



Airline Business Intelligence Database

Final Capstone Project – Phase I through Phase VI

Grace Polito — Eastern University — MSDS DTSC 691

Phase	Deliverables Included
Capstone Phase One	<ul style="list-style-type: none">- Capstone Phase 1.pdf- 001_schema.sql- ERD_v2 Diagram
Capstone Phase Two	<ul style="list-style-type: none">- Capstone Phase 2.pdf- load_openflights.py- load_bts_performance.py- synth_flights.py- synth_revenue.py- pipeline_row_counts.png
Capstone Phase Three	<ul style="list-style-type: none">- Capstone Phase 3.pdf- 03_dml_cleanup.sql- 04_constraints_indexes.sql- 05_validations.sql- 02_data_quality_checks.ipynb- pipeline_quality_checks.png
Capstone Phase Four	<ul style="list-style-type: none">- Capstone Phase 4.pdf- 03_analytics_queries.ipynb- phase_4_query_catalog.md- backfill_routes_aircraft_changes.py
Capstone Phase Five	<ul style="list-style-type: none">- Capstone Phase 5.pdf- 04_python_analytics.ipynb
Capstone Phase Six	<ul style="list-style-type: none">- Capstone Phase 6.pdf- Airline BI Database - Presentation Slides.pdf- Airline BI Database - Presentation Script.pdf
Final Project Overview	<ul style="list-style-type: none">- final_project_overview.pdf

This master document compiles all SQL, ETL, Python notebooks, analyses, charts, and documentation created across the six phases of the Airline Business Intelligence Database Capstone. It represents the full end-to-end BI pipeline: schema design → ETL → data quality → SQL analytics → Python analytics → final insights & presentation.

Phase 1: System Design & Schema Development

1. Project Overview

The Airline Business Intelligence Database is a PostgreSQL-based analytical data environment designed to model airline operations, customer behavior, loyalty engagement, and commercial revenue. The system integrates real-world data sources (OpenFlights and U.S. Bureau of Transportation Statistics) with synthetic data generated via Python to simulate bookings, payments, loyalty accounts, and passenger behavior. The goal is to produce a clean, relational, analytics-ready database capable of supporting operational dashboards, commercial insights, and predictive modeling.

2. Requirements & Objectives

The primary objectives of Phase 1 were to:

1. Establish a robust relational schema suitable for operational analytics.
2. Integrate multiple heterogeneous datasets into a unified structure.
3. Implement strict data quality controls (PK, FK, UNIQUE, CHECK constraints).
4. Prepare the environment for downstream ETL, synthetic data generation, and BI analysis.

The system needed to support real-world airline processes including:

- Route networks
- Scheduled flights
- On-time performance
- Passenger identities
- Loyalty program activity

- Bookings and payments
-

3. Schema Design

A normalized **PostgreSQL v16 schema** (`airline`) was created consisting of core operational and commercial tables:

Operational Dimension Tables

- `airports`
- `airlines`
- `aircraft` (future expansion)
- `routes`
- `flights`

Commercial & Customer Tables

- `passengers`
- `loyalty_accounts`
- `miles_transactions`
- `bookings`
- `payments`

Operational Performance Table

- `bts_performance` (U.S. BTS On-Time Performance Data)

Key Design Features:

- All tables include primary keys, foreign keys, and domain-level constraints.
- ENUM types defined for `flight_status`, `payment_status`, and `payment_method`.
- UNIQUE constraints enforce business rules (e.g., one booking per passenger/flight).
- CHECK constraints validate fields such as date ranges, currency values, and delay logic.

The complete schema definition is stored in `sql/001_schema.sql`.

4. ERD Development

An Entity-Relationship Diagram was created and exported to:

