

1. Pendahuluan

Latar Belakang

Berdasarkan World Health Organization (WHO), penyakit depresi adalah penyakit mental yang paling umum ditemukan di dunia. Pada tahun 2017, lebih dari 300 juta orang memiliki penyakit depresi dengan peningkatan lebih dari 18% dari tahun 2005 hingga 2015 [1]. Peningkatan ini diasosiasikan dengan peningkatan penyakit lain di dunia karena depresi menjadi penyebab utama disabilitas dan memiliki kontribusi besar terhadap keseluruhan penyakit global [2]. Selain itu, depresi juga diestimasi menjadi peringkat kedua sebagai penyebab utama disabilitas pada tahun 2020 [3].

Layanan pengobatan penyakit mental di dunia masih tergolong belum memadai [4]. Penderita penyakit mental di negara berkembang tercatat 76-85% tidak memiliki akses pengobatan yang tepat [4]. Faktanya, belum diketahui adanya uji laboratorium yang dapat diandalkan dalam melakukan diagnosis sebagian besar bentuk penyakit [5]. Hal ini menyebabkan upaya pencegahan lebih diutamakan.

Penelitian pada bidang kesehatan mental secara tradisional masih menggunakan survei, tes kepribadian dan wawancara akademik dalam pengumpulan data [6]. Penelitian di bidang ini masih kekurangan data kuantitatif yang tersedia karena adanya kompleksitas yang terdapat pada penyebab kesehatan mental dan stigma masyarakat yang masih memandang penyakit mental sebagai subjek yang tabu. Sebaliknya, media sosial telah banyak digunakan sebagai sumber data dalam banyak penelitian yang melibatkan hubungan antara penggunaan media sosial dan pola perilaku seperti stress ataupun depresi. Oleh karena itu pemanfaatan media sosial dapat menjadi alternatif baru dalam pencarian informasi mengenai kesehatan mental.

Pada penelitian ini dilakukan analisis sentimen yang bertujuan untuk mendeteksi perilaku depresi pada media sosial. Analisis sentimen atau *sentiment analysis* merupakan suatu pekerjaan yang mengidentifikasi apakah sebuah opini yang disampaikan termasuk kategori positif atau negatif pada sebuah dokumen. Dalam konteks media sosial, *sentiment analysis* berperan dalam menentukan polaritas pada bahasa yang diekspresikan dengan menekankan pada identifikasi kecenderungan perilaku depresi.

Beberapa penelitian serupa telah dilakukan tetapi sebagian besar menggunakan domain *microblog* khususnya twitter seperti model index depresi populasi [7] dan model klasifikasi pembeda depresi, PTSD dan non depresi [8]. Penelitian lain pada *microblog* cina juga dilakukan menggunakan *sentiment analysis* untuk deteksi depresi [6]. Banyak penelitian yang dilakukan pada domain *microblog* khususnya twitter disebabkan penggunaan bahasa yang mendekati kehidupan sehari-hari. Sebaliknya, penelitian pada domain forum belum banyak dilakukan. Forum tidak termasuk dalam *microblog* karena memungkinkan pengguna untuk berbagi informasi atau konten tanpa batasan karakter dan percakapan yang terjadi biasanya berpusat pada topik tertentu. Pada penelitian ini forum yang digunakan adalah Reddit. Reddit merupakan media sosial berorientasi agregasi berita, pemeringkatan konten, dan situs diskusi mengenai topik-topik tertentu. Reddit dipilih karena memiliki kategori berdasarkan topik yang berkaitan dengan kesehatan mental dan kemudahan dalam menemukan pengguna yang telah didiagnosis penyakit kesehatan mental sehingga label data yang diberikan lebih valid. Walaupun begitu, penelitian pada reddit memiliki tantangan karena penggunaan bahasa dan karakter linguistik yang berbeda.

Topik dan Batasannya

Pada penelitian ini penulis melakukan percobaan mengidentifikasi perilaku depresi pada sentimen pengguna media sosial Reddit. Analisis pola penggunaan bahasa pada data depresi dilakukan, akan tetapi fokus utama pada penelitian ini adalah untuk membangun klasifikasi teks yang paling tepat dalam mengidentifikasi teks yang sudah dilabeli depresi dan non-depresi. Proses klasifikasi yang tepat didapatkan dari pemilihan jenis *preprocessing* data terbaik, metode pemilihan fitur dan parameter pada metode klasifikasi yang digunakan. Jenis *preprocessing* data pada penelitian ini terbagi menjadi *stopword removal*, *stemming* dan *lemmatization*. Pada pemilihan fitur, metode yang digunakan adalah *Information Gain* (IG) dan *Categorical Proportional Difference* (CPD) dengan parameter fitur yang berbeda-beda. Parameter jumlah fitur terbaik digunakan oleh IG, sedangkan batas *threshold* digunakan oleh CPD. Sedangkan metode klasifikasi yang digunakan adalah *Multinomial Naïve Bayes* dengan pengaturan parameter *laplace smoothing* berdasarkan beberapa nilai yang telah ditentukan.

Beberapa batasan masalah yang terdapat pada penelitian ini adalah, terbatasnya jumlah data pada dataset sebanyak 659, terbagi menjadi 343 data depresi dan 316 data non-depresi. Pengambilan data depresi dilakukan pada media sosial Reddit dengan cara manual berdasarkan cara yang dianggap paling optimal dalam memperoleh data yang paling tepat, sedangkan data non-depresi yang diperoleh dengan crawling otomatis dengan *library python* berasal dari sub-reddit diluar forum depresi yang berisi banyak sentimen positif. Data non-depresi yang diperoleh dianggap belum cukup ideal, karena tidak termasuk dalam lingkungan yang sama dengan data depresi. Hal ini disebabkan sulitnya mendapatkan data non-depresi yang tidak memiliki karakter yang sama dengan data depresi pada lingkungan atau forum depresi. Batasan lain adalah *splitting* data yang digunakan hanya sekali dan tidak menggunakan *cross validation*.

Tujuan

Penelitian ini bertujuan untuk menganalisis pola penggunaan bahasa pada data depresi dan membuat model dengan performa terbaik yang dapat mengidentifikasi perilaku depresi pada teks media sosial. Penelitian diawali dengan melihat pola penggunaan kata pengganti orang pada dataset depresi. Berikutnya, eksperimen berupa perbandingan jenis *preprocessing*, *feature selection*, dan parameter metode klasifikasi dilakukan untuk mendapatkan performa model terbaik. Penelitian dimulai dengan analisis pengaruh *stopword removal* dan reduksi kata berdasarkan kata dasarnya berupa *stemming* dan *lemmatization*. Selanjutnya perbandingan *Information Gain* (IG) dan *Categorical Proportional Difference* (CPD) dilakukan untuk mendapatkan subset yang menghasilkan performa model terbaik. Terakhir, tuning parameter *smoothing* dalam algoritma klasifikasi dilakukan untuk mendapatkan performa model terbaik.