

MODEL PREDIKTIF BERASASKAN SENTIMEN UNTUK MALAYSIA PASCA COVID-19 PEMBANGUNAN INDUSTRI PELANCONGAN

LIEW SET TENG
MOHD RIDZWAN YAAKUB

Fakulti Teknologi & Sains Maklumat, Universiti Kebangsaan Malaysia

ABSTRAK

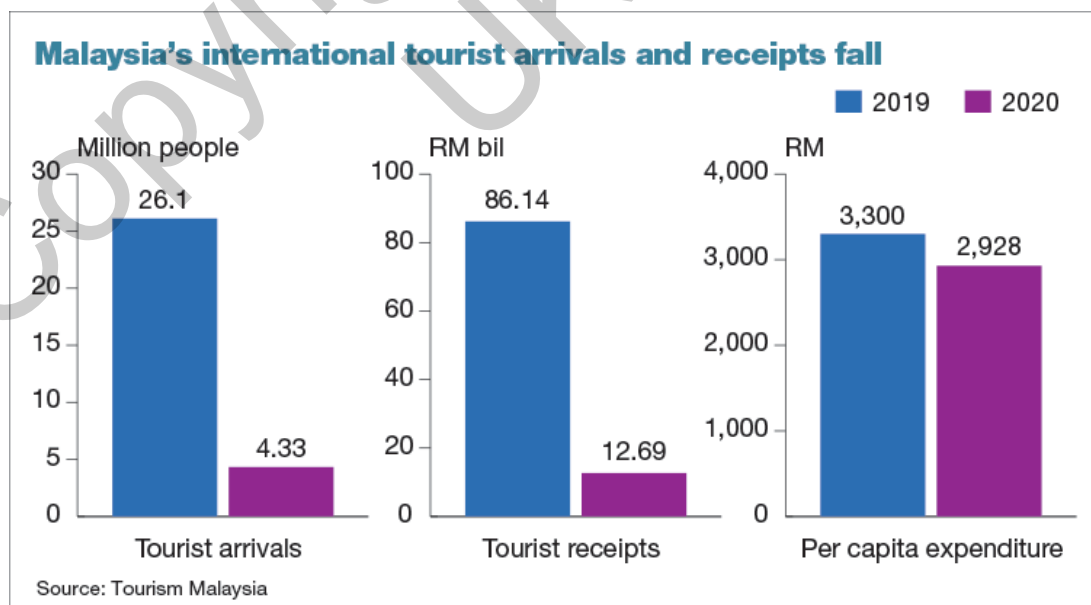
Industri pelancongan merupakan salah satu sektor ekonomi terbesar dan paling pesat berkembang di negara ini kerana ia merupakan penyumbang ketiga terbesar kepada Keluaran Dalam Negara Kasar (KDNK) Malaysia, selepas pembuatan dan komoditi (Hirschmann, 2020). Walau bagaimanapun, wabak COVID-19 telah banyak menjejaskan industri pelancongan di Malaysia bermula dari tahun 2020. Peningkatan kes COVID-19 di Malaysia telah menyebabkan penurunan ketara dalam bilangan pelancong kerana Malaysia telah menutup sempadannya dan melaksanakan Perintah Kawalan Pergerakan (PKP) mengehadkan perjalanan dalam negara. Berdasarkan Motac (Kementerian Pelancongan, Kesenian dan Kebudayaan), ia menyasarkan menerima 24.3 juta ketibaan pelancong antarabangsa dengan pendapatan RM73 bilion dan RM100 bilion untuk pelancongan domestik menjelang 2025 berbanding hanya 4.33 juta ketibaan pelancong antarabangsa dengan pendapatan RM12.7 bilion pada tahun 2020. Oleh itu, pandangan dan respon masyarakat terhadap pembangunan industri pelancongan merupakan agenda penting yang perlu diketahui. Dengan ini, analisis sentimen terhadap pembangunan industri pelancongan pasca COVID-19 di Malaysia telah dijalankan dalam kajian ini. Twitter dipilih sebagai sumber maklumat analisis sentimen kerana ia mengandungi banyak pendapat dan reaksi orang yang berbeza. Kajian ini dijalankan dengan menggunakan klasifikasi dalam teknik pembelajaran mesin terselia dan menggunakan Naïve Bayes dan Mesin Vektor Sokongan (SVM) sebagai pengelas untuk menganalisis sentimen terhadap respon masyarakat terhadap isu ini. Kemudian, pengelas yang mencapai ketepatan tertinggi antara tiga pengelas iaitu SVM digunakan sebagai model ramalan sentimen dalam kajian ini. Hasil analisis sentimen tersebut telah memberi nilai data mengenai respon masyarakat terhadap pembangunan industri pelancongan Malaysia selepas COVID-19 dan boleh dirujuk oleh mereka yang memerlukannya seperti pelabur dan pihak berkuasa berkaitan.

1 PENGENALAN

Industri pelancongan di Malaysia merupakan salah satu sektor penting bagi ekonomi negara kerana ia merupakan penyumbang ketiga terbesar kepada KDNK Malaysia. Industri pelancongan juga memainkan peranan penting dalam ekonomi Malaysia kerana ia boleh membawa masuk banyak pertukaran asing dan mewujudkan banyak peluang pekerjaan. Industri pelancongan merujuk kepada aktiviti orang yang pergi ke tempat lain untuk tujuan riadah, keagamaan, perniagaan atau perubatan. Terdapat 3 jenis pelancongan iaitu pelancongan masuk, pelancongan keluar dan pelancongan domestik. Oleh itu, pembangunan industri pelancongan adalah penting kepada negara.

