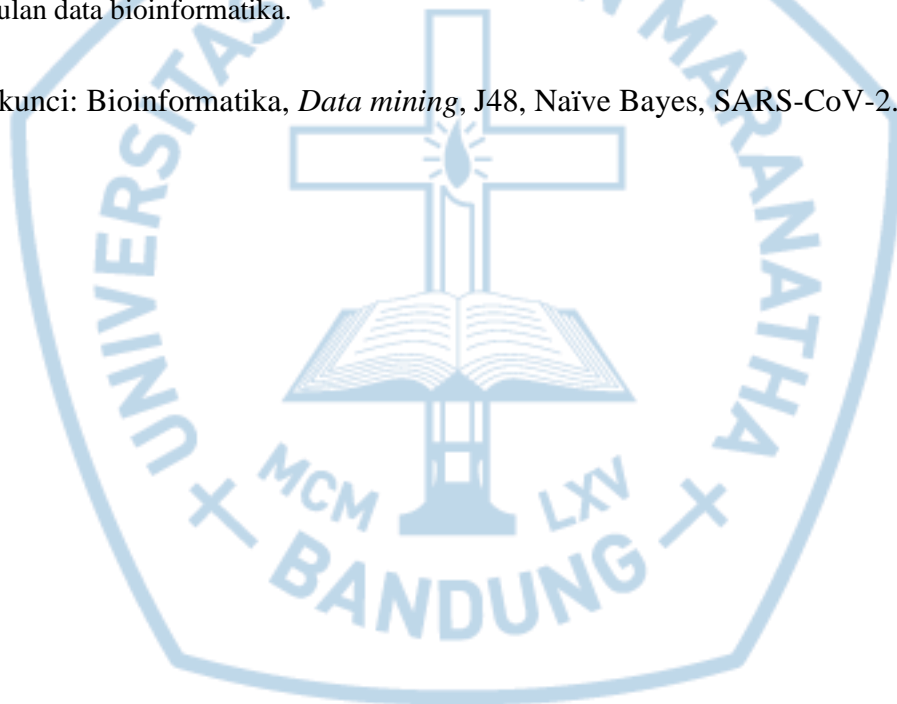


## ABSTRAK

Laporan yang berjudul “Analisis Akurasi Data Protein Virus SARS-CoV-2 dengan menggunakan Metode *Data mining*” ini bertujuan untuk menganalisis hasil penelitian *data mining* pada kumpulan data protein SARS-CoV-2 antara algoritma Naïve Bayes dan J48 dengan pelatihan *10-fold cross-validation*. Proses penelitian ini dimulai dengan pengumpulan data yang disediakan oleh *website* NCBI berupa protein virus SARS-CoV-2. Selanjutnya pada tahap praproses data, kumpulan data tersebut dibersihkan, data yang akan digunakan dari jumlah kumpulan data 1.149.216 menjadi 517.834. Dengan aplikasi WEKA, *data mining* algoritma Naïve Bayes dan J48 akan diterapkan pada kumpulan data yang akan menghasilkan model klasifikasi. Selanjutnya, dengan 500 data uji dari *website* NCBI, model klasifikasi akan memprediksikan kelas protein pada data uji. Pada tahap analisis hasil data, akan dibandingkan kinerja masing-masing algoritma dari hasil detail akurasi juga waktu pembuatan model dan pengujian data. Dapat disimpulkan dengan menggunakan algoritma J48 lebih unggul akurasinya dibanding Naïve Bayes dalam melakukan klasifikasi pada kumpulan data protein SARS-CoV-2. Laporan ini dibuat sebagai pengenalan awal bioinformatika dan proses melakukan *data mining* pada kumpulan data bioinformatika.

Kata kunci: Bioinformatika, *Data mining*, J48, Naïve Bayes, SARS-CoV-2.



## ABSTRACT

*The report entitled “Analisis Akurasi Data Protein Virus SARS-CoV-2 dengan menggunakan Metode Data mining” aims to analyze the result of data mining research on the SARS-CoV-2 protein data set between the Naïve Bayes and J48 algorithms with 10-fold cross-validation training. The research process begins with the collection of data provided by the NCBI website about proteins SARS-CoV-2 virus. At the preprocessing stage, the data set is cleaned, from total data set of 1,149,216 to 517,834. With the WEKA application, the Naïve Bayes and J48 data mining algorithms will be applied to a data set that will produce a classification model. With 500 test data from the NCBI website, the classification model will predict the protein class. At the data analysis stage, the performance of each algorithm will be compared from the detailed results of the accuracy as well as the time for making the model and testing the data. From this research, it can be concluded that using J48 algorithm is more accurate than Naïve Bayes in classifying protein SARS-CoV-2 virus data set. This report was created as an introduction to bioinformatics and the process of data mining on bioinformatics data sets.*

