

**IMPLEMENTASI SYNTHETIC MINORITY
OVERSAMPLING TECHNIQUE (SMOTE) PADA
ALGORITMA EXTREME GRADIENT BOOSTING
(XGBOOST) UNTUK KLASIFIKASI INDEKS STANDAR
PENCEMARAN UDARA (ISPU)
SKRIPSI**



Oleh :

ACHMAD FAUZIHAN BAGUS SAJIWO

20081010069

**PROGRAM STUDI INFORMATIKA
FAKULTAS ILMU KOMPUTER
UNIVERSITAS PEMBANGUNAN NASIONAL "VETERAN"
JAWA TIMUR
2024**

LEMBAR PENGESAHAN SKRIPSI

Judul : IMPLEMENTASI SYNTHETIC MINORITY OVERSAMPLING
TECHNIQUE (SMOTE) PADA ALGORITMA EXTREME
GRADIENT BOOSTING (XGBOOST) UNTUK KLASIFIKASI
INDEKS STANDAR PENCEMARAN UDARA (ISPU)

Oleh : ACHMAD FAUZIHAN BAGUS SAJIWO

NPM : 22081010069

Telah Diseminarkan Dalam Ujian Skripsi Pada :

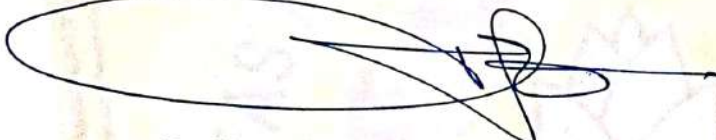
Hari Kamis, Tanggal 04 Juli 2024

Mengetahui

Dosen Pembimbing

Dosen Penguji

1.



Dr. Basuki Rahmat, S.Si. MT.

NIP. 19690723 2021211 002

1.



Eva Yulia Puspaningrum, S.Kom., M.Kom

NIP. 19890705 2021212 002

2.



Achmad Junaidi, S.Kom. M.Kom.

NPT. 3 7811 04 0199 1

2.



Retno Mumpuni, S.Kom., M.Sc

NPT. 172198 70 716054

Menyetujui

Dekan

Koordinator Program Studi

Fakultas Ilmu Komputer,

Informatika



Prof. Dr. Ir. Novirina Hendrasarie, MT.

NIP. 19681126 199403 2 001



Fetty Tri Anggraeny, S.Kom. M.Kom

NIP. 19820211 2021212 005

SURAT PERNYATAAN BEBAS DARI PLAGIASI

Saya, mahasiswa Program Studi Sarjana Informatika Universitas Pembangunan Nasional “Veteran” Jawa Timur, yang bertanda tangan di bawah ini:

Nama : Achmad Fauzihan Bagus Sajiwo

NPM : 20081010069

Menyatakan dengan sesungguhnya bahwa Skripsi/Tugas Akhir yang saya kerjakan berjudul:

“IMPLEMENTASI SYNTHETIC MINORITY OVERSAMPLING TECHNIQUE (SMOTE) PADA ALGORITMA EXTREME GRADIENT BOOSTING (XGBOOST) UNTUK KLASIFIKASI INDEKS STANDAR PENCEMARAN UDARA (ISPU)”

bukan merupakan plagiasi sebagian atau keseluruhan dari Skripsi/Tugas Akhir/Penelitian orang lain dari juga bukan merupakan produk dan software yang saya beli dari pihak lain. Saya juga menyatakan bahwa Skripsi/Tugas Akhir ini secara keseluruhan adalah pekerjaan Saya sendiri, kecuali yang dinyatakan dalam Daftar Pustaka dan tidak pernah diajukan untuk syarat memperoleh gelar di Universitas Pembangunan Nasional “Veteran” Jawa Timur maupun di Institut Pendidikan lain. Bukti hasil pengecekan plagiasi dokumen ini dapat ditelusuri melalui QR Code di bawah.

