

**RANCANG BANGUN SISTEM *DATA CLEANING* UNTUK
MASTER DATA KONSUMEN DI PT XYZ DENGAN
MENERAPKAN METODE *SORTED NEIGHBOURHOOD* DAN
METODE *N-GRAM***

TUGAS AKHIR

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



**UNIVERSITAS
BAKRIE**

**RAHMA MUALIFA
1112001011**

**PROGRAM STUDI INFORMATIKA
FAKULTAS TEKNIK DAN ILMU KOMPUTER
UNIVERSITAS BAKRIE
JAKARTA
2016**

**RANCANG BANGUN SISTEM *DATA CLEANING* UNTUK
MASTER DATA KONSUMEN DI PT XYZ DENGAN
MENERAPKAN METODE *SORTED NEIGHBOURHOOD* DAN
METODE *N-GRAM***

TUGAS AKHIR

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



**UNIVERSITAS
BAKRIE**

**RAHMA MUALIFA
1112001011**

**PROGRAM STUDI INFORMATIKA
FAKULTAS TEKNIK DAN ILMU KOMPUTER
UNIVERSITAS BAKRIE
JAKARTA
2016**

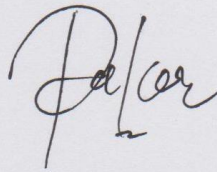
HALAMAN PERNYATAAN ORISINALITAS

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

Nama : Rahma Mualifa

NIM : 1112001011

Tanda Tangan :



Tanggal 10 Agustus 2016

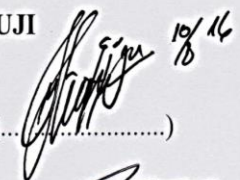
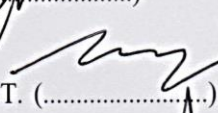
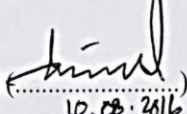
HALAMAN PENGESAHAN

Tugas Akhir ini diajukan oleh:

Nama : Rahma Mualifa
NIM : 1112001011
Program Studi : Informatika
Fakultas : Teknik dan Ilmu Komputer
Judul Skripsi : Rancang Bangun Sistem *Data Cleaning* Untuk Master Data Konsumen di PT XYZ Dengan Menerapkan Metode *Sorted Neighbourhood* dan *N-Gram*

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

DEWAN PENGUJI

Pembimbing : Yusuf Lestanto, S.T., M.Sc. (......)
Penguji 1 : Prof. Dr. Hoga Saragih, S.T., M.T. (......)
Penguji 2 : Gun Gun Gumilar, S.Kom., MMSI (......)
Ditetapkan di : Jakarta
Tanggal : 10 Agustus 2016

UNGKAPAN TERIMA KASIH

Puji dan syukur kehadirat Allah SWT karena atas rahmat-Nya dan karunia-Nya sehingga penulis dapat menyelesaikan Tugas Akhir ini dengan baik. Tugas Akhir dengan judul “Rancang Bangun Sistem *Data Cleaning* Untuk Master Data Konsumen di PT XYZ Dengan Menerapkan Metode *Sorted Neighbourhood* dan Metode *N-Gram*” 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. Oleh karena itu, pada kesempatan ini penulis ingin menyampaikan rasa terima kasih dan penghargaan yang setinggi-tingginya kepada:

1. Hoga Saragih, S.T., M.T., selaku Kepala Program Studi Informatika, yang senantiasa memberikan masukan dan motivasi kepada penulis;
2. Yusuf Lestanto, S.T., M.Sc., selaku dosen pembimbing, yang telah meluangkan waktunya serta memberikan bimbingan, saran, dan perbaikan dalam menyelesaikan penelitian ini;
3. Prof. Dr. Hoga Saragih, S.T., M.T. dan Gun Gun Gumilar, S.Kom., MMSI selaku dosen pembahas dan penguji yang memberikan saran dan perbaikan terhadap penelitian ini;
4. Seluruh Bapak/Ibu dosen Program Studi Informatika UB, yang telah memberikan banyak ilmu, pengetahuan, wawasan kepada penulis selama perkuliahan;
5. Keluarga tercinta, kedua Orang tua penulis (Ariyanto dan Siti Rohmatun) dan saudara kandung penulis (Rahma Maulida dan Rahman Abid) yang telah memberikan dukungan dan doa yang sangat berarti bagi penulis.;

