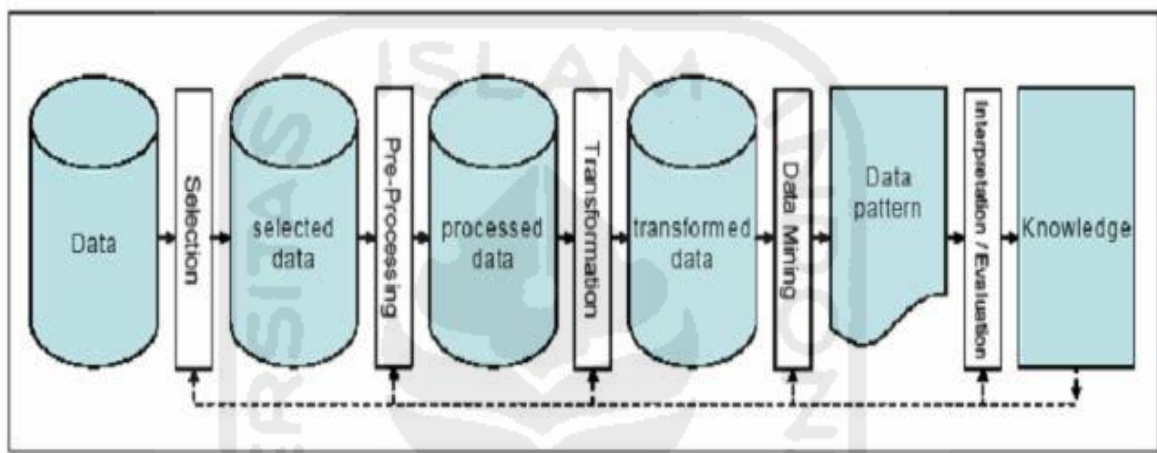
3. Fungsi Asosiasi (*Association*)

Disini data mining juga digunakan untuk mencari hubungan antara karakteristik tertentu .

4. Fungsi pengurutan (*Sequencing*)

Pada Fungsi ini, data mining digunakan untuk mengidentifikasi perubahan pola yang telah terjadi dalam jangka waktu yang tertentu.

Istilah data mining dan knowledge discovery in databases (KDD) sering kali digunakan secara bergantian untuk menjelaskan proses penggalian informasi tersembunyi dalam suatu basis data yang besar. Sebenarnya kedua istilah tersebut memiliki konsep yang berbeda akan tetapi berkaitan satu sama lain. Dan salah satu tahapan dalam keseluruhan proses KDD adalah data mining (Fayyad, 1996). Proses KDD secara garis besar dapat dijelaskan pada Gambar 2.1:



Gambar 2. 1 Fase-fase Dalam Data Mining (Fayyad, 1996)

1. Data Selection

Pemilihan (seleksi) data dari sekumpulan data operasional perlu dilakukan sebelum tahap penggalian informasi dalam KDD dimulai. Data hasil seleksi yang akan digunakan untuk proses data mining, disimpan dalam suatu berkas, terpisah dari basis data operasional.

2. Pre-processing/ Cleaning

Sebelum proses data mining dapat dilaksanakan, perlu dilakukan proses cleaning pada data yang menjadi fokus KDD. Proses cleaning mencakup antara lain membuang duplikasi data, memeriksa data yang inkonsisten, dan memperbaiki kesalahan pada data, seperti kesalahan cetak (tipografi). Juga dilakukan proses Enrichment, yaitu proses “memperkaya” data yang sudah ada dengan data atau informasi lain yang relevan dan diperlukan untuk KDD, seperti data atau informasi eksternal.

3. Transformation

Coding adalah proses transformasi pada data yang telah dipilih, sehingga data tersebut sesuai untuk proses data mining. Proses coding dalam KDD merupakan proses kreatif dan sangat tergantung pada jenis atau pola informasi yang akan dicari dalam basis data.

4. Data mining

Data mining adalah proses mencari pola atau informasi menarik dalam data terpilih dengan menggunakan teknik atau metode tertentu. Teknik, metode, atau algoritma dalam data mining sangat bervariasi. Pemilihan metode atau algoritma yang tepat sangat bergantung pada tujuan dan proses KDD secara keseluruhan.

5. Interpretation/ Evaluation

Pola informasi yang dihasilkan dari proses data mining perlu ditampilkan dalam bentuk yang mudah dimengerti oleh pihak yang berkepentingan. Tahap ini merupakan bagian dari proses KDD yang disebut dengan interpretation. Tahap ini mencakup pemeriksaan apakah pola atau informasi yang ditemukan bertentangan dengan fakta atau hipotesa yang ada sebelumnya.

2.9 Text Mining

Seiring perkembangan teknologi komunikasi, manusia memungkinkan untuk saling berkomunikasi satu sama lain. Salah satu cara untuk berkomunikasi adalah melalui teks. Bidang data mining yang bertujuan untuk menambang informasi dari data teks adalah text mining.

Ungkapan text mining umumnya digunakan untuk menunjukkan system yang menganalisis data yang besar yang berupa teks bahasa alami dan mendeteksi pola penggunaan leksial atau linguistik dalam upaya mengekstrak informasi yang berguna, sebastiani dalam (Sandi 2012). Sementara itu, witten (2005) menyatakan bahwa text mining adalah cabang ilmu baru yang sedang berkembang yang mencoba untuk mengumpulkan informasi yang bermanfaat dari data teks yang natural (*natural language text*). Hal ini biasanya dilakukan dengan mengekstrak informasi dari data teks yang berguna untuk tujuan tertentu. Dengan demikian text mining adalah penggalian informasi yang dilakukan oleh komputer untuk mendapatkan sesuatu yang baru, sesuatu yang tidak diketahui sebelumnya atau menemukan kembali informasi yang diekstrak secara otomatis dari sumber data yang berupa teks yang berbeda-beda.

Proses text mining memerlukan beberapa tahapan, mengingat data teks memiliki karakteristik yang lebih kompleks dari data biasa. Berdasarkan presentasi dari (Auvil 2003) yang menyatakan bahwa secara garis besar sebuah dokumen memiliki karakteristik sebagai berikut:

1. Database teks memiliki ukuran yang besar (*Large textual database*)
2. Memiliki dimensi tinggi, satu kata satu dimensi (*high dimensionaly*)

3. Mengandung frase dan antara frase satu dengan frase lainnya dapat memiliki arti yang berbeda dan saling bergantung satu sama lain (*dependency*)
4. Banyak mengandung kata/kalimat yang ambigu (*Ambiguity*)
5. Mengandung data noise, seperti singkatan. Istilah dan spelling mistake.
6. Mengandung struktur yang tidak baku misalkan singkatan-singkatan pada kata seperti “km di mn?”

Dengan demikian, proses penggalian informasi dari sekumpulan dokumen teks seperti halaman web, *twitter*, dokumen dan lain-lain membutuhkan beberapa proses yang saling terkait. Mengolah dokumen dari yang tidak terstruktur menjadi lebih terstruktur dengan menerapkan beberapa teknik ekstraksi dan penyaringan kata-kata dalam dokumen sekaligus dengan pembobotan tingkat kepentingan kata-kata dengan metode pembobotan. Data hasil pembobotan tersebut kemudian diolah dengan menggunakan teknik data mining sesuai dengan tujuan pengolahan data tersebut. Pada bagian ini akan menjelaskan dasar-dasar teori yang digunakan dalam proses ekstraksi dan pembobotan dokumen dengan metode pembobotan tf-idf.

Bentuk pengubahan dokumen teks yang masih tidak terstruktur menjadi lebih terstruktur adalah dengan mengubahnya menjadi bentuk *spreadsheet*. Kolom menunjuk kepada dokumen dan baris menunjuk kepada kata. Sedangkan sel menunjukkan frekuensi kemunculan suatu kata dalam suatu dokumen.

Agar sebuah dokumen dapat dianalisis dengan menggunakan teknik-teknik data mining, maka dilakukan proses *text transformation* dilakukan dengan mengubah data kata-kata menjadi kata-kata yang lebih baik dan sesuai kemudian dilakukan penghapusan kata-kata yang kurang penting dalam proses selanjutnya dengan proses *stopword list*. Setelah mendapatkan kata-kata unik dan bersih. Maka tahapan feature selection dilakukan dengan melakukan perhitungan jumlah kata yang muncul ataupun dengan teknik statistic lainnya, dimana dimungkinkan dilakukan proses pembobotan kata atau *term weight*. Data *term weight* inilah yang kemudian diolah dengan teknik data mining untuk menghasilkan informasi yang lebih berguna.

1. *Casefolding* adalah mengubah semua huruf dalam dokumen menjadi huruf kecil. Dimana hanya huruf ‘a’ sampai ‘z’ yang diterima. Karakter selain huruf dihilangkan dan dianggap sebagai delimitter (manning, 2008). Dengan demikian data teks yang akan diproses adalah data teks yang semua hurufnya adalah huruf kecil dan juga menghilangkan karakter angka dan symbol lain.
2. Tokenizing

Proses tokenisasi pada data teks adalah melakukan pemecahan sekumpulan karakter (kalimat) menjadi kata-kata (token). Setiap token adalah objek dari suatu tipe, sehingga jumlah token akan lebih banyak daripada tipenya. Contoh, jika di dalam sebuah teks terdapat kata “saya” sebanyak sepuluh buah, maka kesepuluhnya bertipe “saya”.

Hal lain yang berhubungan dengan tokenisasi adalah bagaimana cara memisahkan kata-katanya. Sebagai contoh, karakter *whitespace*. Seperti ada juga karakter-karakter lain seperti *semicolon*, titik, dan symbol-simbol lain yang juga dapat digunakan sebagai pemisah antar kata. Contohnya pada kalimat “siang, sore,malam,” dapat dipisahkan menjadi kata siang, kata sore dan kata malam.

3. Menghapus *Stopword*

Terkadang muncul beberapa kata yang sangat umum pada semua dokumen. Kata-kata tersebut dikenal dengan istilah stop word. *Stopword* adalah kata umum (*common words*) yang biasanya muncul dalam jumlah besar di setiap dokumen dan dianggap tidak memiliki makna. Proses penghapusan *stopword* ini dilakukan untuk setiap dokumen, apabila di dalam dokumen ditemukan kata yang termasuk kedalam daftar *stopword* maka kata tersebut dihapus, sehingga dimensi dokumen menjadi berkurang.

4. *Stemming*

Tahap *stemming* adalah tahap mencari *root* kata dari tiap kata hasil *filtering*. Pada tahap ini dilakukan proses pengembalian berbagai bentukan kata ke dalam suatu representasi yang sama. Tahap ini kebanyakan dipakai untuk teks berbahasa inggris dan lebih sulit diterapkan pada teks berbahasa Indonesia. Hal ini dikarenakan bahasa Indonesia tidak memiliki rumus bentuk baku yang permanen.

2.10 N-Gram

Pada dasarnya, model n-gram adalah model probabilistik yang awalnya dirancang oleh ahli matematika dari Rusia pada awal abad ke-20 dan kemudian dikembangkan untuk memprediksi item berikutnya dalam urutan item. Item bisa berupa huruf / karakter, kata, atau yang lain sesuai dengan aplikasi. Salah satunya, model n-gram yang berbasis kata digunakan untuk memprediksi kata berikutnya dalam urutan kata tertentu. Dalam arti bahwa sebuah n-gram hanyalah sebuah wadah kumpulan kata dengan masing-masing memiliki panjang n kata. Sebagai contoh, sebuah n-gram ukuran 1 disebut sebagai unigram; ukuran 2 sebagai “bigram”; ukuran 3 sebagai "trigram", dan seterusnya.

Pada pembangkitan karakter, *N-gram* terdiri dari substring sepanjang n karakter dari sebuah string, dalam definisi lain n-gram adalah potongan sejumlah n karakter dari sebuah

string. Metode *N-gram* ini digunakan untuk mengambil potongan-potongan karakter huruf sejumlah n dari sebuah kata yang secara kontinuitas dibaca dari teks sumber hingga akhir dari dokumen. Sebagai contoh: kata "TEXT" dapat diuraikan ke dalam beberapa n -gram berikut:

uni-gram : T, E, X, T
bi-gram : TE, EX, XT
tri-gram : TEX, EXT
quad-gram : TEXT, EXT dan seterusnya.

Sedangkan pada pembangkitan kata, metode n -gram ini digunakan untuk mengambil potongan kata sejumlah n dari sebuah rangkaian kata (kalimat, paragraf, bacaan) yang secara kontinuitas dibaca dari teks sumber hingga akhir dari dokumen. Sebagai contoh : kalimat "saya dapat melihat cahaya itu." dapat diuraikan ke dalam beberapa n -gram berikut:

uni-gram : saya, dapat, melihat, cahaya, itu
bi-gram : saya dapat, dapat melihat, itu ada
tri-gram : saya dapat melihat, dapat melihat dia
dan seterusnya.

Salah-satu keunggulan menggunakan n -gram dan bukan suatu kata utuh secara keseluruhan adalah bahwa n -gram tidak terlalu sensitive terhadap kesalahan penulisan yang terdapat pada suatu dokumen (Hanafi 2009).

2.11 Naïve Bayes Classifier

Naïve Bayes adalah sebuah metode yang berdasarkan pada suatu asumsi penyederhanaan dimana dimana nilai atribut secara kondisional saling bebas apabila diberikan nilai output (Santosa, 2007). *Naïve Bayes Classifier* merupakan sebuah metode klasifikasi yang berakar pada teorema Bayes. Ciri utama *Naïve Bayes Classifier* adalah asumsi yang sangat kuat (naif) dari independensi masing-masing kondisi atau kejadian (Natalius, 2010). Secara sederhana NBC mengasumsikan bahwa ada atau tidaknya fitur tertentu dari suatu kelas tidak memiliki keterkaitan dengan keberadaan fitur lainnya.

Keuntungan penggunaan NBC adalah hanya diperlukan sejumlah kecil data pelatihan untuk mengestimasi parameter yaitu rata-rata dan varian dari variable yang diperlukan untuk klasifikasi (Saraswati, 2011). Naïve bayes adalah sebuah algoritma analisa statistic, yang melakukan pengolahan data terhadap data numeric menggunakan probabilitas Bayesian. Klasifikasi-klasifikasi bayes adalah klasifikasi statistic yang dapat memprediksi kelas suati anggota probabilitas. Untuk klasifikasi bayes sederhana yang lebih dikenal sebagai Naïve Bayesian classifier dapat diasumsikan bahwa efek dari suatu nilai atribur sebuah kelas tidak

dipengaruhi nilai dari atribut lainnya. Asumsi ini disebut *class conditional independence* yang diciptakan untuk memudahkan perhitungan, pengertian ini dianggap “naïve”, dalam bahasa lebih sederhana naïve itu mengasumsikan bahwa kemunculan suatu term kata dalam suatu kalimat tidak dipengaruhi kata-kata yang lain, sehingga dalam analisis sentiment kata yang muncul memiliki bobot seluruhnya apakah kalimat tersebut positif ataupun negatif (saputra 2015).

Secara garis besar, cara kerja metode ini dapat direpresentasikan sebagai berikut:

1. Ambil probabilitas positif, dan negatif tiap kata.
2. Hitung rata-rata probabilitas ketiganya
3. Tentukan klasifikasi berdasarkan nilai probabilitas diatas.

Untuk mendapatkan probabilitas dari tiap kata, terlebih dahulu melakukan pembelajaran terhadap setiap kata dan probabilitasnya. Dalam proses pembelajaran ini, diperlukan sebuah training set, yang merupakan sekumpulan kalimat positif, negatif dan netral yang telah diklasifikasikan. Naïve bayes merupakan teknik klasifikasi yang sederhana dan cepat. Teknik ini bekerja dengan baik dengan representasi statistik. Berbeda dengan metode *rule-based*, naïve Bayesian dapat belajar secara incremental. Namun kekurangan dari naïve Bayesian adalah ukuran dari vektor fitur yang dihasilkan cukup besar dan butuh teknik untuk memperkecil ukuran vector tersebut (Darujati 2012).

Untuk kemungkinan kategori bagi dokumen yang diberikan, berikut ini adalah penjelasan dari naïve bayes (Nuraini 2011):

1. Setiap data direpresentasikan sebagai vector berdimensi n yaitu $X = (x_1, x_2, x_3 \dots x_n)$ adalah gambaran dari ukuran yang dibuat di test dari n atribut yaitu $A_1, A_2, A_3 \dots A_m$ dimana m adalah kumpulan kategori yaitu $C_1, C_2, C_3 \dots C_m$. Diberikan data test X yang tidak diketahui kategorinya, maka classifier akan memprediksi bahwa X adalah milik kategori dengan posterior probabilitas tertinggi berdasarkan kondisi X . oleh karena itu, Naïve Bayes Classifier menandai bahwa test X yang tidak diketahui tadi ke kategori C_l jika dan hanya jika dilihat pada persamaan $P(C_l|X) > P(C_j|X)$ untuk $l \leq j \leq m, j \neq i$
Kemudian kita perlu memaksimalkan $P(C_i | X)$ berdasarkan persamaan $P(C_1|X) = \frac{P(X|C_1) \cdot P(C_1)}{P(X)}$
2. P_x adalah konstan untuk semua kategori, hanya $P(X|C_1) \cdot P(C_1)$ yang perlu dimaksimalkan. Jika prior probabilitas mungkin diperkirakan dengan perhitungan $P(C_i) = \frac{S_i}{S}$ dimana S_i adalah jumlah dari data training dari kategori C_i , dan S adalah jumlah total data training.

