



**Klasterisasi Data *Wholesale Customers* Menggunakan
*K-Means++ Clustering***

SKRIPSI

**untuk memenuhi persyaratan dalam menyelesaikan
program sarjana Strata-1 Statistika**

Oleh

ERIEN SYARIF

NIM. 1811017320024

PROGRAM STUDI STATISTIKA

FAKULTAS MATEMATIKA DAN ILMU PENGETAHUAN ALAM

UNIVERSITAS LAMBUNG MANGKURAT

BANJARBARU

JUNI 2023



**Klasterisasi Data *Wholesale Customers* Menggunakan
*K-Means++ Clustering***

SKRIPSI

**untuk memenuhi persyaratan dalam menyelesaikan program sarjana
Strata-1 Statistika**

Oleh

ERIEN SYARIF

NIM. 1811017320024

PROGRAM STUDI STATISTIKA

FAKULTAS MATEMATIKA DAN ILMU PENGETAHUAN ALAM

UNIVERSITAS LAMBUNG MANGKURAT

BANJARBARU

JUNI 2023

SKRIPSI

**KLASTERISASI DATA *WHOLESALE CUSTOMERS* MENGGUNAKAN
*K-MEANS++ CLUSTERING***

Oleh:
ERIEN SYARIF
1811017320024

telah dipertahankan di depan Dosen Penguji pada tanggal 09 Juni 2023

Susunan Dosen Penguji:

Pembimbing I



Yuana Sukmawaty, S.Si., M.Si
NIP 198810152015042002

Penguji I:



Oni Soesanto, S.Si., M.Si
NIP 197301262005011003

Pembimbing II



Irwan Budiman, S.T., M.Kom.
NIP 197703252008121001

Penguji II



I Putu Edy Suardiyana Putra, M.Kom., Ph.D

Banjarbaru, 23 Juni 2023

Mengetahui,
Koordinator Program Studi Statistika



Devi Anggraini, S. Si., M.App.Sci., Ph.D
NIP 198303282005012001

PERNYATAAN

Saya menyatakan bahwa dalam skripsi ini tidak terdapat karya yang pernah diajukan untuk memperoleh gelar sarjana di suatu Perguruan Tinggi dan sepanjang pengetahuan serta kesadaran Saya juga tidak terdapat karya atau pendapat yang pernah ditulis atau diterbitkan oleh orang lain, kecuali yang secara tertulis diacu dalam naskah ini dan disebutkan dalam Daftar Pustaka.

Banjarbaru,

Juni 2023



Erien Syarif

NIM 1811017320024

ABSTRAK

Klasterisasi Data *Wholesale Customers* Menggunakan *K-Means++ Clustering* (Oleh: Erien Syarif; Pembimbing: Yuana Sukmawaty, S.Si., M.Si. dan Irwan Budiman, S.T., M.Kom., 2023; 56 halaman)

Segmentasi pelanggan merupakan proses pengelompokan pelanggan berdasarkan karakteristik dan perilaku pelanggan. Segmentasi dilakukan dengan menganalisis berbagai jenis informasi tentang pelanggan seperti demografis, geografis, dan perilaku. Manfaat segmentasi pelanggan, yaitu mempelajari dan memahami perilaku pelanggan. Penelitian ini bertujuan mengetahui jumlah kluster (k) optimal dari hasil klasterisasi dan menganalisis karakteristik pelanggan dari kluster yang terbentuk. Data yang digunakan adalah data *Wholesale Customers* yang didapatkan dari *UCI Machine Learning Repository*. Data tersebut berisi pengeluaran tahunan pelanggan dari distributor grosir untuk berbagai barang. Metode klasterisasi dalam penelitian ini menggunakan algoritma *K-Means++ Clustering* dan jumlah kluster (k) optimal didapatkan menggunakan metode *Silhouette Coefficient*. Hasil penelitian menunjukkan bahwa k optimal yang didapatkan adalah $k = 2$ dengan nilai *Silhouette Coefficient* sebesar 0.7390. Penelitian ini menghasilkan 2 kluster, yaitu “Pelanggan dengan Pengeluaran Besar” dan “Pelanggan dengan Pengeluaran Kecil”. Jumlah pelanggan pada kluster 1 dan kluster 2 sebanyak 14 dan 426 pelanggan. Kluster 1 memiliki rata-rata pengeluaran lebih besar pada setiap atribut dibandingkan dengan kluster 2. Pengeluaran terbesar pada kluster 1 dan kluster 2 terdapat pada atribut *Milk* dengan rata-rata pengeluaran sebesar 33847.79 dan atribut *Fresh* dengan rata-rata pengeluaran sebesar 11531.44.