Walau bagaimanapun, industri pelancongan berdepan dengan kesan pandemik COVID-19. COVID-19 merupakan penyakit berjangkit yang disebabkan oleh virus SARS-CoV-2. Oleh itu, dunia sedang menghadapi kecemasan kesihatan, sosial dan ekonomi global yang belum pernah terjadi sebelumnya akibat pandemik COVID-19. Perjalanan dan pelancongan adalah antara sektor yang paling terjejas oleh penurunan besar dalam permintaan antarabangsa di tengah-tengah sekatan perjalanan global termasuk banyak sempadan ditutup sepenuhnya, untuk membendung virus itu.

Wabak pandemik yang teruk telah menuntut negara melaksanakan Perintah Kawalan Pergerakan (PKP) dalam usaha membendung peningkatan kes COVID-19 di Malaysia. Bagaimanapun, pergerakan terhad dan sempadan tertutup telah memberi kesan negatif kepada industri pelancongan negara. Sebagai industri berorientasikan rakyat, pelancongan adalah salah satu sektor ekonomi yang paling teruk terjejas. Oleh itu, pelan pemulihan pelancongan yang berkesan adalah penting untuk mengimbangi antara ekonomi dan mata pencarian. Selain itu, peningkatan kes COVID-19 di Malaysia telah menyebabkan banyak lawatan dibatalkan, yang telah menyebabkan penurunan ketara dalam bilangan pelancong ke Malaysia. Dengan ini, analisis sentimen amat berguna untuk dijadikan panduan dan maklumat yang berguna kepada masyarakat yang berminat agar mereka dapat membuat keputusan yang tepat untuk memajukan industri pelancongan.



Rajah 1 Perbandingan ketibaan antarabangsa Malaysia pada tahun 2019 dan 2020

Analisis sentimen ialah proses menentukan sama ada sesuatu teks itu positif, negatif atau neutral. Ia sering digunakan oleh perniagaan untuk menjejaki sentimen dalam data sosial,

mengukur reputasi jenama dan memahami keperluan pelanggan. Analisis sentimen adalah bantuan yang sangat baik untuk manusia kerana ia boleh dilakukan secara automatik, keputusan boleh dibuat berdasarkan jumlah data yang besar dan bukannya gerak hati biasa yang tidak selalu betul.

2 PENYATAAN MASALAH

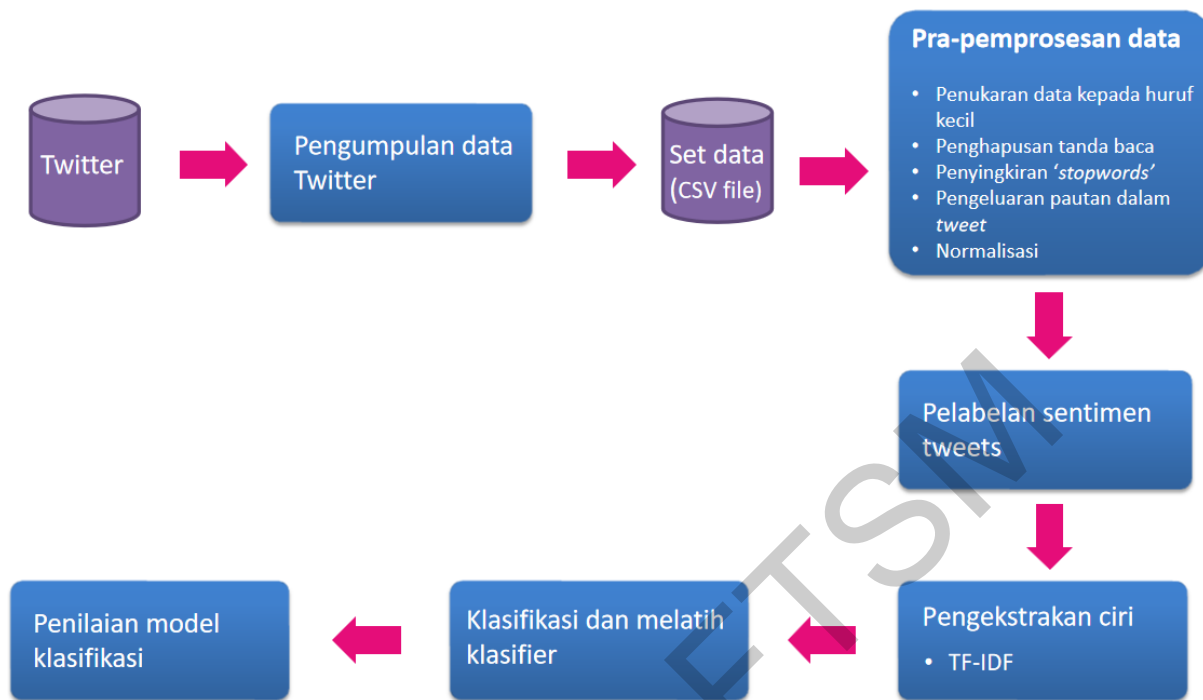
Pada tahun 2020, pembangunan industri pelancongan di Malaysia telah terjejas disebabkan oleh pandemik COVID-19 yang tersebar di seluruh dunia. Oleh itu, pendapatan industri pelancongan telah mencatatkan kemerosotan bilangan pelancong antarabangsa. Jadi, ini telah menyebabkan kemerosotan ekonomi negara (Tourism Malaysia, 2020). Dengan ini, masyarakat dan pihak berkuasa berusaha untuk terus membangunkan industri pelancongan Malaysia. Oleh yang demikian, mereka perlu informasi yang berguna iaitu mengenai sentimen masyarakat terhadap industri pelancongan untuk mengetahui respon dan pandangan masyarakat supaya mereka dapat membuat keputusan yang betul.