Apabila di kemudian hari terbukti bahwa dokumen ini merupakan plagiasi karya orang lain, saya sanggup menerima sanksi sesuai aturan yang berlaku.

Demikian atas perhatiannya disampaikan terima kasih.

Surabaya, 4 Juli 2024

Hormat saya,



Achmad Fauzihan Bagus Sajiwo

NPM. 20081010069



IMPLEMENTASI SYNTHETIC MINORITY OVERSAMPLING TECHNIQUE (SMOTE) PADA ALGORITMA EXTREME GRADIENT BOOSTING (XGBOOST) UNTUK KLASIFIKASI INDEKS STANDAR PENCEMARAN UDARA (ISPU)

Nama Mahasiswa : Achmad Fauzihan Bagus Sajiwo

NPM : 20081010069

Program Studi : Informatika

Dosen Pembimbing : Dr. Basuki Rahmat, S.Si. MT.

Achmad Junaidi, S.Kom, M.Kom.

ABSTRAK

Polusi udara adalah masuknya zat-zat berbahaya ke atmosfer, yang dapat disebabkan oleh tindakan manusia, baik sengaja maupun tidak sengaja, serta oleh peristiwa alam. Menurut *Air Quality Live Index* (AQLI) pada bulan April 2021, DKI Jakarta sebagai ibu kota negara, menempati posisi keenam di dunia dengan kota tingkat kualitas udara yang paling buruk. Untuk menghadapi masalah polusi udara yang terus memburuk, perlu diambil tindakan yang tepat dan efektif. Satu diantaranya adalah melakukan penelitian klasifikasi indeks standar pencemaran udara (ISPU).

Penerapan klasifikasi ISPU membutuhkan metode yang dapat mengolah dan menganalisis pola data dari sensor-sensor yang mengukur tingkat polutan udara. Metode yang digunakan pada penelitian ini adalah *eXtreme Gradient Boosting* (XGBoost). Untuk membantu menyeimbangkan data, pada penelitian ini menggunakan *Synthetic Minority Over-sampling Technique* (SMOTE). Data yang digunakan adalah ISPU DKI Jakarta tahun 2022-2023 yang berasal dari website Satu Data Jakarta: <https://satudata.jakarta.go.id/home>.

Hasil klasifikasi indeks standar pencemaran udara menggunakan algoritma *eXtreme Gradient Boosting* dengan *Synthetic Minority Over-sampling Technique*, didapatkan *accuracy* sebesar 99.63%. Dari perhitungan *confusion matrix* didapatkan nilai *precision*, *recall* dan *f1-score*. Untuk kelas 0 pada *precision* didapatkan nilai sebesar 99%, pada *recall* sebesar 100% dan *f1-score* sebesar 100%. Untuk kelas 1 pada *precision* didapatkan nilai sebesar 100%, pada *recall*

sebesar 99% dan *f1-score* sebesar 100%. Untuk kelas 2 pada *precision* didapatkan nilai sebesar 100%, pada *recall* sebesar 100% dan *f1-score* sebesar 100%.

Kata kunci: ISPU, Klasifikasi, XGBoost, Imbalanced Data, SMOTE

IMPLEMENTATION OF SYNTHETIC MINORITY OVERSAMPLING TECHNIQUE (SMOTE) IN THE EXTREME GRADIENT BOOSTING ALGORITHM (XGBOOST) FOR CLASSIFICATION OF AIR POLLUTION STANDARD INDEX (ISPU)

Nama Mahasiswa : Achmad Fauzihan Bagus Sajiwo

NPM : 20081010069

Program Studi : Informatika

Dosen Pembimbing : Dr. Basuki Rahmat, S.Si. MT.

Achmad Junaidi, S.Kom, M.Kom.

