

ANALISIS SENTIMEN PENGGUNA TWITTER TERHADAP PSBB DI JAKARTA MENGUNAKAN METODE *NAÏVE BAYES CLASSIFIER*

Azini Fauzia Putri¹, Iin Ernawati², Anita Muliawati³

Fakultas Ilmu Komputer

Universitas Pembangunan Nasional Veteran Jakarta

Email : azinifau@gmail.com¹, iinerti@gmail.com², anitamuliawati@upnvj.ac.id³
Jl. Rs. Fatmawati, Pondok Labu, Jakarta Selatan, DKI Jakarta, 12450, Indonesia

Abstrak

PSBB atau Pembatasan Sosial Berskala Besar adalah salah satu tindakan pencegahan dalam penyebaran pandemi COVID-19 yang dilakukan oleh pemerintah. Penerapan PSBB berlangsung hampir di seluruh wilayah Indonesia, salah satunya di Provinsi DKI Jakarta. Setelah PSBB berakhir, kegiatan pembatasan dinamakan PSBB Transisi, di mana adanya kelonggaran dalam beraktivitas dan sejumlah fasilitas umum dibuka dengan memperhatikan protokol kesehatan yang berlaku. Kemudian setelah pemberlakuan PPKM (Pemberlakuan Pembatasan Kegiatan Masyarakat) Jawa dan Bali oleh pemerintah pusat, diberlakukan PSBB ketat di Jakarta yang terus mengalami perubahan peraturan sesuai situasi dan kondisi masyarakat. Pada penelitian ini, dilakukan analisis sentimen masyarakat mengenai PSBB di Jakarta melalui media sosial Twitter dengan metode *Naïve Bayes Classifier*. Data penelitian ini yaitu *tweet* yang didapat dari Twitter menggunakan *keyword* “PSBB DKI Jakarta” yang diambil pada tanggal 1 Februari-31 Maret 2021. Hasil akhir penelitian dengan *oversampling* ini adalah nilai *accuracy* senilai 0.8, nilai *recall* senilai 0.9318, dan nilai *specificity* senilai 0.27.

Kata Kunci: Analisis Sentimen, *tweet*, *Naïve Bayes Classifier*.

1 PENDAHULUAN

Wabah pandemi COVID-19 mulai terjadi pada akhir tahun 2019 di Wuhan, Provinsi Hubei, China. Diduga virus *Corona* bersumber dari hewan kelelawar dan menyebabkan penyakit yang saat ini dikenal sebagai penyakit COVID-19. Virus *Corona* terdeteksi di Indonesia pada awal Maret tahun 2020. Berbagai langkah preventif dan kuratif telah dilakukan pemerintah Indonesia, seperti diselenggarakannya *rapid test*, pelaksanaan PSBB (Pembatasan Sosial Berskala Besar) di berbagai wilayah, dan pemberian bantuan keluarga menengah ke bawah.

Berbagai sentimen muncul mengenai langkah-langkah yang dilakukan pemerintah Indonesia, salah satunya mengenai penerapan PSBB maupun PSBB transisi di Jakarta yang dianggap sudah tidak ketat. Terlihat dari masyarakat yang sering bepergian keluar rumah, padahal pemerintah menganjurkan masyarakat untuk tetap di rumah demi mengurangi penyebaran penyakit COVID-19. Kemudian adanya pemberlakuan PPKM (Pemberlakuan Pembatasan Kegiatan Masyarakat) Jawa dan Bali yang juga diterapkan oleh pemerintah Jakarta dalam bentuk PSBB ketat maupun PSBB yang mengalami perubahan peraturan, seperti pada

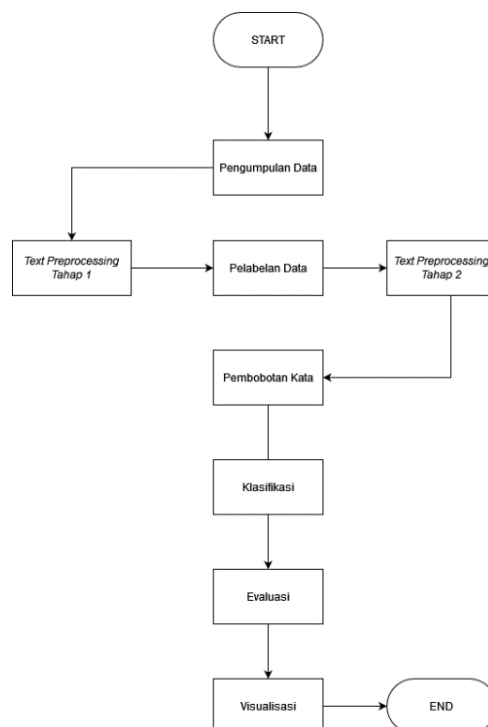
pelaksanaan PSBB ketat pada tanggal 11-25 Januari 2021 & 26 Januari-8 Februari 2021, jam operasional transportasi umum diberlakukan hingga 20.00 WIB. Pemerintah Jakarta menambah jam operasional transportasi umum hingga jam 21.00 WIB pada PSBB terhitung 9 Februari 2021.

