

BAB I

PENDAHULUAN

1.1 Latar Belakang

Kesehatan mental merupakan komponen krusial dari kesejahteraan individu yang menentukan kemampuan seseorang untuk berfungsi secara optimal dalam lingkungan sosial dan profesional [1]. Kondisi ini mencakup kesejahteraan emosional, psikologis, dan sosial yang memengaruhi cara seseorang berpikir, merasa, serta bertindak dalam menghadapi kehidupan sehari-hari. Di skala global, fenomena risiko gangguan mental tidak hanya terjadi pada populasi umum, tetapi juga tercermin secara signifikan pada sektor industri teknologi, di mana prevalensi masalah kesehatan mental menunjukkan angka yang mengkhawatirkan dan memerlukan intervensi proaktif [2]. Jika tidak dikelola dengan baik, tekanan ini dapat memicu gangguan psikologis yang serius.

Di lingkungan profesional, para pekerja sering kali menghadapi tekanan yang signifikan, mulai dari beban kerja yang tinggi hingga tuntutan adaptasi terhadap dinamika industri. Meskipun kesadaran akan kesehatan mental mulai meningkat, banyak individu masih merasa enggan untuk mendiskusikan kondisi mereka secara terbuka karena adanya stigma sosial [1]. Hal ini menyebabkan banyak kasus gangguan mental tidak terdeteksi secara dini. Oleh karena itu, diperlukan pendekatan teknologi yang mampu memberikan prediksi atau klasifikasi risiko kesehatan mental secara otomatis dan objektif guna mendukung langkah diagnosis awal oleh tenaga profesional [3].

Pemanfaatan teknik Artificial Intelligence (AI) dan Machine Learning (ML) telah menunjukkan potensi besar dalam mentransformasi prediksi risiko kesehatan mental menggunakan data perilaku dan lingkungan [2]. Integrasi data survei multi-tahun, seperti yang dilakukan dalam pembersihan dan penggabungan dataset OSMI 2017-2021, terbukti krusial dalam menyediakan data yang valid dan konsisten untuk analisis perilaku kesehatan mental bagi profesional maupun pelajar [4]. Penelitian sebelumnya telah mengevaluasi berbagai algoritma klasifikasi, seperti

Logistic Regression, Support Vector Machine (SVM), dan Neural Networks, untuk memprediksi gangguan mental dengan tingkat akurasi yang bervariasi [3]. Selain itu, penggunaan metode ensemble seperti Gradient Boosting terbukti mampu memberikan performa yang sangat baik dalam menangani dataset kesehatan mental [3], [5].

Namun, terdapat beberapa penelitian yang menjadi dasar urgensi penelitian ini. Pertama, sebagian besar penelitian terdahulu seringkali berfokus pada penggunaan algoritma tunggal atau hanya membandingkan dua hingga tiga model, sehingga belum memberikan gambaran komprehensif mengenai algoritma mana yang paling efisien untuk data kategorikal survei kesehatan mental yang memiliki karakteristik unik. Kedua, masalah ketidakseimbangan kelas (*class imbalance*) pada dataset survei seringkali diabaikan, yang dapat mengakibatkan model memiliki akurasi tinggi namun gagal dalam mendeteksi kelas minoritas (individu yang sebenarnya membutuhkan bantuan). Ketiga, meskipun algoritma sederhana seperti Naive Bayes sering dijadikan baseline, efektivitasnya perlu diuji kembali ketika disandingkan dengan algoritma ensemble modern seperti XGBoost dan AdaBoost dalam konteks data kesehatan mental pekerja yang melibatkan banyak fitur kategorikal [7].

