

**PENERAPAN SISTEM *DATA CLEANING* DALAM MASTER  
DATA DENGAN MENGGUNAKAN ALGORITMA  
*DUPLICATE COUNT STRATEGY*  
(STUDI KASUS: PT XYZ)**

**TUGAS AKHIR**

**Diajukan sebagai salah satu syarat untuk memperoleh gelar Sarjana  
Komputer**



**FILDZAH ADRA ARIFAH**

**NIM 1132001017**

**PROGRAM STUDI INFORMATIKA  
FAKULTAS TEKNIK DAN ILMU KOMPUTER  
UNIVERSITAS BAKRIE  
JAKARTA  
2021**

## **HALAMAN PERNYATAAN ORISINALITAS**

**Tugas Akhir ini adalah hasil karya saya sendiri,  
dan semua sumber baik yang dikutip maupun dirujuk  
telah saya nyatakan dengan benar.**

**Nama : Fildzah Adra Arifah**

**NIM : 1132001017**

**Tanda Tangan :** 

**Tanggal : 1 Juli 2021**

## HALAMAN PENGESAHAN

Tugas Akhir ini diajukan oleh:

Nama : Fildzah Adra Arifah

NIM : 1132001017

Program Studi : Informatika

Fakultas : Teknik dan Ilmu Komputer

Judul Skripsi : Penerapan Sistem *Data Cleaning* dalam Master Data dengan  
Menggunakan Algoritma *Duplicate Count Strategy*  
(Studi Kasus: PT XYZ)

Telah berhasil dipertahankan di hadapan Dewan Penguji dan diterima sebagai  
bagian persyaratan yang diperlukan untuk memperoleh gelar Sarjana Komputer  
pada Program Studi Informatika Fakultas Teknik dan Ilmu Komputer, Universitas  
Bakrie.

### DEWAN PENGUJI

Penguji 1 : Sondang Sibuea, S.Kom., M.Kom.



Penguji 2 : Albert A. Sembiring, S.T, M.T.



Anggota : Prof. Dr. Hoga Saragih, S.T, M.T, IPM.



Anggota : Ihsan Ibrahim, S.T., M.T.



Ditetapkan di : Jakarta

Tanggal : 25 Agustus 2021

## UNGKAPAN TERIMA KASIH

Puji dan syukur atas kehadirat Allah Subhanahu wata'ala atas rahmat-Nya dan karunia-Nya sehingga penulis dapat menyelesaikan Tugas Akhir ini dengan baik. Tugas Akhir dengan judul "**Penerapan Sistem Data Cleaning dalam Master Data dengan Menggunakan Algoritma Duplicate Count Strategy** (Studi Kasus: PT XYZ)" ini ditulis untuk memenuhi salah satu syarat dalam menyelesaikan perkuliahan pendidikan strata satu (S1) pada Program Studi Informatika, Universitas Bakrie.

Banyak pihak yang telah membantu penulis dalam penelitian dan penulisan Tugas Akhir ini, baik itu berupa bimbingan, saran, maupun dukungan secara moril dan materil. Saya menyadari bahwa, tanpa bantuan dan bimbingan dari berbagai pihak, dari masa perkuliahan sampai pada penyusunan Proposal Tugas Akhir ini, sangatlah sulit bagi peneliti untuk menyelesaikannya dan masih banyak kekurangan. Dengan segenap kerendahan hati, melalui kesempatan ini penulis ingin mengucapkan terima kasih kepada:

1. Prof. Dr. Hoga Saragih, S.T., M.T., selaku Kepala Program Studi Informatika dan dosen pembimbing, yang senantiasa memberikan masukan dan motivasi kepada penulis;
2. Ihsan Ibrahim, S.T., M.T., selaku dosen pembimbing, yang telah meluangkan waktunya serta memberikan bimbingan, saran, dan perbaikan dalam menyelesaikan penelitian ini;
3. Sondang Sibuea, S.Kom., M.Kom., selaku dosen pembahas dan penguji yang memberikan saran dan perbaikan terhadap penelitian ini;
4. Albert A. Sembiring, S.T, M.T., selaku dosen pembahas dan penguji yang memberikan saran dan perbaikan terhadap penelitian ini;
5. Seluruh Bapak/Ibu dosen Program Studi Informatika Universitas Bakrie, yang telah memberikan banyak ilmu, pengetahuan, wawasan kepada penulis selama perkuliahan;
6. Keluarga tercinta, kedua Orang tua penulis (Rafles dan Ade Eriani), saudara kandung penulis (Rifdah Adra Shalihah) dan keluarga besar yang telah memberikan dukungan dan doa yang sangat berarti bagi penulis;