Dengan banyaknya perubahan peraturan maupun istilah, muncul berbagai sentimen atau opini dari masyarakat. Oleh karena itu, untuk mengetahui bagaimana opini masyarakat mengenai PSBB di Jakarta, dilakukan penelitian yang memerlukan model klasifikasi yang dapat melakukan analisa sentimen yaitu sentimen positif dan sentimen negatif pengguna media sosial khususnya Twitter yang berisi berbagai pendapat atau komentar terhadap pelaksanaan PSBB di Jakarta. PSBB yang dimaksud yaitu PSBB yang merupakan bentuk penerapan dari PPKM oleh Pemerintah Jakarta yang di antaranya berlangsung dari tanggal 26 Januari-8 Februari 2021, 9-22 Februari 2021, 23 Februari-8 Maret 2021, 9-22 Maret 2021, dan 23 Maret-5 April 2021, dan penulis mengambil data dengan periode 1 Februari-31 Maret 2021.

Analisa dilakukan dengan klasifikasi *tweet* mengenai sentimen masyarakat tentang pelaksanaan PSBB di Jakarta dengan menggunakan metode *Naïve Bayes Classifier*. Data yang digunakan yaitu *tweet* pada tanggal 1 Februari-31 Maret 2021. *Keyword* yang digunakan dalam proses *crawling* yaitu “PSBB DKI Jakarta”.

2 METODOLOGI PENELITIAN

Tahapan penelitian dilakukan berdasarkan *flowchart* yang diilustrasikan pada Gambar 1 sebagai berikut.



Gambar 1: *Flowchart* Tahapan Penelitian

2.1 Pengumpulan Data

Pada tahap pengumpulan data dilakukan menggunakan Twitter API pada Twitter dengan *keyword* “PSBB DKI Jakarta”, di mana proses ini menggunakan *software* R studio dan bahasa pemrograman R. Alasan tidak digunakannya *keyword hashtag* yaitu agar memperoleh data *tweet* terkait PSBB DKI Jakarta meskipun *tweet* tersebut tidak menggunakan *hashtag*. Data

yang digunakan yaitu *tweet* pada tanggal 1 Februari-31 Maret 2021, di mana telah berlangsung pembatasan sosial yaitu tanggal 26 Januari-8 Februari 2021, 9-22 Februari 2021, 23 Februari-8 Maret 2021, 9-22 Maret 2021, dan 23 Maret-5 April 2021.

2.2 Text Preprocessing Tahap 1

Text preprocessing pertama dilakukan dengan menghilangkan *field-field* yang tidak dibutuhkan dari data, menghilangkan data dengan tanggal di luar batasan masalah, menghilangkan *tweet* berduplikat, dan menghilangkan *tweet* yang inkonsisten. Tujuan dilakukannya tahap ini yaitu agar memudahkan *annotator* dalam melakukan pelabelan data sehingga pengerjaan lebih efektif daripada melakukan pelabelan dengan data mentah.

2.3 Pelabelan Data

Labeling atau pelabelan dilakukan dengan memberikan label sentimen positif dan negatif. Proses *labeling* dilakukan secara manual dengan bantuan 3 *annotator* dengan cara membaca setiap *tweet* yang akan diberi label. *Tweet* yang berisi komentar atau pernyataan positif seperti upaya dan tindakan pemerintah dalam menjalankan PSBB, penegakan peraturan, dan masyarakat yang memberikan dukungan dan saran, akan dilambangkan dengan angka 1. Sedangkan *tweet* yang berisi komentar atau pernyataan negatif seperti menjelekkan pemerintah, penutupan perusahaan, dan kritik tidak membangun, akan dilambangkan dengan angka 0.

Dengan menggunakan bantuan manusia (*annotator*) dalam proses *labelling*, tentu akan ada kemungkinan kesalahan, sehingga digunakan metode penilaian bernama *kappa value*. Dengan melakukan metode *kappa value*, maka dapat diketahui apakah dokumen yang dijadikan *dataset* relevan.

2.4 Text Preprocessing Tahap 2

Text preprocessing tahap 2 dilakukan dengan tujuan untuk membersihkan *tweet* sehingga mempermudah dalam memproses data. Pada proses *text preprocessing* tahap 2 menggunakan *software* Google Colab dan bahasa pemrograman Python. *Text Preprocessing* tahap 2 terdiri dari *cleaning*, *case folding*, *tokenization*, normalisasi bahasa, *filtering*, dan *stemming*.

1. *Cleaning*

Pembersihan data dilakukan dengan menghapus karakter-karakter seperti URL, *hashtag*, *username*, dan lainnya.

2. *Case Folding*

Pada proses *case folding*, huruf kapital yang terdapat pada teks *tweet* diubah menjadi *lower case* (huruf tidak kapital). Hal ini dilakukan untuk menghindari adanya perbedaan makna. Contoh pada kata “DKI Jakarta” dengan “dki jakarta” memiliki makna yang sama, yaitu nama provinsi. Dalam pemrograman Python, kata tersebut tidak dianggap sama, sebab Python memiliki *case sensitive* sehingga apabila terdapat kata yang sama namun memiliki perbedaan pada huruf kapital maka dianggap berbeda.

3. *Tokenization*

Setelah melakukan *case folding*, maka proses selanjutnya yaitu melakukan *tokenization* atau tokenisasi. Proses ini dilakukan dengan memecah sekumpulan karakter menjadi satuan kata.

4. Normalisasi Bahasa

Proses normalisasi bahasa dilakukan dengan memanfaatkan daftar kata-kata gaul di mana akan dilakukan pencocokkan antara kata-kata gaul pada data yang telah dilakukan tokenisasi sebelumnya dengan daftar kata-kata gaul (kamus *slangword*) kemudian diubah menjadi kata-kata baku.