3. Diberikan data dengan banyak atribut, ini akan menjadi komputasi yang kompleks untuk mengkomputasi $P(X|C_i)$. Untuk mengurangi komputasi pada saat mengevaluasi $P(X|C_i)$, maka dapat dihitung menggunakan persamaan:

$$P(X|C_i) = \sum_{k=1}^n P(x^k|C_i)$$

dimana X adalah nilai-nilai atribut dalam sampel X dan probabilitas $P(X_1|C_i)$, $P(X_2|C_i)$,....., $P(X_n | C_i)$ dapat diperkirakan dari data training.

2.12 Bag-of-words Model

Bag-of-words adalah model yang mengambil setiap kata dalam sebuah kalimat sebagai fitur, dengan asumsi bahwa setiap fitur memiliki sifat *conditional independence*. Teks direpresentasikan sebagai koleksi dari kata-kata yang tidak berurutan. Setiap fitur mewakili keberadaan satu kata. Semua kata (fitur) dalam merupakan *dictionary*. Tantangan dengan pendekatan ini adalah pilihan kata yang tepat untuk menjadi fitur.

Model *bag-of-words* tidak mampu untuk menangkap relasi antar kata. Sebagai contoh, kalimat *'this movie is excellent'* dan *'this film is great'* akan dianggap sebagai 2 kalimat yang memiliki arti yang berbeda. Bagi kita, sangat jelas kesamaan antara kata *'excellent'* dan *'great'* serta *'movie'* dan *'film'*. Hal ini dapat diatasi dengan menggunakan *semantic handling*, yakni masing-masing pasangan kata tersebut dapat diwakili oleh sebuah fitur tunggal dengan mencari kata sinonimnya, Yessenov dan Misailovic, 2009 dalam (Hilmawan 2014). Model *bag-of-words* dimodifikasi sehingga fitur-fitur yang memiliki makna yang sama hanya memiliki satu kata yang diwakili oleh *cluster* sinonimnya.

2.13 Teknik Pencarian Data

Dalam setiap waktu yang mengakses media sosial seperti Facebook dan Twitter dapat mencapai ratusan ribu bahkan jutaan pengguna, sehingga *database* media sosial mempunyai kapasitas yang sangat besar. Jika melakukan pencarian tanpa suatu teknik yang baik, akan banyak memakan waktu dan tenaga. Untuk itu, diperlukan suatu teknik pencarian yang efektif dan efisien dalam menemukan data sesuai kebutuhan (Saputra 2015).

Search Techniques pada penelitian ini menggunakan salah satu metode pencarian yang dikemukakan oleh Jackie Skinner dari University of Reading yang merupakan Universitas top 1% dunia. Search techniques tersebut memberikan petunjuk tentang opinion mining pada twitter untuk bahasa Indonesia dengan metode Naïve Bayes Classifier:

1. Menggunakan symbol untuk mencari alternative akhiran dan ejaan
2. Menggabungkan konsep dalam pencarian

3. Mencari frase
4. Melakukakan pencarian secara lebih spesifik

Metode dalam *search techniques* yang dipakai dalam penelitian ini adalah menggunakan operator pencarian. Metode ini menggabungkan kata-kata pencarian mencakup sinonim dan dikenal juga sebagai Boolean Searching. Pada metode ini memungkinkan untuk memasukkan banyak kata ataupun konsep dalam pencarian. Sudah banyak peneliti menggunakan Boolean *Searching* dalam teknik pencarian datanya, diantaranya pencarian katalog universitas, OPACs (Online Public Access Catalogs), Katalog perpustakaan online dan terdapat juga strategi pencarian positif dan negatif untuk web.

Pada operator “AND”, misalkan untuk kata kunci pencarian data “bangsat dan brengsek” akan menghasilkan data yang terdapat kata baik itu kata bangsat maupun brengsek. Untuk operator “OR”, misalkan untuk kata kunci pencarian data “bangsat or brengsek” akan menghasilkan data baik itu data yang berisikan bangsat saja, atau brengsek saja maupun data yang berisikan bangsat dan brengsek. Untuk operator “NOT”, misalkan untuk kata kunci pencarian data “bangsat not brengsek” akan menghasilkan data yang berisikan bangsat saja.

2.14 API Twitter

API adalah singkatan dari Application Programming Interface, merupakan aplikasi pendukung dari Jejaring Sosial, salah satu fungsinya yaitu untuk meng-ekstrak data dalam bentuk JSON. API adalah cara yang ditetapkan untuk sebuah program dalam menyelesaikan tugas, biasanya dalam mengambil dan memodifikasi data. Twitter menyediakan API pada hampir setiap fitur yang dapat dimanfaatkan untuk membuat aplikasi, *website*, *widget*, dan proyek lain yang berinteraksi dengan Twitter. Komunikasi antara aplikasi yang dibuat dengan Twitter API dilakukan melalui *Hypertext Transfer Protocol* (HTTP). Twitter API memiliki beberapa komponen yang semuanya bersifat *free* atau gratis.

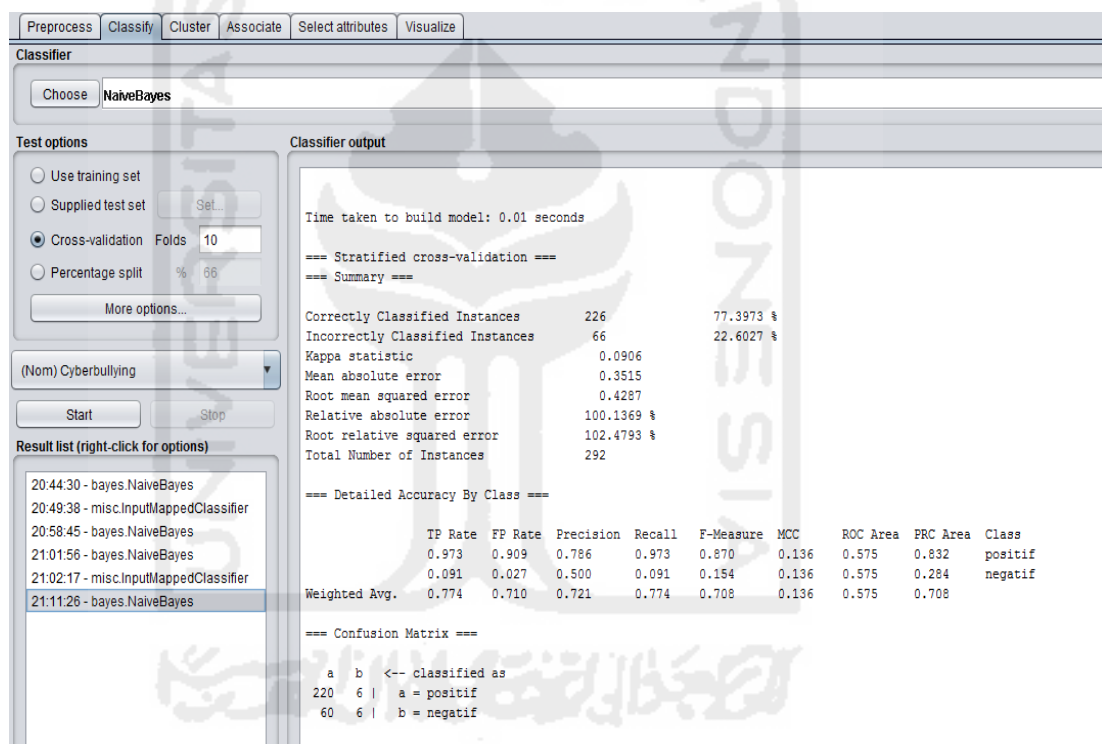
Namun demikian, Twitter memberikan batasan jumlah *tweet* yang bisa didapat dalam satu jam. Beberapa bagian dari Twitter API tidak memerlukan otentikasi dasar penggunaan berupa *username* dan *password*, dan beberapa otentikasi menggunakan *OAuth*. Sebagian besar bagian dari API pada Twitter menggunakan model REST, kecuali untuk API streaming, yang selalu membutuhkan koneksi dengan Internet. Data yang diperoleh dari Twitter dikembalikan kedalam format XML dan JSON. Namun seiring berjalannya waktu, pengembalian data dalam bentuk XML mulai dihapus (Aliandu 2012).

Twitter Search API adalah API yang digunakan untuk menjalankan pencarian terhadap indeks *real time* dari *tweet* terakhir. Terdapat beberapa batasan dengan Twitter Search API, antara lain (Aliandu 2012):

- a. Search API tidak lengkap mengindeks semua *tweet* tetapi indeks 1500 *tweets* terakhir.
- b. Search API tidak dapat digunakan untuk menemukan *tweet* yang umurnya lebih dari seminggu.
- c. Search API mengembalikan maksimal 100 *tweet* yang umurnya lebih dari seminggu.
- d. Query yang terlalu kompleks mungkin tidak akan berhasil
- e. Search tidak mendukung otentifikasi yang berarti semua *Query* adalah anonym.



Grafik tersebut menunjukkan bahwa titik biru adalah kalimat yang mengandung positif bullying sebanyak 77.40%, sementara titik merah adalah kalimat yang mengandung negatif bullying sebanyak 22.60%. Data training tersebut akan diproses dengan metode *Naive Bayes Classifier* dengan validasi data menggunakan metode *10 fold cross validation*, data dibagi mejadi 10 bagian dengan 9/10 bagian digunakan untuk proses *training* dan 1/10 bagian digunakan untuk proses *testing*. Iterasi berlangsung 10 kali dengan variasi data *training* dan *testing* menggunakan kombinasi 10 bagian data. Dari hasil tersebut akan membentuk sebuah pola yang selanjutnya digunakan untuk memprediksi kalimat cyberbullying dan non bullying yang ada pada data baru atau data testing. Proses *Naive Bayes Classifier* pada data training ditunjukkan pada **Gambar 4.2** berikut:



Gambar 4. 2 Evaluasi Naive Bayes Pada Data Training Menggunakan WEKA
 Data yang digunakan pada Data Training sebanyak 292 record, dan didapatkan hasil seperti pada gambar diatas. Sehingga hasil presentasinya adalah:

$$Presentase\ Akurasi = + \frac{\text{banyak prediksi yang benar}}{\text{total banyaknya data}} \times 100\%$$

$$\begin{aligned} \text{Hasilnya} &= (220+60)/(220+6+66+6) \times 100\% \\ &= 77.40\% \end{aligned}$$

Tahap selanjutnya adalah memprediksi data testing menggunakan data training yang diproses sebelumnya. Pada menu *Test Option* pilih *supplied test set* kemudian pilih set untuk memasukkan data uji kemudian di proses dan *save* hasilnya dalam type *.arff* hasil dari prediksi tersebut dapat diketahui melalui menu tools pada WEKA GUI kemudian pilih *Arffviewer*. Hasil prediksi data testing dapat dilihat pada **Gambar 4.3** berikut:

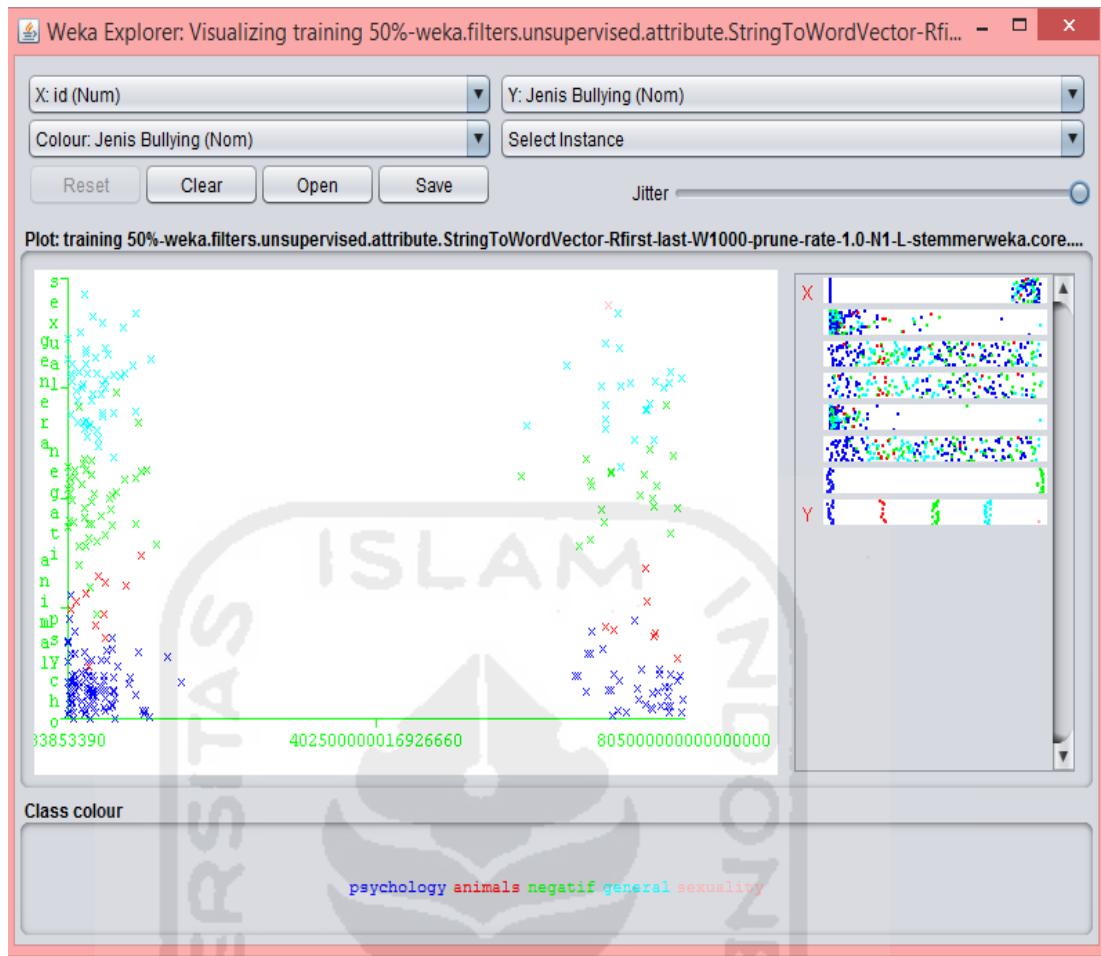
No.	1: id	2: friend	3: follower	4: text	5: prediction margin	6: predicted Cyberbullying	7: Cyberbullying
	Numeric	Numeric	Numeric	Nominal	Numeric	Nominal	Nominal
1	7.75...	125.0	151.0	pron...	0.898794	positif	
2	2.78...	3001.0	3200.0		0.572176	positif	
3	1.53...	137.0	202.0		0.607861	positif	
4	1.52...	460.0	2770.0		0.599858	positif	
5	2.37...	297.0	1750.0		0.605627	positif	
6	1.39...	532.0	1101.0		0.606136	positif	
7	3.32...	4409.0	13276.0		0.271141	positif	
8	5.63...	460.0	467.0		0.60695	positif	
9	2.52...	458.0	1064.0	bad...	0.777891	positif	
10	2.56...	270.0	344.0	das...	0.861038	positif	
11	2.44...	540.0	456.0	das...	0.860625	positif	
12	2.56...	329.0	500.0	das...	0.860967	positif	
13	2.54...	86.0	192.0	das...	0.861197	positif	
14	2.54...	86.0	176.0	das...	0.861197	positif	
15	7.45...	69.0	88.0		0.607776	positif	
16	5.71...	2512.0	16566.0		0.117664	positif	
17	7.88...	1073.0	147.0		0.603781	positif	
18	7.03...	574.0	483.0	das...	0.910553	positif	
19	5.97...	2873.0	5365.0		0.548736	positif	
20	7.59...	7786.0	8455.0		0.333061	positif	
21	7.16...	90.0	98.0		0.607776	positif	
22	7.39...	57.0	59.0		0.607776	positif	
23	2.28...	2514.0	2896.0		0.580774	positif	
24	1.25...	1397.0	236.0		0.601631	positif	
25	7.79...	2178.0	2213.0		0.5904	positif	
26	2.97...	119.0	251.0		0.607861	positif	
27	6.04...	46.0	16.0		0.607776	positif	
28	2.65...	121.0	227.0		0.607861	positif	
29	1.22...	131.0	222.0		0.607171	positif	

Gambar 4.3 Prediksi Data Testing

Jumlah data testing adalah 292 record, setelah melakukan klasifikasi dapat diketahui bahwa sebanyak 96,55% data tweet yang positif mengandung bullying dan negatif bullying hanya 3.44%.