Kata kunci: *K-Means++ Clustering*, Segmentasi Pelanggan, *Silhouette Coefficient*

ABSTRACT

WHOLESALE CUSTOMERS CLUSTERING USING *K-MEANS++ CLUSTERING*
(By Erien Syarif; Supervisor: Yuana Sukmawaty, S.Si., M.Si. and Irwan Budiman, S.T., M.Kom., 2023; 56 pages)

Customer segmentation is the process of grouping customers based on customer characteristics and behavior. Segmentation is done by analyzing different types of information about customers such as demographic, geographic, and behavioral. The benefits of customer segmentation are learning and understanding customer behavior. This study aims to determine the optimal number of clusters (k) from the clustering results and analyze the customer characteristics of the clusters formed. The data used is *Wholesale Customers* obtained from the *UCI Machine Learning Repository*. The data contains the annual spending of customers of wholesale distributors for various goods. The clustering method in this study uses the *K-Means++ Clustering algorithm* and the optimal k is obtained using the *Silhouette Coefficient method*. The results showed that the optimal k obtained was $k = 2$ with a *Silhouette Coefficient* value of 0.7390. This research resulted in 2 clusters, namely "Customers with High Spend" and "Customers with Low Spend". The number of customers in cluster 1 and cluster 2 is 14 and 426 customers. Cluster 1 has a higher average spend on each attribute compared to cluster 2. The largest spending in cluster 1 and cluster 2 is in the *Milk* attribute with an average spending of 33847.79 and the *Fresh* attribute with an average spending of 11531.44.

Keywords: *Customer Segmentation, K-Means++ Clustering, Silhouette Coefficient*

PRAKATA

Alhamdulillah puji syukur dipanjatkan kepada Allah Subhanahu wa Ta'ala yang telah memberikan kesehatan, kekuatan dan kemudahan kepada peneliti, sehingga dapat menyelesaikan penelitian ini. Shalawat serta salam juga tidak lupa kita panjatkan kepada Nabi Muhammad Shallallahu 'alaihi wa sallam serta keluarga dan para sahabat yang telah membawa zaman kegelapan menuju zaman yang penuh cahaya.

Penelitian yang berjudul "**Klasterisasi Data *Wholesale Customers* Menggunakan *K-Means++ Clustering***" ini telah diselesaikan untuk memenuhi salah satu persyaratan dalam menyelesaikan program sarjana di Program Studi Statistika Fakultas MIPA Universitas Lambung Mangkurat.

Penulisan skripsi ini tidak akan tercapai tanpa bantuan dan dukungan dari banyak pihak. Oleh karena itu, penulis ingin mengucapkan terima kasih yang sebesar-besarnya kepada berbagai pihak, yaitu:

1. Orang tua yang telah memberikan segalanya;
2. Ibu Yuana Sukmawaty, S.Si., M.Si selaku dosen pembimbing 1 dan dosen pembimbing akademik serta Bapak Irwan Budiman, S.T., M.Kom. selaku dosen pembimbing 2 yang telah membantu, mengarahkan, memberikan masukan dan membimbing dalam menyelesaikan skripsi;
3. Koordinator Program Studi beserta seluruh jajaran dosen dan staf Program Studi Statistika Fakultas Matematika dan Ilmu Pengetahuan Alam Universitas Lambung Mangkurat (FMIPA ULM);
4. Alpi, Dini, Kethy, Lalu, Umi, Nela, Inal, Nawiw, Aldi, Tomi, Geo, Ardi dan juga teman-teman angkatan 2018 FMIPA ULM;
5. Rekan-rekan pengurus BEM FMIPA ULM Periode 2020 yang telah menjadi tempat untuk mengembangkan diri;
6. Berbagai pihak yang telah turut serta membantu yang tidak dapat disebut satu persatu.

Penelitian ini jauh dari kata sempurna. saran, masukan dan kritik sangat diharapkan untuk perbaikan selanjutnya. Akhir kata penulis berharap semoga skripsi ini dapat memberikan manfaat bagi pembacanya.

Banjarbaru, Juni 2023

Erien Syarif

DAFTAR ISI

COVER JUDUL	i
LEMBAR PENGESAHAN	ii
PERNYATAAN	iii
ABSTRAK	iv
ABSTRACT	v
PRAKATA	vi
DAFTAR ISI	vii
DAFTAR TABEL	ix
DAFTAR GAMBAR	x
DAFTAR LAMPIRAN	xi
BAB I	1
1.1 Latar Belakang	1
1.2 Rumusan Masalah	3
1.3 Tujuan Penelitian	3
1.4 Manfaat Penelitian	3
BAB II	4
2.1 Kajian Penelitian Terdahulu	4
2.2 Kajian Teori	7
2.2.1 <i>Data Mining</i>	7
2.2.2 Normalisasi Data	9
2.2.3 Analisis Kluster	9
2.2.4 <i>K-Means++ Clustering</i>	12
2.2.5 <i>Silhouette Coefficient</i>	14
2.2.6 Segmentasi Pelanggan	15
BAB III	16
3.1 Sumber Data	16
3.2 Atribut Penelitian	16
3.3 Metode Penelitian	17
3.4 Prosedur Penelitian	18
3.5 Alur Penelitian	20
BAB IV	21
4.1 <i>Preprocessing Data</i>	21