5. *Filtering*

Proses *filtering* atau penyaringan merupakan proses penghapusan kata-kata tidak penting pada data. Proses *filtering* dilakukan sebanyak dua kali, yaitu dengan menggunakan *library* NLTK dan daftar *stopword* yang diperoleh dari Rama Prakoso pada situs GitHub yang telah dimodifikasi. Tujuan dilakukan proses *filtering* sebanyak 2 kali yaitu agar proses

penyaringan lebih maksimal sehingga beberapa kata tidak penting yang tidak terhapus dengan *stopword* dari NLTK dapat dihapus dengan daftar *stopword* tambahan yang telah dimodifikasi.

6. Stemming

Tujuan proses *stemming* adalah untuk memperoleh kata dasar dari kata berimbuhan pada awal kata (prefiks) maupun imbuhan pada akhir kata (sufiks). Proses *stemming* dilakukan dengan menggunakan *library* Sastrawi dengan Algoritma StemmerFactory.

2.5 Pembobotan Kata

Pada penelitian ini, TF-IDF (*Term Frequency–Inverse Document Frequency*) digunakan sebagai metode pembobotan kata. Nilai TF diperoleh berdasarkan jumlah kemunculan setiap kata pada tiap dokumen, dan nilai IDF diperoleh berdasarkan jumlah kemunculan kata dalam keseluruhan dokumen.

Untuk menghitung nilai bobot pada suatu kata, dapat menggunakan persamaan (1) yang dideskripsikan sebagai berikut.

$$w_{t,d} = tf_{t,d} \times \log \left(\frac{N}{df_t} \right) \dots\dots\dots (1)$$

Keterangan :

- $w_{t,d}$: Bobot kata (t) pada dokumen (d)
- $tf_{t,d}$: Frekuensi kata (t) pada dokumen (d)
- N : Jumlah dokumen teks.
- df_t : Jumlah dokumen yang mengandung *term* (t).

2.6 Klasifikasi

Metode klasifikasi yang digunakan yaitu metode *Multinomial Naïve Bayes*. Terdapat 2 proses klasifikasi, yaitu proses pelatihan (*training*) dan proses pengujian (*testing*). Tahap *training* dilakukan untuk mendapatkan nilai probabilitas dari setiap kata dari *train data* (data latih). Tahap pengujian dilakukan untuk mendapatkan nilai probabilitas dari *test data* (data uji). Perhitungan *Multinomial Naïve Bayes* dideskripsikan pada persamaan (2).

$$C_{MAP} = \underset{c \in V}{\text{arg max}} P(p) \prod_{t=i}^{|V|} P(W_i|p) \dots\dots\dots (2)$$

Keterangan :

- $P(W_i|p)$: Peluang kemunculan W_i pada kelas p.
- $P(p)$: Peluang kemunculan dokumen yang berada pada kelas p.

2.7 Evaluasi

Evaluasi dilakukan untuk menguji apakah tujuan penelitian sudah tercapai dari penelitian yang dilakukan. Metode evaluasi yang digunakan yaitu *confusion matrix*. *Confusion matrix* terdiri dari informasi perbandingan hasil klasifikasi yang telah dilakukan sistem dengan hasil klasifikasi seharusnya. Tabel 1 mendeskripsikan pemodelan *confusion matrix*.

Tabel 1: Confusion Matrix

		<i>Actual</i>	
		Positive	Negative
<i>Prediction</i>	Positive	True Positive	False Positive
	Negative	False Negative	True Negative

Pada *confusion matrix*, terdapat 4 istilah sebagai hasil dari proses klasifikasi. Penjelasan mengenai istilah-istilahnya dideskripsikan sebagai berikut.

1. *True Positive* (TP) : Data positif yang diprediksi benar.
2. *True Negative* (TN) : Data negatif yang diprediksi benar.
3. *False Positive* (FP) : Data negatif yang diprediksi sebagai data positif.
4. *False Negative* (FN) : Data positif yang diprediksi sebagai data negatif.

Untuk mengukur kinerja algoritma digunakan nilai akurasi, *recall*, dan *specificity* dari model klasifikasi yang telah dibuat. Penjelasan mengenai *performance metrics* dideskripsikan sebagai berikut.

1. *Accuracy*

Accuracy atau akurasi menandakan keakuratan model melakukan klasifikasi dengan benar. Akurasi merupakan perbandingan prediksi benar (positif dan negatif) terhadap keseluruhan data. Nilai *accuracy* dideskripsikan dengan persamaan (3).

$$accuracy = \frac{TP+TN}{TP+TN+FP+FN} \dots \dots \dots (3)$$

2. *Recall*

Recall menandakan keberhasilan model dalam mendapatkan kembali informasi. *Recall* merupakan perbandingan prediksi benar positif dengan keseluruhan data yang benar positif. Nilai *recall* dideskripsikan dengan persamaan (4).

$$recall = \frac{TP}{TP+FN} \dots \dots \dots (4)$$

3. *Specificity*

Specificity merupakan kebenaran dalam prediksi negatif terhadap keseluruhan data negatif. Nilai *specificity* dideskripsikan dengan persamaan (5).

$$specificity = \frac{TN}{TN+FP} \dots \dots \dots (5)$$

3 HASIL DAN PEMBAHASAN

Untuk mengumpulkan data, penulis melakukan *crawling* sebanyak 10 kali dengan tiap data yang diambil memiliki rentang waktu yaitu 30 Januari-6 Februari 2021, 3-10 Februari 2021, 8-15 Februari 2021, 14-20 Februari 2021, 20-28 Februari 2021, 27 Februari-7 Maret 2021, 7-14 Maret 2021, 13-20 Maret 2021, 21-28 Maret 2021, dan 25 Maret-1 April 2021. Tiap hasil *crawling* kemudian digabungkan menjadi 1 *file* sehingga mendapatkan jumlah data (*tweet*) sebanyak 543 data dan disimpan dalam format csv. Hasil pengumpulan data diilustrasikan pada Gambar 2 sebagai berikut.

