

## Turnitin Originality Report

Processed on: 09-Apr-2021 2:30 AM WIB  
 ID: 1553916940  
 Word Count: 1839  
 Submitted: 1

Similarity Index

20%

### Similarity by Source

Internet Sources: 19%  
 Publications: 8%  
 Student Papers: 9%

Implementasi Integrasi Seleksi Data dan  
 Extreme Learning Machine (IDELM) Untuk  
 Klasifikasi DNA Bakteri Patogen By Umi  
 Mahdiyah

1% match (Internet from 12-Oct-2019)

<http://simanis.uin-malang.ac.id/wp-content/uploads/2019/08/Template-SI-MaNI-2019.docx>

1% match (student papers from 02-Jan-2021)

[Submitted to Universitas Gunadarma on 2021-01-02](#)

1% match (Internet from 28-Oct-2016)

<http://english.rejbrand.se/rejbrand/books.asp>

1% match (Internet from 15-Nov-2020)

[http://distan.jogjapro.go.id/wp-content/download/publikasi/review\\_renstra\\_dinas\\_pertanian\\_2012\\_2017.pdf](http://distan.jogjapro.go.id/wp-content/download/publikasi/review_renstra_dinas_pertanian_2012_2017.pdf)

1% match (Internet from 08-Dec-2020)

[http://eprints.undip.ac.id/68473/1/PENUNTUN\\_PRAKTIKUM\\_REKAYASA\\_GENETIKA\\_S1\\_BERBASIS\\_KOMPETENSI.pdf](http://eprints.undip.ac.id/68473/1/PENUNTUN_PRAKTIKUM_REKAYASA_GENETIKA_S1_BERBASIS_KOMPETENSI.pdf)

1% match (Internet from 20-Nov-2020)

<https://artikataku.blogspot.com/2016/08/arti-kata-informasi-dari-contoh-kalimat.html>

1% match (Internet from 18-Dec-2019)

<https://es.scribd.com/document/367272725/Ring-Kasan>

1% match (publications)

[Elvismary Molina de Armas, Maristela Holanda, Daniel de Oliveira, Nalvo F. Almeida, Sérgio Lifschitz. "Chapter 1 A Classification of de Bruijn Graph Approaches for De Novo Fragment Assembly", Springer Science and Business Media LLC, 2020](#)

1% match (student papers from 08-Jan-2018)

[Submitted to Universitas Brawijaya on 2018-01-08](#)

1% match (Internet from 15-Mar-2020)

<http://conferences.uin-malang.ac.id/index.php/SIMANIS/article/cite/64/BibtexCitationPlugin>

1% match (Internet from 11-Nov-2020)

<https://idoc.pub/documents/4-biologi-2007-2878-d47e9woo3yn2>

1% match (Internet from 05-Mar-2021)

<https://iptek.its.ac.id/index.php/limits/article/view/3079>

1% match (Internet from 07-Mar-2021)

<https://firlyamalia.wordpress.com/author/firlyamalia/page/2/>

1% match (student papers from 22-Aug-2013)

[Submitted to iGroup on 2013-08-22](#)

1% match (Internet from 20-Oct-2020)

<https://sinta.ristekbrin.go.id/authors/detail?id=6030679&view=documentsgs>

1% match (Internet from 06-Dec-2014)

<http://tampub.uta.fi/bitstream/handle/10024/66914/978-951-44-8825-2.pdf?sequence=1>

1% match ( )

[Widiarti, Puji. "Perbandingan metode regresi logistik biner dan classification and regression trees \(CART\) untuk klasifikasi diagnosa penyakit diabetes mellitus \(DM\)", 2020](#)

1% match (Internet from 11-Jan-2021)

<http://download.garuda.ristekdikti.go.id/article.php?article=1683255&title=PENANGANAN+PELAKSANAAN+BONGKAR+MUAT+VCM+C2H3CL+DI+KAPAL+MT+GAS+KALIMANTAN+LPG+CARRIER+TYPIC&val=18285>

