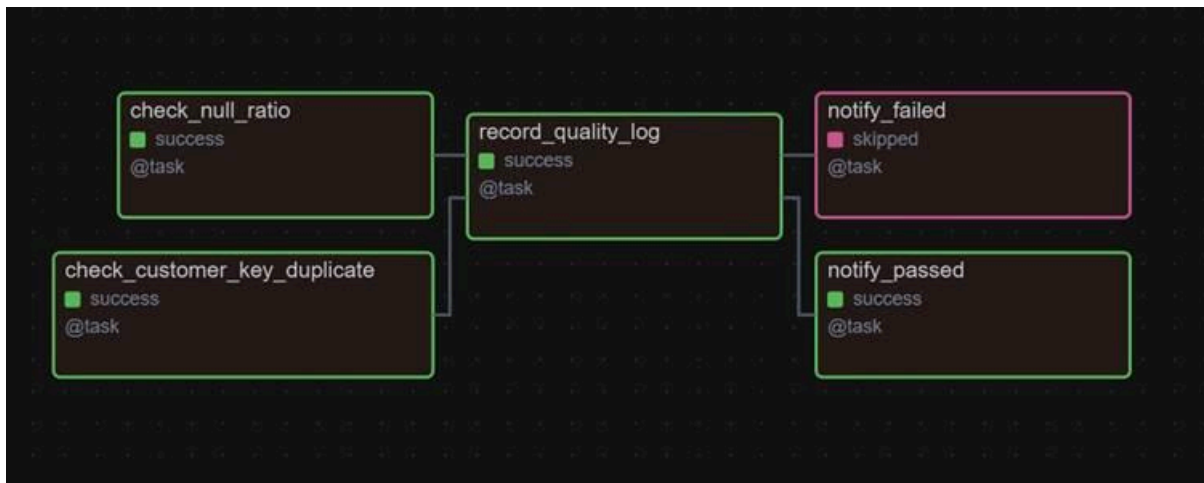


DOKUMENTASI IMPLEMENTASI AIRFLOW

1. Diagram arsitektur



Gambar 1.1 Diagram arsitektur `etl_duckdb_dag.py`



Gambar 1.2 Diagram arsitektur `Dag_data_quality.py`

2. Deskripsi tujuan dari setiap tugas

2.1 Diagram arsitektur `etl_duckdb_dag.py`

Proses ETL diatur dalam sebuah DAG (Directed Acyclic Graph) yang terdiri atas rangkaian tugas (task) yang dijalankan secara berurutan. DAG ini mencerminkan urutan logis proses ETL:

2.1.1 check_data : Memastikan ketersediaan file data mentah (data.csv) dari Google Drive atau lokal.

2.1.2 init: Menginisialisasi struktur tabel di dalam DuckDB, mencakup tabel dimensi (dim_customer, dim_card, dll.) dan tabel fakta (fact_transactions).

2.1.3 extract: Mengekstrak data dari file CSV.

2.1.4 transform: Mengolah data mentah menjadi bentuk tabel dimensi dan fakta yang telah distandarkan.

2.1.5 load: Memasukkan data hasil transformasi ke warehouse DuckDB.

Setelah proses utama selesai, sistem secara otomatis menentukan jalur notifikasi:

- `notify_success` akan mengirim email ke stakeholder jika seluruh proses berhasil.
- `notify_failure` akan aktif jika terjadi kegagalan pada salah satu task dalam pipeline.

2.2 Diagram arsitektur Dag_data_quality.py

DAG ini terdiri dari beberapa task penting yang berjalan secara paralel dan dilanjutkan ke pencatatan dan pengiriman notifikasi. Task-task tersebut adalah:

2.2.1 check_null_ratio : Task ini bertugas menghitung rasio nilai null di setiap kolom dalam tabel fakta `fact_transactions`. Jika rasio nilai null pada suatu kolom melebihi ambang batas tertentu (misalnya 10%), task ini akan gagal dan memicu mekanisme peringatan otomatis.

2.2.2 check_customer_key_duplicate : Fokus pada integritas data di tabel dimensi `dim_customer`, task ini memastikan bahwa `customer_key` bersifat unik dan tidak terjadi duplikasi. Duplikasi kunci pelanggan bisa menyebabkan kesalahan dalam proses join atau analisis downstream.

2.2.3 record_quality_log : Setelah dua pemeriksaan di atas berhasil, hasilnya dicatat dalam log kualitas data. Log ini dapat disimpan dalam bentuk file CSV atau format lainnya, dan menjadi bahan monitoring atau audit historis dari validitas data setiap harinya.

2.2.4 notify_passed dan notify_failed : Dua task ini bertugas mengirimkan notifikasi email kepada pihak terkait. `notify_passed` akan dikirimkan apabila semua pemeriksaan berhasil. Sebaliknya, `notify_failed` akan dikirim jika salah satu pemeriksaan gagal.

Evaluasi Hasil Eksekusi : Berdasarkan visualisasi DAG yang ditampilkan, dapat disimpulkan bahwa semua task utama telah berjalan dengan sukses. Status success berwarna hijau pada semua task menunjukkan bahwa:

- Tidak ada kolom dengan nilai null melebihi threshold
- Tidak ditemukan duplikasi pada `customer_key`
- Log kualitas data telah berhasil dicatat
- Notifikasi keberhasilan berhasil dikirimkan
- Task notifikasi kegagalan `notify_failed` otomatis di-skip karena tidak relevan

Hal ini menandakan bahwa kualitas data saat ini berada dalam kondisi yang baik dan dapat digunakan dengan aman untuk analisis atau visualisasi lebih lanjut.

3. Informasi penjadwalan dan dependensi

DAG dijadwalkan untuk berjalan setiap hari (@daily) dengan pengaturan catchup=False, artinya hanya akan mengeksekusi untuk hari-hari berjalan dan tidak akan memproses data backlog. Ketergantungan antar task didefinisikan secara eksplisit menggunakan TaskFlow API dan operator >>.

4. Pengaturan pemantauan dan peringatan

Pengawasan sistem dilakukan dengan pendekatan kombinasi:

4.1 Airflow Web UI memberikan gambaran visual dari alur dan status masing-masing task.

4.2 Sistem email notifikasi dikonfigurasi untuk mengirim peringatan ke tiga akun email UGM jika terjadi kegagalan, serta konfirmasi keberhasilan proses.

DAG tambahan **dag_data_quality** digunakan untuk memantau kualitas data, mengevaluasi duplikasi customer_key, dan rasio nilai null, serta mencatat hasil monitoring ke dalam file **log_data_quality.csv**.

5. Prosedur pemulihan kegagalan

Jika terjadi kegagalan pada salah satu task, Airflow akan:

- a. Menandai task sebagai gagal,
- b. Menjalankan task notify_failure untuk memberi tahu tim melalui email,
- c. Melakukan retry otomatis jika sudah dikonfigurasi (retries),
- d. Memberikan akses ke log kesalahan yang dapat diperiksa melalui UI Airflow.