user_id	status_id	created_at	screen_name	text	source	display_text	reply_to_text_id
68930552	1.3581E+18	2/6/2021 16:10	mediaindi	Pemprov DKI masih	TweetDeck	231	NA
1363462284	1.3581E+18	2/6/2021 13:57	SatpolPP	Data Pengawasan dan	Twitter for Android	280	NA
1363462284	1.357E+18	2/3/2021 14:15	SatpolPP	Data Pengawasan dan	Twitter for Android	275	NA
1363462284	1.3579E+18	2/6/2021 3:15	SatpolPP	Jum'at (5/2) , Beberapa	Twitter for Android	277	NA
1363462284	1.356E+18	1/31/2021 22:20	SatpolPP	Data Pengawasan dan	Twitter for Android	280	NA

Gambar 2: Hasil Pengumpulan Data

Setelah didapatkan data, dilakukan *text preprocessing* tahap 1. Tahapan yang dilakukan pada proses ini adalah sebagai berikut.

1. Menghilangkan *field-field* yang tidak dibutuhkan.
Pada proses *crawling*, diperoleh *field* yang tidak terpakai seperti *user_id*, *status_id*, *source*, *display*, dan *reply*. Maka *field* yang tidak terpakai akan dihilangkan sehingga tersisa *field* tanggal *tweet* (*created_at*), nama akun (*screen_name*), dan isi *tweet* (*text*).
2. Menghilangkan *tweet* yang bertanggal selain tanggal 1 Februari-31 Maret 2021.
Pada penelitian ini digunakan data *tweet* yang dimulai dari tanggal 1 Februari-31 Maret 2021, sehingga *tweet* yang bertanggal selain 1 Februari-31 Maret 2021 akan dihapus dan hasilnya diurutkan mulai dari tanggal paling awal. Setelah dihilangkan, tersisa 537 data saat ini.
3. Menghilangkan *tweet* yang yang berduplikat
Dikarenakan proses *crawling* dilakukan beberapa kali dalam tiap minggunya, tentu menjadi salah satu faktor yang menyebabkan *tweet* berduplikat.
4. Menghilangkan *tweet* yang termasuk data inkonsisten
Data inkonsisten yang dimaksud yaitu *tweet* yang memiliki struktur yang hampir sama, namun memiliki perbedaan isi yang tidak terlalu berpengaruh pada isi teks. Pada Gambar 3, selain adanya perbedaan *link*, juga terdapat perbedaan isi *tweet* yaitu pada tanggal periode, di mana jika tanggal tersebut dihapus, maka tidak terlalu berpengaruh dengan isi teks. Penghapusan *tweet* dilakukan pada *tweet* yang waktunya lebih baru dibanding waktu *tweet* sebelumnya.

created_at	screen_name	text
2/2/2021 12:11	SatpolPP_DKI	Data Pengawasan dan Penindakan Pelanggaran PSBB oleh Satpol PP DKI Jakarta periode 11 Januari s/d 02 Februari 2021 (berdasarkan Pergub No.3 tahun 2021). - Salam sehat dan selalu patuhi Protokol Kesehatan. @aniesbaswedan @BangAriza @ISCLounge @DKIJakarta #JAKI #dkijakarta AG https://t.co/A9kNmBLIX4
3/5/2021 0:15	SatpolPP_DKI	Data Pengawasan dan Penindakan Pelanggaran PSBB oleh Satpol PP DKI Jakarta periode 11 Januari s/d 04 Maret 2021 (Berdasarkan Pergub No.3 tahun 2021) - Salam sehat dan selalu patuhi Protokol Kesehatan @aniesbaswedan @ArizaPatria @ISCLab @DKIJakarta #JAKI #DKIJakarta AG https://t.co/ESN932bWwY2

Gambar 3: Contoh *Tweet* Data Inkonsisten

Setelah dilakukan penghapusan *tweet* yang berduplikat dan data inkonsisten, data yang tersisa sebanyak 275 data dan siap untuk dilakukan tahap selanjutnya yaitu pelabelan data atau *labelling* data. Contoh hasil pelabelan data dideskripsikan pada Tabel 2.

Tabel 2: Contoh Hasil Pelabelan *Tweet*

Text	Pihak 1	Pihak 2	Pihak 3	Hasil Akhir
Pemerintah Provinsi DKI Jakarta mengeluarkan sejumlah aturan tentang denda bagi pelanggar aturan protokol kesehatan pada masa pandemi Covid-19. #TempoGrafis https://t.co/e9PqpUfOwo https://t.co/SgC2zvwWIC	1	1	1	1
Pandemi Covid-19 Kian Darurat, 911 Perusahaan di DKI Jakarta Ditutup Sementara Selama Masa PSBB https://t.co/q3aBjKq0Sn	1	0	0	0
@Sopianhaidi @imadya @PDI_Perjuangan @dprddkijakarta Saat Pemprov DKI Tetapkan PSBB Sampai Kerahkan Preman Kejar Warga Yg Langgar Prokes Tapi Gubernur & Wagubnya Malah Hadiri Kerumunan Terjadilah Pembangkangan Warga Dan Meledaklah Warga Yg Terpapar Covid19	0	0	0	0

Hasilnya *tweet* dengan label positif sebanyak 222 data dan *tweet* dengan label negatif sebanyak 53 data. Dikarenakan pelabelan menggunakan *annotator* lebih dari 2 orang, maka digunakan persamaan *fleiss kappa*. Tabel 3 mendeskripsikan sampel yang digunakan dalam perhitungan nilai *kappa*.