6. Tim Magang (Mario Joel, Aulia Syarifuddin, Mega Silviana, Celina Maya Cantika, M. Fadil). Terima kasih telah memberikan semangat, motivasi, dukungan, suka cita dan kebersamaan selama ini;
7. Teman VMG (Destalia Dianing Putri, Septy Dwi Aryani, Eka Juniar, Andining Tyas, Diti Puspa Permata). Terima kasih telah memberikan semangat, motivasi, dukungan, suka cita dan kebersamaan selama ini;
8. Tim Seperjuangan (Sawitri Sadanti, Stefanny Uliarta, Sarah Putri, Rien Pratama, Rahmad Pratama Dita, Fazz Faidurrahman, Evi Margaretha, Rizky Akbarie, Chandra Setiawan). Terima kasih telah menjadi teman yang selalu memberikan semangat, motivasi, dukungan dan suka cita selama lebih dari 4 tahun.
9. Teman-teman TIF 2011 senasib seperjuangan. Terima kasih sudah menemani dan bekerja sama selama lebih dari 4 tahun masa studi di UB;
10. Seluruh pihak yang terlibat dalam penyusunan Tugas Akhir ini yang tidak dapat penulis sebutkan satu persatu;

Dengan segala keterbatasan yang ada, penulis menyadari bahwa penyusunan tugas akhir ini masih jauh dari kesempurnaan. Untuk itu, saran dan kritik akan selalu diterima agar penulis dapat memperbaiki setiap kekurangan untuk kesempurnaan dimasa mendatang.

Akhirnya, penulis menyampaikan ucapan terima kasih dan semoga Allah SWT membalas segala kebaikan serta melimpahkan berkat dan rahmat-Nya kepada semua pihak yang telah membantu selama ini. Penulis berharap semoga Tugas Akhir ini berguna dan bermanfaat bagi kita semua.

Jakarta, 10 Agustus 2016

Rahma Mualifa

HALAMAN PERNYATAAN PERSETUJUAN PUBLIKASI

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

Nama : Rahma Mualifa
NIM : 1112001011
Program Studi : Informatika
Fakultas : Teknik dan Ilmu Komputer
Jenis Tugas Akhir : Rancang Bangun

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:

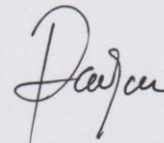
Rancang Bangun Sistem *Data Cleaning* Untuk Master Data Konsumen di PT XYZ Dengan Menerapkan Metode *Sorted Neighbourhood* dan *N-Gram*

beserta perangkat yang ada (jika diperlukan). Dengan Hak Bebas Royalti Noneksklusif ini Universitas Bakrie berhak menyimpan, mengalih media/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 : 10 Agustus 2016

Yang menyatakan



Rahma Mualifa

Rancang Bangun Sistem *Data Cleaning* Untuk Master Data Konsumen di PT XYZ Dengan Menerapkan Metode *Sorted Neighbourhood* dan *N-Gram*

Rahma Mualifa

ABSTRAK

Penelitian ini membahas tentang rancang bangun sistem *data cleaning* untuk dapat mendeteksi duplikasi data yang ada pada master data konsumen Divisi *Consumer Care* PT XYZ. Metode yang digunakan dalam penelitian ini untuk mendeteksi duplikasi data adalah dengan menerapkan pendekatan metode *Sorted Neighbourhood* (SNM) dan *N-Gram*. Sistem *data cleaning* ini bertujuan membantu *user* untuk dapat mempermudah menemukan duplikasi data. Selain itu, sistem ini juga dapat membantu *user* untuk dapat merapikan format penulisan telepon dan fax yang ada pada master data konsumen Divisi *Consumer Care* PT XYZ. Sistem yang akan dibangun adalah sistem *web based* dengan menggunakan bahasa pemrograman C#. Hasil dari sistem *data cleaning* yang dibangun kemudian akan diuji coba kepada *user* dan dinilai seberapa efektif metode SNM dan N-Gram dalam mendeteksi duplikasi data dengan menghitung nilai *recall* dan *precision* terhadap hasil proses deteksi duplikasi data.