4.2.2 Klasifikasi jenis Cyberbullying

Data yang digunakan sama dengan data pada klasifikasi cyberbullying dan non bullying, bedanya hanya pada pelabelan untuk data training. Pelabelan ini dilakukan merujuk pada peneliti terdahulu yang mengklasifikasikan kata bullying yang banyak digunakan di Indonesia. Jumlah data training juga sebanyak 292 record dengan jumlah atribut adalah 8. Proses klasifikasi ini juga sama dengan proses diatas. **Gambar 4.4** dibawah adalah tampilan proses data training untuk jenis bullying:



Gambar 4. 4 Grafik Data Training Jenis Cyberbullying Menggunakan WEKA

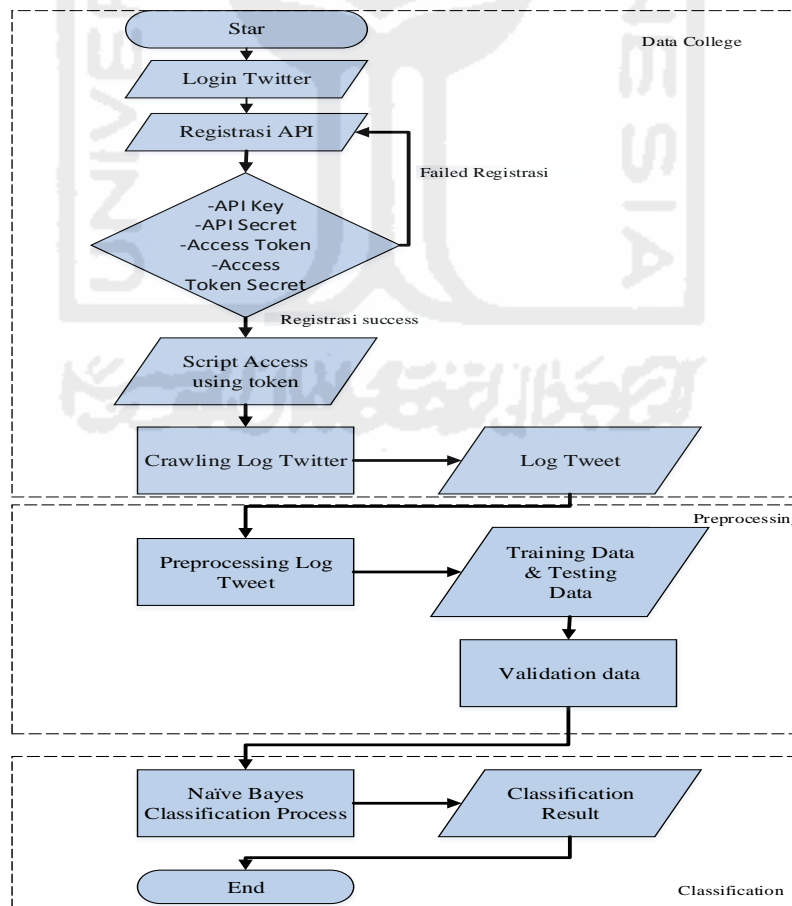
Berdasarkan Grafik tersebut diatas dapat diketahui bahwa titik berwarna biru adalah kalimat jenis bullying yang *related psychology* seperti kata makian “bodoh”, “goblok” dan lain-lain, sementara untuk titik yang berwarna merah adalah jenis bullying yang *related animals* seperti makian dengan kata “anjing” dan semacamnya. Untuk titik yang berwarna hijau adalah kalimat untuk yang mengandung negatif bullying, untuk yang berwarna biru adalah jenis bullying secara *general* seperti kata “setan”, “iblis” dan lain-lain, dan untuk yang berwarna *pink* adalah jenis bullying *related sexuality* seperti umpatan kata “banci” dan semacamnya.

Data training yang telah diproses dilakukan klasifikasi untuk dijadikan relasi pengujian data testing. Dari hasil klasifikasi data training dan validasi *10 cross validation folds* akan memprediksi jenis bullying pada data testing. **Gambar 4.5** dibawah adalah proses klasifikasi untuk data training jenis cyberbullying:

Bab 3 Metodologi Penelitian

3.1 Metodologi Penelitian

Penelitian ini terdiri dari tiga tahap yang pertama adalah teknik pengumpulan *Log Data*. *Log data Tweet* data dari jejaring sosial *Twitter* di *crawling* menggunakan API *Twitter* yang telah disediakan oleh *developer Twitter*. kedua adalah *preprocessing* data. *Log Data* yang telah di *crawling* dari *twitter* menghasilkan data mentah yang tidak terstruktur, *preprocessing* atau pembersihan data dilakukan agar data menjadi terstruktur dan memudahkan pada saat analisis, Ketiga adalah Klasifikasi, dari *log twitter* yang telah dibersihkan selanjutnya akan dirubah dalam bentuk *vector* untuk kemudian diklasifikasikan dengan metode *Naïve Bayes Classification* (NBC) menggunakan *Machine Learning WEKA*. *Flowchart* penelitian dapat dilihat pada **Gambar 3.1**:



Gambar 3. 1 Alur Penelitian

3.2 Perangkat Pendukung Penelitian

Untuk mendukung penelitian ini dibutuhkan beberapa perangkat keras maupun perangkat lunak diantaranya adalah:

1. Perangkat keras:

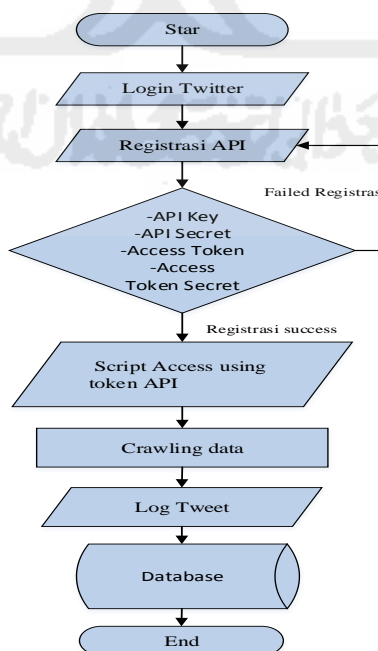
- a. Prosesor Intel(R) Core(TM) i5-3210M CPU @ 2.40GHz (4 CPUs)
- b. Hard Disk Drive 500 GB
- c. RAM 4 GB

2. Perangkat Lunak:

- a. Sistem operasi Windows 8.1 Pro 64-bit
- b. Notepad ++ v6.6.9 sebagai editor
- c. Anaconda 4.0.0
- d. WEKA 3.8.0
- e. Akun Twitter
- f. API Twitter
- g. Sastrawi master untuk stemming data
- h. Stopword Tala
- i. JsonLint untuk validasi file *Json*
- j. Convert csv untuk mengubah file *Json* ke CSV atau ke Excel

3.3 Pengumpulan Data *Log Tweet*

Proses pengambilan *Log Tweet* dapat dilihat pada alur dibawah ini:

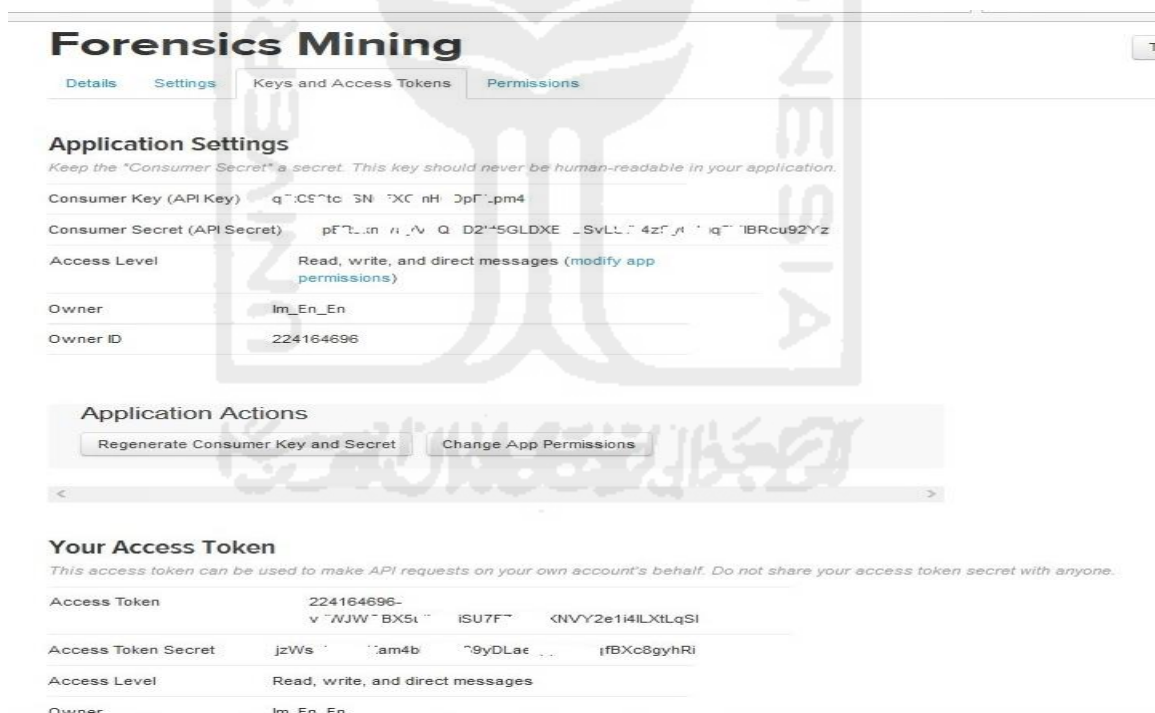


Gambar 3. 2 Teknik Pengumpulan Data

Proses pengumpulan data diawali dengan melakukan *Login* pada akun Twitter. Setelah *Login* lakukan daftar aplikasi untuk mendapatkan *access tokens* berupa *consumer_key*, *consumer_secret*, *access_token*, dan *access_secret*, hal ini dibutuhkan agar bisa mengakses *Twitter Search API*. Untuk bisa melakukan komunikasi antara *token Access* dan *Twitter Search API* maka dibuat sebuah *script* sebagai media untuk *crawling*. Meskipun demikian terdapat beberapa batasan dengan *Twitter Search API* diantaranya adalah:

- a. Hanya dapat melakukan indeks pada 1500 *tweets* terakhir.
- b. Data yang umurnya lebih dari seminggu tidak dapat di *crawling*
- c. Pencarian data yang umurnya lebih dari seminggu maksimal 100 tweet
- d. Pencarian dengan Query yang kompleks kemungkinan tidak berhasil

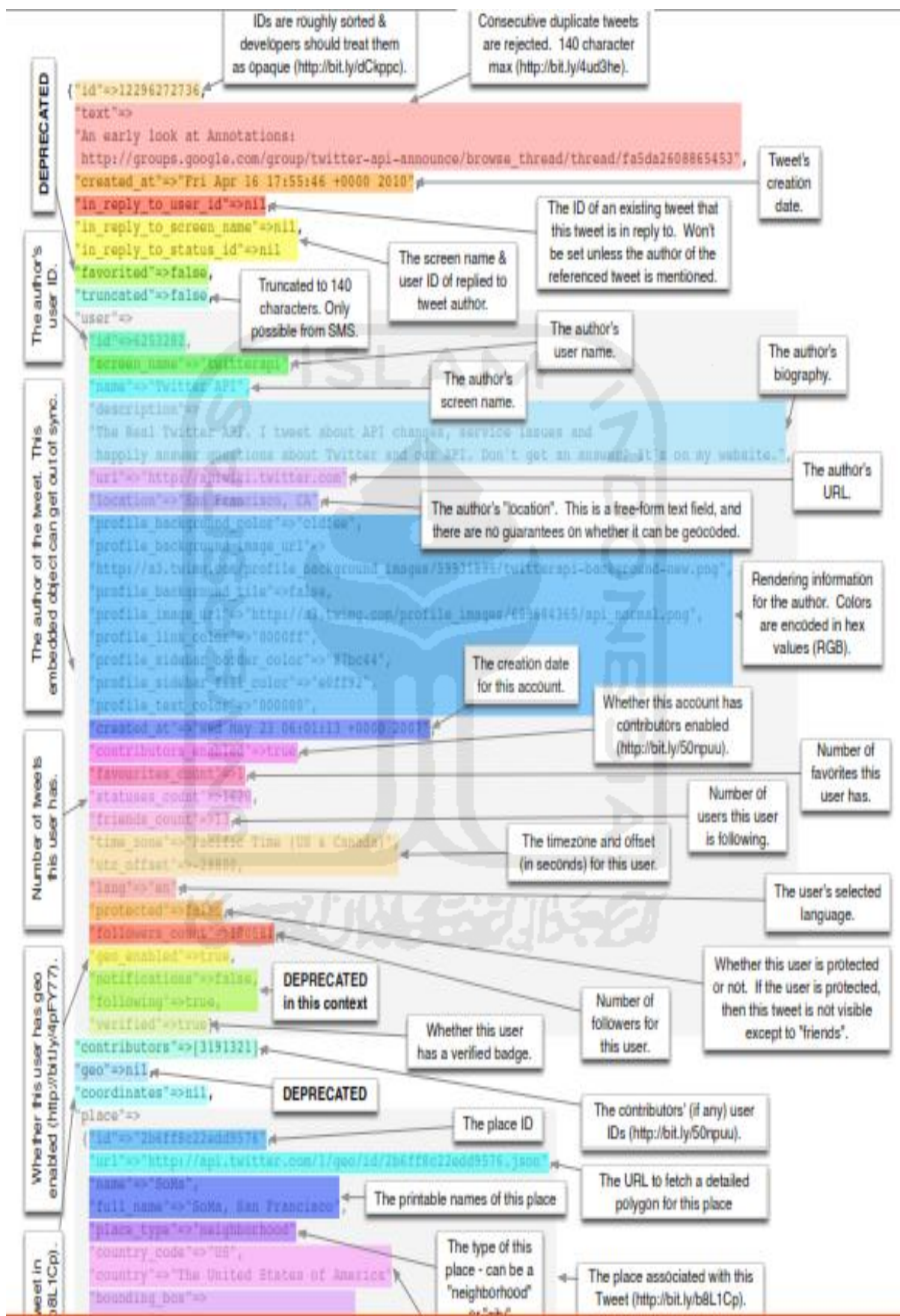
Dengan mengetahui batasan tersebut diatas maka dapat dilakukan pencarian data sesuai akses yang diberikan. Selanjutnya, setelah *token access* didapatkan maka dapat dilakukan pencarian data menggunakan *script* yang telah dibuat sebelumnya sesuai Query yang diinginkan. Pengumpulan data dilakukan secara random pada periode November-Desember 2016 sebanyak 1000 data, namun setelah melakukan filter data menjadi 583 data. **Gambar 3.3** dibawah adalah contoh daftar aplikasi untuk mendapatkan akses token.



Gambar 3.3 Access Token

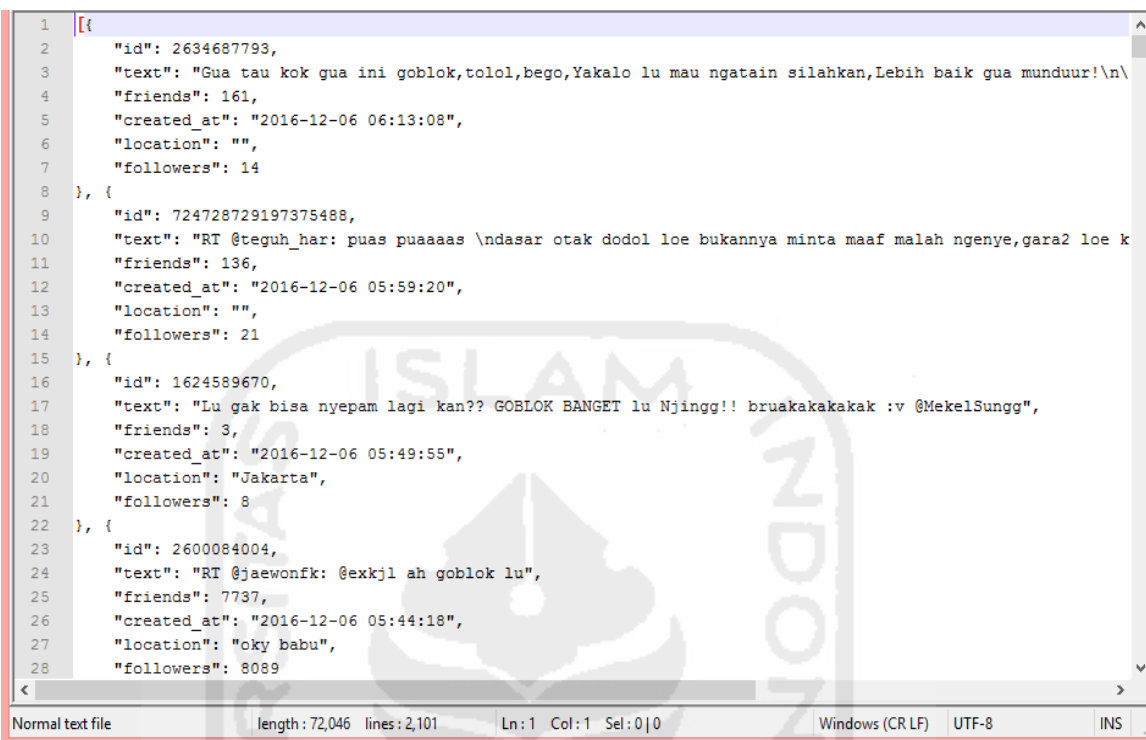
Untuk satu akun *Twitter* banyak objek yang bisa menjadi sumber Informasi, seperti ID unik, *text*, waktu pembuatan akun twitter, ID user yang *me-reply tweet author*, *biography user*, lokasi, Negara, bahkan sampai type lokasi user dapat diketahui, maka untuk memudahkan penelitian ini, dilakukan parsing data pada *script* sehingga saat *crawling* hanya

sebagian objek yang di ekstrak. Sebelum dilakukan parsing struktur data pada twitter terlihat seperti pada **Gambar 3.4** (Raffi Krikorian 2010).