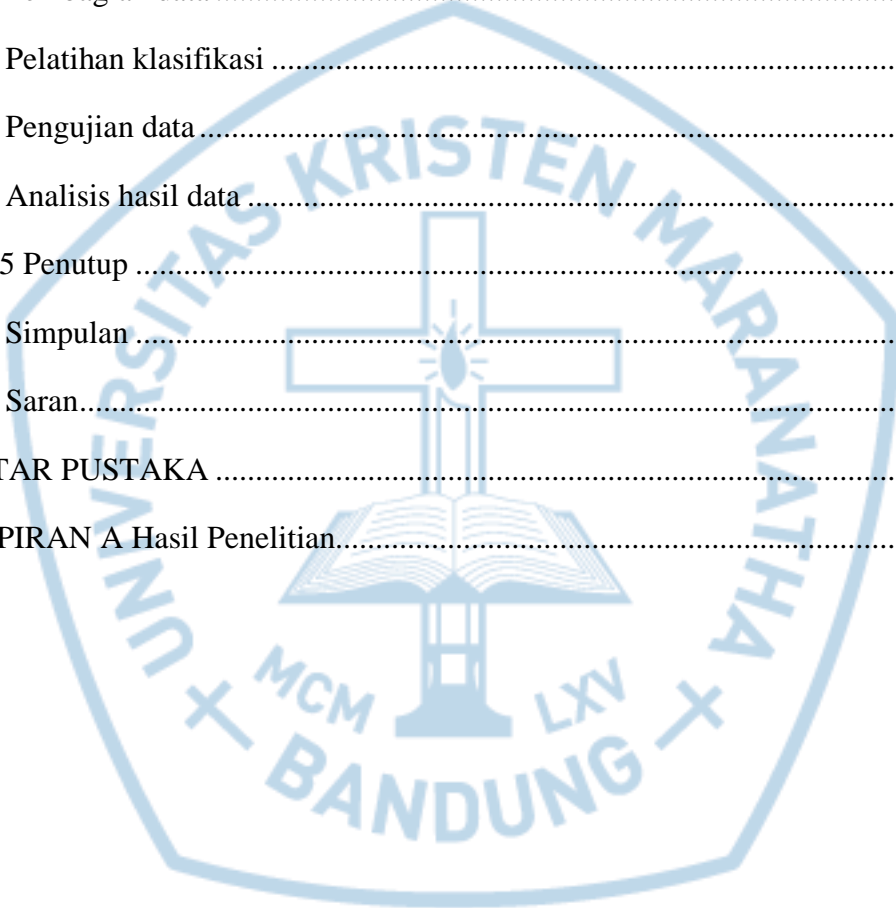
*Keywords: Bioinformatics, Data mining, J48, Naïve Bayes, SARS-CoV-2.*



## DAFTAR ISI

LEMBAR PENGESAHAN .....	i
PERNYATAAN ORISINALITAS LAPORAN PENELITIAN .....	ii
PERNYATAAN PUBLIKASI LAPORAN PENELITIAN.....	iii
PRAKATA.....	iv
ABSTRAK .....	v
ABSTRACT .....	vi
DAFTAR ISI.....	vii
DAFTAR GAMBAR .....	ix
DAFTAR TABEL.....	x
DAFTAR SINGKATAN .....	xi
DAFTAR ISTILAH .....	xii
BAB 1 PENDAHULUAN .....	1
1.1 Latar Belakang .....	1
1.2 Rumusan Masalah .....	2
1.3 Tujuan Pembahasan .....	2
1.4 Ruang Lingkup.....	2
1.5 Sumber Data.....	3
1.6 Sistematika Penyajian .....	3
BAB 2 KAJIAN TEORI .....	4
2.1 Bioinformatika .....	4
BAB 3 METODOLOGI PENELITIAN.....	10
3.1 Pengumpulan data .....	11
3.2 Praproses data.....	11
3.3 Pembagian data .....	11

3.4 Pelatihan klasifikasi .....	11
3.5 Pengujian data .....	11
3.6 Analisis hasil data .....	11
BAB 4 Hasil dan pembahasan.....	12
4.1 Pengumpulan data .....	12
4.2 Praproses data.....	13
4.3 Pembagian data .....	15
4.4 Pelatihan klasifikasi .....	16
4.5 Pengujian data .....	19
4.6 Analisis hasil data .....	21
BAB 5 Penutup .....	23
5.1 Simpulan .....	23
5.2 Saran.....	23
DAFTAR PUSTAKA .....	24
LAMPIRAN A Hasil Penelitian.....	A-1