Kata kunci: *Data cleaning*, Deteksi Duplikasi Data, *Sorted Neighbourhood*, *N-gram*

Rancang Bangun Sistem *Data Cleaning* Untuk Master Data Konsumen di PT XYZ Dengan Menerapkan Metode *Sorted Neighbourhood* dan *N-Gram*

Rahma Mualifa

ABSTRACT

This research discusses data cleaning system to detect data duplication are exists in customer master data at Consumer Care Division PT XYZ. Method will be used in this research to detect data duplication is by implementing Sorted Neighbourhood Method (SNM) and N-Gram. This system aims to assist user to find duplicate data easier. Moreover, it also can assist user to fix the format number both phone and fax number within customer master data at Consumer Care Division PT XYZ. System will be developed as web-based system by using C# programming language. The result of this system development will be tested by user and rated its effectiveness through implementation of SNM and N-Gram. To compute effectiveness will be obtained recall and precision value to determine how effective this system in detecting duplication data.

Kata kunci: *Data cleaning, Duplicate Detection, Sorted Neighbourhood, N-gram*

DAFTAR ISI

HALAMAN PERNYATAAN ORISINALITAS.....	i
HALAMAN PENGESAHAN.....	ii
UNGKAPAN TERIMA KASIH.....	iii
HALAMAN PERNYATAAN PERSETUJUAN PUBLIKASI.....	v
ABSTRAK	vi
ABSTRACT.....	vii
DAFTAR ISI.....	viii
DAFTAR GAMBAR	xi
DAFTAR TABEL.....	xiii
DAFTAR RUMUS	xv
DAFTAR LAMPIRAN.....	xvi
DAFTAR SINGKATAN	xvii
1 Pendahuluan	1
1.1 Latar Belakang Masalah	1
1.2 Perumusan Masalah.....	3
1.3 Tujuan Penelitian.....	3
1.4 Manfaat Penelitian.....	3
1.5 Batasan Masalah	4
1.6 Sistematika Penulisan	4
2 Tinjauan Pustaka	6
2.1 Penelitian Terkait.....	6
2.2 Pengertian Sistem	8
2.3 <i>Data Cleaning</i>	9
2.4 Metode <i>Data Cleaning</i>	10
2.4.1 Algoritma Deteksi Duplikasi Data	10
2.4.2 Metode <i>Sorted Neighbourhood</i> Sebagai Metode Untuk Deteksi Duplikasi Data.....	11
2.4.3 Algoritma Perhitungan Kemiripan Antar String	12

2.4.4	Algoritma Pendekatan <i>N-Gram</i> Sebagai Algoritma Perhitungan Kemiripan Antar <i>String</i>	14
2.5	Sistem Berbasis <i>Web</i>	15
2.6	Bahasa Pemrograman ASP.NET	16
2.7	Model <i>Waterfall</i>	16
2.8	Pemrograman Berorientasi Objek	18
2.9	<i>Unified Modelling Language</i> (UML)	19
3	Metodologi Penelitian	23
3.1	Tahap Analisis dan Definisi Kebutuhan.....	23
3.1.1	Prosedur Yang Sedang Berjalan.....	24
3.1.2	Master Data Konsumen Divisi <i>Consumer Care</i> di PT XYZ..	25
3.1.3	Struktur Organisasi.....	27
3.1.4	Bisnis Proses.....	30
3.1.5	Sistem <i>Data Cleaning</i> Yang Diajukan	30
3.1.6	Definisi Kebutuhan Sistem.....	31
3.2	Tahap Perancangan Sistem.....	32
3.3	Tahap Implementasi	32
3.4	Tahap Pengujian	32
3.5	Tahap Pemeliharaan	33
4.	Implementasi dan Pembahasan	34
4.1	Perancangan Sistem.....	34
4.1.1	Perancangan Alur Algoritma Deteksi Duplikasi Data	34
4.1.2	Perancangan <i>Database</i>	44
4.1.3	UML (<i>Unified Modelling Language</i>).....	45
4.1.3.1	<i>Use Case Diagram</i>	45
4.1.3.2	<i>Activity Diagram</i>	51
4.1.3.3	<i>Class Diagram</i>	55
4.2	Implementasi	56
4.2.3	Implementasi Sistem	56
4.2.4	Implementasi GUI (<i>Graphical User Interface</i>).....	57
4.3	Pengujian	62