3 OBJEKTIF KAJIAN

Matlamat kajian ini adalah untuk mengumpul serta menganalisis respon dan sentimen masyarakat terhadap pasca COVID-19 pembangunan industri pelancongan Malaysia. Selain itu, mengkaji pengelas yang paling sesuai untuk melatih dan mengira kebarangkalian ketepatan. Akhir sekali, menghasilkan data analisis yang dalam bentuk mudah difahami dan mampu menjadi rujukan kepada masyarakat.

4 METOD KAJIAN

Penggunaan model pembangunan yang bersesuaian adalah penting bagi memastikan kajian berjalan lancar dan memperoleh hasil yang berkualiti dan tepat. Rajah 1 menunjukkan proses pembangunan model analisis sentimen dalam kajian ini. Kajian ini mempunyai fasa pengumpulan data, fasa pra-pemprosesan data, fasa pelabelan sentimen *tweets*, fasa pengekstrakan ciri, fasa klasifikasi dan melatih klasifier dan fasa penilaian model klasifikasi.



Rajah 2 Proses pembangunan model analisis sentimen

4.1 Fasa Pengumpulan Data

Fasa ini amat penting dalam kajian ini kerana data berkualiti yang digunakan dalam proses analisis sentimen dapat menghasilkan keputusan yang tepat. Bagi memperoleh data tersebut, Twitter dipilih sebagai sumber untuk pengumpulan data yang disebabkan oleh Twitter mempunyai banyak maklumat yang orang ramai berikan setiap hari tentang produk, personaliti, acara, komuniti. Oleh itu, banyak data pandangan dan sentimen berbeza yang ditulis daripada orang yang berbeza boleh dikumpulkan. Pustaka Tweepy digunakan untuk mendapatkan semula set data mentah daripada aplikasi Twitter menggunakan API Twitter. Akhirnya, data yang dikumpul akan disimpan dalam format fail *comma-separated values* (CSV). Fail CSV memudahkan proses pengubahsuaian dan manipulasi maklumat.

4.2 Fasa Pra-pemrosesan Data

Fasa ini melibatkan proses pembersihan data yang dikumpul. Proses ini ialah teknik perlombongan data untuk menukar data mentah kepada maklumat yang berguna dan mudah difahami. Langkah pra-pemrosesan yang dijalankan ialah menukar karakter kepada huruf kecil, menghapuskan sebarang perkataan yang kurang bermakna dalam analisis sentimen seperti pautan laman web (URL), nama pengguna (@-mention), hashtag (#), RT (*retweet*), simbol, nombor serta emoji. Selepas itu, penyingkiran *stopwords* untuk menapis perkataan yang paling biasa dalam

bahasa kerana ia tidak membawa makna yang ketara kepada ayat itu. Seterusnya, proses normalisasi iaitu kaedah yang bertanggungjawab untuk mengelompokkan bentuk kata yang berbeza dalam bentuk akar, yang mempunyai makna yang sama. Proses-proses ini adalah amat penting dalam pra-pemprosesan data bagi meningkatkan ketepatan keputusan analisis.

	tweet
0	RT @sarahav8n: From Malaysia to South Korea fo...
1	@cook4beginners @ERosson1982 They had shrimp c...
2	Malaysia is open again, and this team went to ...
3	Emirates signs MoC with Malaysia Tourism Board...
4	MY #MalaysiaAirlines will expand its internati...

Rajah 3 Set data asal

	cleaned_tweet
0	malaysia south korea stray kid conce somewhat ...
1	shrimp chip lived malaysia singapore father la...
2	malaysia open team went test travel
3	emirate sign moc malaysia tourism board emirat...
4	expand international network new direct flight...

Rajah 4 Hasil pra-pemprosesan data

4.3 Fasa Pelabelan Sentimen Tweets

Sebelum analisis sentimen dilaksanakan, proses pelabelan sentimen bagi setiap *tweet* akan dilakukan pada fasa ini. Untuk mengklasifikasikan data kepada kelas positif, negatif dan neutral, VADER (*Valence Aware Dictionary for Sentiment Reasoning*) dan Textblob akan digunakan dan hasil kedua-dua penganalisis sentimen ini juga akan dibandingkan bagi memilih teknik yang lebih sesuai.

