

BAB I

PENDAHULUAN

1.1 Latar belakang masalah

Pada era informasi sekarang ini, media penyimpanan dan teknologi retrieval di internet semakin banyak digunakan, sehingga data yang disimpan bertambah dari jumlah maupun dimensi [1]. Pada umumnya, data yang dimiliki sering kali mengandung *noise*, *missing value*, *redundant* dan *irrelevant* feature yang bisa membuat pengolahan data kurang efisien. Salah satu solusi yang bisa ditawarkan untuk menghasilkan data yang berkualitas adalah dengan memilih *feature-feature* penting saja untuk diproses.

Feature subset selection (FSS) merupakan masalah umum yang terdapat pada *data mining* mengenai bagaimana memilih features yang relevan dan menghilangkan features yang redundan (berulang) dari feature set yang asli [2]. Semakin kecil ukuran feature subset, makin optimal feature subset yang didapat. Jika feature dipilih secara tepat maka dapat meningkatkan akurasi dan performansi dari proses klasifikasi [3]. Karena itu, *feature selection* menjadi sangat perlu dilakukan untuk mendukung *preprocessing* untuk *machine learning*.

Feature selection untuk preprocess data sudah cukup banyak diterapkan. Feature selection bisa dibagi dalam tiga jenis, yaitu *filter model*, *wrapper model* dan *embedded model*. Contoh-contoh penerapan model filter antara lain Mutual Information [6] dan Relief algoritma [5]. Mutual information (MI) yang menggunakan pendekatan information measure sangat efektif untuk mengevaluasi relevansi dari setiap variabel input namun gagal untuk mengeliminasi feature yang redundan, begitu pula pada relief algoritma yang menggunakan pembobotan untuk melihat feature mana yang tidak relevan namun tidak mempertimbangkan redundancy data.

Salah satu teknik *feature subset selection* yang dapat digunakan adalah *Mutual Information and Redundancy-Synergy Coefficient*. Metode ini merupakan salah satu model berpendekatan *information measure* yang merupakan pengembangan dari metode mutual information, namun sudah dilengkapi dengan cara mengurangi feature yang redundan. Pada penerapan mutual information sebelumnya [6], memungkinkan untuk menghilangkan feature yang tidak relevan tapi tidak bisa mengurangi feature yang redundan.

Terkadang feature redundan dihilangkan namun menghilangkan feature yang sebenarnya relevan [3]. Pada metode ini akan dihitung feature subset mana yang memiliki nilai informasi yang tinggi. Setelah itu diterapkan Redundancy-Synergy coefficient untuk mengukur tingkat *redundancy* data tanpa mengurangi tingkat relevansi data. Jika pada paper sebelumnya [2] yang diamati lebih menekankan pada kecepatan proses *preprocessing* dan tingkat akurasi umum, maka pada tugas akhir ini lebih diamati hasil tingkat ketepatan hasil klasifikasi masing-masing nilai kelas (diukur dengan precision dan recall) serta waktu pembentukan model saat dilakukan klasifikasi dengan algoritma klasifikasi

Naives Bayes. Algoritma klasifikasi tersebut dipilih karena menggunakan keseluruhan dataset saat pemrosesan dan tidak terdapat preprocess sebelumnya. Selain itu algoritma klasifikasi Naives bayes juga merupakan algoritma klasifikasi yang cepat dalam prosesnya.

Hipotesis awal dari pengerjaan tugas akhir ini adalah, dengan melakukan *preprocessing* menggunakan metode mutual information dan *redundancy-synergy coefficient* dapat meningkatkan nilai *precision* dan *recall* serta waktu pembentukan model saat dilakukan proses klasifikasi.

1.2 Perumusan masalah

Berdasarkan latar belakang yang telah diuraikan di atas, maka dapat dirumuskan beberapa masalah antara lain:

1. Bagaimana melakukan reduksi dimensi data menggunakan *Mutual Information and Redundancy-Synergy coefficient* agar pengolahan data untuk proses selanjutnya lebih baik?
2. Bagaimana performansi reduksi data setelah dilakukan *feature selection* dengan *Mutual Information and Redundancy-Synergy coefficient* dinilai dari parameter uji seperti *precision* dan *recall* serta waktu pembentukan model saat menggunakan algoritma klasifikasi?

1.3 Batasan Masalah

Adapun beberapa batasan masalah dalam pengerjaan TA ini, antara lain:

1. Dataset yang digunakan adalah data dengan tipe atribut diskrit
2. Untuk analisis hasil, parameter yang digunakan untuk mengevaluasi hasil adalah *precision* dan *recall* serta waktu pembentukan model saat hasil *feature selection* diterapkan pada algoritma klasifikasi.
3. Tidak menangani *data cleaning (missing value, noisy data)*, *data integration* dan *data transformation*.
4. Evaluasi dataset hasil reduksi menggunakan algoritma klasifikasi naives bayes

1.4 Tujuan

Tujuan dari pengerjaan tugas akhir ini antara lain adalah:

1. Mereduksi dimensi data menggunakan metode *Mutual Information and Redundancy-Synergy coefficient* sehingga akan lebih efisien saat dilakukan pemrosesan selanjutnya, yaitu proses klasifikasi dengan *naives bayes*
2. Menganalisis performansi hasil preprocessing menggunakan *Mutual Information and Redundancy-Synergy coefficient* menggunakan algoritma klasifikasi dengan parameter-parameter pengukuran *precision*, *recall* dan waktu pembentukan model saat dilakukan proses klasifikasi.