7. Fuaidah Hanani dan Sairam Salim, yang telah membantu dalam proses penelitian, penulisan dan memberikan dukungan juga saran;
8. Muhammad Fadhil Zulmardhiya, yang bersedia menjadi tempat curhat mengenai suka duka dalam penyusunan tugas akhir ini, juga telah memberikan doa, dukungan, motivasi, semangat dan kasih sayang hingga saat ini;
9. Seluruh teman Informatika seangkatan, senior dan junior. Terima kasih atas kebersamaan selama masa studi di kampus;
10. Seluruh pihak yang terlibat dalam penyusunan Tugas Akhir ini yang tidak dapat penulis sebutkan satu persatu;

Penulis berharap semoga Allah Subhanahu wata'ala membalas kebaikan seluruh pihak yang telah membantu penulis dalam menyelesaikan Tugas Akhir ini. Penulis berharap Tugas Akhir ini bermanfaat bagi pihak-pihak terkait ke depannya.

Jakarta, 1 Juli 2021



Fildzah Adra Arifah

## HALAMAN PERNYATAAN PERSETUJUAN PUBLIKASI

Sebagai sivitas akademik Universitas Bakrie, saya yang bertanda tangan di bawah ini:

Nama : Fildzah Adra Arifah  
NIM : 1132001017  
Program Studi : Informatika  
Fakultas : Teknik dan Ilmu Komputer  
Jenis Tugas Akhir : Skripsi

Demi pengembangan ilmu pengetahuan, menyetujui untuk memberikan kepada Universitas Bakrie **Hak Bebas Royalti Noneksklusif (Non-exclusive Royalty-Free Right)** atas karya ilmiah saya yang berjudul:

### PENERAPAN SISTEM DATA CLEANING DALAM MASTER DATA DENGAN MENGGUNAKAN ALGORITMA DUPLICATE COUNT STRATEGY (STUDI KASUS: PT XYZ)

Beserta perangkat yang ada (jika diperlukan). Dengan Hak Bebas Royalti Noneksklusif ini Universitas Bakrie berhak menyimpan, mengalihmedia/formatkan, mengelola dalam bentuk pangkalan data (*database*), merawat, dan mempublikasikan tugas akhir saya selama tetap mencantumkan nama saya sebagai penulis/pencipta dan sebagai pemilik Hak Cipta untuk kepentingan akademis.

Demikian pernyataan ini saya buat dengan sebenarnya.

Dibuat di : Jakarta  
Pada tanggal : 24 Agustus 2021

Yang menyatakan



( Fildzah Adra Arifah )

**PENERAPAN SISTEM *DATA CLEANING* DALAM MASTER  
DATA DENGAN MENGGUNAKAN ALGORITMA  
*DUPPLICATE COUNT STRATEGY* (STUDI KASUS: PT XYZ)**

Fildzah Adra Arifah<sup>1</sup>

---

**ABSTRAK**

Adanya permasalahan berupa duplikasi data dalam sebuah master data, peneliti melakukan penerapan sistem untuk *data cleaning* yang dapat mendeteksi duplikasi data pada master data konsumen Divisi *Consumer Care* PT XYZ. Dalam penelitian ini digunakan algoritma untuk mendeteksi duplikasi data yaitu dengan menerapkan pendekatan metode *Duplicate Count Strategy* (DCS++) dan *N-Gram*. Sistem *data cleaning* diperuntukan bagi *Sales Admin* agar mempermudah dalam menemukan duplikasi data. Selain itu, sistem ini juga dibuat untuk merapikan format penulisan telepon dan fax yang ada pada master data konsumen Divisi *Consumer Care* PT XYZ. Penerapan ini dibangun dengan menggunakan bahasa pemrograman C#. Hasil dari penerapan sistem *data cleaning* yang dibangun akan dinilai seberapa efektif metode DCS++ dan N-Gram dengan menghitung nilai *recall* dan *precision* dalam mendeteksi duplikasi data.