Berdasarkan permasalahan tersebut, penelitian ini bertujuan untuk membandingkan kinerja dari enam algoritma utama, yaitu Support Vector Machine (SVM), Random Forest, Naive Bayes, XGBoost, Logistic Regression, dan AdaBoost, dalam mengklasifikasikan tingkat kesehatan mental pekerja. Dengan menerapkan alur kerja Knowledge Discovery in Databases (KDD) yang sistematis, penelitian ini bertujuan untuk menemukan model klasifikasi yang paling reliabel dan akurat. Menerapkan teknik pra-pemrosesan data yang tepat dan integrasi teknik Synthetic Minority Over-sampling Technique (SMOTE) untuk menangani ketidakseimbangan data, diharapkan penelitian ini dapat menemukan model yang paling reliabel. Hasil penelitian ini diharapkan dapat menjadi sistem pendukung keputusan bagi organisasi dalam melakukan intervensi dini terhadap risiko kesehatan mental pekerja secara lebih akurat dan objektif.

1.2 Rumusan Masalah

Berdasarkan latar belakang yang telah diuraikan, rumusan masalah dalam penelitian ini dinyatakan sebagai berikut:

1. Bagaimana cara mengatasi tingginya tingkat tekanan kerja yang berpotensi memicu gangguan kesehatan mental namun sulit terdeteksi dini akibat stigma sosial di organisasi?
2. Bagaimana efektivitas performa algoritma Machine Learning seperti Support Vector Machine (SVM), Naive Bayes, Logistic Regression, Random Forest, AdaBoost, dan XGBoost dengan model klasifikasi dalam meningkatkan akurasi serta reliabilitas prediksi kesehatan mental pekerja?
3. Bagaimana pengaruh penerapan teknik Synthetic Minority Over-sampling Technique (SMOTE) dalam menangani kendala ketidakseimbangan data (class imbalance) untuk meminimalisir bias model dan meningkatkan reliabilitas prediksi?
4. Algoritma klasifikasi manakah yang paling optimal untuk diimplementasikan sebagai alat deteksi dini terhadap risiko gangguan kesehatan mental di lingkungan profesional?

1.3 Batasan Masalah

Batasan masalah dalam penelitian ini ditetapkan untuk memfokuskan ruang lingkup penelitian, yaitu sebagai berikut:

1. Penelitian ini menggunakan dataset sekunder "*Mental Health in Tech Survey*" yang bersumber dari Kaggle, dengan fokus pada responden yang bekerja di perusahaan teknologi.
2. Target klasifikasi dibatasi pada variabel treatment, yaitu memprediksi apakah seorang individu membutuhkan penanganan kesehatan mental atau tidak.
3. Metode klasifikasi yang digunakan meliputi Random Forest, Support Vector Machine, Naive Bayes, Ada Boost, Logistic Regression, dan

XGBoost.

4. Batasan pengolahan data mencakup pembersihan data (data cleaning), penanganan nilai yang hilang (missing values), Label Encoding untuk variabel kategorikal, dan teknik Resampling untuk menangani ketidakseimbangan data (imbalanced data).
5. Evaluasi hasil klasifikasi dibatasi pada penggunaan Confusion Matrix, tingkat akurasi (Accuracy), Precision, Recall, dan F1-Score.
6. Seluruh proses komputasi, mulai dari eksplorasi data hingga pengujian model, dilakukan menggunakan bahasa pemrograman Python versi 3.x di lingkungan Google Colab.

1.4 Tujuan Penelitian

Tujuan utama dari penelitian ini adalah untuk merancang dan mengimplementasikan model klasifikasi yang optimal guna mendeteksi tingkat kesehatan mental pekerja di lingkungan profesional. Secara terperinci, tujuan penelitian ini adalah sebagai berikut:

1. Merancang sistem klasifikasi berbasis Machine Learning yang mampu mengidentifikasi risiko gangguan kesehatan mental secara objektif untuk meminimalisir keterlambatan penanganan akibat stigma sosial di organisasi.
2. Mengevaluasi dan membandingkan efektivitas kinerja algoritma Machine Learning yang mencakup model konvensional (Support Vector Machine, Naive Bayes, Logistic Regression) dan model ensemble learning (Random Forest, AdaBoost, XGBoost) dalam memprediksi tingkat kesehatan mental pekerja.
3. Menganalisis pengaruh penerapan Synthetic Minority Over-sampling Technique (SMOTE) dalam menangani kendala ketidakseimbangan data (class imbalance), guna meminimalisir bias pada model klasifikasi serta meningkatkan reliabilitas prediksi pada kelas minoritas.
4. Mengidentifikasi algoritma klasifikasi yang menunjukkan kinerja paling superior berdasarkan metrik evaluasi (Accuracy, Recall, Precision, F1-

Score, dan ROC-AUC), untuk kemudian direkomendasikan sebagai alat bantu keputusan (decision support system) dalam upaya intervensi kesehatan mental yang proaktif di lingkungan profesional.

