

# Tugas Kelompok Data Mining

## Kelas C

Kelompok : 4

Anggota kelompok:

1. Imam Mubarak.A (60900123064)
2. Habil Nasruddin (60900123072)

### Step by Step Proses Pengolahan Data

Nama Dataset : penjualan\_barang.csv

Sumber URL : <https://www.kaggle.com/datasets/bejopamungkas/transaksi-pembelian-penjualan-sembako?select=penjualan+barang.csv>

Pengolahan data ini dilakukan melalui empat tahap utama: Collecting, Integrasi, Analisis, dan Validasi. Setiap tahap menghasilkan file output tersendiri agar proses dapat ditelusuri ulang (data lineage) dan setiap perubahan dataset terdokumentasi dengan baik. Dataset yang digunakan adalah penjualan\_barang.csv, berisi catatan transaksi penjualan produk.

#### 1. Collecting Data (Pengumpulan Data)

```
D: > Semester 5 > Data Mining > Tahap_2_DataMining > 📡 Collecting_Data.py
 1 import pandas as pd
 2
 3 # === Tahap 1: Collecting Data ===
 4 try:
 5     # Membaca dataset asli
 6     df = pd.read_csv("penjualan_barang.csv")
 7
 8     # Simpan hasil tahap collecting
 9     df.to_csv("data_collecting.csv", index=False)
10
11     print("✓ Tahap 1 (Collecting) selesai → data_collecting.csv berhasil dibuat")
12 except Exception as e:
13     print("✗ Error di tahap Collecting:", e)
```

Tahap awal adalah mengumpulkan data mentah dari sumber yang tersedia.

- Sumber data: file penjualan\_barang.csv.
- Dataset ini berisi informasi transaksi penjualan dengan atribut antara lain:
  - tanggal\_transaksi → tanggal transaksi penjualan.
  - id\_barang → kode unik barang.

- nama\_barang → nama produk yang dijual.
- jumlah → jumlah unit barang yang dibeli.
- harga\_satuan → harga satuan barang.
- total\_harga → hasil perkalian jumlah × harga satuan.

Langkah yang dilakukan:

- Membaca file CSV menggunakan library pandas.
- Menyimpan ulang data mentah ke file data\_collecting.csv tanpa perubahan.

Tujuan:

- Menyediakan salinan awal dataset sebelum dilakukan pembersihan.
- Jika ada kesalahan pada tahap berikutnya, data mentah ini bisa menjadi referensi ulang.

Output: data\_collecting.csv.

## 2. Integrasi Data (Data Integration)

```

1  import pandas as pd
2
3  # === Tahap 2: Integrasi Data ===
4  try:
5      df = pd.read_csv("data_collecting.csv")
6
7      # Hapus baris duplikat
8      df_integration = df.drop_duplicates()
9
10     # Normalisasi nama kolom
11     df_integration.columns = [col.strip().lower().replace(" ", "_") for col in df_integration.columns]
12
13     # Simpan hasil integrasi
14     df_integration.to_csv("data_integration.csv", index=False)
15
16     print("✓ Tahap 2 (Integrasi) selesai → data_integration.csv berhasil dibuat")
17 except Exception as e:
18     print("✗ Error di tahap Integrasi:", e)

```

Setelah data dikumpulkan, tahap berikutnya adalah integrasi dan standarisasi data agar lebih konsisten.

Langkah yang dilakukan:

- Menghapus baris duplikat untuk menghindari penghitungan ganda.
- Menstandarkan nama kolom: semua huruf kecil dan spasi diganti dengan underscore (\_).
- Mengecek konsistensi tipe data dasar (misalnya, kolom jumlah & harga tetap numerik).

Tujuan:

- Memastikan dataset rapi, konsisten, dan siap dipakai untuk analisis.
- Tahap integrasi penting terutama jika data berasal dari berbagai sumber

(misalnya beberapa file penjualan dari cabang berbeda).

Output: data\_integration.csv.