Gambar 3. 4 Map Twitter

Namun setelah dilakukan parsing objek yang diekstrak menjadi beberapa objek diantaranya adalah *created_at*, *text*, *Location*, *followers*, *friends* dan *id_str*, seperti terlihat pada **Gambar 3.5** dibawah:



```
1 [{"id": 2634687793,
2   "text": "Gua tau kok gua ini goblok,tolol,bego,Yakalo lu mau ngatain silahkan,Lebih baik gua munduur!\n",
3   "friends": 161,
4   "created_at": "2016-12-06 06:13:08",
5   "location": "",
6   "followers": 14
7 }, {
8   "id": 724728729197375488,
9   "text": "RT @teguh_har: puas puaaaaas \ndasar otak dodol loe bukannya minta maaf malah ngenye,gara2 loe k
10  "friends": 136,
11  "created_at": "2016-12-06 05:59:20",
12  "location": "",
13  "followers": 21
14 }, {
15  "id": 1624589670,
16  "text": "Lu gak bisa nyepam lagi kan?? GOBLOK BANGET lu Njingg!! bruakakakakak :v @MekelSungg",
17  "friends": 3,
18  "created_at": "2016-12-06 05:49:55",
19  "location": "Jakarta",
20  "followers": 8
21 }, {
22  "id": 2600084004,
23  "text": "RT @jaewonfk: @exkjl ah goblok lu",
24  "friends": 7737,
25  "created_at": "2016-12-06 05:44:18",
26  "location": "oky babu",
27  "followers": 8089
28 }
```

Gambar 3.5 Hasil Ekstrak Data Setelah di Parsing

Proses pengumpulan data ini menggunakan beberapa *library* dalam pengambilan data *tweet* diantaranya adalah *library OAuth*, *Twitter REST API v1.1*, *Jsonpickle* dan *Tweepy*. *Library OAuth* digunakan untuk proses otentikasi sedangkan *Twitter REST API* digunakan untuk mengirimkan pesan kepada Twitter dan menerima *status update*, *Jsonpickle* digunakan untuk hasil pencarian data agar data tersebut dalam bentuk format *Json* sementara *Tweepy* digunakan untuk menghubungkan pemrograman yang digunakan ke Twitter. Untuk teknik pencarian data menggunakan operator pencarian. Metode ini menggabungkan kata-kata pencarian mencakup sinonim dan dikenal juga sebagai *Boolean Searching*. Pada metode ini memungkinkan untuk memasukkan banyak kata ataupun konsep dalam pencarian. Metode ini mengindikasikan hasil yang didapat berdasarkan operator “AND”, “OR”, dan “NOT”.

Pada operator “AND”, misalkan untuk kata kunci pencarian data dengan kata bullying “Bangsat and Bajingan” akan menghasilkan data yang terdapat kata bangsat saja maupun bajingan saja. Untuk operator “OR”, misalkan untuk kata kunci pencarian data kata bullying “Bangsat or Bajingan” akan menghasilkan data yang berisikan bangsat saja, atau bajingan saja maupun data yang berisikan bajingan dan bangsat. Untuk operator “NOT”, misalkan

untuk kata kunci pencarian data kata bullying “Bangsat not Bajingan” akan menghasilkan data yang berisikan kata bangsat saja.

Keywords dalam pencarian data menggunakan kata yang sering digunakan untuk melakukan bullying misalnya kata, “bangsat”, “dasar monyet” dan lain-lain. Metode ini mengikuti peneliti terdahulu yang pada penelitiannya menganalisis *Gender Bullying* sehingga kata kunci hanya berdasar pada LGBT seperti kata “gay” dan “bitch” (Sanchez 2011). Sementara untuk *keywords* pencarian juga mengacu pada penelitian terdahulu yang mana penelitian tersebut terdapat kata-kata bullying yang banyak digunakan di Indonesia untuk melakukan bullying pada jejaring sosial (Margono, 2011). Meskipun *keywords* pencarian adalah kata-kata bullying namun tidak semua maksud dari kata tersebut adalah untuk membullying, misalnya pada kalimat berikut “siapa yang kamu sebut babi” dan “babi kamu”. Walaupun kedua kalimat ini menggunakan kata “babi” namun tidak semuanya bermakna untuk membullying seseorang.

Implementasi pengambilan data dilakukan dengan membuat *file script parsing* yang bertugas melakukan *cron job*, untuk penelitian ini menggunakan *python*. **Gambar 3.6** merupakan kode program proses *query* pengaksesan data ke Twitter yang telah di parsing.

```
import tweepy,sys,jsonpickle

consumer_key = 'qfkC99tcnSN0FXQmHCDpRLpm4'
consumer_secret = 'pER9knwayVPQQD2kt5GLDXEHDSvLSfb4zSyGJgq5YIBRcu92Yz'

#inisialisasi
qry='bangsat'
maxTweets = 300
tweetsPerQry = 100
fName='Parse_Data_tweet.json'

#Proses Parsing
parseTweets=[]

#getData
auth = tweepy.AppAuthHandler(consumer_key,consumer_secret)
api = tweepy.API(auth, wait_on_rate_limit=True,wait_on_rate_limit_notify=True)
if (not api):
    sys.exit('Autentikasi gagal, cek "Consumer Key" & "Consumer Secret" Twitter anda')
parseTweets = []
#inisialisasi
sinceId=None;max_id=-1;tweetCount=0

print("Mulai mengunduh maksimum {0} tweets".format(maxTweets))
with open(fName,'w') as f:
    while tweetCount < maxTweets:
```

Gambar 3. 6 Script Program Pengumpulan Data


```

try:
if (max_id <= 0):
if (not sinceId):
    new_tweets=api.search(q=qry,count=tweetsPerQry)
    else:
        new_tweets=api.search(q=qry,count=tweetsPerQry,since_id=sinceId)
else:
    if (not sinceId):
        new_tweets=api.search(q=qry,count=tweetsPerQry,max_id=str(max_id - 1))
    else:
        new_tweets=api.search(q=qry,count=tweetsPerQry,max_id=str(max_id
1),since_id=sinceId)
if not new_tweets:
    print("Tidak ada lagi Tweet ditemukan dengan Query="{0}"".format(qry));break
for tweet in new_tweets:
# f.write(jsonpickle.encode(tweet._json,unpicklable=False)+'\n')
userTweet = { }
userTweet = tweet.user
parseTweets.append({
    "id" : userTweet.id,
    "location" : userTweet.location,
    "text" : tweet.text,
    "created_at" : tweet.created_at,
    "followers": userTweet.followers_count,
    "lang": userTweet.lang,
    "friends": userTweet.friends_count
})

f.write(jsonpickle.encode(parseTweets,unpicklable=False))

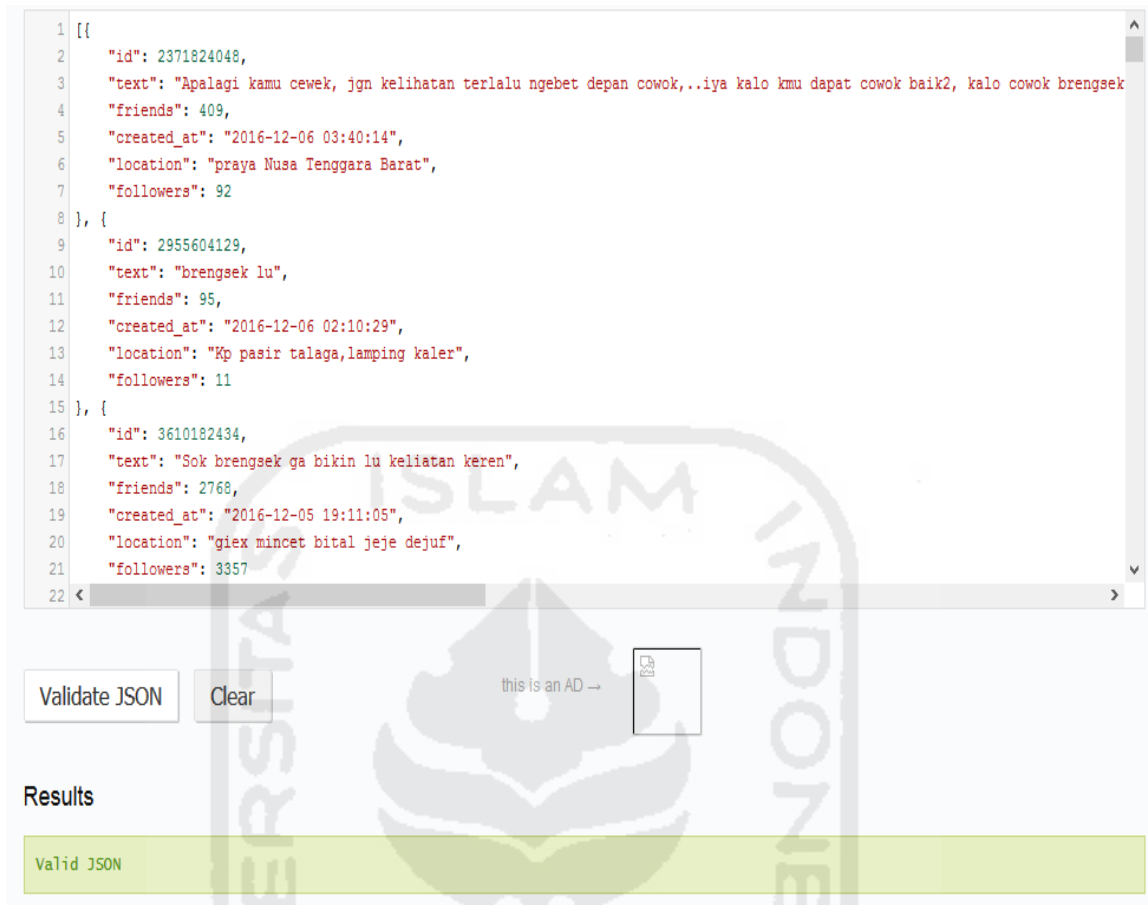
tweetCount+=len(new_tweets)
sys.stdout.write("\r");sys.stdout.write("Jumlah Tweets yang tersimpan: %.0f"
%tweetCount);sys.stdout.flush()
max_id=new_tweets[-1].id
except tweepy.TweepError as e:
    print("some error : " + str(e));break
print ("\nSelesai! {0} tweets tersimpan di "{1}"".format(tweetCount,fName))

```

Gambar 3.6 Script Program Pengumpulan Data (Lanjutan)

Proses pengambilan data *tweet* dilakukan dengan memanggil fungsi *search* dari *library twitter*. Namun sebelum proses pencarian dilakukan, terlebih dahulu dideklarasikan variabel *\$consumerKey*, *\$consumerSecret*, *\$accessToken*, *\$accessTokenSecret*. Variabel *\$consumerKey* dan *\$consumerSecret* berisi *OAuth setting* aplikasi yang didaftarkan ke Twitter. Variabel *\$accessToken* dan *\$accessTokenSecret* merupakan akses token untuk mengakses Twitter. Hasil yang diperoleh dalam proses ini yaitu berupa file dalam bentuk file *Json*. Untuk memastikan file tersebut file *Json* maka peneliti melakukan verifikasi secara

online menggunakan *JsonLint*, **Gambar 3.7** dibawah adalah contoh verifikasi *Json* menggunakan *Jsonlint*:



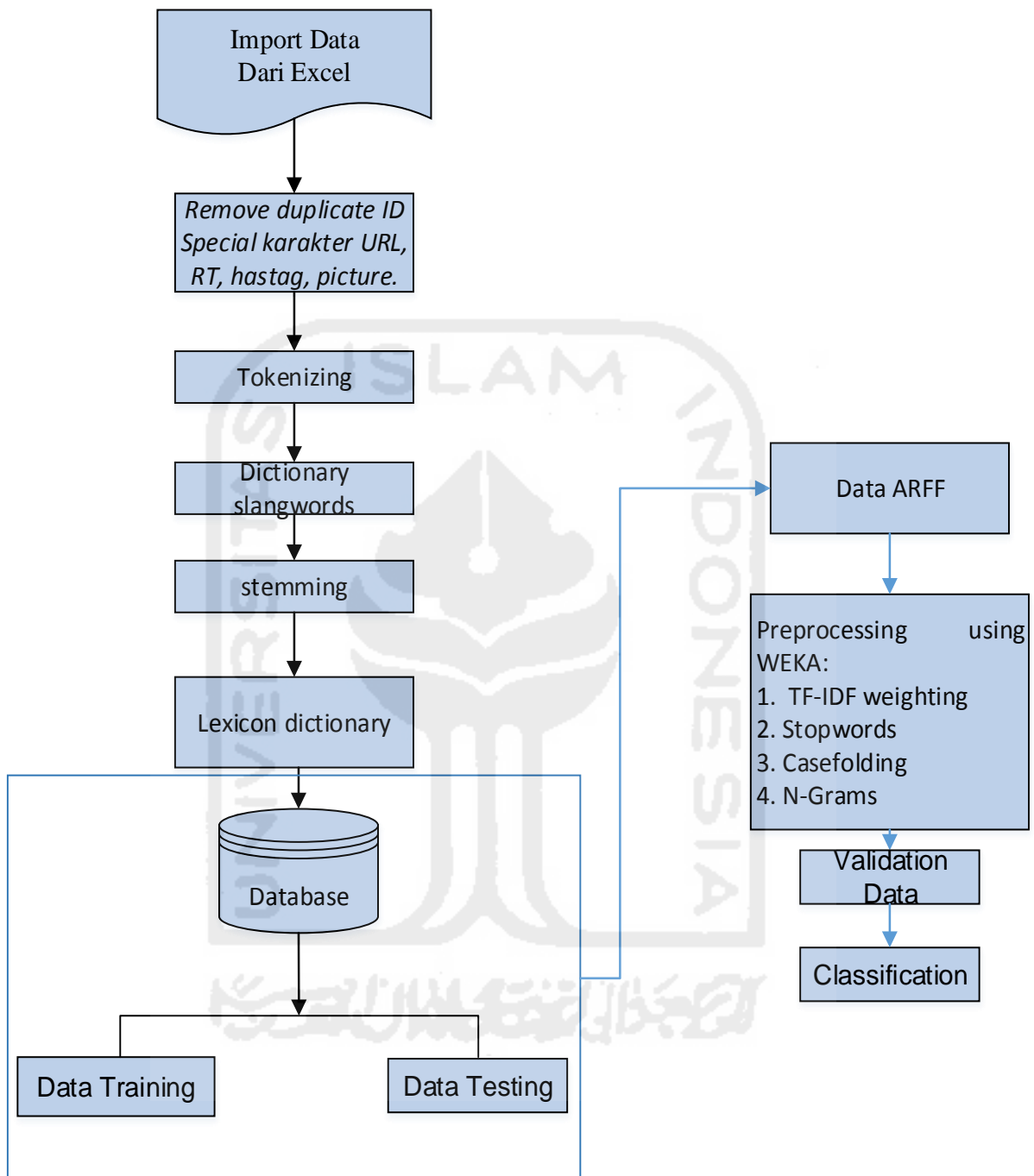
Gambar 3. 7 Hasil Validasi file Json Menggunakan JsonLint

Aplikasi ini tersedia di <http://jsonlint.com/>, setelah melakukan verifikasi selanjutnya data di *convert* dalam bentuk *.csv* atau excel agar lebih mudah dalam *cleansing* data, pada penelitian ini menggunakan *convert .csv* secara online dan tersedia di <http://www.convertcsv.com/json-to-csv.htm>, selanjutnya data disimpan ke dalam Database untuk diolah lebih lanjut.

3.4 Preprocessing Data

Data yang telah diubah dalam bentuk *.csv* selanjutnya dilakukan pembersihan atau *preprocessing*. Hal ini dilakukan agar mendapatkan data terstruktur yang mudah diolah baik secara manual maupun menggunakan *machine learning*. *preprocessing* ini terdiri dari penghapusan ID duplicate, penghapusan URL, penghapusan karakter khusus *hashtag*, *RT* dan gambar, Normalisasi dari kata tidak baku menjadi kata baku, *tokenizing* yaitu memecah kalimat menjadi kata, *casefolding* yang merupakan perubahan semua kalimat dalam bentuk huruf kecil, *stopword* yaitu membuang akhiran kata *possessive* seperti *-kah*, *-lah*, *-pun*,

Stemming atau pencarian akar kata, dan *N-grams*. Alur dari *preprocessing* ini dapat dilihat pada **Gambar 3.8** berikut:



Gambar 3.8 Teknik Preprocessing Data

3.4.1 Menghapus Special Karakter

Pada tahap ini dilakukan proses penghapusan karakter-karakter yang dapat mengganggu proses analisis, baik data training maupun data testing. Karakter yang dihapus adalah Duplicate ID, URL, RT, Gambar, *hastag* dan *special character* lainnya seperti tanda baca koma, kurung dll. Proses penghapusan dilakukan secara manual menggunakan Excel dengan

metode *fine* dan *replace* (Saputra 2015). **Tabel 3.1** menunjukkan penghapusan dengan *Fine* dan *Replace*:

Tabel 3. 1 Tabel Find dan Replace

Yang dihapus	Find	Replace	Keterangan
URL	http*[spasi]	[spasi]	Link di depan
URL	[spasi]http*[spasi]	[spasi]	Link ditengah
URL	[spasi]http*	[spasi]	Link dibelakang
Gambar	Pic.Twitter*[spasi]	[spasi]	Gambar didepan
Gambar	[spasi]pic.Twitter*	[spasi]	Gambar dibelakang
Gambar	[spasi]pic.twitter*[spasi]	[spasi]	Gambar ditengah
@	@*[spasi]	[spasi]	Akun didepan
@	[spasi]@*	[spasi]	Akun dibelakang
@	[spasi]@* [spasi]	[spasi]	Akun ditengah
#	#[spasi]	[spasi]	Hastag didepan
#	[spasi]#*	[spasi]	Hastag ditengah
#	[spasi]#* [spasi]	[spasi]	Hastag dibelakang

3.4.2 Normalisasi kalimat

Normalisasi kalimat di perlukan untuk menyetarakan kata pada kalimat. Adapun langkah-langkah dalam Normalisasi adalah sebagai berikut:

a. Tokenizing

Pada tahap ini dilakukan secara manual pada Excel dengan cara mengganti spasi menjadi koma. *Tokenizing* merupakan pemotongan string input berdasarkan tiap kata menyusunnya. Sebagai contoh pada *tweet*:

“Kamu fikir nemu dompet orang di jalan itu berkah Tolol” akan di pecah menjadi kata kamu, kata fikir, kata nemu, kata dompet, kata orang dan seterusnya. Setelah dilakukan tokenizing selanjutnya akan mudah dilakukan normalisasi kalimat dari kata tidak baku menjadi baku atau kata *slang* menjadi baku dengan merujuk pada Kamus Besar Bahasa Alay (KBBA).

b. Kamus KBBA

Komentar yang diberikan seseorang tidak semuanya bahasa baku, banyak sekali yang menggunakan bahasa gaul, misalnya: “gue”, “loe” dan lain-lain, serta tidak jarang pula yang menggunakan potongan kata, misalnya: “yg”, “brp”, “bgm” dan lain-lain. Kata yang tidak dinormalisasi lebih dahulu akan dikenali oleh *machine learning* sebagai kata yg berbeda, misalnya: ‘semoga’ dan ‘smoga’ yang seharusnya memiliki makna yang sama akan menjadi

beda makna dikarenakan penulisannya yang berbeda. Untuk itu dilakukan normalisasi kata dari yang tidak baku menjadi baku. Untuk normalisasi ini menggunakan bantuan kamus KBBA. Dibawah ini adalah Contoh tabel dari kata tidak baku menjadi kata baku:

Tabel 3. 2 Contoh kata tidak baku menjadi baku

Kata Tidak Baku	Kata Baku
Brp	Berapa
Sm	Sama
Njir	Anjing
Syg	Sayang
Klw, Low, Klo	Kalau
Kamuh, Kamyu, ello, elu	Kamu
Aj	Saja
Nyet	Monyet

c. Penggunaan Rumus

Rumus yang digunakan untuk mengganti kata tidak baku menjadi kata baku adalah
`=IF(ISNA(VLOOKUP('data training'!$1:$1048576,KBBA!$1:$1048576,2,FALSE))=TRUE,'data training'!$1:$1048576,VLOOKUP('data training'!$1:$1048576,KBBA!$1:$1048576,2,FALSE)).`

Keterangan:

Data Training = tabel data training berisi kalimat yang dipecah menjadi kata

KBBA = nama table yang berisi kata tidak baku dan kata baku

ISNA = rumus yang digunakan untuk mengatasi output berupa #N/A yang artinya Not Available, sehingga tidak perlu dilakukan penghapusan satu persatu.