ABSTRACT

Air pollution is the entry of harmful substances into the atmosphere, which can be caused by human actions, whether intentional or unintentional, as well as by natural events. According to the Air Quality Live Index (AQLI) in April 2021, DKI Jakarta, as the nation's capital, is in sixth place in the world with the city with the worst air quality level. To deal with the problem of air pollution which continues to worsen, appropriate and effective action needs to be taken. One of them is conducting research on the classification of the air pollution standard index (ISPU).

Implementing the ISPU classification requires a method that can process and analyze data patterns from sensors that measure air pollutant levels. The method used in this research is eXtreme Gradient Boosting (XGBoost). To help balance the data, this research used Synthetic Minority Over-sampling Technique (SMOTE). The data used is the DKI Jakarta ISPU for 2022-2023 which comes from the Satu Data Jakarta website: <https://satudata.jakarta.go.id/home>.

The results of the standard air pollution index classification using the eXtreme Gradient Boosting algorithm with Synthetic Minority Over-sampling Technique, obtained an accuracy of 99.63%. From the confusion matrix calculations, precision, recall and f1-score values are obtained. For class 0, the precision value is 99%, the recall is 100% and the f1-score is 100%. For class 1, the precision score was 100%, the recall was 99% and the f1-score was 100%. For

class 2, the precision score was 100%, the recall was 100% and the f1-score was 100%.

Keywords: *ISPU, Classification, XGBoost, Imbalanced Data, SMOTE*

KATA PENGANTAR

Puji syukur kita panjatkan kehadiran Allah Swt. yang telah memberikan rahmat dan hidayah-Nya sehingga penulis dapat menyelesaikan laporan skripsi yang berjudul “Implementasi Synthetic Minority Oversampling Technique (SMOTE) Pada Algoritma Extreme Gradient Boosting (XGBoost) Untuk Klasifikasi Indeks Standar Pencemaran Udara (ISPU)” ini tepat pada waktunya. Penulis berusaha dengan semaksimal mungkin dalam penyusunan laporan ini baik dari hasil bimbingan bersama dosen pembimbing, maupun diskusi dengan teman-teman penulis.

Surabaya, 04 Juli 2024

Penulis

Achmad Fauzihan Bagus Sajiwo

UCAPAN TERIMA KASIH

Dalam kesempatan ini, penulis mengucapkan banyak terima kasih kepada semua pihak yang telah membantu dalam pelaksanaan kegiatan perkuliahan maupun penyusunan laporan ini kepada :

1. Orang tua penulis yang selalu memberikan doa serta dukungan baik secara moril ataupun materil kepada penulis
2. Bapak Prof. Dr. Ir. Akhmad Fauzi, M.MT., selaku Rektor Universitas Pembangunan Nasional “Veteran” Jawa Timur.
3. Ibu Prof. Dr. Ir. Novirina Hendrasarie, MT., selaku Dekan Fakultas Ilmu Komputer, Universitas Pembangunan Nasional “Veteran” Jawa Timur.
4. Ibu Fetty Tri Anggraeny, S.Kom, M.Kom. selaku Koordinator Program Studi Informatika Universitas Pembangunan Nasional “Veteran” Jawa Timur.
5. Bapak Andreas Nugroho Sihananto, S.Kom., M.Kom., selaku Koordinator Skripsi Program Studi Informatika Universitas Pembangunan “Veteran” Jawa Timur.
6. Bapak Dr. Basuki Rahmat, S.Si. MT., selaku dosen pembimbing 1 yang telah membimbing penulis dari awal penyusunan laporan sampai penanda tangan laporan ini.
7. Bapak Achmad Junaidi, S.Kom, M.Kom., selaku dosen pembimbing 2 yang telah membimbing penulis dari awal penyusunan laporan sampai penanda tangan laporan ini.
8. Ibu Henni Endah Wahanani, ST. M.Kom. selaku dosen wali Program Studi Informatika Fakultas Ilmu Komputer Universitas Pembangunan Nasional “Veteran” Jawa Timur.
9. Seluruh dosen dan staff Program Studi Informatika dan staff Fakultas Ilmu Komputer Universitas Pembangunan Nasional “Veteran” Jawa Timur.
10. Badan Legislatif Mahasiswa Fakultas Ilmu Komputer Universitas Pembangunan Nasional “Veteran” Jawa Timur (BLM FASILKOM UPN “Veteran” Jawa Timur) dan Badan Eksekutif Mahasiswa Fakultas Ilmu Komputer Universitas Pembangunan Nasional “Veteran” Jawa Timur (BEM FASILKOM UPN “Veteran” Jawa Timur) yang telah menjadi