4.3.1	Pengujian <i>White Box</i>	62
4.3.1.1	Pengujian Algoritma SNM	62
4.3.1.2	Pengujian Algoritma N-Gram	66
4.3.2	Pengujian <i>Black Box</i>	72
4.3.2.1	Pengujian <i>Functionality</i>	72
4.3.2.2	Pengujian <i>Usability</i>	74
4.3.2.3	Pengujian <i>Compatibility</i>	75
4.3.2.4	Pengujian <i>Performance</i>	76
4.3.3	Evaluasi Data.....	77
5.	Simpulan dan Saran.....	81
5.1	Simpulan.....	81
5.2	Saran	83
	Daftar Pustaka	84
	Lampiran – Lampiran.....	89

DAFTAR GAMBAR

Gambar 2.1 Bagian-Bagian Komponen dari Suatu Sistem dapat Mengendalikan Operasinya Sendiri	8
Gambar 2.2 <i>Waterfall Model</i>	17
Gambar 3.1 Struktur Organisasi Divisi <i>Consumer Care</i> PT XYZ.....	27
Gambar 3.2 Gambar Bisnis Proses Deteksi Data Kembar Pada Master Data Konsumen Divisi <i>Consumer Care</i> PT XYZ	30
Gambar 4.1 <i>Flowchart</i> Algoritma Deteksi Duplikasi Data	35
Gambar 4.2 <i>Flowchart</i> Tahap <i>Pra-cleaning</i>	38
Gambar 4.3 <i>Flowchart</i> Tahap Tokenisasi	40
Gambar 4.4 <i>Flowchart</i> Tahap Pemecahan Kata Berdasarkan Nilai N-Gram	42
Gambar 4.5 <i>Flowchart</i> Perhitungan Nilai Kemiripan Antar <i>Record</i>	43
Gambar 4.6 Rancangan <i>Database</i> Sistem <i>Data Cleaning</i> Pada <i>Database</i> Master Data Konsumen PT XYZ	45
Gambar 4.7 <i>Use Case</i> Sistem <i>Data Cleaning</i>	46
Gambar 4.8 <i>Activity Diagram</i> <i>Login</i>	51
Gambar 4.9 <i>Activity Diagram</i> <i>Detection Duplication</i>	52
Gambar 4.10 <i>Activity Diagram</i> <i>Import Data</i>	52
Gambar 4.11 <i>Activity Diagram</i> <i>View Duplicate Detection Result</i>	53
Gambar 4.12 <i>Activity Diagram</i> <i>Fixing Contact Number</i>	53
Gambar 4.13 <i>Activity Diagram</i> <i>Save Fix Result</i>	54
Gambar 4.14 <i>Activity Diagram</i> <i>Export Data</i>	54
Gambar 4.15 <i>Class Diagram</i>	55
Gambar 4.16 Tampilan <i>Login</i>	57
Gambar 4.17 Halaman Utama Sistem <i>Data Cleaning</i>	58
Gambar 4.18 Tampilan <i>View Clean Data</i>	59
Gambar 4.19 Tampilan <i>Drop Down List View Data</i> dan <i>Area</i>	59
Gambar 4.20 Halaman <i>Fix Contact Number</i>	60
Gambar 4.21 Tampilan Hasil <i>Fix Contact Number</i>	60
Gambar 4.22 Tampilan Halaman <i>Import Data</i>	61

Gambar 4.23 Tampilan Pengisian Halaman <i>Import Data</i>	61
Gambar 4.24 Gambar Grafik Hasil Pengujian <i>Usability</i> Pada <i>User Sistem Data Cleaning</i> PT XYZ	74
Gambar 4.25 Gambar Grafik Hasil Pengujian $D_{small} = 2500$ Data.....	79
Gambar 4.26 Gambar Grafik Hasil Pengujian $D_{large} = 25.000$ Data.....	80