4.4 Fasa Pengekstrakan Ciri

Fasa pengekstrakan ciri dilakukan untuk mendapatkan ciri-ciri penting bagi set data yang digunakan. Ciri utama ialah TF-IDF (kekerapan dokumen terbalik kekerapan jangka) yang merupakan ukuran statistik yang menilai sejauh mana sesuatu perkataan berkaitan dengan

dokumen dalam koleksi dokumen. Ia boleh dilakukan dengan mendarabkan metrik istilah kekerapan (TF) dan kekerapan dokumen songsang (IDF). TF ialah kiraan kekerapan perkataan yang muncul dalam dokumen, iaitu bilangan kali perkataan muncul dalam dokumen dan IDF ialah berapa biasa atau jarang sesuatu perkataan muncul dalam keseluruhan set dokumen. Dalam fasa ini, perpustakaan *Scikit-Learn* Python mengandungi kelas *TfidfVectorizer* yang boleh digunakan untuk menukar ciri teks kepada vektor ciri TF-IDF.

4.5 Fasa Klasifikasi dan Melatih Klasifier

Fasa ini memerlukan pengelas untuk mengenal pasti kategori set data yang diberikan, dan pengelas ini digunakan terutamanya untuk meramalkan output bagi kategori data. Kajian ini dijalankan dengan menggunakan teknik pembelajaran mesin terselia untuk perlombongan data. Pendekatan ini melatih model dengan sumber data berlabel. Model terlatih kemudiannya boleh membuat ramalan untuk output mempertimbangkan data input baru yang tidak berlabel. Pengelas Naïve Bayes dan Mesin Vektor Sokongan (SVM) akan digunakan untuk ramalan. Ketiga-tiga model klasifikasi akan dilatih dengan data latihan dan diuji dengan data ujian.

4.6 Fasa Penilaian Model Klasifikasi

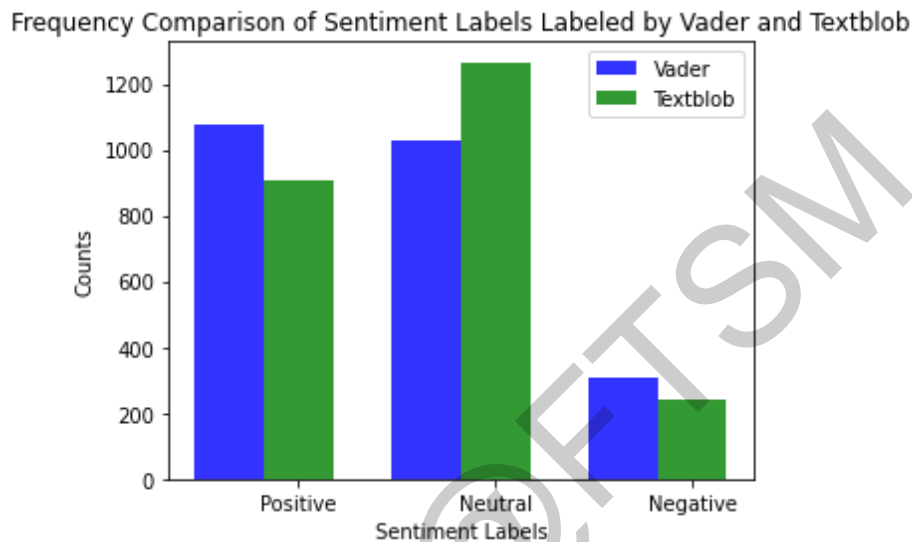
Fasa ini bermatlamat untuk menilai serta menguji hasil analisis sentimen bagi mengenal pasti ketepatan klasifikasi iaitu klasifikasi Naïve Bayes dan klasifikasi Mesin Vektor Sokongan (SVM). Matriks konfusi digunakan untuk menunjukkan model prestasi daripada data ujian. Laporan klasifikasi boleh menghasilkan nilai ketepatan (*Accuracy*), dapatan (*Precision*), kejituan (*Recall*) dan ukuran-f1 (F1) yang memberikan maklumat yang lebih terperinci tentang model prestasi. Selain itu, nilai AUC (*Area Under The Curve*) adalah ukuran keupayaan klasifikasi untuk membezakan antara kelas. Lebih tinggi AUC, lebih baik prestasi model dalam membezakan antara kelas positif dan negatif. Seterusnya, model yang mencapai prestasi terbaik dalam hipotesis ujian kemudiannya dipilih sebagai model akhir dan digunakan untuk membuat ramalan ke atas data.

5 HASIL KAJIAN

Terdapat beberapa hasil kajian berdasarkan analisis yang telah dilakukan iaitu pelabelan sentimen *tweets*, penilaian model klasifikasi, analisis sentimen dan papan pemuka analisis sentimen.

5.1 Hasil Pelabelan Sentimen Tweets

Rajah 5 menunjukkan hasil pelabelan sentimen *tweets* yang dilakukan oleh VADER dan Textblob yang divisualisasikan dalam bentuk graf.