1% match (Internet from 06-Apr-2021)

<https://ejournal3.undip.ac.id/index.php/geodesi/article/download/19330/18338>

1% match (student papers from 29-Apr-2019)

[Submitted to UIN Sultan Syarif Kasim Riau on 2019-04-29](#)

1% match (student papers from 29-Mar-2021)

[Submitted to Universitas Sumatera Utara on 2021-03-29](#)

1% match (Internet from 02-May-2020)

<http://ipi.portalgaruda.org/index.php?id=250309&ipp=5&mod=profile&page=7&ref=author>

1% match ( )

[Kusumaningrum, Ageng Pramesthi. "Optimasi Parameter Support Vector Machine menggunakan Genetic Algorithm untuk Klasifikasi Microarray Data", 2017](#)

1% match (Internet from 07-Dec-2018)

<http://unsri.portalgaruda.org/?id=591985&journal=9626&mod=profile&ref=author>

1% match (Internet from 25-Sep-2020)

<https://worldwidescience.org/topicpages/l/linear+learning+machines.html>

Prosiding [Seminar Nasional Integrasi Matematika dan Nilai Islami Vol.3, No.1](#), September 2019, pp. xx~xx p-ISSN: 2580-4596; e-ISSN: 2580-460X [Implementasi Integrasi Seleksi Data dan Extreme Learning Machine \(IDELM\) Untuk Klasifikasi DNA Bakteri Patogen](#) Umi Mahdiyah1, Lilia Sinta Wahyuniar2 1,2Jurusan Teknik Informatika, Universitas Nusantara PGRI Kediri umimahdiyah@gmail.com, li2asint@gmail.com Info Artikel ABSTRAK Riwayat Artikel: PCR adalah salah satu metode untuk mendeteksi keberadaan mikroba dalam tubuh. Dibanding metode lain, metode ini tergolong akurat, Diterma: cepat, dan bisa diandalkan. Keuntungan PCR lainnya adalah sekuens Direvisi: DNA dari mikroba atau strain infeksi yang baru ditemukan. Diterbitkan: Selanjutnya data DNA yang ditemukan tersebut dilakukan pencarian masuk pada jenis DNA apa, sehingga dibutuhkan metode yang Kata Kunci: optimal. Sekuensing DNA dapat dimanfaatkan [untuk menentukan identitas](#) maupun [fungsi gen atau fragmen DNA](#) lainnya dengan cara Pertama membandingkan sekuens-nya dengan sekuens DNA lain yang sudah Kedua diketahui. [Integrasi Seleksi data dan Extreme Learning Machine ini](#) Ketiga dipilih sebagai metode untuk klasifikasi DNA karena data DNA Keempat merupakan data yang besar serta karakteristik datanya yang Kelima kebanyakan adalah data yang imbalance. Pada proses penelitian data yang akan diolah terlebih dahulu diuraikan fragmennya, selanjutnya dilakukan ekstraksi fitur, kemudian dilakukan proses klasifikasi. Hasil dari pengklasifikasian tersebut diperoleh precision, recall, gmean, dan akurasi masing-masing sebesar 0,917, 0,824, 0,793, dan 0,943. [Copyright © 2019 SIMANIS. All rights reserved. Korespondensi: Penulis, Jurusan Matematika, UIN Maulana Malik Ibrahim Malang, Jl. Gajayana No. 50 Malang, Jawa Timur, Indonesia 65144 penulis@gmail.com](#) 1. **PENDAHULUAN (10 PT)** DNA merupakan unsur yang sangat penting dan mendasar pada setiap organisme, DNA sendiri dituliskan dalam bentuk sebuah urutan symbol-simbol yang mengkodekan ciri dari suatu organisme tersebut. [Urutan tersebut dikenal sebagai sequence DNA, yang merupakan informasi paling mendasar suatu gen atau genom karena mengandung instruksi yang dibutuhkan untuk pembentukan tubuh makhluk