

Nama : Imron Bagus Sajiwo

Npm : 5230311037

Prodi : Sistem Informasi A

1. Apa perbedaan data cleaning & data reduction berikan penjelasan dan tuliskan contohnya.

Jawab:

A. Data Cleaning (Pembersihan Data)

Data cleaning adalah proses mengidentifikasi dan memperbaiki kesalahan atau inkonsistensi dalam data sehingga dapat meningkatkan kualitas data untuk analisis lebih lanjut. Kesalahan yang sering diperbaiki mencakup data yang hilang, data duplikat, data tidak relevan, atau data yang tidak sesuai format.

Contoh:

ID Pelanggan	Nama	Usia	Email	Kota
101	Andi	25	andi@email.com	Jakarta
102	Budi	30	NULL	Bandung
103	Siti	-5	siti@email.com	Surabaya
104	Budi	30	budi@email.com	Bandung

- Data "NULL" pada email Budi diperbaiki dengan mengisi data yang valid.
- Usia -5 adalah nilai tidak valid dan perlu diperbaiki atau dihapus.
- Data duplikat untuk Budi perlu dihapus atau dikoreksi.

B. Data Reduction (Reduksi Data)

Data reduction adalah proses menyederhanakan atau mengurangi ukuran dataset tanpa kehilangan informasi yang signifikan. Ini dilakukan untuk meningkatkan efisiensi penyimpanan dan kecepatan analisis data.

Contoh:

Dataset awal memiliki 10 fitur (kolom), tetapi setelah analisis, hanya 5 fitur yang paling relevan dipilih.

Sebelum reduksi:

ID	Nama	Usia	Email	Kota	Penghasilan	Status	Pendidikan	Hobi	Tanggal Daftar
1	Andi	25	andi@email.com	Jakarta	10 Juta	Menikah	S1	Sepak bola	6/12/2021

Setelah reduksi (hanya fitur penting yang dipilih):

ID	Nama	Usia	Kota	Penghasilan
1	Andi	25	Jakarta	10 Juta

2. Carilah data dari web yang sesuai dengan program studi saudara, jelaskan isinya, dan apa kontribusi data tersebut terhadap program studi saudara, lakukan pengolahan awal agar data tersebut siap diolah, apa saja yang saudara lakukan. berikan penjelasan caranya.

Jawab:

A. Isi Dataset

Dataset ini berisi **253.680 entri** dengan **22 kolom** yang mencatat indikator kesehatan terkait penyakit jantung dari survei BRFSS (Behavioral Risk Factor Surveillance System) tahun 2015.

Beberapa Kolom Utama dalam Dataset:

- HeartDiseaseorAttack – Indikator apakah seseorang pernah mengalami serangan jantung atau penyakit jantung.
- HighBP (Tekanan Darah Tinggi) – Menunjukkan apakah seseorang memiliki tekanan darah tinggi.
- HighChol (Kolesterol Tinggi) – Menunjukkan apakah seseorang memiliki kadar kolesterol tinggi.
- BMI (Indeks Massa Tubuh) – Mengukur apakah seseorang memiliki berat badan ideal, kurang, atau berlebih.
- Smoker (Perokok) – Indikator apakah seseorang pernah merokok minimal 100 batang dalam hidupnya.
- Stroke – Riwayat apakah seseorang pernah mengalami stroke.
- Diabetes – Indikator apakah seseorang memiliki diabetes.
- PhysActivity (Aktivitas Fisik) – Menunjukkan apakah seseorang melakukan aktivitas fisik secara rutin.
- Fruits & Veggies (Konsumsi Buah dan Sayur) – Menunjukkan kebiasaan makan sehat seseorang.

- HvyAlcoholConsump (Konsumsi Alkohol Berlebihan) – Menunjukkan apakah seseorang mengonsumsi alkohol dalam jumlah berlebihan.
- AnyHealthcare (Akses Layanan Kesehatan) – Menunjukkan apakah seseorang memiliki akses ke layanan kesehatan.
- NoDocbcCost (Kesulitan Bertemu Dokter karena Biaya) – Indikator seseorang tidak bisa menemui dokter karena kendala biaya.
- GenHlth, MentHlth, PhysHlth – Penilaian terhadap kondisi kesehatan umum, kesehatan mental, dan kesehatan fisik seseorang.
- Demografi (Sex, Age, Education, Income) – Informasi tentang jenis kelamin, usia, tingkat pendidikan, dan pendapatan seseorang.