IF = adalah fungsi (kondisi jika benar, jika salah).

VLOOKUP = rumus yang berfungsi untuk mencari kolom pertama dalam satu rentang sel, kemudian mengembalikan nilai apapun yang ada pada baris yang sama.

Penjelasan rumus:

Apabila kata pada workseet pada table Data training \$1 sampai \$1048576 tidak terdapat satupun pada workseet KBBA dari rentang \$1 sampai \$1048576, maka kata tidak diubah

menjadi kolom ke-2 atau kolom B pada workseet KBBA, melainkan akan dikeluarkan output berupa #N/A, karena dituliskan rumus ISNA, maka kata akan dikembalikan seperti semula atau kata tidak terjadi perubahan. Dan apabila kata tersebut terdapat pada workseet KBBA pada kolom \$1 sampai \$1048576, maka kata tersebut akan diubah menjasi kata yang terdapat pada kolom berikutnya atau ke-2 pada workseet KBBA.

3.4.3 Stemming

Proses ini adalah tahap mencari akar kata dari tiap kata hasil *filtering*. Proses ini mengubah kata menjadi kata dasarnya dengan menghilangkan imbuhan-imbuhan pada kata dalam dokumen atau mengubah kata kerja menjadi kata benda. *Stem* (akar kata) adalah kata inti setelah imbuhan dihilangkan (awalan dan akhiran). Misalnya kata "perancangan" dan "merancang" akan diubah menjadi sebuah kata yang sama, yaitu "rancang". Proses stemming sangat tergantung kepada bahasa dari kata yang akan di-stem.

Penelitian ini menggunakan Sastrawi Master. Sastrawi master adalah *library* php sederhana yang menyediakan *stemming* kata bahasa Indonesia. kamus kata dasar yang digunakan Sastrawi berasal dari kateglo.com dengan sedikit perubahan dan masing-masing mempunyai lisensi Sastrawi dan lisensi kateglo. Sastrawi dapat diunduh secara gratis di alamat <https://github.com/sastrawi/sastrawi>.


Untuk melakukan proses ini peneliti menggunakan bahasa pemrograman python. Namun terlebih dahulu install Library master sastrawi pada python, kemudian buat *script* pada *console* python seperti berikut:

```
# import StemmerFactory class
from Sastrawi.Stemmer.StemmerFactory import StemmerFactory
# create stemmer
factory = StemmerFactory()
stemmer = factory.create_stemmer()
# stemming process
sentence = 'kamu memang politisi kampung yang sedang membela tuannya'
output = stemmer.stem(sentence)
print(output)
# hasil stemming
# "kamu memang politis kampung yang sedang bela tuan"
```

Gambar 3. 9 Script Stemming

3.4.4 Penggunaan Lexicon

Penggunaan lexicon pada prosesnya sama dengan tahap normalisasi cleansing, yaitu *preprocessing* data menggunakan excel. Tetapi perbedaannya adalah pada workseet yang berisi kamus, pada normalisasi kolom A berisi kata tidak baku dan pada kolom B berisi kata baku, sedangkan pada proses pemanfaatan lexicon ini, kolom A berisi kamus bullying, pronoun atau kata ganti orang kedua dan ketiga misanya, “kamu”, “kau”, “anda” dan kamus negasi seperti “bukan” dan “tidak”, sementara kolom B berisi kata bullying yang saya ubah menjadi “badword” untuk kamus bullying dan kata pronoun untuk kamus pronoun dan kata negasi untuk kamus negasi. **Gambar 3.10** dibawah adalah Contoh kamus lexicon.



	A	B	C
16	menteri	pronoun	
17	presiden	pronoun	
18	walikota	pronoun	
19	gubernur	pronoun	
20	mereka	pronoun	
21	bapak	pronoun	
22	ibu	pronoun	
23	kalian	pronoun	
24	perempuan	pronoun	
25	tidak	negasi	
26	bukan	negasi	
27	sableng	badword	
28	gila	badword	
29	jelek	badword	
30	edan	badword	
31	tolol	badword	
32	bego	badword	
33	geblek	badword	
34	goblok	badword	
35	dongo	badword	
36	bodoh	badword	
37	buta	badword	
38	tuli	badword	
39	sableng	badword	

Gambar 3. 10 Kamus Lexicon

Rumus yang digunakan sama dengan rumus pada proses Normalisasi kata tidak baku menjadi kata baku. **Gambar 3.11** dibawah merupakan contoh data sebelum penggunaan lexicon:

E	F	G	H	I	J	K	L	M	N	O	P	Q	R	S	T	U
396	brengsek	kamu	hebat	juga	gitu	saja	pakai	tanya								
1565	brengsek	kamu	hebat	juga	gitu	saja	pakai	tanya								
941179	brengsek	kamu	hebat	juga	gitu	saja	pakai	tanya								
122	kamu	terlalu	brengsek	buat	sayang	sama	dia	ohseh								
8	brengsek	juga	diko	main	sendiri	saja	kamu	modus	sempat	dasar	brengsek					
138	dasar	brengsek														
216	dasar	brengsek														
2825	arti	kamu	juga	alien	gayung	geblek										
97	jakarta	bekasi	tiga	seperdua	jam	dasar	geblek									
278	geblek	otak	kamu	balik												
514	geblek	kamu														
1253	saya	baca	dan	kamu	geblek											
735	tai	kontok	kamu	omong	apa	geblek	kutil									
1671	geblek	memang	jidat	kamu	jong	yang	jendol									
146	apa	apa	pc	gaya	banget	kamu	pc	orang	geblek	malah	nanya	saya	memang	saya	pacar	galer
1239	maksimal	jek	tidak	mungkin	positif	geblek	kamu									
214	dasar	gila														
1915	dasar	cewek	gila	memang												
3474	gila	kamu	allahu													
8892	kamu	gila														
322	malique	respect	gila	kamu												
521	bodor	gila	pokok	kamu	asik	cinta	kamu									
8646	ngapain	kamu	gila													
685	gila	baik	banget	kamu												

Gambar 3. 11 Data Sebelum Penggunaan Lexicon

Setelah menggunakan lexicon data tersebut menjadi seperti pada **Gambar 3.12**

berikut:

E	F	G	H	I	J	K	L	M	N	O	P	Q	R	S	T	U	
158	396	badword	pronoun	hebat	juga	gitu	saja	pakai	tanya								
159	1565	badword	pronoun	hebat	juga	gitu	saja	pakai	tanya								
160	941179	badword	pronoun	hebat	juga	gitu	saja	pakai	tanya								
161	122	pronoun	terlalu	badword	buat	sayang	sama	pronoun	ohseh								
162	8	badword	juga	diko	main	sendiri	saja	pronoun	modus	sempat	dasar	badword					
163	138	dasar	badword														
164	216	dasar	badword														
165	2825	arti	pronoun	juga	badword	gayung	badword										
166	97	jakarta	bekasi	tiga	seperdua	jam	dasar	badword									
167	278	badword	otak	pronoun	balik												
168	514	badword	pronoun														
169	1253	saya	baca	dan	pronoun	badword											
170	735	badword	kontok	pronoun	omong	apa	badword	kutil									
171	1671	badword	memang	jidat	pronoun	jong	yang	jendol									
172	146	apa	apa	pc	gaya	banget	pronoun	pc	orang	badword	malah	nanya	saya	memang	saya	pacar	galer
173	1239	maksimal	jek	negasi	mungkin	positif	badword	pronoun									
174	214	dasar	badword														
175	1915	dasar	pronoun	badword	memang												
176	3474	badword	pronoun	allahu													
177	8892	pronoun	badword														
178	322	malique	respect	badword	pronoun												
179	521	bodor	badword	pokok	pronoun	asik	cinta	pronoun									
180	8646	ngapain	pronoun	badword													
181	685	badword	baik	banget	pronoun												

Gambar 3. 12 Data Setelah Penggunaan Lexicon

3.4.5 Data training dan Data testing

Setelah penggunaan lexicon pada data, maka data dapat dibagi menjadi dua bagian yaitu data training dan data testing. Metode pembagian data ini dibagi seimbang yaitu 50% data training dan 50% untuk data testing karena data yang tidak seimbang klasifikasi yang dibangun

memiliki kecenderungan untuk mengabaikan *minority class* (Buntoro 2016). Selanjutnya data di ubah ke format ARFF, Namun untuk data training sebelum di ubah ke format ARFF terlebih dahulu dilakukan pelabelan secara manual berdasarkan pattern atau pola yang mengindikasikan bahwa kalimat tersebut mengandung bullying. Tabel dibawah menunjukkan pola bahwa suatu kalimat mengandung bullying jika terdapat unsur sebagai berikut (Yin 2009):

Tabel 3. 3 Pola Cyberbullying

BadWord!	Pronoun
Kamu	BadWord
...	BadWord	Pronoun	...
Pronoun	BadWord
Pronoun	BadWord	...

Tetapi untuk kalimat yang menggunakan kata negasi diikuti kata BadWord maka kalimat tersebut menjadi negatif bullying. Demikian juga kalimat yang mengandung unsur pertanyaan disertai pronoun dan BadWord maka kalimat tersebut bernilai negatif bullying. Tabel dibawah merupakan pola kalimat bullying yang disertai negasi dan Question.

Tabel 3. 4 Pola Negasi

Negasi	Badword
Question	Pronoun	BadWord	...

Setelah dilakukan pelabelan secara manual untuk data training, selanjutnya adalah mengubah file menjadi ARFF. Proses perubahan data bisa dilakukan secara manual maupun otomatis. Perubahan secara manual dilakukan dengan cara data diubah ke.txt terlebih dahulu kemudian menambahkan @relation untuk nama datanya, @attribute berupa text type data string, @attribute @@class@@ {positif,negatif} merupakan kelas atribut berupa positif, dan negatif kemudian @data yang berisi datanya yang ditambahkan single quote dan dilabeli “pos” untuk kalimat positif, “neg” untuk kalimat negatif lalu file .txt di save dengan ekstension ARFF. Sementara pengubahan secara otomatis dilakukan dengan cara mengubah file ke bentuk .CSV, lalu buka tools WEKA, open file, setelah data terbuka save as kembali dengan mengubah type data .CSV menjadi ARFF.

3.4.6 Pengolaan Data Menggunakan WEKA

Weka adalah aplikasi data mining open source berbasis Java. Aplikasi ini dikembangkan pertama kali oleh Universitas Waikato di Selandia Baru sebelum menjadi bagian dari Pentaho. Weka terdiri dari koleksi algoritma machine learning yang dapat digunakan untuk melakukan generalisasi / formulasi dari sekumpulan data sampling. Walaupun kekuatan Weka terletak pada algoritma yang makin lengkap dan canggih, kesuksesan data mining tetap terletak pada faktor pengetahuan manusia implementornya. Tugas pengumpulan data yang berkualitas tinggi dan pengetahuan pemodelan dan penggunaan algoritma yang tepat diperlukan untuk menjamin keakuratan formulasi yang diharapkan (susanto 2012). Penelitian ini menggunakan WEKA versi 3.8.0 Gambar dibawah adalah *Interface* dari WEKA 3.8:



Gambar 3. 13 Interface WEKA

Preprocessing menggunakan WEKA dilakukan dengan cara sebagai berikut:

a. Mengubah Data ke dalam bentuk Vektor

Pada tahap ini, data yang berupa kalimat yang sudah dilabeli dengan positif dan negatif akan diubah kedalam bentuk vector. Adapun caranya adalah pada aplikasi WEKA open file dan arahkan ke file .arff yang sudah diolah sebelumnya, setelah itu tekan tombol choose pada filter dan pilih StringToWordVector. Kemudian filters, Unsupervised, attribute kemudian StringToWordVector lalu Apply. Hal ini dilakukan pada data training maupun data testing.

Pada bentuk vector ini, masing-masing token mewakili satu attribute, contoh untuk data training 50%, data yang sudah diubah kedalam vector dengan jumlah data positif bullying sebanyak 226 dan data negatif bullying sebanyak 66 sehingga totalnya menjadi 292 data training.

b. Pembobotan TF-IDF

Proses pengubahan data teks menjadi data vektor dilakukan dengan membaca kata satu persatu dan menghitung nilai tf-idf. Nilai tf-idf adalah kemunculan kata (*term frequency*) dalam kalimat dikalikan log jumlah dokumen/*record* dibagi jumlah dokumen/*record* yang mengandung kata yang dimaksud.

Pada penelitian ini Pembobotan TF-IDF dan preprocessing dengan menggunakan WEKA dilakukan dengan cara klik text box yang berisikan StringToWordVector. Setelah muncul gambar (weka.gui.genericobjecteditor), lakukan pilihan sesuai dengan preprocessing yang akan dilakukan, seperti *casefolding*, *Token N-gram*, *Penggunaan Stopword* dan penghapusan emoticon lalu apply.

c. Stopword Removal

Stopwords removal adalah proses menghilangkan kata-kata yang umum digunakan dan tidak mempunyai informasi yang berharga pada suatu konteks. Kamus stopwords yang digunakan berasal dari (Tala 2003) yang diunduh disitus <http://hikaruyuuki.lecture.ub.ac.id/kamus-kata-dasar-dan-stopword-list-bahasa-indonesia/>.

Contohnya dapat dilihat pada table dibawah:

Tabel 3. 5 Tabel Stopword Tala

Ada
adalah
adanya
adapun
agak
agaknya
agar
akan
akankah
akhir
akhiri
akhirnya
aku
akulah
amat
amatlah
sampaikan
sana
sangat
sangatlah
satu
...dan seterusnya

Untuk menggunakan Stopword Tala Bahasa Indonesia dilakukan dengan cara melakukan klik pada tulisan “weka-3-8-0” kemudian memilih stopwords yang akan digunakan. Untuk mengubah semua huruf kecil dengan memilih “true” pada lowercasetoken. Untuk melakukan normalisasi panjang dokumen terhadap seluruh data dengan memilih “normalize all data” pada normalizeDocLength.

d. N-Gram

Penelitian ini mengimplementasikan tokenisasi N-Gram yang tidak terikat dengan satu aturan bahasa apapun, Tokenisasi menggunakan N-Gram adalah tahap pemrosesan dimana teks input dibagi menjadi unit-unit kecil yang disebut *token* sepanjang n karakter. Dalam bahasa Indonesia, frasa dengan satu kesatuan arti memiliki maksimal 3 kata, pembagian *token* dibagi menjadi Unigram, Bigram, Trigram dan N-Gram, berikut contoh pemecahan pada kalimat “orang pada buta semua pendukung semu”.

Unigram: yaitu *token* yang terdiri dari hanya satu kata, menghasilkan: “orang”, “pada”, “buta”, “semua”, “pendukung”, “semu”.

Bigram: yaitu *token* yang terdiri dari dua kata, menghasilkan: “orang pada”, “pada buta”, “buta semua”, “semua pendukung”, “pendukung semu”.

Trigram: yaitu *token* yang terdiri dari tiga kata, menghasilkan: “orang pada buta”, “pada buta semua”, “buta semua pendukung”, “semua pendukung semu”.

Proses N-gram pada penelitian ini juga menggunakan *machine learning* WEKA.