hidup](#)[1]. PCR(Polymerase Chain Reaction) [2] adalah salah satu metode untuk mendeteksi keberadaan mikroba dalam tubuh. Dibanding metode lain, metode ini tergolong akurat, cepat, dan bisa diandalkan. Keuntungan PCR lainnya adalah sekuens DNA dari mikroba atau strain infeksi yang baru ditemukan. Selanjutnya data DNA yang ditemukan tersebut dilakukan pencarian masuk pada jenis DNA apa, sehingga dibutuhkan metode yang optimal. Sekuensing DNA dapat dimanfaatkan [untuk menentukan identitas](#) maupun [fungsi gen atau fragmen DNA](#) lainnya [dengan cara membandingkan](#) sekuens-nya [dengan sekuens DNA lain yang sudah diketahui](#). Perbandingan yang banyak dilakukan adalah dengan proses pensejajaran 2 atau lebih sequence DNA. Proses pensejajaran ini membutuhkan waktu yang cukup lama. Sehingga dalam penelitian ini dibahas URL Prosiding: <http://conferences.uin-malang.ac.id/index.php/SIMANIS> 2 [Prosiding SI MaNIIs \(Seminar Nasional Integrasi Matematika dan Nilai Islami\)](#) pengkalsifikasian DNA dengan menggunakan Machine Learning[3]. Dalam proses kalsifikasi dengan machine learning tidak mensejajarkan, tetapi lebih ke pembuatan model dari masing- masing jeis DNA sehingga jika ada DNA baru yang diujikan dapat dikenali. Dalam penelitian ini dilakukan pengelompokan DNA dengan metode klasifikasi, yakni metode [IDELM\(Integrasi Seleksi Data dan Extreme Learning Machine\)](#) [4]. Metode [IDELM](#) dipilih dikarenakan data DNA yang sangat besar dan merupakan data dengan tipe imbalance data, maka perlu dilakukan seleksi data terlebih dahulu. Seleksi data ini dilakukan untuk mengatasi masalah imbalanced data. Konsep dari IDELM sendiri adalah melakukan proses seleksi data dan klasifikasi secara bersamaan. Seleksi data dilakukan untuk memilih mana fitur yang perlu digunakan dan mana fitur yang sudah terwakili dengan fitur yang terpilih sebelumnya. Proses seleksi data dan klasifikasi dilakukan secara bersamaan [untuk menghindari masalah inkonsistensi](#) antara [proses seleksi data dan](#) klasifikasi. Sehingga dengan menerapkan metode IDELM pada proses klasifikasi untuk pengelompokan DNA ini diharapkan hasil yang dicapai lebih optimal. 2. METODE PENELITIAN (10 PT) Langkah-langkah dari [penelitian ini adalah sebagai berikut: 1\) Pengumpulan Data Dalam](#) tahap [ini](#) juga dilakukan proses [pengumpulan data dengan](#) cara melakukan pengambilan data berupa FASTA dari DNA pada web NCBI. Data yang telah digunakan [dalam penelitian ini dapat dilihat pada tabel 1. Tabel 1. Data Penelitian No](#) Bakteri Patogen Bakteri Nonpatogen 1 Bartonella bacilliformis Bifidobacterium animals 2 Bordetella pertussis Bifidobacterium bifidum 3 Borrelia recurrentis Bifidobacterium breve 4 Haemophilus influenzae Bifidobacterium adolescentis 5 Haemophilus ducreyii Bifidobacterium longum 6 Streptococcus salivarius Lactobacillus delbrueckii 7 Streptococcus Pyogenes Lactobacillus acidophilus 8 Streptococcus mutans Lactobacillus brevis 9 Streptococcus agalactiae Lactobacillus fermentum 10 Brucella abortus Lactobacillus reuteri Selanjutnya data disiapkan sedemikian