wadah bagi penulis untuk belajar berorganisasi dan menambah relasi yang dapat berguna bagi masa depan penulis.

11. Teman-teman penulis yang senantiasa bersama dari semester 1 sampai sekarang.

Akhir kata, penulis berharap laporan ini dapat memberikan manfaat serta menjadi bahan referensi bagi penulis maupun pembaca. Penulis menyadari bahwa laporan yang ditulis ini masih jauh dari kata sempurna. Oleh karena itu, kritik dan saran yang membangun penulis butuhkan demi masa depan bersama.

DAFTAR ISI

LEMBAR PENGESAHAN.....	ii
SURAT PERNYATAAN ANTI PLAGIAT.....	iii
ABSTRAK	iv
ABSTRACT	vi
KATA PENGANTAR	viii
UCAPAN TERIMA KASIH	ix
DAFTAR ISI	xi
DAFTAR TABEL.....	xiv
DAFTAR GAMBAR	xvi
DAFTAR KODE PROGRAM	xviii
BAB I PENDAHULUAN	1
1.1. Latar Belakang	1
1.2. Rumusan Masalah.....	4
1.3. Tujuan Penelitian	4
1.4. Manfaat Penelitian	5
1.5. Batasan Masalah.....	5
BAB II TINJAUAN PUSTAKA	6
2.1. Penelitian Terdahulu	6
2.2. Pencemaran Udara	9
2.3. Indeks Standar Pencemaran Udara	10
2.4. Klasifikasi	12
2.5. eXtreme Gradient Boosting.....	12
2.6. Imbalance Dataset.....	14
2.7. Synthetic Minority Over-sampling Technique.....	14

BAB III METODOLOGI	16
3.1. Kebutuhan Perangkat	16
3.1.1. Perangkat Keras (<i>Hardware</i>)	16
3.1.2. Perangkat Lunak (<i>Software</i>)	16
3.2. Sumber Data Penelitian.....	16
3.3. Studi Pustaka	18
3.4. Tahapan Penelitian	18
3.5. Identifikasi Masalah.....	19
3.6. Pengumpulan Data.....	19
3.7. Preprocessing Data	19
3.8. Pembagian Data.....	22
3.9. Proses Klasifikasi	22
3.10. Evaluasi Model.....	23
BAB IV HASIL DAN PEMBAHASAN	24
4.1. Pengumpulan Data.....	24
4.2. Implementasi Program	31
4.2.1. Library	31
4.2.2. Preprocessing Data.....	32
4.2.2.1. Penggabungan Data.....	32
4.2.2.2. Transformation Data.....	33
4.2.2.3. Data Cleaning.....	35
4.2.2.4. Normalisasi Data.....	41
4.2.2.5. Feature Selection.....	43
4.2.2.6. Balancing Data.....	44
4.2.3. Pembagian Data	46
4.2.4. Klasifikasi XGBoost	47

4.3. Pembahasan Pengujian.....	47
4.3.1. Tanpa Balancing Data	48
4.3.2. Teknik SMOTE.....	55
4.3.3. Teknik Random Oversampling	62
4.3.4. Teknik Random Undersampling	70
4.3.5. Perbandingan Akurasi Model	77
4.4. Evaluasi Model.....	78
BAB V KESIMPULAN	83
5.1. Kesimpulan	83
5.2. Saran	83
DAFTAR PUSTAKA	85