4.1.1	Menghapus Atribut yang Tidak Terlibat dalam Klasterisasi.....	21
4.1.2	Mendeteksi <i>Missing Value</i> , Duplikasi Data, dan <i>Outlier</i>	22
4.2	Transformasi Data.....	23
4.3	Pengujian Asumsi Multikolinieritas.....	24
4.4	Klasterisasi Menggunakan <i>K-Means++ Clustering</i>	25
4.5	Evaluasi Hasil Klaster Menggunakan <i>Silhouette Coefficient</i>	32
4.6	Analisis Hasil Klaster	33
BAB V	36
5.1	KESIMPULAN	36
5.2	SARAN	36
DAFTAR PUSTAKA	37
LAMPIRAN	40
RIWAYAT HIDUP	59

DAFTAR TABEL

Tabel 1. Kajian Penelitian Terdahulu	4
Tabel 2. Nilai Silhouette Coefficient.....	15
Tabel 3. Atribut Penelitian.....	16
Tabel 4. Data Wholesale Customers	21
Tabel 5. Data Tanpa Atribut Channel dan Region	22
Tabel 6. Hasil Normalisasi Data	24
Tabel 7. Data Tanpa Atribut Grocery	26
Tabel 8. Perhitungan untuk Mendapatkan C2	27
Tabel 9. Centroid pertama.....	28
Tabel 10. Hasil Perhitungan Jarak dan Pengelompokkan Data Iterasi ke-1	28
Tabel 11. Centroid iterasi ke-1	29
Tabel 12. Nilai Selisih Centroid Iterasi ke-1 dan Centroid Pertama.....	30
Tabel 13. Hasil Perhitungan Jarak dan Pengelompokkan Data Iterasi ke-2	30
Tabel 14. Centroid Iterasi ke-2.....	31
Tabel 15. Nilai Selisih Centroid Iterasi ke-2 dan Centroid Iterasi ke-1.....	31
Tabel 16. Centroid Iterasi ke-9.....	31
Tabel 17. Hasil Klaster yang Terbentuk.....	31

DAFTAR GAMBAR

Gambar 1. Proses KDD	7
Gambar 2. Flowchart Alur Penelitian	20
Gambar 3. Grafik Box-Plot Data	23
Gambar 4. Grafik dan Nilai Silhouette Coefficient dari Hasil Kluster Data dengan Normalisasi Data	33
Gambar 5. Jumlah Pelanggan pada Setiap Kluster	34
Gambar 6. Perbandingan Rata-Rata Pengeluaran Pelanggan Setiap Kluster	34
Gambar 7. Pengeluaran Pelanggan untuk Setiap Atribut	35

DAFTAR LAMPIRAN

Lampiran 1 Tampilan Website UCI Machine Learning Repository	40
Lampiran 2 Data <i>Wholesale Customers</i>	41
Lampiran 3 Input <i>Google Colaboratory</i> untuk <i>Preprocessing Data</i>	42
Lampiran 4 Input <i>Google Colaboratory</i> untuk Transformasi	43
Lampiran 5 Hasil Transformasi Menggunakan Logaritma Natural	44
Lampiran 6 Hasil Transformasi Menggunakan Akar Kuadrat.....	45
Lampiran 7 Input <i>Google Colaboratory</i> untuk <i>Z-Score Normalization</i>	46
Lampiran 8 Hasil <i>Z-Score Normalization</i> Data Tanpa Transformasi	47
Lampiran 9 Hasil <i>Z-Score Normalization</i> Data Transformasi Logaritma Natural	48
Lampiran 10 Hasil <i>Z-Score Normalization</i> Data Transformasi Akar Kuadrat.....	49
Lampiran 11 Input <i>Google Colaboratory</i> untuk Asumsi Korelasi.....	50
Lampiran 12 Nilai Korelasi Data dengan Sebelum dan Sesudah Melakukan Normalisasi Data	51
Lampiran 13 Nilai Korelasi Data dengan Transformasi Logaritma Natural serta Data dengan Transformasi Logaritma Natural dan dilanjutkan dengan Normalisasi Data	52
Lampiran 14 Nilai Korelasi Data dengan Transformasi Akar Kuadrat serta Data dengan Transformasi Akar Kuadrat dan dilanjutkan dengan Normalisasi Data	53
Lampiran 15 Input <i>Google Colaboratory</i> untuk Klasterisasi	54
Lampiran 16 Nilai <i>Silhouette Coefficient</i> untuk Semua Atribut.....	55
Lampiran 17 Nilai <i>Silhouette Coefficient</i> Tanpa Atribut " <i>Grocery</i> "	56
Lampiran 18 Nilai <i>Silhouette Coefficient</i> Tanpa Atribut " <i>Detergents_Paper</i> "	57
Lampiran 19 Tabel Hasil Clustering	58