rupa (dibuat dalam bentuk fragmen DNA) untuk diklasifikasikan. Pengumpulan data ini dilakukan oleh ketua peneliti. Adapun proses penyiapan data adalah: a. Penguraian Data Menjadi Beberapa Fragmen Pada proses penguraian data DNA menjadi beberapa fragmen digunakan aplikasi MetaSim[5]. Gambar 1. Simulator MetaSim Simulator ini menguraikan DNA menjadi masing-masing fragmennya. b. Ekstraksi Fitur dengan n-mers Selanjutnya dari data tersebut dilakukan ekstraksi fitur dengan menggunakan n-mers. Metode ini digunakan untuk mengetahui banyaknya kemunculan substring tertentu pada sebuah string[6]. Artinya untuk data DNA yang tersusun dari 4 jenis basa (A,C,G,T), sehingga jika n=3 maka akan ada 43 yaitu membentuk 64 substring, sedangkan jika n=4 maka akan ada 44 yaitu 256 substring, dalam hal ini n yang digunakan adalah 3, 4, dan 5. Pada setiap data dicari pula nilai rata-rata dan standar deviasinya untuk menjadi bagian dari fitur yang digunakan. Selanjutnya setiap data diberikan label sesuai dengan jenisnya, patogen dilabeli dengan 1 dan nonpatogen dilabeli dengan 0. c. Normalisasi Data Proses normalisasi data pada penelitian ini adalah menggunakan normalisasi minmax dengan rumus dapat dilihat pada persamaan 1. Prosiding [Seminar Nasional Integrasi Matematika dan Nilai Islami](#) Vol. 3, No. 1, September 2019: xx – xx p-ISSN: 2580-4596; e-ISSN: 2580-460X 
$$x_{norm} = \frac{(x_{max} - x_{min})}{(m_{max} - m_{min})} \times (m_{max} - m_{min}) + m_{min}$$
 (1) d. Dalam UIN Sultan Syarif Kasim Riau on 2019-04-29">[penelitian ini, data training dan data testing dibagi](#) menggunakan UIN Sultan Syarif Kasim Riau on 2019-04-29">[k-fold cross validation](#) yang diilustrasikan pada tabel 2. Data training dikelompokkan menjadi 5 kelompok, jika 1 kelompok untuk data testing maka sisanya untuk training. Sehingga berdasar ilustrasi pada tabel 2 proses training dan testing sebanyak 5 kali. Tabel 2. Ilustrasi 5-fold cross validation Keterangan:  $x_{norm}$  = [data hasil normalisasi](#)  $x_{min}$  = [nilai minimum dari data per kolom](#)  $x_{max}$  = [nilai maximum dari data per kolom](#)  $m_{min}$  = [adalah batas minimum yang kita berikan](#)  $m_{max}$  = [adalah batas maximum yang kita berikan](#) Pembagian data training dan testing [Kelompok 1 Kelompok 1 Kelompok 1 Kelompok 1 Kelompok 2 Kelompok 2 Kelompok 2 Kelompok 2 Kelompok 3 Kelompok 3 Kelompok 3 Kelompok 3 Kelompok 4 Kelompok 4 Kelompok 4 Kelompok 4 Kelompok 5 Kelompok 5 Kelompok 5 Kelompok 5 Kelompok 5](#) Keterangan: : Data Training : Data Testing 2) Proses Training Tahap desain dan perancangan sistem [akan menerjemahkan syarat kebutuhan ke sebuah perancangan perangkat lunak sebelum dibuat coding](#). Pada [proses ini](#) dibuat perancangan arsitektur perangkat lunak, struktur data, dan algoritma prosedural. pada tahap desain dan perancangan sistem ini dilakukan dengan diskusi bersama antara ketua peneliti dan anggota. 3) Pengujian Sistem (Testing & Integration) Pada tahap ini dilakukan pengujian terhadap program yang telah dibuat dengan cara melakukan uji coba terhadap semua fungsi pada sistem. Tahap pengujian sistem juga dilakukan sendiri oleh ketua peneliti, pengujian dilakukan dengan menggunakan berbagai macam data DNA, sehingga dapat diketahui seberapa baik kemampuan sistem yang dibuat. Untuk pengujian kemampuan dilakukan analisis confusion matriks, krena data yang digunakan termasuk dalam imbalanced data. Imbalanced data merupakan kasus khusus dalam Machine Learning. Ukuran evaluasi memainkan peran penting dalam machine learning. Ukuran tersebut digunakan untuk mengevaluasi dan mengarahkan algoritma pembelajaran. Jika pilihan ukuran mengabaikan kelas minoritas, maka algoritma pembelajaran tidak akan mampu menangani masalah imbalanced data dengan baik. Ukuran yang umum digunakan untuk dalam penelitian biasanya adalah tingkat klasifikasi keseluruhan yaitu akurasi. Namun pada imbalanced dataset, tingkat