DAFTAR TABEL

Tabel 2.1 Nilai Ambang Batas	10
Tabel 2.2 Indeks Standar Pencemaran Udara	11
Tabel 3.1 Data ISPU	17
Tabel 4.1 Contoh Data ISPU Tahun 2021	24
Tabel 4.2 Contoh Data ISPU Tahun 2022	27
Tabel 4.3 Contoh Data ISPU Tahun 2023	29
Tabel 4.4 Contoh Hasil Penggabungan Data	32
Tabel 4.5 Contoh Data Sebelum Transformasi	33
Tabel 4.6 Contoh Data Setelah Transformasi	34
Tabel 4.7 Contoh Data Sebelum dibersihkan.....	35
Tabel 4.8 Contoh Data yang Missing Value dihapus	36
Tabel 4.9 Contoh Data yang Missing Value diisi Median.....	37
Tabel 4.10 Contoh Data yang Missing Value diisi Mean.....	38
Tabel 4.11 Hasil Penghapusan Kolom yang Tidak dDigunakan	41
Tabel 4.12 Hasil Normalisasi Data yang Missing Value dihapus	42
Tabel 4.13 Hasil Normalisasi Data yang Missing Value diisi Median dan Mean ..	42
Tabel 4.14 Perbandingan Correlation Matrix	43
Tabel 4.15 Hasil Feature Selection.....	44
Tabel 4.16 Hasil Balancing Data Missing Value dihapus	45
Tabel 4.17 Hasil Balancing Data Missing Value diisi Median dan Mean.....	46
Tabel 4.18 Confusion Matrix Missing Value Hapus.....	49
Tabel 4.19 Confusion Matrix Missing Value Mean.....	51
Tabel 4.20 Confusion Matrix Missing Value Median.....	53
Tabel 4.21 Confusion Matrix Missing Value Hapus dengan Teknik SMOTE.....	56

Tabel 4.22	Confusion Matrix Missing Value Mean dengan Teknik SMOTE.....	58
Tabel 4.23	Confusion Matrix Missing Value Median dengan Teknik SMOTE	61
Tabel 4.24	Confusion Matrix Missing Value Hapus dengan Teknik ROS	63
Tabel 4.25	Confusion Matrix Missing Value Mean dengan Teknik ROS	66
Tabel 4.26	Confusion Matrix Missing Value Median dengan Teknik ROS	68
Tabel 4.27	Confusion Matrix Missing Value Hapus dengan Teknik RUS	71
Tabel 4.28	Confusion Matrix Missing Value Mean dengan Teknik RUS	73
Tabel 4.29	Confusion Matrix Missing Value Median dengan Teknik RUS	76
Tabel 4.30	Perbandingan Accuracy, Precision, Recall dan F1-Score	78
Tabel 4.31	Hasil Validasi K-Fold.....	81