Tabel 3: Sampel Data Untuk *Kappa Value*

Text	Pihak 1	Pihak 2	Pihak 3	F ₁	F ₀
Pemerintah Provinsi DKI Jakarta mengeluarkan sejumlah aturan tentang denda bagi pelanggar aturan protokol kesehatan pada masa pandemi Covid-19. #TempoGrafis https://t.co/e9PgpUfOwo https://t.co/SgC2zvwWIC	1	1	1	3	0
Pandemi Covid-19 Kian Darurat, 911 Perusahaan di DKI Jakarta Ditutup Sementara Selama Masa PSBB https://t.co/q3aBjKq0Sn	1	0	0	1	2
@Sopianhaidi @imadya @PDL_Perjuangan @dprddkijakarta Saat Pemprov DKI Tetapkan PSBB Sampai Kerahkan Preman Kejar Warga Yg Langgar Prokes Tapi Gubernur & Wagubnya Malah Hadiri Kerumunan Terjadilah Pembangunan Warga Dan Meledaklah Warga Yg Terpapar Covid19	0	0	0	0	3

Field F₁ berisi jumlah *annotator* yang memberikan label positif (1) dan *field* F₀ berisi jumlah *annotator* yang memberikan label negatif (0). Persamaan *kappa value* dideskripsikan pada persamaan (6).

$$Kappa = \frac{P_0 - P_e}{1 - P_e} \dots\dots\dots (6)$$

Keterangan :

Kappa : Koefesien dari nilai kesepakatan di mana 0 untuk persetujuan secara kebetulan, 1 untuk persetujuan total.

P₀ : Proporsi frekuensi pengamatan.

P_e : Peluang kesepakatan antar pengamat.

Untuk mendapatkan nilai P₀ dan P_e dibutuhkan nilai P_i dan P_j. Dari sampel data di atas, maka didapat nilai P_i pada data pertama yang dideskripsikan pada persamaan (7).

$$P_i = \frac{1}{n(n-1)} (\sum_{j=1}^k n_{ij}^2 - n_{ij}) \dots\dots\dots (7)$$

$$P1 = \frac{1}{3(3-1)} * (3^2 + 0^2 - 3)$$

$$P1 = \frac{1}{6} * 6 = 1$$

Untuk nilai P_i pada sampel kedua dideskripsikan sebagai berikut.

$$P2 = \frac{1}{3(3-1)} * (1^2 + 2^2 - 3)$$

$$P2 = \frac{1}{6} * 2 = 0,3333$$

Untuk nilai P_i pada sampel ketiga dideskripsikan sebagai berikut.

$$P3 = \frac{1}{3(3-1)} * (0^2 + 3^2 - 3)$$

$$P3 = \frac{1}{6} * 6 = 1$$

Lakukan perhitungan P_i untuk seluruh data, sehingga dihasilkan nilai P_i untuk keseluruhan data yang dideskripsikan sebagai berikut.

$$\sum_{i=1}^N P_i = P1 + P2 + P3 + \dots$$

$$\sum_{i=1}^N P_i = 1 + 0,3333 + 1 + \dots$$

$$\sum_{i=1}^N P_i = 258,3333$$

Kemudian lakukan perhitungan \bar{P}_0 dengan persamaan (8) sehingga dihasilkan perhitungan yang dideskripsikan sebagai berikut.

$$\bar{P}_0 = \frac{1}{N} * \sum_{i=1}^N P_i \dots\dots\dots (8)$$

$$\bar{P}_0 = \frac{1}{275} * 258,3333 = 0,939394$$

Berdasarkan data yang diperoleh, diketahui jumlah label positif (F_1) sebanyak 659 label dan jumlah label negatif (F_0) sebanyak 166 label. Sehingga untuk mendapatkan nilai P_j untuk kategori positif dan negatif dapat diperoleh dengan persamaan (9) dan perhitungannya yang dideskripsikan sebagai berikut.

$$P_j = \frac{n.j}{Nn} = \frac{1}{Nn} \sum_{i=1}^N n_{ij} \dots\dots\dots (9)$$

Keterangan :

$n.j$ = $\sum_{i=1}^N n_{ij}$ = Total jumlah label untuk kategori.

N : Jumlah data.

n : Jumlah *annotator*.

$$P_1 = \frac{1}{275 * 3} * 659 = \frac{659}{825} = 0,798788$$

$$P_0 = \frac{1}{275 * 3} * 166 = \frac{166}{825} = 0,201212$$

Setelah mendapatkan nilai P_j pada kedua kategori, maka perhitungan \bar{P}_e dengan menggunakan persamaan (10) yang dideskripsikan sebagai berikut.

$$\bar{P}_e = \sum_{j=1}^k P_j^2 \dots\dots\dots (10)$$

$$\bar{P}_e = P_1^2 + P_0^2$$

$$\bar{P}_e = 0,798788^2 + 0,201212^2$$

$$= 0,638062 + 0,040486$$

$$= 0,678548$$

Setelah mendapat nilai \bar{P}_0 dan \bar{P}_e , maka perhitungan nilai kappa dengan menggunakan persamaan (6) yang dideskripsikan sebagai berikut.