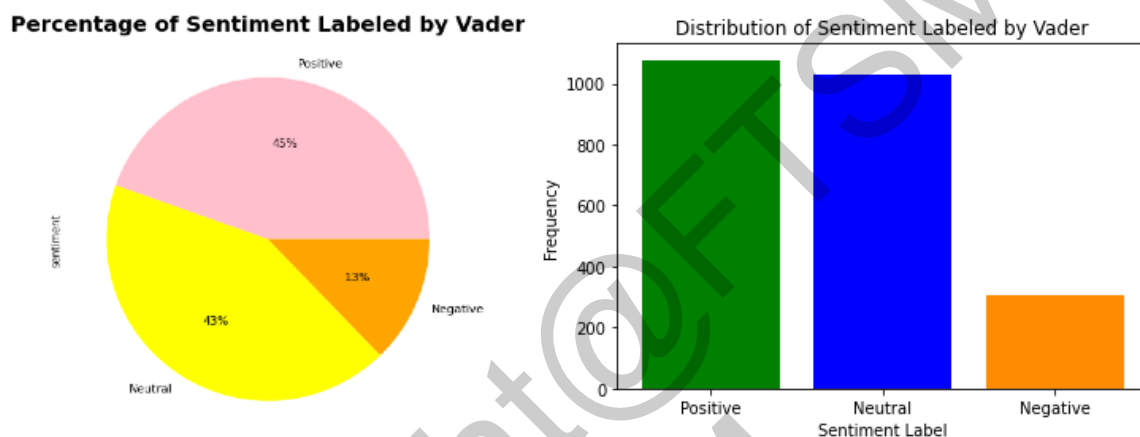
Rajah 5 Perbandingan bilangan label sentimen yang dilabelkan oleh Vader dan Textblob

Berdasarkan rajah 5, Vader telah melabelkan 1077 *tweets* sebagai sentimen positif yang lebih tinggi berbanding Textblob yang hanya mempunyai 909 *tweets* berlabel positif. Bagi sentimen neutral, Textblob mempunyai bilangan *tweets* tertinggi iaitu 1265 *tweets* yang telah melepasi bilangan *tweets* yang dilabel neutral oleh Vader (1030 *tweets*). *Tweets* negatif merekodkan jumlah terendah iaitu hanya 308 *tweets* yang dilabelkan oleh Vader dan 241 *tweets* yang dilabelkan oleh Textblob.

Dengan mengambil kira perkara ini, Vader telah dipilih sebagai penganalisis sentimen untuk melabelkan sentimen pada set data. Hal ini kerana Vader boleh melabelkan sentimen set data dengan lebih tepat dan lebih sesuai digunakan dalam analisis sentimen berbanding Textblob. Ini kerana Textblob telah melabelkan *tweets* yang bersentimen neutral lebih banyak daripada sentimen positif dan negatif. Jadi, set data yang dilabelkan oleh Textblob tidak sesuai untuk digunakan dalam analisis sentimen kerana ia boleh menyebabkan ralat ramalan oleh pengelas terlatih.

5.2 Analisis Sentimen

Hasil analisis sentimen yang dilakukan oleh VADER pada set data divisualisasikan dalam bentuk graf bar dan carta pai untuk memudahkan pengguna memahami maklumat tersebut. Keputusan analisis sentimen dikeluarkan dalam carta pai dan carta bar dalam rajah 6 yang menggambarkan peratusan dan bilangan *tweets* yang bersentimen positif, neutral dan negatif. Set data yang digunakan dalam kajian ini mendapati 45% *tweets* adalah positif, 43% adalah neutral dan 13% adalah negatif.



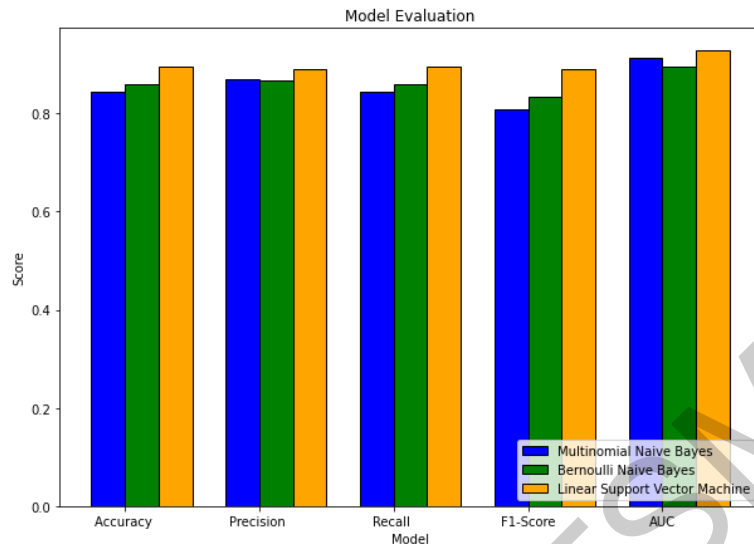
Rajah 6 Output analisis sentimen

5.3 Penilaian Model Klasifikasi

Prestasi pengelasan dalam pengelasan set data yang diproses dinilai untuk membandingkan prestasi model untuk memilih model yang paling sesuai sebagai model ramalan dalam kajian ini. Jadual 1 menunjukkan penilaian bagi tiga pengelasan yang digunakan dalam pengelasan set data dan rajah 7 menunjukkan perbandingan penilaian pengelasan.