klasifikasi keseluruhan tidak lagi menjadi ukuran yang cocok, karena kelas minoritas tidak berpengaruh pada akurasi dibandingkan dengan kelas mayoritas. Oleh karena itu, ukuran lainnya telah dikembangkan untuk menilai kinerja classifier untuk data yang imbalanced. Berbagai ukuran yang umum didefinisikan berdasarkan confusion matrix. [Confusion matrix untuk klasifikasi biner ditunjukkan](#) dalam [Tabel 3. Tabel 3 Confusion matrix untuk klasifikasi biner True Nilai Sebenarnya False Prediksi True TP \(True Positive\) FP \(False Positive\) False FN \(False Negative\) TN \(True Negative\)](#) Judul makalah harus ringkas dan jelas, menggambarkan hasil penelitian (Penulis Pertama) [4 Prosiding SI MaNis \(Seminar Nasional Integrasi Matematika dan Nilai Islami\)](#) Diantara berbagai kriteria evaluasi, ukuran yang paling relevan dengan data yang imbalanced yaitu precision, recall, sensitifity, specificity, dan geometric mean (G-mean)[7]. Precision dalam artikel ini adalah berapa persen bakteri yang benar patogen dari keseluruhan bakteri yang diprediksi patogen. Sedangkan recall dalam penelitian ini adalah ketepatan berapa persen bakteri yang diprediksi patogen dibandingkan keseluruhan bakteri yang sebenarnya patogen. TP precision ? . TP ? FP TP recall ? . TP ? FN Ukuran Geometric mean (G-mean) ini digunakan ketika performa dari kedua kelas yang bersangkutan dan diharapkan tinggi secara bersamaan. Geometric mean telah digunakan beberapa peneliti untuk mengevaluasi classifier pada dataset yang imbalanced . [G-mean](#) mengindikasikan [keseimbangan antara kinerja klasifikasi pada kelas mayoritas dan minoritas](#). Ukuran G-mean diambil berdasarkan sensitifity (akurasi dari data positif) dan specificity (akurasi data negatif). sensitifity=recall TP specificity ? 1 ? FP ? TN G ? mean ? sensitifity ? specificity 3. [HASIL DAN PEMBAHASAN \(10 PT\) Hasil penelitian dapat dilihat pada](#) tabel 4. Dari tabel 4 ditunjukkan seberapa besar nilai precision, recall, g-mean, dan akurasi dari pengklasifikasian data patogen dan non patogen. 3-mers Precision Recall 0,950 Tabel 4. Hasil Penelitian 1000bp 0,905 G-mean Akurasi 0,829 0,970 4-mers 0,900 0,783 0,774 0,930 5-mers 0,850 0,680 0,725 0,890 Rata-rata 0,900 0,789 2000bp 0,776 0,930 3-mers 0,950 0,905 0,829 0,970 4-mers 0,900 0,818 0,793 0,940 5-mers 0,900 0,750 0,757 0,920 Rata-rata 0,917 0,824 0,793 0,943 Dari tabel sangat terlihat akurasi yang dimiliki cukup baik, dapat dilihat dari rata-rata precision, recall, G-mean, dan akurasi, masing masing pada sequence DNA yang panjang fragmennya 1000bp adalah 0,900, 0,789, 0,776, dan 0,930. Sedangkan untuk yang ukuran fragmennya 2000bp masing-masing nilai precision, recall, G-mean, dan akurasinya adalah 0,917, 0,824, 0,793, dan 0,943. Perhitungan akurasi yang digunakan tidak hanya perhitungan akurasi biasa karena data DNA termasuk dalam kategori data yang berkarakter imbalance. Hal tersebut terjadi karena meskipun sama-sama menggunakan 10 data masing - masing untuk data patogen dan non patogen saat dilakukan penguraian fragmen dengan panjang fragmen 1000bp dan 2000bp banyak fragmen yang terbentuk memiliki perbandingan yang jauh kberbeda nata data patogen dan non patogen, perbandingannya sekitar 1:10. Dari tabel di atas juga dapat dilihat pada masing masing panjang fragmen yang berbeda, jika dilihat dari proses ekstraksi fitur sequence DNA tersebut nika nili n pada n-mers semakin besar maka akurasi yang di dapat semakin sedikit. Sedangkan jika dilihat dari panjangnya fragmen semakin panjang fragmen yang dibuat, dalam penelitian ini nilai akurasinya semakin bagus. Prosiding [Seminar Nasional Integrasi Matematika dan Nilai Islami](#) Vol. 3, No. 1, September 2019: xx – xx p-ISSN: 2580-4596; e-ISSN: 2580-460X 5 4. KESIMPULAN (10 PT) Dari uraian yang telah dituliskan dapat disimpulkan bahwa klasifikasi sequence DNA bakteri patogen non patogen memiliki performansi yang cukup