1.5 Metodologi penyelesaian masalah

Dalam pengerjaan tugas akhir ini, digunakan metodologi penyelesaian masalah sebagai berikut:

1. Identifikasi Masalah

Permasalahan yang muncul dan melatarbelakangi dibuatnya tugas akhir ini antara lain:

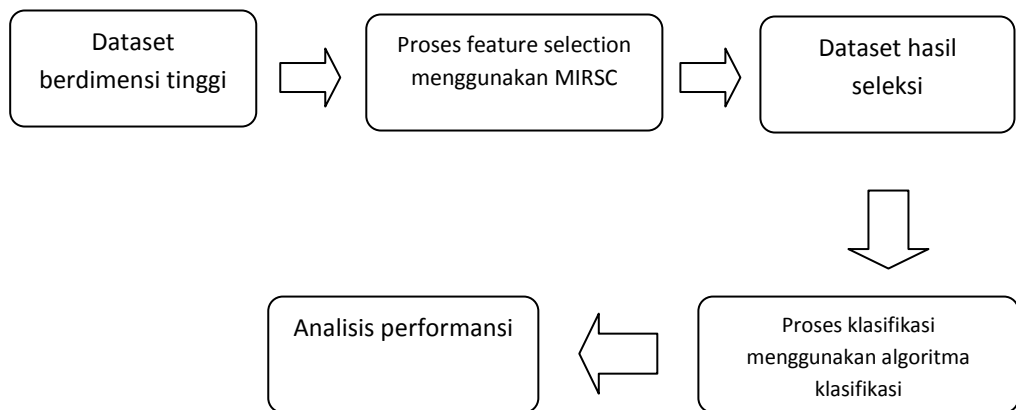
- ❖ Bagaimana memilih feature subset yang minimal agar dataset yang berdimensi besar bisa lebih efisien ketika diproses di task klasifikasi.
- ❖ Bagaimana menghilangkan feature redundan dan irrelevant pada sebuah dataset berdimensi besar, namun tetap memiliki nilai informasi yang sama dengan dataset asli

2. Study literatur

Melakukan studi literatur yang berhubungan dengan feature selection, klasifikasi, mutual information dan materi-materi yang berhubungan dengan pengerjaan TA melalui sumber-sumber literatur maupun sumber online.

3. Desain Penelitian

Pada pengerjaan tugas akhir ini, alur pengerjaan yang akan dilakukan adalah sebagai berikut:



Gambar 1.1 Alur pengerjaan

Keterangan:

1. Dipilih sebuah dataset berdimensi tinggi untuk dilakukan preprocessing

2. Kemudian dataset diproses menggunakan metode mutual information dan redundancy-synergy coefficient agar dapat menghilangkan data yang berulang dan tidak relevan
3. Kemudian dari proses tersebut didapat dataset hasil reduksi
4. Dataset hasil reduksi diproses menggunakan algoritma klasifikasi, yaitu naives bayes
5. Setelah itu, hasil klasifikasi akan dibandingkan dengan hasil pemrosesan dataset sebelum di preproses agar dapat dilihat perbandingan *precision*, *recall* serta waktu pembentukan model saat dilakukan klasifikasi.

Aplikasi diimplementasikan menggunakan Matlab, proses klasifikasi akan diuji menggunakan Weka sebagai alat bantu. Algoritma klasifikasi yang digunakan untuk pengujian adalah Naives Bayes

4. Analisis

Hipotesis awal dari pengerjaan Tugas akhir ini adalah, dengan melakukan *preprocessing* menggunakan metode mutual information dan redundancy-synergy coefficient dapat meningkatkan *precision*, *recall* dan waktu pembentukan model saat dilakukan proses klasifikasi.

1.5 Sistematika Penulisan

Penulisan tugas akhir ini dibagi ke dalam 5 bab, antara lain :

1. **Bab 1 Pendahuluan**
Berisi latar belakang masalah, perumusan masalah, batasan masalah, tujuan masalah, metodologi penyelesaian masalah dan sistematika penulisan
2. **Bab 2 Landasan Teori**
Berisi landasan teori tentang KDD, metode *mutual information* dan *redundancy-synergy coefficient* dan klasifikasi.
3. **Bab 3 Perancangan Aplikasi**
Berisi analisa dan perancangan aplikasi yang akan dibuat dengan bahasa permodelan DFD
4. **Bab 4 Analisa Hasil**
Berisi penjelasan mengenai implementasi hasil perancangan, uji coba terhadap sistem, dan analisa perangkat lunak yang dibangun.
5. **Bab 5 Kesimpulan dan Saran**
Berisi kesimpulan dan saran dari hasil penelitian untuk pengembangan lebih lanjut.