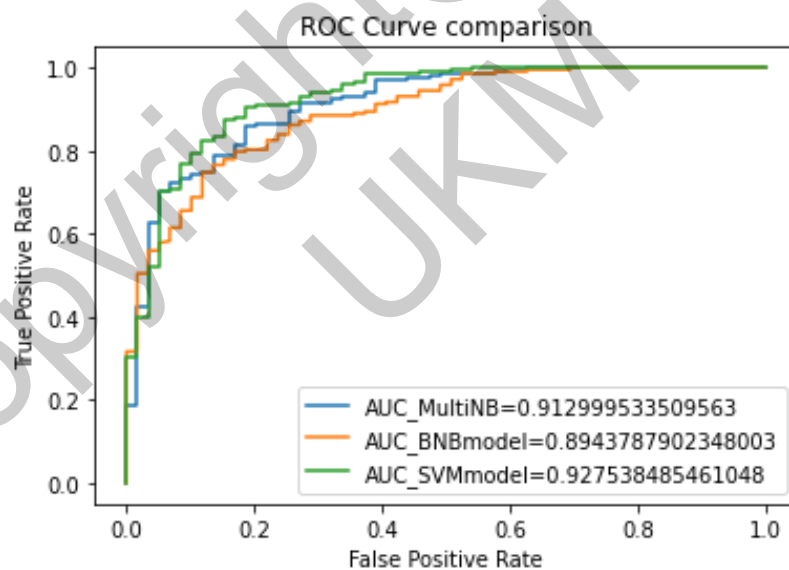
Jadual 1 Penilaian pengelasan

Model	Accuracy	Precision	Recall	F1	AUC
Multinomial	0.8448	0.8703	0.8448	0.8072	0.9130
Bernoulli	0.8592	0.8672	0.8592	0.8348	0.8944
Support Vector Machine (SVM)	0.8953	0.8914	0.8953	0.8903	0.9275



Rajah 7 Carta bar perbandingan prestasi pengelas

Berdasarkan rajah 7, model SVM mencapai skor ketepatan (*Accuracy*), dapatan (*Precision*), kejituan (*Recall*), ukuran-f1 (F1) dan AUC tertinggi berbanding model Multinomial dan Bernoulli.



Rajah 8 Perbandingan Lengkung ROC

Rajah 8 menunjukkan lengkung ROC diplot dengan TPR melawan FPR di mana TPR berada pada paksi-y dan FPR berada pada paksi-x. Dengan ini, nilai yang lebih kecil pada paksi-x plot menunjukkan positif palsu yang lebih rendah dan negatif benar yang lebih tinggi. Walau bagaimanapun, nilai yang lebih besar pada plot paksi-y menunjukkan positif benar yang lebih

tinggi dan negatif palsu yang lebih rendah. Kesimpulannya, semakin banyak lengkung itu memeluk sudut kiri atas plot, semakin baik model itu mengklasifikasikan data ke dalam kategori.

Berdasarkan rajah 8, model SVM mencapai prestasi yang paling baik antara ketiga-tiga model yang diuji kerana lengkungnya adalah yang paling memeluk sudut kiri atas plot disebabkan AUC yang tertinggi.

Dalam kajian ini, keputusan ujian hipotesis dibuat seperti berikut:

Jika:

nilai-p > 0.05: Gagal menolak hipotesis nol bahawa model mempunyai prestasi min yang sama dan sebarang perbezaan yang diperhatikan dalam ketepatan min adalah kebarangkalian kebetulan statistik.

nilai-p <= 0.05: Menolak hipotesis nol bahawa model mempunyai prestasi min yang sama, yang bermaksud perbezaan itu mungkin ketara.

Jadual 2 Ujian Hipotesis bagi model SVM & BNB

	SVM	BNB
Min Ketepatan	0.839	0.811
Nilai-p	0.011	
statistik-t	3.930	
Hipotesis	Perbezaan antara prestasi min adalah ketara	

Berdasarkan jadual 2, nilai-p dalam kes ini ialah 0.011 iaitu kurang daripada nilai alfa (0.05), oleh itu, menolak hipotesis nol, menunjukkan bahawa sebarang perbezaan yang diperhatikan antara algoritma adalah ketara.

Jadual 3 Ujian Hipotesis bagi model SVM & MultiNB

	SVM	MultiNB
Min Ketepatan	0.839	0.809
Nilai-p	0.006	
statistik-t	4.595	
Hipotesis	Perbezaan antara prestasi min adalah ketara	

Mengikut jadual 3, nilai-p adalah 0.006 iaitu kurang daripada nilai alfa (0.05), oleh itu, hipotesis nol juga ditolakkan, menunjukkan bahawa sebarang perbezaan yang diperhatikan antara algoritma adalah ketara.

Jadual 4 Ujian Hipotesis bagi model MultiNB & BNB

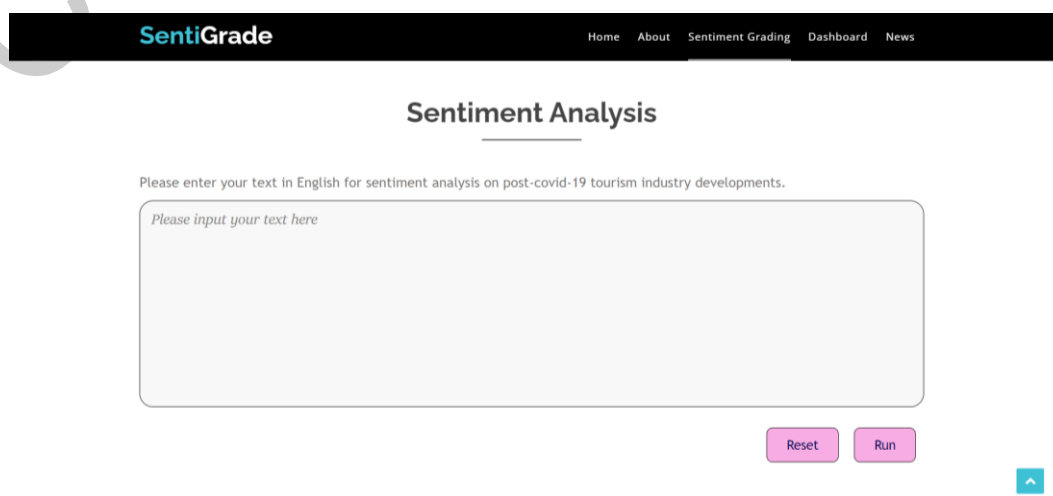
	MultiNB	BNB
Min Ketepatan	0.809	0.811
Nilai-p	0.791	
statistik-t	0.280	
Hipotesis	Model mempunyai prestasi min yang sama	