**Kata kunci:** *Data cleaning*, Deteksi Duplikasi Data, *Duplicate Count Strategy*, *N-gram*

---

<sup>1</sup> Mahasiswa Program Studi Informatika, Universitas Bakrie

**IMPLEMENTATION OF DATA CLEANING SYSTEM IN MASTER DATA  
USING DUPLICATE COUNT COUNT STRATEGY ALGORITHM (CASE  
STUDY: PT XYZ)**

Fildzah Adra Arifah<sup>2</sup>

---

**ABSTRACT**

*The existence of a problem in the form of duplication of data in a master data, the researchers implemented a system for data cleaning that can detect duplication of data in the consumer master data of the Consumer Care Division of PT XYZ. In this study, an algorithm is used to detect duplication of data by applying the Duplicate Count Strategy (DCS++) and N-Gram method approaches. The data cleaning system is intended for Sales Admins to make it easier to find duplicate data. In addition, this system was also created to tidy up the format of writing telephone and faxes in the consumer master data of the Consumer Care Division of PT XYZ. This application is built using the C# programming language. The results of the implementation of the data cleaning system that was built will be assessed how effective the DCS++ and N-Gram methods are by calculating the recall and precision values in detecting duplication of data.*

**Keywords:** Data Cleaning, Duplicate Detection, Duplicate Count Strategy, N-Gram

---

<sup>1</sup> Student of Informatics Program, Bakrie University

## DAFTAR ISI

<b>HALAMAN PERNYATAAN ORISINALITAS .....</b>	ii
<b>HALAMAN PENGESAHAN.....</b>	iii
<b>UNGKAPAN TERIMA KASIH .....</b>	iv
<b>ABSTRAK.....</b>	vii
<b>ABSTRACT.....</b>	viii
<b>DAFTAR ISI.....</b>	ix
<b>DAFTAR GAMBAR.....</b>	xi
<b>DAFTAR TABEL .....</b>	xii
<b>DAFTAR RUMUS .....</b>	xiii
<b>BAB I.....</b>	1
<b>PENDAHULUAN .....</b>	1
1.1 <b>Latar Belakang Masalah.....</b>	1
1.2 <b>Perumusan Masalah.....</b>	2
1.3 <b>Tujuan Penelitian .....</b>	2
1.4 <b>Manfaat Penelitian .....</b>	3
1.5 <b>Batasan Masalah.....</b>	3
1.6 <b>Sistematika Penulisan .....</b>	4
<b>BAB II .....</b>	5
<b>TINJAUAN PUSTAKA .....</b>	5
2.1 <b>Penelitian Terkait .....</b>	5
2.2 <b>Deteksi Duplikasi .....</b>	7
2.3 <b>Data Cleaning .....</b>	8
2.4 <b>Algoritma Deteksi Duplikasi Data .....</b>	9
2.5 <b>Algoritma Duplicate Count Strategy++ .....</b>	9
<b>BAB III.....</b>	11
<b>METODOLOGI PENELITIAN .....</b>	11
3.1 <b>Metode Perancangan dan Pengembangan.....</b>	11
3.1.1 Pengamatan dan Perencanaan .....	11
3.1.2 Analisis Kebutuhan Aplikasi .....	11
3.1.3 Perancangan dan Pembangunan.....	11
3.1.4 Implementasi.....	23
3.1.5 Pengujian.....	23
3.2 <b>Objek Penelitian.....</b>	23

<b>3.3 Jenis Penelitian.....</b>	23
<b>3.4 Metode Pengumpulan Data.....</b>	24
<b>BAB IV .....</b>	25
<b>IMPLEMENTASI DAN PEMBAHASAN.....</b>	25
<b>4.1 Perancangan Sistem .....</b>	25
<b>4.1.1 Alur Algoritma Deteksi Duplikasi Data .....</b>	25
<b>4.1.2 Rancangan Database .....</b>	36
<b>4.1.3 UML (<i>Unified Modelling Language</i>) .....</b>	36
<b>4.2 Implementasi.....</b>	48
<b>4.2.1 Implementasi Sistem .....</b>	48
<b>4.2.2 Implementasi <i>Graphical User Interface (GUI)</i> .....</b>	49
<b>4.3 Pengujian.....</b>	51
<b>4.3.1 Evaluasi Data .....</b>	51
<b>4.3.2 Pengujian <i>White Box</i> .....</b>	57
<b>BAB V .....</b>	65
<b>SIMPULAN DAN SARAN.....</b>	65
<b>5.1 Simpulan .....</b>	65
<b>5.2 Saran .....</b>	66
<b>DAFTAR PUSTAKA .....</b>	67
<b>Lampiran 1 – <i>Requirement Elicitation</i>.....</b>	69
<b>Lampiran 2 – <i>Template Import Data</i>.....</b>	72