1.5 Manfaat Penelitian

Manfaat dari penelitian ini diharapkan dapat memberikan kontribusi signifikan baik secara teoritis maupun praktis dalam bidang data science dan kesehatan mental. Secara teoritis, penelitian ini memperkaya literatur mengenai efektivitas integrasi metode SMOTE untuk menangani ketidakseimbangan data serta optimasi hyperparameter pada algoritma Naive Bayes dan XGBoost untuk kasus klasifikasi data psikologis yang kompleks. Secara praktis, penelitian ini memberikan manfaat bagi organisasi maupun perusahaan sebagai instrumen deteksi dini yang reliabel dan terukur untuk memantau kondisi kesehatan mental karyawan secara objektif, sehingga kebijakan intervensi dapat dirumuskan secara tepat sasaran guna memitigasi risiko gangguan psikologis di lingkungan produktif. Selain itu, bagi peneliti selanjutnya, hasil analisis komparatif yang mencakup evaluasi mendalam melalui metrik akurasi, presisi, recall, hingga ROC Curve dan AUC Score ini dapat dijadikan acuan metodologis dalam memilih teknik preprocessing termasuk Label Encoding dan MinMaxScaler serta algoritma yang paling optimal untuk mengolah dataset survei dengan karakteristik data yang tidak seimbang.

1.6 Sistematika Penulisan

Sistematika penulisan skripsi ini disusun sebagai berikut:

BAB I PENDAHULUAN, berisi tentang latar belakang masalah yang menggambarkan kerangka pada penelitian yang akan dilakukan, diikuti dengan rumusan masalah, tujuan penelitian, manfaat penelitian, dan sistematika penelitian.

BAB II TINJAUAN PUSTAKA, berisi tentang tinjauan pustaka yang terdiri dari studi literatur yang membahas penelitian-penelitian terdahulu yang bersumber dari jurnal-jurnal untuk membandingkan penelitian yang akan dilakukan, kemudian dasar teori yang membahas teori-teori yang relevan terhadap alur penelitian sehingga dapat memperkuat penelitian.

BAB III METODE PENELITIAN, berisi tentang alur penelitian secara teknis mulai dari pengumpulan dataset survei hingga tahapan evaluasi model. Penjelasan difokuskan pada prosedur pra-pemrosesan data yang meliputi pembersihan data (data cleaning), transformasi fitur melalui label encoding, dan seleksi enam fitur utama yang paling berpengaruh. Bab ini juga merinci strategi pemodelan menggunakan pembagian data 90:10 serta penggunaan metrik evaluasi seperti Accuracy, Recall, dan ROC-AUC untuk mengukur ketangguhan algoritma.

BAB IV HASIL DAN PEMBAHASAN, menyajikan fase inti penelitian yang berisi analisis hasil pembersihan data dan implementasi enam algoritma klasifikasi yang telah dirancang. Pembahasan berfokus pada evaluasi kinerja model melalui tabel komparasi, kurva ROC, dan confusion matrix untuk mengidentifikasi algoritma dengan tingkat presisi dan sensitivitas tertinggi. Di bagian akhir bab, dilakukan sintesis hasil yang membandingkan capaian akurasi penelitian ini terhadap beberapa penelitian terdahulu guna memvalidasi kontribusi ilmiah dari model yang diusulkan.

BAB V PENUTUP, berisi kesimpulan yang diperoleh dari penelitian serta saran untuk penelitian selanjutnya.