$$Kappa = \frac{\bar{P}_0 - \bar{P}_e}{1 - \bar{P}_e}$$

$$Kappa = \frac{0,939394 - 0,678548}{1 - 0,678548}$$

$$Kappa = \frac{0,260846}{0,321452} = 0,811461$$

Berdasarkan pada skala *kappa value*, dengan hasil perhitungan nilai *kappa* sebesar 0,811461, maka hasil termasuk dalam kategori *almost perfect* yang berarti sangat cukup baik. Tahapan selanjutnya yaitu *text preprocessing* tahap 2 yang terdiri dari *cleaning*, *case folding*, *tokenization*, normalisasi bahasa, *filtering*, dan *stemming*. Hasil *text preprocessing* tahap 2 dideskripsikan pada Tabel 4.

Tabel 4: Text Preprocessing Data Sampel

Data Sebelum Text Preprocessing Tahap 2	Data Sesudah Text Preprocessing Tahap 2	Label
Pemerintah Provinsi DKI Jakarta mengeluarkan sejumlah aturan tentang denda bagi pelanggar aturan protokol kesehatan pada masa pandemi Covid-19. #TempoGrafis https://t.co/e9PqpUfOwo https://t.co/SgC2zwwWIC	perintah provinsi dki jakarta keluar atur denda langgar atur protokol sehat pandemi covid	Positif
Pandemi Covid-19 Kian Darurat, 911 Perusahaan di DKI Jakarta Ditutup Sementara Selama Masa PSBB https://t.co/q3aBjKq0Sn	pandemi covid kian darurat usaha dki jakarta tutup psbb	Negatif
@Sopianhaidi @imadya @PDI_Perjuangan @dprddjakarta Saat Pemprov DKI Tetapkan PSBB Sampai Kerahkan Preman Kejar Warga Yg Langgar Prokes Tapi Gubernur & Wagubnya Malah Hadiri Kerumunan Terjadilah Pembangkangan Warga Dan Meledaklah Warga Yg Terpapar Covid19	perintah provinsi dki tetap psbb kerah preman kejar langgar protokol sehat gubernur wagubnya hadir kerumun bangkang ledak papar covid	Negatif

Selanjutnya dilakukan pembobotan kata dengan TF-IDF pada data. Perhitungan TF-IDF pada sampel dideskripsikan pada Tabel 5.

Tabel 5: Perhitungan TF-IDF Data Sampel

Term	tf			df	IDF	w		
	D1	D2	D3			D1	D2	D3
perintah	1	0	1	2	0.176	0.176	0	0.176
provinsi	1	0	1	2	0.176	0.176	0	0.176
dki	1	1	1	3	0	0	0	0
jakarta	1	1	0	2	0.176	0.176	0.176	0
keluar	1	0	0	1	0.477	0.477	0	0
atur	2	0	0	1	0.477	0.954	0	0
denda	1	0	0	1	0.477	0.477	0	0
langgar	1	0	1	2	0.176	0.176	0	0.176
protokol	1	0	1	2	0.176	0.176	0	0.176
sehat	1	0	1	2	0.176	0.176	0	0.176
pandemi	1	1	0	2	0.176	0.176	0.176	0
covid	1	1	1	3	0	0	0	0
kian	0	1	0	1	0.477	0	0.477	0
darurat	0	1	0	1	0.477	0	0.477	0
usaha	0	1	0	1	0.477	0	0.477	0
tutup	0	1	0	1	0.477	0	0.477	0
psbb	0	1	1	2	0.176	0	0.176	0.176
tetap	0	0	1	1	0.477	0	0	0.477
kerah	0	0	1	1	0.477	0	0	0.477
preman	0	0	1	1	0.477	0	0	0.477
kejar	0	0	1	1	0.477	0	0	0.477
gubernur	0	0	1	1	0.477	0	0	0.477
wagubnya	0	0	1	1	0.477	0	0	0.477
hadir	0	0	1	1	0.477	0	0	0.477
kerumun	0	0	1	1	0.477	0	0	0.477
bangkang	0	0	1	1	0.477	0	0	0.477
ledak	0	0	1	1	0.477	0	0	0.477
papar	0	0	1	1	0.477	0	0	0.477

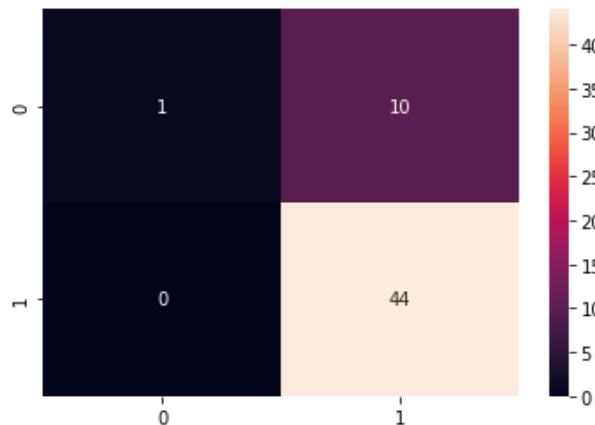
Setelah didapatkan nilai bobot, maka selanjutnya dilakukan proses klasifikasi dengan melakukan prediksi label dari data uji (*data testing*) yang diperoleh dari data latih (*data training*) kemudian hasilnya dibandingkan dengan label sebenarnya dari data uji tersebut. Pembagian *data training* dan *data testing* menggunakan perbandingan 80 : 20 yang dideskripsikan pada Tabel 6.