Cara penggunaan N-Gram pada WEKA adalah pilih NGramTokenizer dengan cara klik pada tombol “choose” pada tokenizer, kemudian pilih Ngramtokenizer. Selanjutnya memecah kata dengan mengubah angka pada NgramMaxSize dan NgramMinSize yang terdapat pada gambar. Untuk unigram NgramMaxSize diubah menjadi 1 dan NgramMinSize menjadi 2. Kemudian Ngram dengan mengubah angka Ngrammaxsize menjadi 3 dan Ngrammaxsize menjadi 1. Selanjutnya mengubah delimiter, dan menghapus emoticon.

e. Validation Data

Penerapan untuk classifier akan diuji sesuai dengan pilihan yang ditetapkan dan sesuai kebutuhan penelitian. Ada beberapa test option yang bias dipilih pada WEKA sebelum melakukan klasifikasi yaitu:

1. Use training set

Pengetesan dilakukan dengan menggunakan data training itu sendiri.

2. Supplied test set

Pengetesan dilakukan dengan menggunakan data lain. Dengan menggunakan option inilah, bisa dilakukan prediksi terhadap data tes.

3. Cross-validation

Pada cross-validation, akan ada pilihan berapa fold yang akan digunakan. Nilai default-nya adalah 10. Mekanisme-nya adalah sebagai berikut : Data training dibagi menjadi k buah subset (subhimpunan). Dimana k adalah nilai dari fold. Selanjutnya, untuk tiap dari subset, akan dijadikan data tes dari hasil klasifikasi yang dihasilkan dari k-1 subset lainnya. Jadi, akan ada 10 kali tes. Dimana, setiap datum akan menjadi data tes sebanyak 1 kali, dan menjadi data training sebanyak k-1 kali. Kemudian, error dari k tes tersebut akan dihitung rata-ratanya.

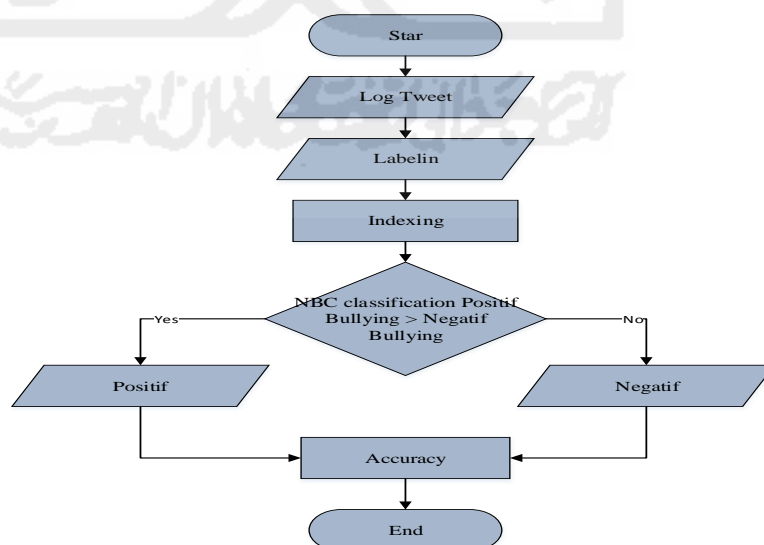
4. Percentage split

Hasil klasifikasi akan dites dengan menggunakan k% dari data tersebut. k merupakan masukan dari user.

Dalam penelitian ini, digunakan Cross-validation dengan nilai default 10 sehingga disebut juga *10 cross validation folds*. Tujuan penggunaan model ini adalah untuk menentukan pola untuk data testing terhadap data training tujuannya untuk membatasi masalah *overfitting* dan mengetahui bahwa model ini mengeneralisasi data pada data testing untuk mendapatkan hasil klasifikasi.

3.5 Klasifikasi

Dalam menentukan akurasi dengan menggunakan metode Naïve bayes, dilakukan berdasarkan probabilitas kemunculan kata. Gambar 3.14 dibawah adalah alur dari metode Naive Bayes Classifier:



Gambar 3. 14 Flowchart Naive Bayes Classifier

Data teks yang digunakan adalah data bersih yang telah melalui preprocessing, selanjutnya diberi label secara manual pada data training, setelah itu dilakukan pembobotan TF-IDF, dan validasi data menggunakan *10 fold cross validation* lalu klasifikasi menggunakan Naïve Bayes, untuk teks yang positif cyberbullying akan di klasifikasikan ke class positif bullying dan teks yang negatif bullying akan ke class negatif bullying. Demikian pula untuk jenis cyberbullying akan diklasifikasikan ke class masing-masing Seperti jenis bullying yang *related psychology* akan diklasifikasikan ke *class related psychology* dan seterusnya untuk jenis cyberbullying yang lain. Keseluruhan proses ini dilakukan pada *machine learning* WEKA.

Penentuan probabilitas positif bullying dan negatif bullying secara manual dapat dilihat pada Contoh, untuk jumlah data training adalah 292, untuk data positif bullying sebanyak 226 dan negatif bullying sebanyak 66:

- Probabilitas data positif $P(Y=\text{positif})=226/292 = 0,77$
- Probabilitas data negatif $P(Y=\text{negatif})=60/292 = 0,20$
- 292 merupakan jumlah seluruh data (226+60).

Proses selanjutnya yaitu set data testing pada WEKA, data testing ini juga merupakan data bersih, dari pola yang telah di proses pada data training akan mengikuti pola untuk data testing sehingga hasil klasifikasi dapat diprediksi. Untuk mengetahui hasil klasifikasi keseluruhan dapat dilihat pada hasil klasifikasi data training kemudian tambahkan pada hasil prediksi.

Bab 4 Hasil dan Pembahasan

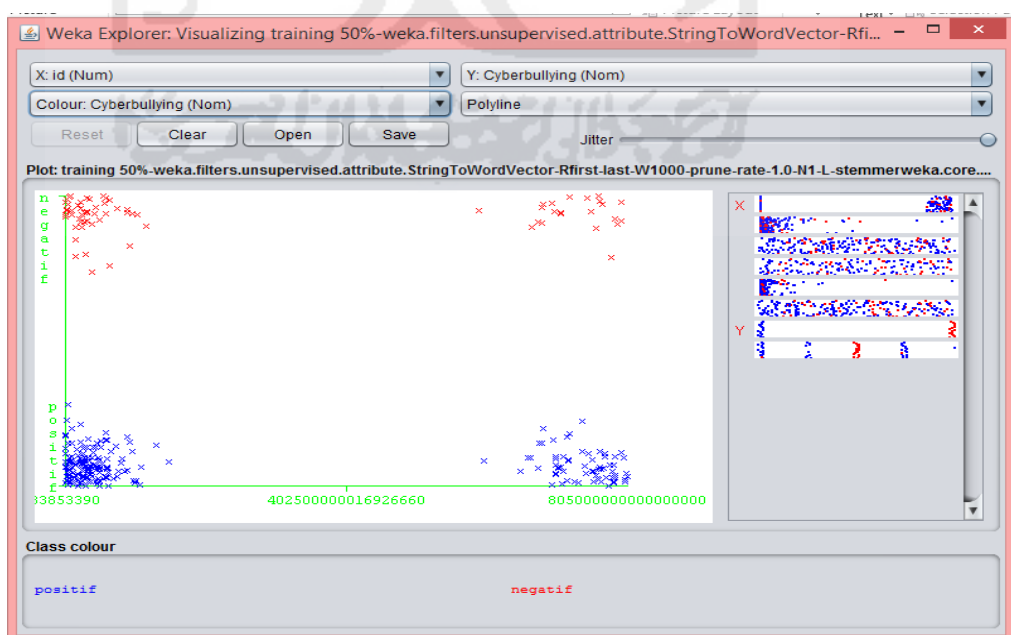
4.1 Deskripsi Penelitian

Pada bagian ini melakukan evaluasi data berdasarkan pembagian *data set*. Jumlah data yang digunakan sebanyak 583 data tweet. Untuk membandingkan akurasi dalam mengklasifikasikan data, peneliti melakukan beberapa pembagian data training dan data testing secara random diantaranya. Dari perbandingan tersebut dapat diketahui yang mana lebih akurat dalam klasifikasi data dengan menggunakan *Naïve Bayes Classifier*.

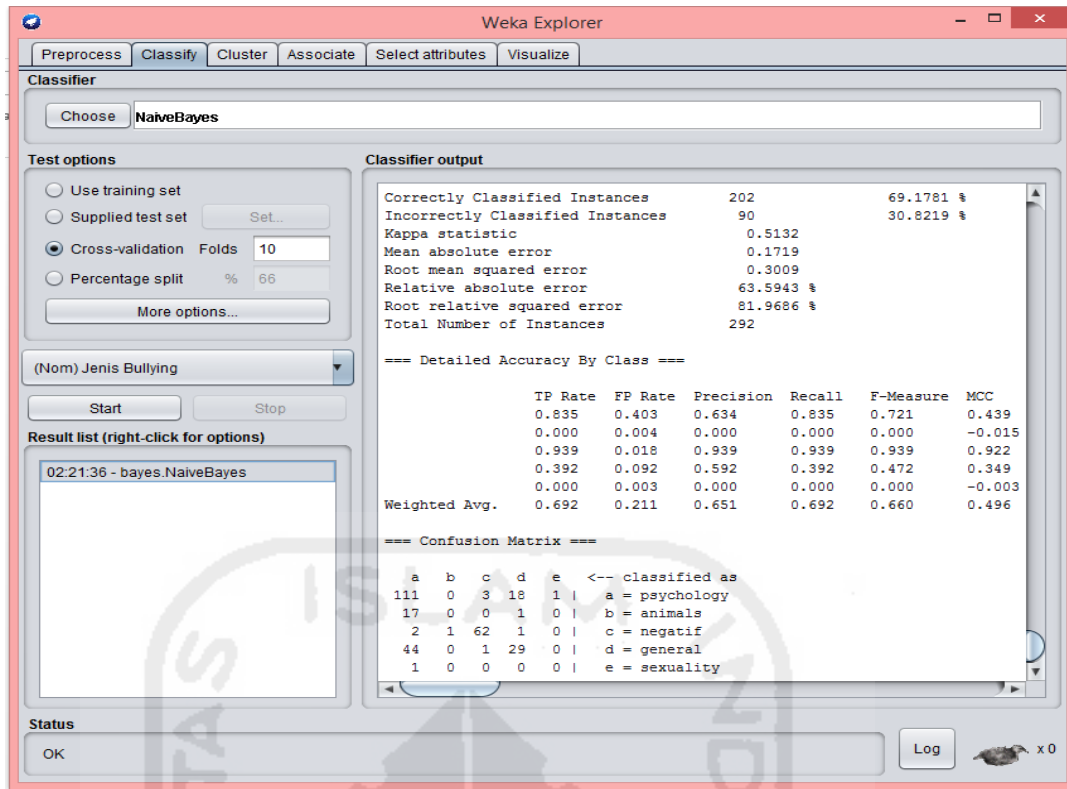
4.2 Evaluasi Hasil Penelitian

4.2.1 Klasifikasi Cyberbullying dan Non-Bullying

Data training yang sudah di bersihkan kemudian diubah dalam format *.arff* untuk selanjutnya akan diolah menggunakan aplikasi WEKA untuk mengetahui hasil klasifikasi cyberbullying dan non bullying pada *Log tweet*. Dataset dibagi seimbang yaitu 50% data training dan 50% untuk data testing karena data yang tidak seimbang klasifikasi yang dibangun memiliki kecenderungan untuk mengabaikan *minority class*. Setelah diproses menggunakan WEKA maka akan menampilkan grafik seperti pada **Gambar 4.1** yang menunjukkan hasil pengklasifikasian sesuai atribut masing masing dan pengaruh atribut lainnya.



Gambar 4.1 Grafik Klasifikasi Data Training Menggunakan WEKA



Gambar 4.5 Klasifikasi Data Training Jenis Cyberbullying

Gambar diatas menunjukkan keakuratan dalam klasifikasi jenis bullying 69,17%. Secara manual dapat dirumuskan sebagai berikut:

$$\text{Presentase Akurasi} = + \frac{\text{banyak prediksi yang benar}}{\text{total banyaknya data}} \times 100\%$$

Hasil klasifikasi

$$\begin{aligned}
 &= \\
 &= (111+0+62+29+0)/(111+0+3+18+17+0+0+1+0+2+1+62+1+0+44+0+1+29+0+0+1+0) \\
 &*100\% \\
 &=202/291 \times 100\% \\
 &= 69,17\%
 \end{aligned}$$

Proses selanjutnya adalah klasifikasi untuk data testing. Seperti sebelumnya, pilih menu *Test Option* dan klik *Supplied test set* lalu pilih data testing pada penyimpanan data lalu proses. Hasil prediksi data tersebut di *save* dalam bentuk *.arff* dan hasilnya dapat dilihat pada WEKA GUI menu *Tools* dan pilih *ArffViewer*. **Gambar 4.6** dibawah adalah hasil dari prediksi jenis cyberbullying yang telah diklasifikasi.

File Edit View

prediksi cyberbullying 50%.arff prediksi jenis cyberbullying 50%.arff

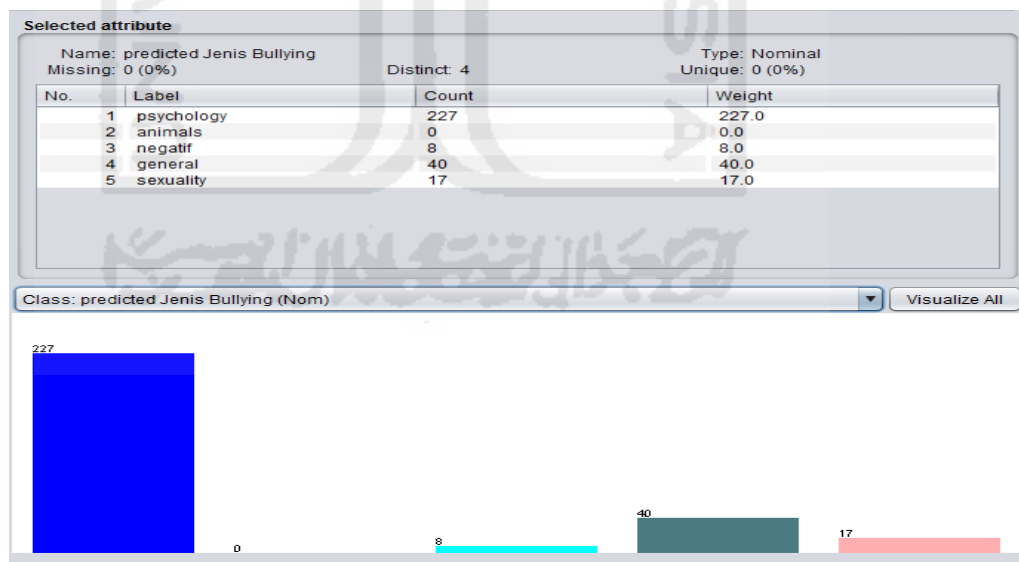
Relation: training 50% jenis bullying-weka.filters.unsupervised.attribute.StringToWordVector-Rfirst-last-W1000-prune-rate-1.0-N1-L-stemmerw...

No.	1: id	2: friend	3: follower	4: text	5: prediction margin	6: predicted Jenis Bullying	7: Jenis Bullying
	Numeric	Numeric	Numeric	Nominal	Numeric	Nominal	Nominal
1	8.0E...	289.0	329.0		0.340678	psychology	
2	1.49...	6408.0	11309.0	gen...	-0.232862	negatif	
3	7.91...	34.0	12.0		0.341005	psychology	
4	7.69...	3056.0	12050.0		-0.329035	general	
5	7.46...	45.0	33.0		0.325707	psychology	
6	7.17...	547.0	320.0	das...	-0.136495	general	
7	3.12...	31.0	1629.0	das...	-0.119261	general	
8	2.81...	340.0	402.0	das...	-0.139865	general	
9	2.76...	1848.0	600.0	das...	-0.092479	general	
10	7.62...	5.0	43.0	das...	-0.158412	general	
11	7.42...	60.0	7.0	das...	-0.151084	general	
12	3.11...	100.0	85.0	gen...	-0.306231	general	
13	3.01...	794.0	561.0	gen...	-0.273337	general	
14	6.01...	359.0	410.0	gen...	-0.293891	general	
15	3.02...	457.0	59.0	gen...	-0.303241	general	
16	4.82...	247.0	119.0	gen...	-0.303136	general	
17	4.91...	117.0	135.0		-0.951104	sexuality	
18	1.94...	4188.0	4775.0		0.366778	psychology	
19	3.05...	620.0	449.0	naik...	0.026061	psychology	
20	2.20...	202.0	560.0	gen...	-0.287861	general	
21	2.38...	281.0	229.0	gen...	-0.295618	general	
22	4.10...	224.0	137.0	gen...	-0.303136	general	
23	1.34...	332.0	643.0	gen...	-0.28443	general	
24	4.48...	233.0	354.0		0.341921	psychology	
25	2.16...	195.0	487.0	gen...	-0.287861	general	
26	1.70...	2446.0	2515.0		0.351744	psychology	
27	1.33...	63.0	79.0	gen...	-0.306231	general	
28	7.23...	152.0	65.0	gen...	-0.907222	sexuality	
29	2.00...	248.0	420.0	gen...	-0.303136	general	

Data testing untuk prediksi jenis bullying sebanyak 292 record. Dari hasil prediksi...

Gambar 4. 6 Prediksi Klasifikasi Jenis Cyberbullying

Data testing untuk prediksi jenis bullying sebanyak 292 record. Dari hasil prediksi klasifikasi tersebut diketahui jenis cyberbullying yang banyak digunakan. Gambar 4.8 dibawah menunjukkan hasil secara grafik:



Gambar 4. 7 Grafik Prediksi Jenis Cyberbullying

Gambar diatas menunjukkan bahwa jenis cyberbullying yang *related psychology* sebanyak 77.73%, dan untuk jenis cyberbullying *related animals* terprediksi 0.00%, sementara kalimat yang mengandung negatif bullying teridentifikasi sebanyak 2.75%, dan untuk jenis cyberbullying *general* teridentifikasi

sebanyak 13.69% dan jenis cyberbullying *related sexuality* teridentifikasi sebanyak 5.82%.