## DAFTAR GAMBAR

Gambar 4.1 <i>Flowchart</i> Algoritma Deteksi Duplikasi Data .....	26
Gambar 4.2 <i>Flowchart</i> Pra-cleaning.....	28
Gambar 4.3 <i>Flowchart</i> Tokenisasi.....	31
Gambar 4.4 <i>Flowchart</i> Pemecahan Kata dengan Nilai <i>N-Gram</i> .....	34
Gambar 4.5 <i>Flowchart</i> Perhitungan Nilai Kemiripan Antar-record .....	35
Gambar 4.6 <i>Database</i> Sistem <i>Data Cleaning</i> PT XYZ.....	36
Gambar 4.7 <i>Use Case</i> Sistem <i>Data Cleaning</i> .....	37
Gambar 4.8 <i>Activity Diagram</i> <i>Login</i> .....	42
Gambar 4.9 <i>Activity Diagram</i> <i>Duplication Detection</i> .....	43
Gambar 4.10 <i>Activity Diagram</i> <i>Import Data</i> .....	44
Gambar 4.11 <i>Activity Diagram</i> <i>Duplicate Detection Result</i> .....	44
Gambar 4.12 <i>Activity Diagram</i> <i>Fix Contact Number</i> .....	45
Gambar 4.13 <i>Activity Diagram</i> <i>Fix Result</i> .....	45
Gambar 4.14 <i>Activity Diagram</i> <i>Export Data</i> .....	46
Gambar 4.15 <i>Class Diagram</i> .....	47
Gambar 4.16 Tampilan <i>Login</i> .....	49
Gambar 4.17 Halaman Sistem <i>Data Cleaning</i> .....	49
Gambar 4.18 <i>View Clean Data</i> .....	50
Gambar 4.19 <i>Import Data</i> .....	51
Gambar 4.20 Grafik $D_{small} = 2700$ Data.....	54
Gambar 4.21 Grafik $D_{large} = 27.000$ Data .....	56

## DAFTAR TABEL

Tabel 4.1 Rincian Karakter yang Dihilangkan pada Proses Pendekripsi .....	27
Tabel 4. 2 Contoh Sebelum Proses Pra-cleaning .....	28
Tabel 4.3 Contoh Setelah Melakukan Tahap Pra-Cleaning.....	28
Tabel 4.4 Contoh Setelah Proses Tokenisasi .....	29
Tabel 4.5 Setelah Token Diurutkan dan Digabungkan.....	30
Tabel 4.6 Aktor dan Sistem dalam <i>Use Case Login</i> .....	38
Tabel 4.7 Aktor dan Sistem dalam <i>Use Case Duplicate Detection</i> .....	39
Tabel 4.8 Aktor dan Sistem dalam <i>Use Case View Duplicate Detection Result</i> .....	39
Tabel 4.9 Aktor dan Sistem dalam <i>Use Case Import Data</i> .....	40
Tabel 4.10 Aktor dan Sistem dalam <i>Use Case Fix Contact Number</i> .....	41
Tabel 4.11 Deskripsi Aksi Aktor dan Respon Sistem Untuk <i>Use Case Fix Result</i> ....	41
Tabel 4.12 Aktor dan Sistem dalam <i>Use Case Export Data</i> .....	42
Tabel 4.13 $D_{small} = 2700$ Menggunakan Nilai <i>Token</i> , <i>N-Gram</i> dan <i>Threshold</i> .....	53
Tabel 4.14 $D_{large} = 27.000$ Menggunakan Nilai <i>Token</i> , <i>N-Gram</i> dan <i>Threshold</i> .....	55
Tabel 4.15 Pengujian Algoritma DCS++ – Fungsi <i>cleanData()</i> .....	57
Tabel 4.16 Pengujian Algoritma DCS++ – Fungsi <i>ChopTheWords()</i> .....	58
Tabel 4.17 Pengujian Algoritma DCS++ – Fungsi <i>gramming()</i> .....	60

## DAFTAR RUMUS

Rumus 4.1 <i>N-Gram</i> dalam Menghitung Kemiripan Antar- <i>string</i> .....	32
Rumus 4.2 Rumus <i>precision</i> .....	52
Rumus 4.3 Rumus <i>recall</i> .....	52
Rumus 4.4 Rumus <i>f-measure</i> .....	52

## DAFTAR LAMPIRAN

Lampiran 1 <i>Requirement Elicitation</i> .....	69
Lampiran 2 <i>Template Import Data</i> .....	71