Tabel 6: Pembagian Data Latih dan Data Uji

	Label		Jumlah
	Positif	Negatif	
Data Latih	178	42	220
Data Uji	44	11	55
Total	222	53	275

Dengan menggunakan *Multinomial Naïve Bayes*, didapatkan hasil klasifikasi yang dideskripsikan pada Tabel 7 berupa *confusion matrix*.

Tabel 7: Confusion Matrix Hasil Klasifikasi



Dari Tabel 7, diketahui informasi yaitu nilai TN (*True Negative*) sebesar 1, nilai TP (*True Positive*) sebesar 44, nilai FN (*False Negative*) sebesar 0, dan nilai FP (*False Positive*) sebesar 10. Untuk mengukur kinerja algoritma digunakan nilai akurasi, *recall*, dan *specificity* dengan hasil sebagai berikut.

$$accuracy = \frac{TP + TN}{TP + TN + FP + FN} = \frac{44 + 1}{44 + 1 + 10 + 0} = 0,818$$

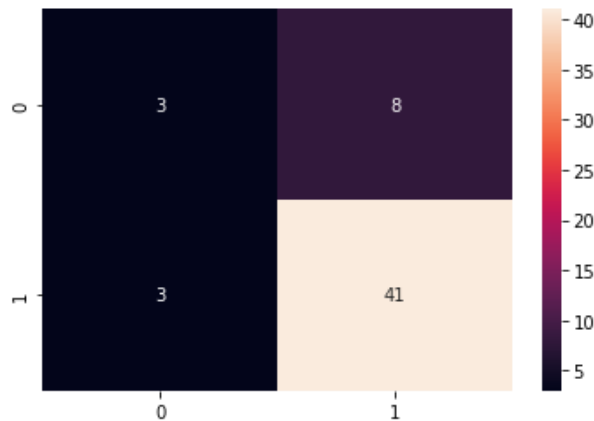
$$recall = \frac{TP}{TP + FN} = \frac{44}{44 + 0} = 1,0$$

$$specificity = \frac{TN}{TN + FP} = \frac{1}{1 + 10} = 0,09$$

Nilai *recall* dan *specificity* memiliki perbedaan nilai yang jauh. Hal ini disebabkan pada data latih, jumlah kelas positif yaitu sebanyak 178 data lebih banyak dibanding jumlah kelas negatif sebanyak 42 data, sehingga jumlah data antar kelas tidak seimbang.

Untuk mengatasi ketidakseimbangan pada data, maka solusi yang digunakan yaitu *oversampling* dengan menggunakan SMOTE (*Synthetic Minority Over-Sampling Technique*). Metode SMOTE akan menambah jumlah data kelas minor (kelas negatif) sehingga jumlahnya setara dengan kelas mayor (kelas positif) dengan cara membangkitkan data buatan. Hasil klasifikasi model setelah dilakukan *oversampling* dideskripsikan pada Tabel 8.

Tabel 8: *Confusion Matrix* Hasil Klasifikasi Dengan *Oversampling*



Dari Tabel 8, diketahui informasi yaitu nilai TN sebesar 3, nilai TP sebesar 41, nilai FN sebesar 3, dan nilai FP sebesar 8. Hasil penghitungan nilai *accuracy*, *recall*, dan *specificity* dideskripsikan sebagai berikut.

$$accuracy = \frac{TP + TN}{TP + TN + FP + FN} = \frac{41 + 3}{41 + 3 + 8 + 3} = 0,8$$

$$recall = \frac{TP}{TP + FN} = \frac{41}{41 + 3} = 0,9318$$

$$specificity = \frac{TN}{TN + FP} = \frac{3}{3 + 8} = 0,27$$

Dari hasil evaluasi yang telah dilakukan, nilai akurasi tidak jauh berbeda dengan sebelumnya yaitu sebesar 0,8. Nilai *recall* menurun menjadi 0,9318 dari sebelumnya sebesar 1,0. Namun terjadi peningkatan terhadap nilai *specificity* yaitu sebesar 0,27 dari nilai sebelumnya yaitu 0,09.

Berdasarkan jumlah kata-kata yang sering muncul pada *tweet*, visualisasi kata-kata dalam bentuk *wordcloud* sentimen positif mengenai PSBB di Jakarta diilustrasikan pada Gambar 4.



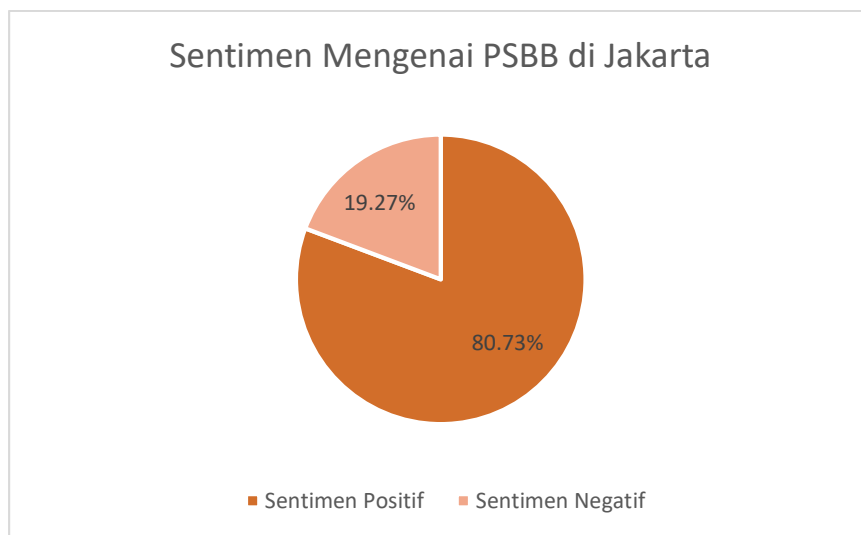
Gambar 4: *Wordcloud* Sentimen Positif PSBB di Jakarta