4.2.3 Hasil Klasifikasi

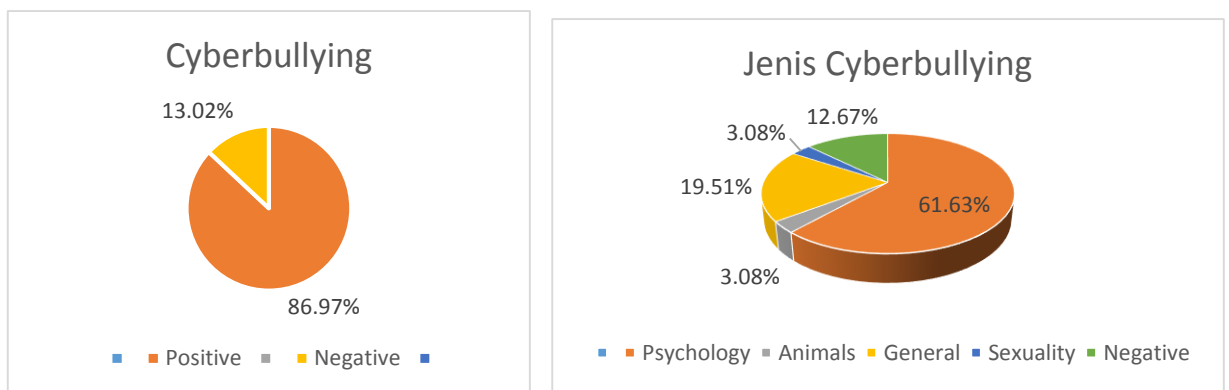
Hasil klasifikasi dapat dilihat pada tabel 4.1 dibawah ini:

Tabel 4. 1 Hasil Presentase Klasifikasi

Klasifikasi	Data Training NBC		Data prediction NBC		Presentase	
	Positif	Negatif	Positif	Negatif	Positif	Negatif
Cyberbullying	77.40%	22.60%	96.55%	3.44%	86.97%	13.02%
Psychology	45.56%		77.73%		61.63%	
Animals	6.16%		0.00%		3.08%	
General	25.34%	22.60%	13.69%	2.75%	19.51%	12.67%
Sexuality	0.34%		5.82%		3.08%	

Perbandingan akurasi pada data training dan prediksi dari data testing dapat diketahui dari tabel diatas bahwa untuk klasifikasi cyberbullying dan non bullying yang teridentifikasi mengandung kalimat bullying pada data training sebanyak 77.40% dan negatif bullying sebanyak 22.60%, untuk prediksi klasifikasi cyberbullying diketahui positif bullying yang teridentifikasi sebanyak 96.55% dan negatif bullying hanya 3.44%.

Untuk data training jenis cyberbullying dapat diketahui bahwa jenis bullying yang *related psychology* sebesar 45.56%, *related animlas* 6.16%, *general bullying* 25.34% dan *sexuality* 0.34%, sementara yang terdeteksi untuk negatif bullying 22.60%. disini terlihat ada perbedaan deteksi negatif bullying sekitar 1.37% untuk klasifikasi cyberbullying dan klasifikasi untuk jenis cybebullying. Untuk hasil prediksi teridentifikasi *related psychology* 77.73%, *related animals* 0.00%, *general bullying* 13.69% dan *sexuality* 5.82% dan terdeteksi negatif 2.75%. untuk presentase keseluruhan dapat dilihat pada **Gambar 4.8** dibawah ini:



Gambar 4. 8 Grafik Cyberbullying dan Jenis Cyberbullying

Secara akumulatif presentase Cyberbullying dan Jenis Cyberbullying dapat dihitung dengan menjumlahkan positif cyberbullying pada data training dan data *prediction*, demikian pula untuk negatif bullying pada data training dan data *prediction* sehingga didapatkan hasil akumulatif keseluruhan dari Cyberbullying baik data training maupun data prediction. Untuk positif cyberbullying diperoleh 86.97%, negatif bullying 12.67%. Untuk jenis cyberbullying related psychology 61.63%, related animals 3.08%, General Bullying 19.51% dan sexuality 3.08%.

Cyberbullying di Indonesia pada periode November sampai Desember 2016, dengan jumlah data sebanyak 583 adalah 86.97%, untuk jenis cyberbullying yang banyak digunakan adalah related psychology sebanyak 61.63% itu berarti kata umpatan dan makian yang ditujukan untuk membullying banyak yang menyerang secara psikologi seperti kata “goblok”, “idiot”, “tolol”, “sarap” dan lain-lain. Untuk bullying yang related animal atau memaki seseorang dengan sebutan binatang sebanyak 3.08% seperti kata “anjing”, “babi” dan lain-lain. Untuk general bullying sebanyak 19.51%, kata makian yang umum digunakan seperti kata makian “setan”, “keparat”, “bajingan” dan lain-lain. Dan makian yang related sexuality sebanyak 3.08% seperti kata makian “banci”, “lonte” dan lain-lain.

Penelitian terkait cyberbullying di Indonesia, dibandingkan dengan penelitian ini secara keseluruhan terletak pada objek dan sampling data yang digunakan. Penelitian lainnya untuk sampling banyak menggunakan kuesioner pada siswa SMP dan SMA pengguna sosial media dengan melakukan wawancara mendalam, pengamatan dan dokumentasi, ada juga dengan mengambil sampel pertemanan pada jejaring sosial kemudian melakukan analisis dengan teori pendekatan psikologis, ilmu komunikasi, ilmu sosial dan Informatika. Sementara penelitian ini, mengambil sampel data langsung dari *database* Twitter kemudian melakukan analisis secara menyeluruh pada data yang diperoleh dengan Algoritma Data Mining Naïve Bayes Classifier.

Penggunaan Algoritma Naïve Bayes Classifier untuk penelitian ini kurang maksimal dibandingkan penelitian lainnya. Sebagai contoh klasifikasi berita secara otomatis akurasi 90.23%, klasifikasi untuk informasi kemacetan lalu lintas melalui twitter akurasi 93.58%, klasifikasi untuk menilai kelayakan kredit pada BCA Finance akurasi 92.54%. Sementara untuk penelitian ini akurasi untuk klasifikasi konten yang memuat cyberbullying yaitu 86.97% jadi akurasinya lebih rendah dibanding kasus lainnya namun cukup untuk memberi informasi tingginya konten yang mengandung bullying pada Twitter.

Bab 5 Kesimpulan dan Saran

5.1 Kesimpulan

Berdasarkan penelitian yang telah dilakukan maka kesimpulan yang dapat ditarik adalah sebagai berikut:

1. Teknik pengumpulan bukti digital pada Jejaring Sosial Twitter dilakukan dengan cara, melakukan registrasi aplikasi pada link yang telah disediakan oleh Twitter untuk mendapatkan *Access Token*. Setelah akses token didapatkan kemudian buat *script* sebagai media *crawling* antara akun *Twitter* dan *API Search Twitter*, *Access Token* yang telah didapatkan di input pada *script* agar mendapat akses untuk melakukan *crawling* data, pembuatan *script* bisa menggunakan bahasa pemrograman apa saja seperti java, PHP dan lain-lain, namun untuk penelitian ini peneliti menggunakan Python dari Anaconda 4.0 untuk membuat *script* tersebut, selanjutnya melakukan *Crawling* data untuk mendapatkan *Log Twitter* dan di simpan dalam bentuk file *Json*.
2. Klasifikasi bukti digital *Cyberbullying* pada Jejaring Sosial Twitter menggunakan *Naïve Bayes Classifier* dilakukan dengan cara: pertama adalah mengumpulkan *Log tweet* dari database Twitter. Kedua, lakukan pembersihan data atau preprocessing data baik secara manual maupun otomatis dengan *Machine Learning* yang sesuai, dalam hal ini peneliti menggunakan WEKA. Ketiga, lakukan pembobotan term Frequency dan Invers Document Frequency (TF-IDF) untuk menghitung bobot masing-masing teks, kemudian lakukan validasi data, penelitian ini menggunakan *10 fold cross validation* dan terakhir adalah lakukan klasifikasi menggunakan *Naïve Bayes*. Adapun Hasil klasifikasi bukti digital *log tweet* cyberbullying adalah untuk teks yang positif bullying sebanyak 86.97%, negatif bullying 13.02%. untuk Jenis cyberbullying yang digunakan related *psychology* sebanyak 61.63%, related animals sebanyak 3.08%, general bullying sebanyak 19.57% dan sexuality sebanyak 3.08% dan negatif bullying 12.67%. hasil tersebut diatas untuk data periode November sampai desember 2016 dan dapat disimpulkan bahwa pelaku cyberbullying di Indonesia untuk periode tersebut cukup tinggi.

5.2 Saran

Berdasarkan hasil pengujian tersebut maka penulis memberikan saran sebagai berikut:

1. Hasil *crawling* data tidak semuanya berbahasa Indonesia, terdapat gabungan bahasa daerah dan bahasa asing, maka perlu pengembangan pada *script* agar dapat melakukan *filter* untuk *crawl* bahasa Indonesia saja.
2. Proses *preprocessing* menjadi lambat karena *Library stemming* yang digunakan masih dilakukan perbaikan pada hasil *stemming*, sehingga pada data yang lebih banyak butuh waktu yang lebih lama untuk mengubah data tersebut menjadi terstruktur, untuk penelitian selanjutnya Perlu dikembangkan kamus aplikasi untuk melakukan *stemming* dan *stopwords* khusus untuk analisis *cyberbullying*, *harassment* dan *hate speech* agar akurasi untuk klasifikasi lebih baik dan lebih cepat. Selain itu data yang digunakan diperbanyak.



Daftar Pustaka

- Akaichi, J. 2013. *Sosial Networks 'Facebook' Statutes Updates Mining For Sentiment Classification*. SosialCom/PASSAT/BigData/EconCom/ BioMedCom, Le Bardo, Tunisia.
- Akbar, A.U., 2013. Implikasi Hukum Kebebasan Berpendapat di Jejaring Sosial Dalam Terwujudnya Delik Penghinaan, Skripsi, Fakultas Hukum Universitas Hasanuddin.
- Aliandu,P., 2012, Analisis Sentimen Tweet Berbahasa Indonesia di Twitter, *Tesis*, Program Studi S2 Ilmu Komputer, Fakultas Matematika Dan Ilmu Pengetahuan Alam, Universitas Gajah Mada.
- Anugroho,P., 2016. Klasifikasi Email Spam dengan Metode Naïve Bayes Classifier Menggunakan Java Programming. Politeknik Elektronik Negeri Surabaya.
- Auvil., Loretta., Searsmith., Duane. *IUsing Text Mining for Spam Filtering*, Automated Learned Group National Center for Supercomputing Applications University of Illinois.
- Berson, alex dkk. (2000). *Building data mining application for CRM*. Mc Graw–Hill. United states of America.
- Buntoro, A.G. 2016, Analisis Sentimen Hatespeech Pada Twitter Dengan Metode Naïve Bayes Classifier dan Support Vector Machine, *Jurnal Dinamika Informatika* volume 5, No 2.
- Ciptohartono,C.C. 2014, Algoritma Klasifikasi Naïve Bayes untuk Menilai Kelayakan. Tersedia di <http://eprints.dinus.ac.id/5439/1/13671.pdf>.
- CHFI Module XLVIII. 2011. *Investigating Sosial Networking Websites for Evidence*. (n.d.)
- Darujati,C., Gumelar, B.A. 2012, Pemanfaatan Teknik Supervised Untuk Klasifikasi Teks Bahasa Indonesia, *Jurnal Link* Vol. 16 No.1.
- Dinakar, K., Jones, B., Havasi, C., Lieberman, H., Picard, R. 2012, *Common Sense Reasoning for Detection, Prevention, and Mitigation of Cyberbullying*. ACM Transactions on Interactive Intelligent Systems. Vol. 2. No. 3.
- Firman, M dan Ngazis, A. N. 2012. Cyberbullying Ancaman Bagi Anak di Internet. (<http://fokus.news.viva.co.id/news/read/279625-cyberbullying-efek-samping-internet-bagi-anak>), (di unduh 4 oktober 2015).
- Fayyad, U., Piatetsky-Shapiro, G., Smyth, P. 1996. *From Data Mining to Knowledge Discovery in Databases*. *AI Magazine*

- Hanafi, A. 2009. Pengenalan Bahasa Suku Bangsa Indonesia Berbasis Teks Menggunakan Metode N-gram. IT TELKOM.
- Hamzah, Andi dan Marsita, Boedi D. 2012. Aspek-aspek Pidana dibidang Komputer. Jakarta : Sinar Grafika.
- Hafilizara, M., Adisantoso, J. 2014. Metode Smoothing dalam Naïve Bayes untuk Klasifikasi Email Spam. <http://repository.ipb.ac.id/handle/123456789/74528>.
- Haryati, 2014. Cyberbullying Sisi Lain Dampak Negatif Internet. Pusatian dan Pengembangan Aplikasi Informatika dan Informasi dan Komunikasi Publik. Jakarta Pusat.
- Hinduja, S., Patchin J. W. 2015, Statistic Cyberbullying <http://cyberbullying.org/statistics/>
- Hidayatullah, A.F., SN,A. 2014. Analisis Sentimen dan Klasifikasi Kategori Terhadap Tokoh Publik pada Twitter. Seminar Nasional Informatika UPN "Veteran" Yogyakarta.
- Hosseinmardi, H., Mattson, S.A., Rafiq, R.I., Han, R., Lv, Q., Mishra, A. 2015. *Detection of Cyberbullying Incidents on the Instagram Sosial Network*. arXiv: 1503.03909 v1 (cs.SI).
- Hilmawan, B.L. 2014. Aplikasi Mobile Untuk Analisis Sentimen Pada Google Play. Tesis. Fakultas Matematika dan Ilmu Pengetahuan Alam, Universitas Gajah Mada.
- JsonLint tersedia di <http://jsonlint.com/>
- Kansara, K. B., Shekocar, N. M. 2015. *A Framework for Cyberbullying Detection in Sosial Network*. International Journal of Current Engineering ang Technology. Vol. 5, No. 1. E-ISSN: 227-4106.
- Kayarkar,P,V., Ricchariya,P., Motwani, A. 2014. *Mining Frequent Sequences for Emails in Cyber Forensics Investigation*. International Journal of Computer Applications.
- Margono, H., Yi Xun., Raikundalia, G.K. 2014. *Mining Indonesia Cyber Bullying Patterns in Sosial*. Proceedings of the Thirty-Seventh Australasian Computer Science Conference (ACSC), Auckland, New Zealand.
- Manning, C. D., Raghavan, P., & Schütze, H., 2008, *Introduction to Information Retrieval*, Cambridge University Press, Cambridge.
- Machsun, R., 2016, Fenomena Cyberbullying Pada Remaja. Jurnal Ilmu Perpustakaan, Informasi, Kearsipan Khizanah Al-Hikmah, 35-44.
- Nuh, M. 2012. Digital Forensic : Panduan Praktis Investigasi Komputer. Salemba Infotek.
- Nuraini, I., Susanto, B., Proboyekti, U. 2011, Implementasi Naïve Bayes Classifier Pada Program Bantu Penentuan Buku Referensi Matakuliah, ti.ukdw.ac.id.

- Natalius, S. 2010. Metoda Naïve Bayes Classifier dan Penggunaannya pada Klasifikasi Dokumen. Program Studi Sistem dan Teknologi Informasi, Sekolah Teknik Elektro dan Informatika, ITB.
- Nalini, K., Sheela, D.L.J. 2014. *A survey on Datamining in Cyber Bullying*. International Journal on Recent and Innovation Trends in Computing and Communication. ISSN: 2321-8169 Volume: 2 Issue: 7 1865 – 1869.
- Nurjanah, S. 2014. Pengaruh Penggunaan Media Sosial Facebook Terhadap Perilaku Cyberbullying Pada Siswa Sman 12 Pekanbaru. Skripsi. Universitas Riau.
- Orebaugh, A., Allnut, J. 2009. *Classification of Instant Messaging Communications for Forensics Analysis*. The International Journal of Forensic Computer Science, IJoFCS 22-28.
- Kayarkar, V.P., Nirt, R.P., Motwani, A. 2014. *Mining Frequent Sequences for Email in Cyber Forensics Investigation*. International Journal of Computer Application (0975-8887), Volume 85-No 17.
- Rahayu, F. S. 2012. Cyberbullying Sebagai Dampak Negatif Penggunaan Teknologi Informasi. Jurnal Sistem Informasi, Vol 8, No 1, e-ISSN: 2502-6631
- Rodiyansyah, S.F., Winarko, E. 2012. *Klasifikasi Posting Twitter Kemacetan Lalu Lintas Kota Bandung Menggunakan Naïve Bayesian Classification*. IJCCS, Vol.6, No.1, January 2012, pp. 91~100.
- Satalina., Dina. 2014. Kecenderungan Perilaku Cyberbullying Ditinjau Dari Tipe Kepribadian Ekstrovert Dan Introvert. Ejournal UMM, Vol. 02, No. 02, ISSN: 2301-8267.
- Sanchesz, H., Kumar, S. 2011. *Twitter Bullying Detection*, Dept of Computer Science UC Santa Cruz.
- Saraswati, N.W.S. 2011. *Text Mining dengan Metode Naïve Bayes dan Support Vector Machines untuk Sentiment Analysis*. Tesis. Universitas Udayana Denpasar.
- Singhal, P., Bansal, A. 2013. *Improved Textual Cyberbullying Detection Using Data Mining*. International Journal of Information and Computation Technology. ISSN 0974-2239 Volume 3, Number 6, pp. 569-576.
- Susanto, S., Suryadi, D. 2010. Pengantar data mining, menggali pengetahuan dari bongkahan data. C.V. Andi Offset, Yogyakarta.
- Sucahyo, Y. G. 2013. Data Mining Menggali Informasi yang Terpendam. IlmuKomputer.com .(Online)<http://wsilfi.staff.gunadarma.ac.id/Downloads/files/4413/yudho-datamining.pdf> (diakses 10 Desember 2016).
- Santosa, B. 2007. Data Mining, Teknik Pemanfaatan Data Untuk Keperluan Bisnis, Teori dan Aplikasi. Graha Ilmu. Yogyakarta.