## DAFTAR GAMBAR

Gambar 3.1 Tahapan Penelitian .....	10
Gambar 4.1 Pengumpulan data tahap 1.....	12
Gambar 4.2 Pengumpulan data tahap 2.....	12
Gambar 4.3 Pengumpulan data tahap 3.....	13
Gambar 4.4 Sampel Kumpulan Data Protein SARS-CoV-2 Sebelum Dibersihkan .....	14
Gambar 4.5 Tabel Kumpulan Data SARS-CoV-2 Setelah Dibersihkan.....	14
Gambar 4.6 Konversi Data Latih Melalui WEKA.....	15
Gambar 4.7 Hasil <i>Summary</i> Klasifikasi Naïve Bayes.....	16
Gambar 4.8 Hasil <i>Summary</i> Klasifikasi J48 .....	16
Gambar 4.9 Hasil Detail Akurasi Berdasarkan Kelas Klasifikasi Naïve Bayes ...	17
Gambar 4.10 Hasil Detail Akurasi Berdasarkan Kelas Klasifikasi J48.....	17
Gambar 4.11 Hasil Confusion Matrix Klasifikasi Naïve Bayes .....	19
Gambar 4.12 Hasil Confusion Matrix Klasifikasi J48 .....	19
Gambar 4.13 Mengganti Kelas Prediksi .....	20
Gambar 4.14 Atribut Data Uji yang Harus Diganti .....	20
Gambar 4.15 Memuat Data Uji.....	20
Gambar 4.16 Memuat Model Klasifikasi dan Pilihan Proses untuk Prediksi .....	21
Gambar 4.17 Sampel Hasil Prediksi .....	22

## DAFTAR TABEL

Tabel 4.1 Perbandingan Akurasi Algoritma Naïve Bayes dan J48 ..... 21



## DAFTAR SINGKATAN

<i>ARFF</i>	<i>Attribute Relation File Format</i>
<i>CART</i>	<i>Classification and Regression Trees</i>
<i>CBIIT</i>	<i>Center for Biomedical Informatics &amp; Information Technology</i>
<i>CSV</i>	<i>Comma-Separated Value</i>
<i>DDBJ</i>	<i>DNA Database of Japan</i>
<i>DNA</i>	<i>Deoxyribonucleic Acid</i>
<i>EMBL</i>	<i>European Molecular Biology Laboratory</i>
<i>Envelope protein</i>	<i>The CoV E protein</i>
<i>ID3</i>	<i>Iterative Dichotomiser 3</i>
<i>KDD</i>	<i>Knowledge Discovery in Database</i>
<i>KNN</i>	<i>K-Nearest Neighbour</i>
<i>Membrane glycoproteins</i>	<i>membrane proteins</i>
<i>NCBI</i>	<i>National Center for Biotechnology Information</i>
<i>Nucleocapsid phosphoprotein</i>	<i>The multifunctional nucleocapsid protein</i>
<i>ORF10 protein</i>	<i>Open Reading Frame 10 protein</i>
<i>ORF1a polyprotein</i>	<i>Open Reading Frame 1a polyprotein</i>
<i>ORF1ab polyprotein</i>	<i>Open Reading Frame 1ab polyprotein</i>
<i>ORF3a protein</i>	<i>Open Reading Frame 3a protein</i>
<i>ORF6 protein</i>	<i>Open reading Frame 6 protein</i>
<i>ORF7b protein</i>	<i>Open Reading Frame 7 protein</i>
<i>ORF8 protein</i>	<i>Open Reading Frame 8 protein</i>
<i>ORF10 protein</i>	<i>Open Reading Frame 10 protein</i>

## DAFTAR ISTILAH

Bioinformatika	berupa konsep dari biologi yang memandang segi molekul dan menerapkan teknik informatika untuk memahami informasi yang berkaitan dengan molekul dalam skala besar.
<i>Data mining</i>	merupakan metode untuk menggali suatu pola dari kumpulan data yang besar dan kompleks dengan menggunakan algoritma.
DNA	molekul kompleks yang berisi informasi yang diperlukan untuk membangun dan memelihara organisme.
J48	singkatan dari weka.classifiers.trees.j48 yang menggunakan algoritma Quinlan yaitu C4.5.
K-fold Cross-Validation	salah satu pelatihan klasifikasi di WEKA.
Naïve Bayes	klasifikasi yang menekankan kepada bidang statistika yaitu probabilitas.
RNA	molekul yang mempunyai peran dalam pengkodean, pengurai, regulasi dan ekspresi gen.