DAFTAR GAMBAR

Gambar 2.1	Flowchart XGBoost	13
Gambar 2.2	Ilustrasi Cara Membuat Titik Data Sintetis Pada SMOTE	15
Gambar 3.1	Diagram Alur Penelitian	19
Gambar 3.2	Diagram Alur Tahap Preprocessing Data	20
Gambar 3.3	Ilustrasi Cara Kerja XGBoost.....	23
Gambar 4.1	Bar Plot Nilai Rata-Rata Parameter Setiap Kategori Tahun 2021	25
Gambar 4.2	Bar Plot Nilai Rata-Rata Parameter Setiap Stasiun Tahun 2021.....	26
Gambar 4.3	Bar Plot Nilai Rata-Rata Parameter Setiap Kategori Tahun 2022	27
Gambar 4.4	Bar Plot Nilai Rata-Rata Parameter Setiap Stasiun Tahun 2022.....	28
Gambar 4.5	Bar Plot Nilai Rata-Rata Parameter Setiap Kategori Tahun 2023	29
Gambar 4.6	Bar Plot Nilai Rata-Rata Parameter Setiap Stasiun tahun 2023	30
Gambar 4.7	Tipe Data Awal Pada Dataset.....	31
Gambar 4.8	Bar Plot Jumlah Kategori Kualitas Udara	34
Gambar 4.9	Tipe Data Setelah Transformasi	35
Gambar 4. 10	Pengecekan Missing Value Sebelum dihapus	39
Gambar 4.11	Pengecekan Missing Value Setelah dibersihkan	39
Gambar 4.12	Pengecekan Jumlah Setiap Kategori Sebelum dibersihkan.....	40
Gambar 4.13	Pengecekan Jumlah Setiap Kategori Missing Value Hapus.....	40
Gambar 4.14	Pengecekan Jumlah Setiap Kategori Missing Value	41
Gambar 4.15	Correlation Matrix	43
Gambar 4.16	Hasil Pengujian Missing Value dihapus	48
Gambar 4.17	Hasil Akurasi dan MSE Missing Value dihapus	48
Gambar 4.18	Hasil Pengujian Missing Value diisi Mean	50
Gambar 4.19	Hasil Akurasi dan MSE Missing value diisi Mean.....	50

Gambar 4.20 Hasil Pengujian Missing Value diisi Median	52
Gambar 4.21 Hasil Akurasi dan MSE Missing Value diisi Median	53
Gambar 4.22 Hasil Pengujian Missing Value dihapus dengan Teknik SMOTE	55
Gambar 4.23 Hasil Akurasi dan MSE Missing Value Hapus Teknik SMOTE	55
Gambar 4.24 Hasil Pengujian Missing Value diisi Mean dengan Teknik SMOTE	57
Gambar 4.25 Hasil Akurasi dan MSE Missing Value Mean Teknik SMOTE.....	58
Gambar 4.26 Hasil Pengujian Missing Value Median Teknik SMOTE.....	60
Gambar 4.27 Hasil Akurasi dan MSE Missing Value Median Teknik SMOTE.....	60
Gambar 4.28 Hasil Pengujian Missing Value Hapus dengan Teknik ROS	62
Gambar 4.29 Hasil Akurasi dan MSE Missing Value Hapus dengan Teknik ROS	63
Gambar 4.30 Hasil Pengujian Missing Value diisi Mean dengan Teknik ROS.....	65
Gambar 4.31 Hasil Akurasi dan MSE Missing Value Mean Teknik ROS	65
Gambar 4.32 Hasil Pengujian Missing Value diisi Median dengan Teknik ROS....	67
Gambar 4.33 Hasil Akurasi dan MSE Missing Value Median Teknik ROS	68
Gambar 4.34 Hasil Pengujian Missing Value dihapus dengan Teknik RUS	70
Gambar 4.35 Hasil Akurasi dan MSE Missing Value Hapus dengan Teknik RUS	70
Gambar 4.36 Hasil Pengujian Missing Value diisi Mean dengan Teknik RUS.....	72
Gambar 4.37 Hasil Akurasi dan MSE Missing Value Mean Teknik RUS	73
Gambar 4.38 Hasil Pengujian Missing Value diisi Median dengan Teknik RUS....	75
Gambar 4.39 Hasil Akurasi dan MSE Missing Value Median Teknik RUS	75
Gambar 4.40 Perbandingan Hasil Akurasi	77

DAFTAR KODE PROGRAM

Kode Program 4.1 Library Pada Klasifikasi XGBoost.....	31
Kode Program 4.2 Pemanggilan Dataset Setiap Tahun.....	32
Kode Program 4.3 Penggabungan Dataset	32
Kode Program 4.4 Penggunaan Teknik SMOTE	44
Kode Program 4.5 Penggunaan Teknik Random Oversampling	45
Kode Program 4.6 Penggunaan Teknik Random Undersampling	45
Kode Program 4.7 Pembagian Data Training dan Testing	46
Kode Program 4.8 Pembuatan Model	47