Row No.	HeartDiseas...	HighBP	HighChol	CholCheck	BMI	Smoker	Stroke	Diabetes	PhysActivity	Fruits	Veggies
9	1	1	1	1	30	1	0	2	0	1	1
10	0	0	0	1	24	0	0	0	0	0	1
11	0	0	0	1	25	1	0	2	1	1	1
12	0	1	1	1	34	1	0	0	0	1	1
13	0	0	0	1	26	1	0	0	0	0	1
14	0	1	1	1	28	0	0	2	0	0	1
15	0	0	1	1	33	1	1	0	1	0	1
16	0	1	0	1	33	0	0	0	1	0	0
17	0	1	1	1	21	0	0	0	1	1	1
18	0	0	0	1	23	1	0	2	1	0	0
19	0	0	0	0	23	0	0	0	0	0	1
20	0	0	1	1	28	0	0	0	0	0	0
21	1	1	1	1	22	0	1	0	0	1	0
22	0	1	1	1	38	1	0	0	0	1	1
23	0	0	0	1	28	1	0	0	0	0	1
24	0	1	0	1	27	0	0	2	1	1	1
25	0	1	1	1	28	1	0	0	0	1	1
26	0	0	0	1	19	0	0	0	1	1	1

Penjelasan:

penjelasan diolah datanya saya pertama – tama memukan Rapid miner kemudian saya memberikan operator read CSV kemudian saya import data tersebut kemudian saya menyambungkan datanya lalu saya running untuk pengolahan awal.

3. Setelah pengolahan awal saudara lakukan, terapkan Operator PCA pada data saudara gunakan beberapa pilihan jumlah atributnya. Jelaskan langkah-langkahnya dan tuliskan contoh hasilnya.

Jawab:

Contoh Hasil sesudah di PCA

ExampleSet (PCA)

Open in: Turbo Prep, Auto Model, Interactive Analysis

Filter (253,680 / 253,680 examples):

Row No.	pc_1	pc_2	pc_3	pc_4	pc_5	pc_6	pc_7	pc_8	pc_9	pc_10	pc_11
1	-4.525	-0.821	-0.405	1.118	-0.136	1.093	-1.495	-0.831	-0.228	-0.459	0.078
2	-0.324	-5.452	1.643	-3.090	0.711	-2.418	1.157	-1.052	3.055	-1.658	-0.452
3	-4.636	-2.211	-2.060	1.694	-0.891	0.725	-0.195	1.222	1.168	-2.296	1.187
4	0.423	0.602	-0.508	-0.879	2.067	0.184	-0.612	-0.106	-0.318	0.532	-0.526
5	0.228	1.231	-0.902	-0.767	1.448	-0.103	-0.700	0.636	0.729	-0.792	-0.366
6	0.839	1.959	0.244	-0.624	-0.988	-0.740	-0.377	-0.030	-0.414	-1.147	0.303
7	-0.877	-0.187	1.388	0.994	-0.062	1.870	-0.094	-0.192	0.573	-0.348	1.111
8	-1.753	0.762	0.189	-1.108	0.799	0.871	-1.006	0.073	-0.088	-0.799	-0.585
9	-6.443	0.105	-2.530	0.393	-1.243	0.315	-0.057	-0.648	-0.399	-1.006	-0.170
10	0.312	-0.952	1.270	-0.185	1.171	1.078	0.486	-0.423	-0.700	0.466	1.005
11	0.592	1.670	0.009	-0.296	-0.638	-0.698	0.113	-0.847	-1.123	-0.094	-0.319
12	-3.589	0.196	-1.534	-0.047	0.563	0.961	-1.091	-1.382	-0.313	-0.380	1.171
13	-0.100	-0.949	0.092	0.310	-0.300	1.757	-0.124	-0.751	-0.520	0.055	1.141
14	-2.840	1.202	0.316	0.999	1.550	0.267	-0.958	-0.596	1.116	0.057	0.590
15	-3.448	-2.960	-2.369	0.055	-2.695	-0.186	1.895	3.075	1.000	0.582	-1.966
16	0.944	-0.253	1.181	1.896	0.341	0.645	0.228	1.030	1.168	-0.008	-0.611
17	-0.315	0.880	-0.766	-1.259	1.930	0.070	-0.694	0.605	0.304	-0.647	-0.569
18	0.310	-0.174	2.482	0.218	-0.082	0.495	0.670	0.370	-0.896	-0.148	-1.296

ExampleSet (253,680 examples, 0 special attributes, 15 regular attributes)

Penjelasan:

Pertama tama Saya Import Data dari Read CSV kemudian saya memberikan Select Atributies Untuk Memilih Data yang akan dianalisis tetapi saya lupa dalam memilih data nya jadinya saya melakukan all Data jadinya data nya dianalisi semua, kemudian saya memberikan noralmalize untuk menormalkan data kemudian saya memberikan PCA tersebut agar data dapat dianalisis, setelah itu saya merunning data tersebut.