baik. Dari hasil yang ada dapat dilihat bahwa rata-rata seluruh hasil melebihi 0,75, artinya akurasinya sekira 75%. 5. UCAPAN TERIMAKASIH (10 PT) Termakasih saya sampaikan kepada instansi UNP Kediri yang telah memberikan support yang sangat optimal. Serta diucapkan pula terimakasih pada pihak KEMENRISTEK Dikti telah memberikan bantuan dana untuk terlaksananya penelitian ini. Terimakasih pula kepada seluruh pihak yang telah mendukung dan membantu terlaksananya penelitian ini. DAFTAR PUSTAKA (10 PT) [1] Polanski A, Kimmel M. 2007. Bioinformatics. Berlin (DE): Springer. [2] [Alberts, Bruce; Johnson, Alexander; Lewis, Julian; Morgan, David; Raff, Martin; Roberts, Keith; Walter, Peter \(2014\). Molecular Biology of the Cell, Sixth Edition. Garland Science. United State of America](#) [3] [Wang J, Zaki M, Toivonen H, Shasha D. 2005. Data Mining in Bioinformatics. London \(UK\): Springer.](#) [4] Mahdiyah, Umi, Imah, E. M., Irawan, M. I.. 2017. [Integrating Data Selection And Extreme Learning Machine To Predict Protein-Ligand Binding Site, Contemporary Engineering Science, vol. 9.](#) [5] Richter DC, et al. 2009. User manual for MetaSim V0.9.5 [Internet]. [diunduh 2014 Juni 5]. Tersedia pada: [ab.informatik.uni.tuebingen.de/software/metasim/download/V095/manual.pdf](#) [6] [Erbert M, Rechner S, Müller-Hannemann M. Gerbil. 2017. a fast and memory-efficient k-mer counter with GPU-support. Algorithms for Molecular Biology. Springer Nature. doi:10.1186/s13015-017-0097-9](#) [7] Tharwat, A., 2018. Classification Assessment Methods, Applied Computing and Informatics [8] [Huang, G., Zhu, Q. dan Siew, C., \(2006a\), "Extreme Learning Machine: Theory and Applications", Neurocomputing, Vol. 70, 489-501](#)