Jadual 4 menunjukkan bahawa nilai-p ialah 0.791 yang jauh lebih besar daripada nilai alfa (0.05), oleh itu, menyebabkan gagal untuk menolak hipotesis nol, menunjukkan model mempunyai prestasi min yang sama.

Kesimpulannya, model SVM mencapai ketepatan dan prestasi tertinggi antara ketiga-tiga model. Oleh itu, ia dipilih sebagai model ramalan dalam kajian ini untuk analisis sentimen.

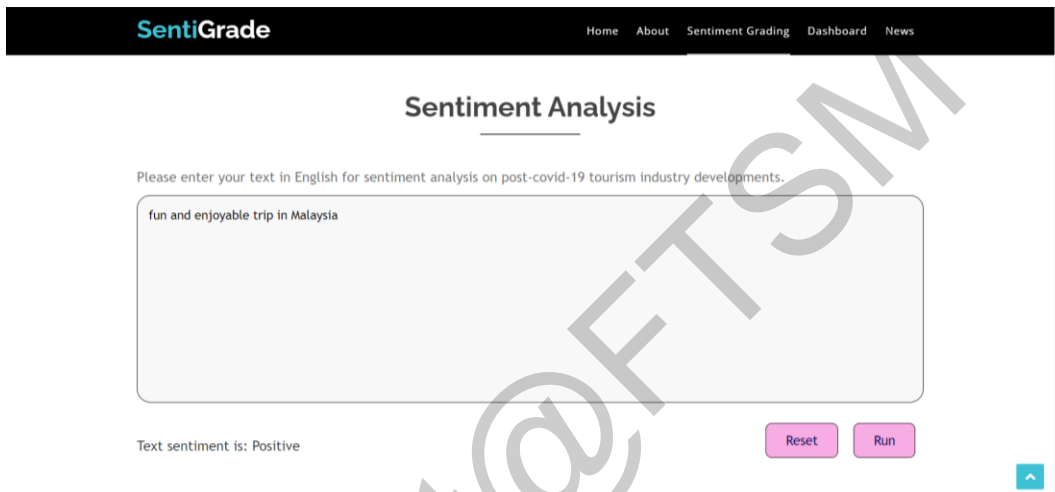
5.4 Papan Pemuka Analisis Sentimen

Papan pemuka analisis sentimen dibangunkan dengan menggunakan bahasa pengaturcaraan PHP dan HTML, berdasarkan hasil kajian yang diperolehi dengan bahasa pengaturcaraan Python. Perisian yang digunakan ialah Sublime Text Editor.