Berdasarkan Gambar 4, kata dalam sentimen positif yang banyak muncul selain kata kunci “PSBB DKI Jakarta” yaitu “perintah”, “provinsi”, dan “panjang”, di mana pemerintah provinsi DKI Jakarta terus berupaya menekan penyebaran penyakit COVID-19 dengan memperpanjang pemberlakuan PSBB secara berkala. Adapun visualisasi *wordcloud* untuk sentimen negatif yang diilustrasikan pada Gambar 5.



Gambar 5: Wordcloud Sentimen Negatif PSBB di Jakarta

Berdasarkan Gambar 5, kata dalam sentimen negatif yang banyak muncul selain kata kunci yaitu “langgar”, di mana banyak terjadi pelanggaran yang dilakukan oleh cafe dengan beroperasi di luar jam operasi yang telah dibatasi oleh pemerintah. Visualisasi mengenai perbandingan jumlah sentimen positif dan negatif diilustrasikan pada Gambar 6.



Gambar 6: Perbandingan Sentimen Pada PSBB di Jakarta

Berdasarkan Gambar 6, dapat dilihat jumlah sentimen positif lebih banyak dibandingkan dengan sentimen negatif dengan masing-masing persentase sebesar 80,73% untuk sentimen positif dan 19,27% untuk sentimen negatif.

4 KESIMPULAN

4.1 Kesimpulan

Berdasarkan hasil penelitian yang dihasilkan, dapat disimpulkan hal-hal sebagai berikut.

1. Opini masyarakat pengguna Twitter terhadap PSBB (Pembatasan Sosial Berskala Besar) di Jakarta cenderung positif. Terlihat pada hasil *wordcloud* sentimen positif, kata yang banyak muncul yaitu “perintah”, “provinsi”, dan “panjang”, di mana pemerintah provinsi DKI Jakarta terus berupaya menekan penyebaran penyakit COVID-19 dengan memperpanjang pemberlakuan PSBB secara berkala.

2. Metode *Naïve Bayes Classifier* dilakukan dengan menggunakan metode *Multinomial Naïve Bayes*. Setelah data dilakukan *text preprocessing*, pelabelan data, dan pembobotan kata, data dibagi menjadi data latih dan data uji. Kemudian dilakukan proses klasifikasi yang hasilnya akan dievaluasi menggunakan *confusion matrix*.
3. Hasil evaluasi klasifikasi sentimen terhadap PSBB di Jakarta dengan menggunakan metode *Naïve Bayes Classifier* dengan 80% data latih dan 20% data uji serta dilakukan *oversampling* diperoleh nilai akurasi sebesar 0.8, nilai *recall* sebesar 0.9318, dan nilai *specificity* sebesar 0.27.

4.2 Saran

Saran yang dapat digunakan sehingga pengembangan penelitian ini ke depannya agar lebih baik yaitu sebagai berikut.

1. Pengambilan data menggunakan *keyword* yang bervariasi seperti “PSBB Jakarta” atau “PSBB DKI” dan menggunakan atribut lokasi *tweet* yaitu di Jakarta, sehingga data yang terkumpul lebih banyak dan penggunaanya berasal dari masyarakat Jakarta yang menggunakan Twitter.
2. Data bisa diperoleh dari berbagai media sosial ataupun jejaring sosial lainnya, sehingga lebih banyak opini yang terkumpul.
3. Menambah periode pengumpulan data untuk waktu selanjutnya.
4. Pelabelan data diharapkan menggunakan ahli sehingga hasil pelabelan bisa lebih tepat.

Referensi

Alrajak, M. Suyudi. (2020). Analisis Sentimen Terhadap Pelayanan PT. PLN di Jakarta pada Media Sosial *Twitter* Menggunakan Metode *K-Nearest Neighbor* (K-NN). Skripsi. FIK, Informatika, Universitas Pembangunan Nasional Veteran Jakarta, Jakarta.

Analisis. (2016). Diakses pada Januari 19, 2022, dari KBBI: <https://kbbi.kemdikbud.go.id/entri/analisis>.

Eka Fahreza H. (2017). PENERAPAN *SYNTHETIC MINORITY OVERSAMPLING TECHNIQUE* TERHADAP DATA TIDAK SEIMBANG PADA MODEL *MULTIVARIATE ADAPTIVE REGRESSION SPLINE*. Skripsi. Universitas Hasanuddin.

I Gusti Naufhal Daffa Adnyana. (2021). *ANALISIS SENTIMEN TENTANG UU CIPTA KERJA MENGGUNAKAN ALGORITMA NAÏVE BAYES*. Skripsi. FIK, Informatika, Universitas Pembangunan Nasional Veteran Jakarta.

Yuli Nurhasinah. (2021). Beda PSBB dan Karantina Wilayah, Apa Saja? [Halaman web]. Diakses dari <https://indonesiabaik.id/infografis/beda-psbb-dan-karantina-wilayah-apa-saja>.