DAFTAR TABEL

Tabel 2.1 Perbandingan Metode dari Beberapa Penelitian Terkait	7
Tabel 2.2 Tabel Simbol-Simbol Pada <i>Use Case Diagram</i>	20
Tabel 2.3 Tabel Simbol-Simbol Pada <i>Activity Diagram</i>	21
Tabel 2.4 Tabel Simbol-Simbol Pada <i>Class Diagram</i>	22
Tabel 3.1 Kamus Data Konsumen Divisi <i>Consumer Care</i> PT XYZ	25
Tabel 3.2 Tabel <i>Role</i> dan Deskripsi Kerja Divisi <i>Consumer Care</i> PT XYZ.....	28
Tabel 4.1 Tabel Rincian Karakter Yang Akan Dihilangkan Pada Proses Deteksi Duplikasi Data	36
Tabel 4.2 Tabel Sebelum Dilakukan Proses <i>Pra-cleaning</i>	37
Tabel 4.3 Tabel Contoh Setelah Melewati Tahap <i>Pra-Cleaning</i>	37
Tabel 4.4 Tabel Contoh Setelah Proses Tokenisasi	39
Tabel 4.5 Tabel Setelah Token Diurutkan dan Digabungkan.....	39
Tabel 4.6 Tabel <i>Clean</i>	44
Tabel 4.7 Deskripsi Aksi Aktor dan Respon Sistem Untuk <i>Use Case Login</i>	46
Tabel 4.8 Deskripsi Aksi Aktor dan Respon Sistem Untuk <i>Use Case Duplicate Detection</i>	47
Tabel 4.9 Deskripsi Aksi Aktor dan Respon Sistem Untuk <i>Use Case View Duplicate Detection Result</i>	48
Tabel 4.10 Deskripsi Aksi Aktor dan Respon Sistem Untuk <i>Use Case Import Data</i>	48
Tabel 4.11 Deskripsi Aksi Aktor dan Respon Sistem Untuk <i>Use Case Fixing Contact Number</i>	49
Tabel 4.12 Deskripsi Aksi Aktor dan Respon Sistem Untuk <i>Use Case Save Fix Result</i>	50
Tabel 4.13 Deskripsi Aksi Aktor dan Respon Sistem Untuk <i>Use Case Edit Export Data</i>	50
Tabel 4.14 Tabel Hasil Pengujian Algoritma SNM – Fungsi <i>cleanData()</i>	63
Tabel 4.15 Tabel Hasil Pengujian Algoritma SNM – Fungsi <i>ChopTheWords()</i> ..	64
Tabel 4.16 Tabel Hasil Pengujian Algoritma SNM – Fungsi <i>gramming()</i>	67

Tabel 4.17 Tabel Hasil Pengujian <i>Functionality</i>	73
Tabel 4.18 Tabel Hasil Pengujian <i>Usability</i>	75
Tabel 4.19 Tabel Hasil Pengujian <i>Compatibility</i>	75
Tabel 4.20 Tabel Hasil Pengujian <i>Performance</i>	76

DAFTAR RUMUS

Rumus 2.1 Rumus <i>N-Gram</i> Untuk Menghitung Kemiripan Antar <i>String</i>	14
Rumus 4.1 Rumus <i>precision</i> , <i>recall</i> , dan <i>f-measure</i>	78

DAFTAR LAMPIRAN

Lampiran 1 – Wawancara	89
Lampiran 2 – <i>Requirement Elicitation</i>	93
Lampiran 3 – <i>Software Requirement Specification</i>	98
Lampiran 4 – <i>Template Import Data</i>	128
Lampiran 5 – Hasil Pengujian Sistem	129
Lampiran 6 – Hasil Percobaan Data	135

DAFTAR SINGKATAN

IDE	Integrated Development Environment
IGASIS	Intra-Governmental Access To Shared Information System
KDD	Knowledge Discovery in Databases
OOP	Object Oriented Programming
SNM	Sorted Neighbourhood Method
UML	Unified Modelling Language