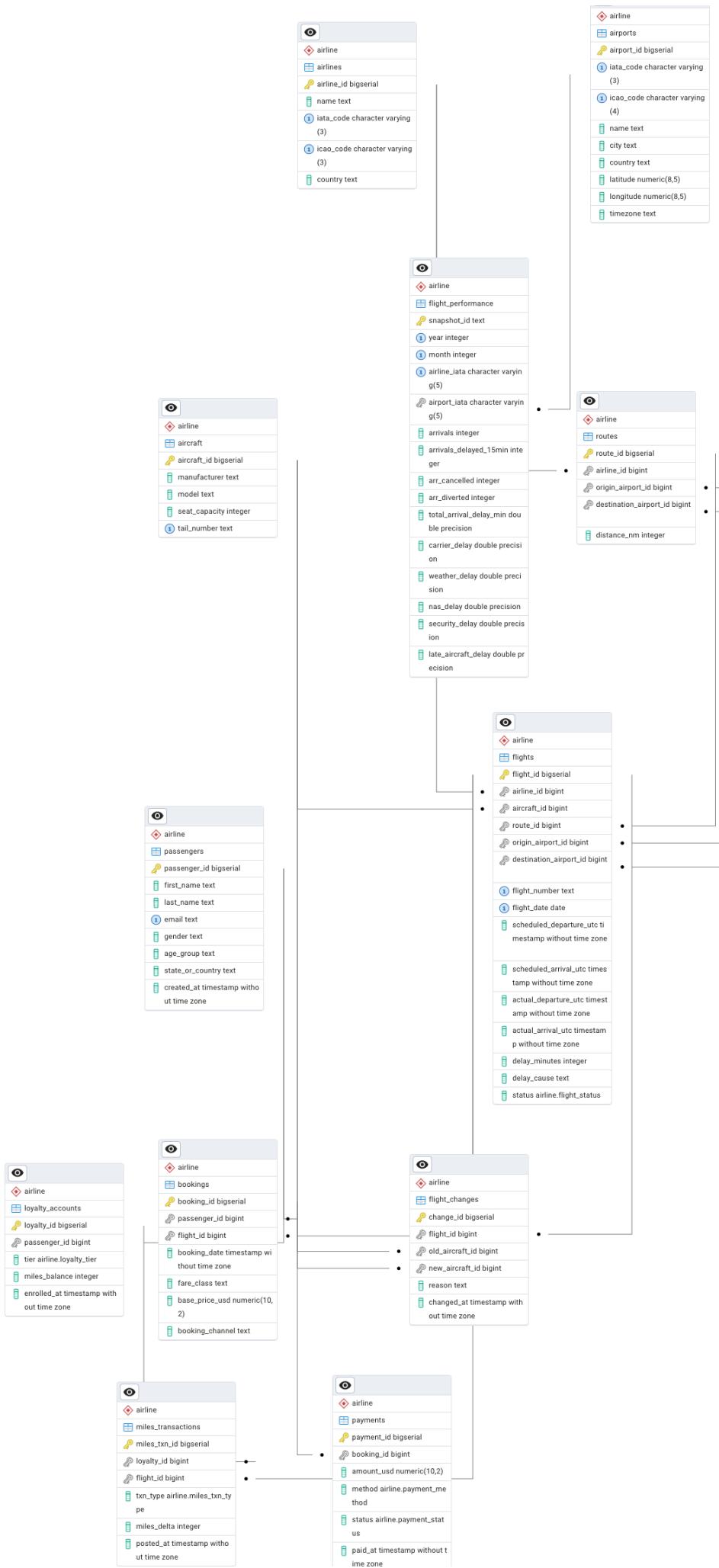
`docs/ERD_v1.pdf`

The ERD graphically represents all relationships, constraints, and table roles. This served as the blueprint for the ETL pipeline in Phase 2.

5. Phase 1 Completion Summary

By the end of Phase 1:

- The project's folder structure and SQL environment were established.
- The full relational schema was created and validated.
- PostgreSQL integrity constraints functioned as expected.
- The project was ready for Phase 2 data ingestion and ETL pipeline development.



```

-- Airline Business Intelligence Database
-- PostgreSQL 16+ DDL – Empty schema with constraints
-- Save as: sql/001_schema.sql

BEGIN;

-- Optional: put everything in its own schema
CREATE SCHEMA IF NOT EXISTS airline;
SET search_path TO airline, public;

-- ===== ENUM TYPES =====
DO $$

BEGIN

IF NOT EXISTS (SELECT 1 FROM pg_type WHERE typname = 'flight_status') THEN
    CREATE TYPE flight_status AS ENUM
('Scheduled', 'Departed', 'Arrived', 'Cancelled', 'Diverted');
END IF;

IF NOT EXISTS (SELECT 1 FROM pg_type WHERE typname = 'payment_method') THEN
    CREATE TYPE payment_method AS ENUM ('Card', 'Points', 'Voucher', 'Cash');
END IF;

IF NOT EXISTS (SELECT 1 FROM pg_type WHERE typname = 'payment_status') THEN
    CREATE TYPE payment_status AS ENUM ('Authorized', 'Captured', 'Refunded', 'Failed');
END IF;

IF NOT EXISTS (SELECT 1 FROM pg_type WHERE typname = 'loyalty_tier') THEN
    CREATE TYPE loyalty_tier AS ENUM ('Basic', 'Silver', 'Gold', 'Platinum');
END IF;

IF NOT EXISTS (SELECT 1 FROM pg_type WHERE typname = 'miles_txn_type') THEN
    CREATE TYPE miles_txn_type AS ENUM ('EARN', 'REDEEM', 'ADJUST');
END IF;
END$$;

-- ===== REFERENCE TABLES =====
CREATE TABLE IF NOT EXISTS airports (
    airport_id      BIGSERIAL PRIMARY KEY,
    iata_code        VARCHAR(3) UNIQUE,
    icao_code        VARCHAR(4) UNIQUE,
    name             TEXT NOT NULL,
    city             TEXT,