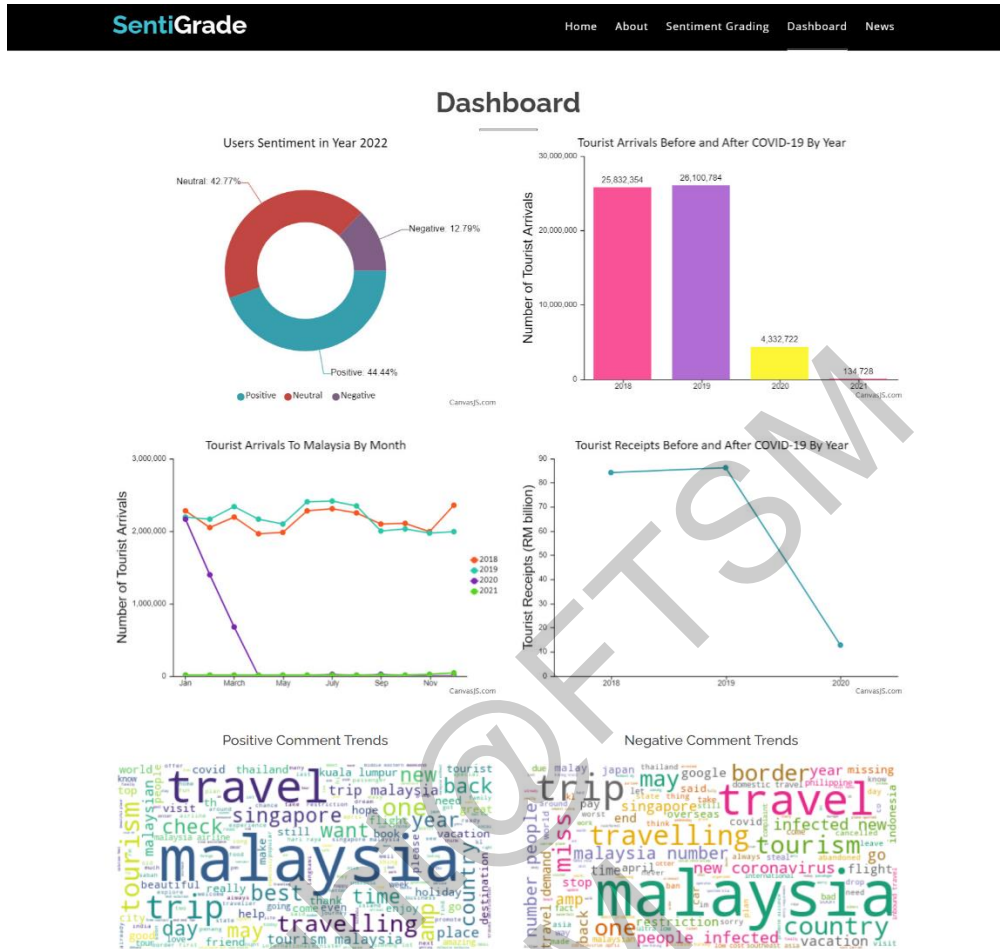
Rajah 9 Antara muka analisis sentimen

Rajah 9 merupakan antara muka bagi pengguna untuk melakukan analisis sentimen. Dengan ini, pengguna boleh melakukan analisis sentimen dengan mengisi teks atau ayat dalam bahasa Inggeris di ruangan yang disediakan dan kemudian menekan butang 'Run' untuk menjalankan analisis dan mendapatkan keputusan analisis.



Rajah 10 Antara muka keputusan analisis sentimen

Rajah 10 merupakan antara muka keputusan analisis sentimen. Keputusan analisis akan dipaparkan kepada pengguna dan pengguna dapat mengetahui sentimen tentang teks/ayat yang dimasukkan oleh pengguna sama ada positif atau negatif. Pengguna boleh menekan butang 'Reset' untuk mengosongkan teks dan output yang dipaparkan.



Rajah 11 Antara muka hasil analisis

Dalam rajah 11, hasil keputusan analisis yang diperoleh dalam analisis sentimen serta maklumat lain yang penting tentang tajuk kajian ini turut dipaparkan dalam bentuk yang mudah difahami bagi dijadikan rujukan kepada pengguna yang memerlukannya.

6 KESIMPULAN

Pengkaji telah berjaya membangunkan analisis sentimen terhadap pembangunan industri pelancongan pasca COVID-19 di Malaysia yang menepati objektif kajian yang dibincangkan pada awal kajian. Model SVM menghasilkan keputusan analisis sentimen yang paling tepat pada set data yang dikumpul daripada Twitter. Selain itu, papan pemuka juga dibangunkan dalam bentuk yang mudah difahami dan senang digunakan oleh pengguna.

Terdapat beberapa limitasi dalam kajian ini iaitu tarikh *tweets* yang telah dikumpul dalam kajian ini tertumpu pada tahun 2022 walaupun tarikh tersebut telah ditetapkan semasa dalam proses mengekstrak *tweets* daripada Twitter. Jadi, ini telah mengehadkan bilangan *tweets* yang boleh digunakan sebagai set data dalam analisis sentimen.

Cadangan penambahbaikan kajian ini pada masa hadapan ialah mengumpul lebih banyak *tweets* untuk digunakan dalam analisis sentimen. Sejumlah besar data yang digunakan untuk melatih pengelas boleh meningkatkan prestasi model dengan menghasilkan keputusan analisis yang lebih tepat. Seterusnya, menggunakan teknik lain seperti penandaan *Part-of-speech* (POS) untuk mengetahui aspek terpenting dalam data yang dikumpul. Di samping itu, menggunakan data profil pengguna yang dikumpul untuk melakukan analisis bagi mendapatkan informasi tentang pengguna yang telah diberi pandangan.

Dengan ini, kajian ini diharapkan dapat memberi maklumat dan rujukan yang berguna mengenai pandangan masyarakat tempatan dan luar negara terhadap pembangunan industri pelancongan Malaysia pasca COVID-19 kepada mereka yang memerlukan dan akhirnya dapat membawa impak positif kepada negara.

7 RUJUKAN

- Unwto.org. 2021. Glossary of tourism terms | UNWTO. <https://www.unwto.org/glossary-tourism-terms> [13 Oktober 2021]
- Shashank Gupta. 2018. Sentiment Analysis: Concept, Analysis and Applications. <https://towardsdatascience.com/sentiment-analysis-concept-analysis-and-applications-6c94d6f58c17> [13 Oktober 2021]
- Inés Roldós, 2020. NLP, Machine Learning & AI, Explained. <https://monkeylearn.com/blog/nlp-ai/> [13 Oktober 2021]
- Bruno Stecanella. 2019. what-is-tf-idf. <https://monkeylearn.com/blog/what-is-tf-idf/> [13 Disember 2021]
- Google Developers. 2020. Classification: ROC Curve and AUC | Machine Learning Crash Course | Google Developers. <https://developers.google.com/machine-learning/crash-course/classification/roc-and-auc> [28 Mei 2022]
- iSixSigma. 2022. HYPOTHESIS TESTING. <https://www.isixsigma.com/dictionary/hypothesis-testing/> [28 Mei 2022]
- Brownlee, J., 2020. Hypothesis Test for Comparing Machine Learning Algorithms. Machine Learning Mastery. <https://machinelearningmastery.com/hypothesis-test-for-comparing-machine-learning-algorithms/> [28 Mei 2022]

Liew Set Teng (A176374)
Mohd Ridzwan Yaakub
Fakulti Teknologi & Sains Maklumat,
Universiti Kebangsaan Malaysia