- Susanto, E. 2012. Data Mining Menggunakan WEKA (online)
- Sandi, F. R., 2012, Klasifikasi Posting Twitter Kemacetan Lalu Lintas Kota Bandung Menggunakan Naïve Bayesian Classification, Tesis, Program Studi S2 Ilmu Komputer, FMIPA, UGM.
- Satyawati, D.A.I., Purwani, S.P.M. 2014. Pengaturan Cyberbullying Dalam Undang-Undang Nomor 11 tahun 2008 Tentang Informasi dan Transaksi Elektronik. Bagian Hukum Pidana Fakultas Hukum Universitas Udayana
- Setty, S., Jadit, R., Shaikh, S., Mattikallis, C., Mudenagudi, V. 2014. *Classification of Facebook News Feeds and Sentiment Analysis*. Hubli, India.
- Saputra, N. 2015. Analisis Sentiment Berbasis Lexicon dan Emoticon. Program Pascasarjana Fakultas Teknik UGM Yogyakarta.
- Sastrawi Master Tersedia di <https://github.com/sastrawi/sastrawi>
- Taibah, S. A. 2013. Urgensi Kriminalisasi Cyberbullying di Indonesia, Skripsi, Fakultas Hukum, Universitas Indonesia.
- Tala Stopwords Tersedia di <http://hikaruyuuki.lecture.ub.ac.id/kamus-kata-dasar-dan-stopword-list-bahasa-indonesia/>.
- Yurnalita. 2016. Cyberbullying pada Jejaring Sosial Twitter (Analisis Semiotika Trending Topic), Fakultas Ilmu Sosial dan Ilmu Politik, Universitas Syiah Kuala Darussalam, Banda Aceh.
- Yin, D., Xue, Z., Hong, L. 2009. *Detection of Harassment on Web 2.0*. Departement of Computer Science and Engineering Lehigh University.
- www.kominfo.go.id 2013

Lampiran

DAFTAR KAMUS BESAR BAHASA ALAY (KBBA)

Kata Alay	Kata Baku
anjir	Anjing
lu	Kamu
serah	terserah
setres	Stres
lo	Kamu
gk	Tidak
aja	Saja
bodo	bodoh
stress	stres
ntar	sebentar
nying	anjing
bener	Benar
ga	Tidak
udah	sudah
gue	Saya
gw	Saya
njink	anjing
ama	Sama
gomballl	gombal
tp	Tapi
udah	sudah
kyanya	seperti
sono	Sana
yg	Yang
org	orang
hdp	hidup
lyk	Layak
elu	kamu
td	Tadi
blg	bilang
kagak	Tidak
gmpang	gampang
lgi	Lagi
gag	tidak
da	Ada
tv	Tivi
adlah	adalah
lo	kamu
amp	sampai

skrg	sekarang
ngmgn	ngomongin
pny	punya
sdh	sudah
tak	tidak
pasca	sesudah
aje	Saja
diem	diam
udh	sudah
dimn	dimana
idup	hidup
ntu	Itu
jgn	jangan
lg	Lagi
udh	sudah
kek	kayak
gt	begitu
ngmg	ngomong
de	adek
loe	kamu
cewe	cewek
ga	tidak
aja	Saja
klo	kalau
aje	Saja
nie	Ini
jd	Jadi
gk	tidak
ngak	tidak
td	Tadi
bcr	bicara
jwb	jawab
gak	tidak
terjd	terjadi
tdk	tidak
doank	Saja
sono	sana
cm	Cuma
muke	muka
dah	
dh	sudah
jing	anjing
pake	pakai
gw	saya
jg	Juga

loe	kamu
ngga	tidak
bgmn	bagaimana
gua	saya
hurup	huruf
dpt	dapat
jodo	jodoh
lbh	lebih
bca	baca
liat	lihat
malem	malam
buar	Biar
njs	najis
chabai	cabe
kacian	kasihan
kg	tidak
kaga	tidak
kbeh	Saja
inget	Ingat
pe	sampai
byr	Bayar
sm	Sama
trima	Terima
pntas	Pantas
klwr	Keluar
ente	Kamu
ame	Sama
kebo	Kerbau
spt	seperti
ava	Apa
dr	Dari
emg	memang
drpd	daripada
njir	Anjing
anjay	Anjing
kgk	Tidak
tmn	Teman
syg	sayang
ptng	potong
aj	Saja
jng	Jangan
skg	sekarang
kyk	seperti
uda	Sudah
nyet	monyet

udeh	0	Sudah
mulu		Selalu
laknak		Laknat
blm		Belum
anjeeenngg		anjing
bangsad		bangsat
ngriyok		ngeroyok
ape		Apa
kw		palsu
bangke		bangkai
full		penuh
smua		semua
ae		Saja
trs		terus
knp		kenapa
kpan		kapan
kgalah		tidaklah
brani		berani
bukn		bukan
nyebut		menyebut
cino		Cina
jowo		jawa
hongkon		hongkong
mulutx		mulutnya
msih		masih
pdhl		padahal
plg		paling
trus		terus
km		kamu
jga		Juga
ngpain		ngapain
mrasa		merasa
skolah		sekolah
msh		masih
kalo		kalau
sampe		sampai
byk		banyak
ngebales		membalas
gmn		gimana
ngtwain		menertawakan
mw		mau
ngsh		ngasih
nga		tidak
wong		orang

jdi	Jadi
ud	sudah
nyang	yang
mang	memang
dn	Dan
keknya	sepertinya
hrs	harus
sobat	teman
kl	kalau
dianggep	dianggap
nape	kenapa
pd	pada
dng	dengan
ngerti	mengerti
emang	memang
tar	sementara
brisik	berisik
ndut	gendut
ngeselin	menjengkelkan
nyebelin	menyebalkan
idungnya	hidungnya
fine	baiklah
send	kirim
message	pesan
cwe	cewek
srp	sarap
mbe	kambing
jirr	anjing
anjrit	anjing
like	suka
tamvan	tampan
pcr	pacar
wkt	ketika
tlnga	telinga
w	saya
mkn	makan
ngulurkan	mengulurkan
tanganny	tangannya
kayanya	sepertinya
beneran	benar
mgkn	mungkin
kebaca	terbaca
bezakan	membedakan
maen	main
ngaca	berkaca

ngilangg	menghilang
pasca	setelah
seneng	senang
orok	bayi
kemane	kemana
ketauan	ketahuan
nggak	tidak
vinterlah	pintarlah
ngga	tidak
anjir	anjing
ngerusak	merusak
dr	Dari
jgn	jangan
cem	macam
ude	sudah
idup	hidup
gt	begitu
klakuan	kelakuan
elo	kamu
gt	begitu
tar	sebentar
guoblog	goblok
ae	Saja
pulak	pula
njingg	anjing
pny	punya
bgt	banget
lsg	langsung
sblm	sebelum
anj	anjing
bangkek	bangkai
trllu	terlalu
deket	dekat
anjyr	anjing
ngaku	mengaku
nyett	monyet
bapake	bapaknya
ush	usah
ati	Hati
taun	tahun
asu	anjing
asuu	anjing
trmsk	termasuk
anjiinq	anjing
ngasi	memberi

pict	gambar
it	Itu
vangke	bangkai
dmna	dimana
jlek	jelek
bnget	banget
dri	Dari
ndak	tidak
dgn	dengan
klu	kalau
tu	Itu
sapa	siapa
bsk	besok
diem	diam
dapet	dapat
karna	karena
ampe	sampai
bsa	Bisa
bg	bagi
denger	dengar
jkt	jakarta
sm	sama
kejer	kejar
ngapa	kenapa
tivu	Tipu
dtg	datang
deketin	dekat
krn	karena
nanggep	tanggap
temen	teman
ngandelin	mengandalkan
nyerah	menyerah
saiton	setan
bhs	bahasa
sempet	sempat
begoo	bego
pgn	ingin
dpn	depan
anggep	anggap
with	dengan
ga	tidak
brensex	brensek
gw	saya
gak	tidak
kagak	tidak

ngatas	mengatas
betpa	betapa
taikkk	Tai
bngt	banget
emng	memang
kta	Kata
tkt	takut
gtu	Gitu
sma	sama
macem	seperti
prcuma	percuma
nunggy	menunggu
gilak	Gila
mo	mau
kturunan	keturunan
mslh	masalah
hr	Hari
pantes	pantas
segitunye	segitunya
enggak	tidak
ngatain	mengatakan
jalanin	jalankan
ngrasa	merasa
lonely	kesepian
lag	Lagi
nyesel	menyesal
pinter	pintar
bkn	bukan
jt	Juta
ky	kayak
istifar	istigfar
ngira	Kira
ttg	tentang
brpa	berapa
mn	mana
ad	Ada
pke	pakai
ista	istana
pda	pada
aq	saya
spti	seperti
meperlakukan	memperlakukan
boong	bohong
ngajakin	mengajak
br	baru

keliatan	kelihatan
dijanjiin	dijanjakan
activitynya	aktifitasnya
at	atau
ilang	hilang
ngedem	ngadem
mz	Mas
pnya	punya
rbu	Ribu
mantab	mantap
dsar	dasar
goblokkkk	goblok
laporin	laporkan
asliny	aslinya
curutt	curut
ngerasain	merasakan
goblo	goblok
guah	saya
dipamerin	dipamerkan
diliatin	diperlihatkan
dapetin	mendapatkan
gede	besar
pantesin	pantasan
cuman	cuma
gublok	goblok
ngehargain	menghargai
nurutin	menurutkan
babik	babi
hny	hanya
maluin	malu
yakiiin	yakin
bedain	membedakan
nyantai	santai
kate	Kata
gelitikin	gelitik
ngeladenin	meladeni
jir	anjing
males	malas
ane	saya
mcam	seperti
dIm	dalam
kemaren	kemarin
naek	Naik
tetep	tetap
mnt	minta

jelekk	jelek
ngadu	mengadu
diemin	diamkan
smpe	sampai
kmu	kamu
trusss	terus
mikiran	memikirkan
ngekang	mengekang
kesian	kasihan
ditanyain	ditanyakan
sanggub	sanggup
tanggap	tanggap
du	dulu
rante	rantai
taik	Tai
crut	curut
blh	boleh
perhatiin	perhatikan
kayak	seperti
kecium	tercium
keparatt	keparat
mksd	maksud
jatohnya	jatuhnya
yang	yang
jebolin	menjebol
tauu	Tau
kluar	keluar
syng	sayang
lgian	lagian
jdiin	jadikan
bs	Bisa
ptar	putar
wktu	waktu
sgala	segala
skap	sikap
prnh	pernah
nyakitin	menyakiti
anjer	anjing
nyimpen	menyimpan
ig	instagram
cmn	cuma
temuin	temukan
ngikutin	mengikuti
bales	balas
engga	tidak

ceramahin	ceramahi
am	sama
ngajak	mengajak
tangkap	tangkap
aer	Air
ngamuk	mengamuk
skr	sekarang
gwa	saya
eh	
sok	
eh	
ah	
ta	
kok	
woe	
deh	
hh	
heh	
nya	
ya	
ak	
ms	masa
nyari	Cari
mah	
kan	
balasan	balasan
soalnya	karena
balesin	membalas
bhaha	
wkwk	
sih	
nandingi	menandingi
bu	Ibu
gimana	bagaimana
anjeng	anjing
WKWKWKWKW	
an	
usrnm	username
puter	putar
dijadiin	dijadikan
pikirin	pikirkan
si	
die	Dia
tuh	
wkw	

dibanggain	dibanggakan
awalnye	awalnya
taunye	taunya
pea	peak
godain	menggoda
orngnye	orangnya
diajarin	diajarkan
ptsin	putusan
ye	
jawabannye	jawabannya
yuk	
nye	
emak	Ibu
non	bukan
faham	paham
ngurus	mengurus
mainin	memainkan
ny	
lah	
nih	
dpnya	Dp
bner	Benar
cah	
dateng	Dating
huft	
smw	Semua
kaya	Kayak
abis	Habis
ha	
gih	
in	
dih	
ih	
cinak	Cina
ngatasnamakan	mengatasnamakan
haa	
hahaha	
bitch	Jalang
fak	
woy	
doang	Saja
laper	Lapar
yeeh	
ngasih	beri
ala	

ngeledekin	meledek
ank	anak
mh	
dskolain	disekolahkan
brsykr	bersyukur
dkit	sedikit
duit	uang
ngmng	bicara
smbgrn	sembarangan
diketawaiin	ketawa
yeay	
sip	
oh	
cie	
yaa	
bangett	banget
nyengir	tertawa
pengen	ingin
suapin	suap
ehehe	
wkkk	
ngeliat	melihat
et	
loh	
aw	
ditawarin	ditawarkan
LOL	
bah	
se	
HAHAHAHAHAHA	
ngomongin	omong
ahox	ahok
hahah	
nutup	tutup
nhe	
slma	selama
slalu	selalu
ni	
tu	
pala	kepala
layanin	melayani
na	
amittt	amit
siain	siakan
sialen	sial

nnti	nanti
dya	Dia
iyuh	
anjink	anjing
cowo	cowok
kesel	jengkel
wa	
ditinggalin	ditinggalkan
bapaknye	bapak
matiin	mati
elah	
bacanya	baca
pah	
nemu	temu
nun	
noh	
aelah	
dikatain	kata
masi	masih
hii	
alah	
huh	
Helloooooo	
yah	
woi	
Haha	
bruh	
bro	
ka	
shitttt	
tae	Tai
buset	
liatin	lihat
huhuhu	
dikatainnya	kata
orgnya	orang
gilaan	gila
bangett	banget
kedengeran	dengar
ngelawak	lawak
fto	foto
drmn	darimana
kerennn	keren
galakk	galak
idung	hidung

nti	nanti
apaan	apa
tida	tidak
lagy	lagi
dong	
ngadem	adem
idih	
balikin	balik
ngambil	ambil
wkaka	
heleh	
balesnya	balas
ngapah	apa
kzl	jengkel
puaaaas	puas
laah	
jiaah	
iyalah	lya
tuk	untuk
wkwwkwj	
to	
mbl	mobil
nyumpahi	sumpah
au	
asyuuuu	Asu
wkwkwk	
kuy	kamu
cma	cuma
nyepam	spam
disalahin	salah
wooi	
opp	
utk	untuk
pdhal	padahal
jaat	jahat
ehe	
didatenging	datang
hih	
tong	
chabe	cabe
chabean	cabe
bini	istri
ngeganggu	ganggu
kau	kamu
huehue	

boy	teman
temenan	
eyy	
Hahahaha	
It	
goblog	goblok
kesenangan	senang
uke	
so	
ngevote	vote
bosen	bosan
oiya	
nongol	muncul
tempelin	tempel
yha	
ahhh	
lebayyyyyyy	lebay
banyakin	banyak
tkng	tukang
Ahahahayyy	
spah	sumpah
aduh	
bantuin	bantu
wkwkwkw	
hehe	
capa	siapa
lahhh	
ahh	
wkwkw	
twit	tweet
ngurusin	urus
tmpt	tempat
nrima	terima
sdg	sedang
tereakin	teriak
halah	
blagu	belagu
lampiasin	lampias
huhu	
hi	
hu	
ho	
wkkw	
he	
yaoi	

jkw	jokowi
dtng	datang
interest	tertarik
ellooo	kamu
kaleee	kali
gemesin	gemes
yak	
yaaa	
keh	
HAAA	
ngeshare	share
wkakakak	
lyaaay	
heheheh	
wkwkwkwk	
wkek	
huf	
pade	pada
lepasin	lepas
geblekk	geblek
ksempatan	kesempatan
dibantuin	bantu
banggsat	bangsat
yo	
nu	
luuuu	kamu
nah	
polotikus	politikus
g	tidak
woooi	
bodooh	bodoh
wkwkkkkk	
jawabanye	jawab
only	hanya
taekk	Tai
gilaa	gila
hsilnya	hasil
gin	gini
pun	
wkwkkwk	
nang	nangis
sich	
katain	kata
wow	
cuihhh	

hah	
masukin	masuk
cilakak	celaka
kbnykn	banyak
hahaa	
wakakakak	
erti	ngerti
makanya	maka
bajing	bajingan
tuch	
wkk	
dibebasin	bebas
hei	
datengin	datang
kepercayaan	percaya
kepercayaan	percaya
tok	
naikin	naik
lot	lemot
gwe	saya
hahahahah	
bhahahah	
aduin	adu
emaknya	ibu
kwwkkwwkk	
gegara	gara
ngomong	omong
cwek	cewek
ehh	
eeh	
hahahhahaa	
rrrr	
hahhhaaha	
bruakakakakak	
iblissss	iblis
luh	kamu
kampung	kampungan
cakeppp	cakep
ente	kamu
jhaaa	
maap	maaf
ihh	
hyung	
paan	apa
asa	rasa

ko	
oppa	
cepat	cepat
tajem	tajam
ckckck	
cuk	
napa	kenapa
now	sekarang
doain	doa
niatin	niat
gueh	saya