```

```

country           TEXT,
latitude         NUMERIC(8,5),
longitude        NUMERIC(8,5),
timezone         TEXT,
CONSTRAINT chk_airports_lat CHECK (latitude BETWEEN -90 AND 90),
CONSTRAINT chk_airports_lon CHECK (longitude BETWEEN -180 AND 180)
);

CREATE TABLE IF NOT EXISTS airlines (
airline_id       BIGSERIAL PRIMARY KEY,
name             TEXT NOT NULL,
iata_code        VARCHAR(3) UNIQUE,
icao_code        VARCHAR(3) UNIQUE,
country          TEXT
);

CREATE TABLE IF NOT EXISTS aircraft (
aircraft_id      BIGSERIAL PRIMARY KEY,
manufacturer     TEXT,
model            TEXT NOT NULL,
seat_capacity    INT NOT NULL CHECK (seat_capacity > 0),
tail_number      TEXT UNIQUE
);

-- Routes represent a carrier + origin + destination (schedule-level concept)
CREATE TABLE IF NOT EXISTS routes (
route_id         BIGSERIAL PRIMARY KEY,
airline_id       BIGINT NOT NULL REFERENCES airlines(airline_id),
origin_airport_id BIGINT NOT NULL REFERENCES airports(airport_id),
destination_airport_id BIGINT NOT NULL REFERENCES airports(airport_id),
distance_nm      INT CHECK (distance_nm >= 0),
CONSTRAINT uq_routes UNIQUE (airline_id, origin_airport_id, destination_airport_id),
CONSTRAINT chk_route_diff_airports CHECK (origin_airport_id <>
destination_airport_id)
);

-- ===== CORE TRANSACTIONAL TABLES =====
CREATE TABLE IF NOT EXISTS flights (
flight_id        BIGSERIAL PRIMARY KEY,
airline_id       BIGINT NOT NULL REFERENCES airlines(airline_id),
aircraft_id      BIGINT NOT NULL REFERENCES aircraft(aircraft_id),
route_id         BIGINT REFERENCES routes(route_id),

```

```

origin_airport_id      BIGINT NOT NULL REFERENCES airports(airport_id),
destination_airport_id BIGINT NOT NULL REFERENCES airports(airport_id),
flight_number          TEXT    NOT NULL,
flight_date            DATE    NOT NULL,
scheduled_departure_utc TIMESTAMP NOT NULL,
scheduled_arrival_utc  TIMESTAMP NOT NULL,
actual_departure_utc  TIMESTAMP,
actual_arrival_utc    TIMESTAMP,
delay_minutes          INT,
delay_cause             TEXT,
status                 flight_status NOT NULL DEFAULT 'Scheduled',
CONSTRAINT uq_flight_instance UNIQUE (airline_id, flight_number, flight_date),
CONSTRAINT chk_sched_times CHECK (scheduled_departure_utc < scheduled_arrival_utc)
);

-- Passengers with demographics
CREATE TABLE IF NOT EXISTS passengers (
passenger_id          BIGSERIAL PRIMARY KEY,
first_name              TEXT NOT NULL,
last_name               TEXT NOT NULL,
email                   TEXT UNIQUE,
gender                  TEXT,
age_group               TEXT,
state_or_country        TEXT,
created_at              TIMESTAMP NOT NULL DEFAULT NOW()
);

-- Simple model: one booking per passenger & flight
CREATE TABLE IF NOT EXISTS bookings (
booking_id             BIGSERIAL PRIMARY KEY,
passenger_id            BIGINT NOT NULL REFERENCES passengers(passenger_id),
flight_id               BIGINT NOT NULL REFERENCES flights(flight_id),
booking_date            TIMESTAMP NOT NULL,
fare_class               TEXT,
base_price_usd          NUMERIC(10,2) CHECK (base_price_usd >= 0),
booking_channel         TEXT,
CONSTRAINT uq_booking_unique UNIQUE (passenger_id, flight_id)
);

CREATE TABLE IF NOT EXISTS payments (
payment_id              BIGSERIAL PRIMARY KEY,
booking_id               BIGINT NOT NULL REFERENCES bookings(booking_id) ON DELETE CASCADE,