### 3. Analisis Data (Data Analysis)

```
1 import pandas as pd
2
3 # === Tahap 3: Analisis Data ===
4 try:
5     df_integration = pd.read_csv("data_integration.csv")
6
7     # Statistik deskriptif (tanpa datetime_is_numeric biar kompatibel semua versi pandas)
8     desc_stats = df_integration.describe(include="all")
9
10    # Jumlah missing values
11    missing = df_integration.isnull().sum()
12
13    # Buat ringkasan analisis
14    analysis = pd.DataFrame({
15        "missing_values": missing,
16        "data_type": df_integration.dtypes.astype(str)
17    })
18
19    # Gabungkan dengan statistik deskriptif
20    analysis = analysis.join(desc_stats.transpose(), how="left")
21
22    # Simpan hasil analisis
23    analysis.to_csv("data_analysis.csv")
24
25    print("✓ Tahap 3 (Analisis) selesai → data_analysis.csv berhasil dibuat")
26 except Exception as e:
27     print("✗ Error di tahap Analisis:", e)
```

Pada tahap ini, dilakukan analisis eksplorasi awal untuk memahami isi dataset.

Langkah yang dilakukan:

- Menghitung statistik deskriptif: rata-rata, median, minimum, maksimum, dan standar deviasi pada kolom numerik.
- Mengidentifikasi jumlah missing values pada setiap kolom.
- Menyimpan tipe data (int, float, object) setiap kolom untuk memastikan kesesuaian.
- Menyimpan ringkasan analisis ke file CSV.

Contoh hasil analisis:

- Rata-rata kolom jumlah menunjukkan barang paling sering dibeli dalam jumlah kecil (1–5 unit).
- Kolom harga\_satuan bervariasi, ada produk murah dan produk premium.
- Terdapat missing values pada total\_harga karena beberapa baris belum terisi dengan benar.

Tujuan:

- Memberikan gambaran umum dataset.
- Menentukan masalah kualitas data yang harus ditangani di tahap validasi.

Output: data\_analysis.csv.

### 4. Validasi Data (Data Validation)

```

1 import pandas as pd
2
3 # === Tahap 4: Validasi Data ===
4 try:
5     df_integration = pd.read_csv("data_integration.csv")
6
7     # Hapus baris yang semua kolomnya kosong
8     df_validation = df_integration.dropna(how="all")
9
10    # Isi missing values: numerik → median, kategorikal → mode
11    for col in df_validation.columns:
12        if df_validation[col].dtype in ["int64", "float64"]:
13            median_val = df_validation[col].median()
14            df_validation[col] = df_validation[col].fillna(median_val)
15        else:
16            if not df_validation[col].mode().empty:
17                mode_val = df_validation[col].mode()[0]
18                df_validation[col] = df_validation[col].fillna(mode_val)
19
20    # Simpan hasil validasi
21    df_validation.to_csv("data_validation.csv", index=False)
22
23    print("✓ Tahap 4 (Validasi) selesai → data_validation.csv berhasil dibuat")
24 except Exception as e:
25     print("✗ Error di tahap Validasi:", e)

```

Tahap validasi bertujuan memastikan dataset final bersih dan bebas masalah.

Langkah yang dilakukan:

- Menghapus baris yang seluruh kolomnya kosong.
- Mengisi nilai kosong (imputasi):
  - Numerik → diganti dengan median (agar tidak bias terhadap nilai ekstrem).
  - Kategorikal → diganti dengan mode (nilai yang paling sering muncul).
- Memastikan tidak ada lagi nilai NaN setelah proses validasi.

Tujuan:

- Menghasilkan dataset final yang bisa langsung dipakai untuk analisis lanjutan, machine learning, atau visualisasi.
- Menjamin kualitas data agar hasil analisis tidak menyesatkan.

Output: data\_validation.csv.

## Kesimpulan

Proses pengolahan data dari penjualan\_barang.csv telah mengikuti alur standar:

1. Dataset mentah disimpan ulang di data\_collecting.csv.
2. Dataset dibersihkan dan disatukan di data\_integration.csv.
3. Hasil analisis eksplorasi disimpan di data\_analysis.csv.
4. Dataset final yang sudah tervalidasi disimpan di data\_validation.csv.

Dengan menyimpan output di setiap tahap, proses ini memenuhi prinsip data lineage, sehingga setiap perubahan dapat dilacak dan dikembalikan ke versi sebelumnya jika diperlukan.

## **Tabel Ringkasan**

<b>Tahap</b>	<b>File Hasil</b>	<b>Perubahan Utama</b>
Collecting	data_collecting.csv	Salinan dataset mentah dari penjualan_barang
Integrasi	data_integration.csv	Hapus duplikat, standarisasi nama kolom
Analisis	data_analysis.csv	Statistik deskriptif, missing values, tipe data
Validasi	data_validation.csv	Isi missing values (median/mode), data final