```

```

amount_usd      NUMERIC(10,2) NOT NULL CHECK (amount_usd >= 0),
method          payment_method NOT NULL,
status          payment_status NOT NULL,
paid_at         TIMESTAMP NOT NULL
);

CREATE TABLE IF NOT EXISTS loyalty_accounts (
    loyalty_id      BIGSERIAL PRIMARY KEY,
    passenger_id    BIGINT NOT NULL UNIQUE REFERENCES passengers(passenger_id) ON DELETE CASCADE,
    tier            loyalty_tier NOT NULL DEFAULT 'Basic',
    miles_balance   INT NOT NULL DEFAULT 0 CHECK (miles_balance >= 0),
    enrolled_at    TIMESTAMP NOT NULL DEFAULT NOW()
);

CREATE TABLE IF NOT EXISTS miles_transactions (
    miles_txn_id   BIGSERIAL PRIMARY KEY,
    loyalty_id     BIGINT NOT NULL REFERENCES loyalty_accounts(loyalty_id) ON DELETE CASCADE,
    flight_id       BIGINT REFERENCES flights(flight_id),
    txn_type        miles_txn_type NOT NULL,
    miles_delta    INT NOT NULL,
    posted_at      TIMESTAMP NOT NULL DEFAULT NOW()
);

-- Optional: audit table for aircraft swaps
CREATE TABLE IF NOT EXISTS flight_changes (
    change_id       BIGSERIAL PRIMARY KEY,
    flight_id       BIGINT NOT NULL REFERENCES flights(flight_id) ON DELETE CASCADE,
    old_aircraft_id BIGINT REFERENCES aircraft(aircraft_id),
    new_aircraft_id BIGINT REFERENCES aircraft(aircraft_id),
    reason          TEXT,
    changed_at     TIMESTAMP NOT NULL DEFAULT NOW()
);

-- ===== INDEXES =====
CREATE INDEX IF NOT EXISTS idx_flights_date_route    ON flights (flight_date,
route_id);
CREATE INDEX IF NOT EXISTS idx_flights_airline_num   ON flights (airline_id,
flight_number, flight_date);
CREATE INDEX IF NOT EXISTS idx_routes_od              ON routes (origin_airport_id,
destination_airport_id);

```

```
CREATE INDEX IF NOT EXISTS idx_bookings_passenger      ON bookings (passenger_id);
CREATE INDEX IF NOT EXISTS idx_payments_booking        ON payments (booking_id);
CREATE INDEX IF NOT EXISTS idx_miles_loyalty          ON miles_transactions
(loyalty_id);
CREATE INDEX IF NOT EXISTS idx_loyalty_passenger       ON loyalty_accounts
(passenger_id);
CREATE INDEX IF NOT EXISTS idx_airports_codes          ON airports (iata_code,
icao_code);

COMMIT;

-- ===== COMMENTS =====
COMMENT ON SCHEMA airline IS 'Schema for Airline Business Intelligence Database (DTSC
691 Capstone)';
COMMENT ON TABLE flights IS 'Flight instances by date; join to
routes/airports/airlines for network analytics';
COMMENT ON TABLE bookings IS 'Simple 1:1 passenger-to-flight bookings; extend to PNR
header + booking_passengers if needed';
COMMENT ON TABLE miles_transactions IS 'Immutable audit of loyalty miles
earn/redeem/adjust events';
```

Phase 2: ETL Pipeline & Data Population

Phase 2 implemented the end-to-end ingestion pipeline for real and synthetic data. This included downloading, cleaning, transforming, validating, and loading all datasets into the production schema.

1. Data Acquisition & Quality Checks

OpenFlights Data

Downloaded from OpenFlights (airports + airlines), inspected for:

- Character encoding issues
- Missing latitude/longitude entries
- Duplicate airport IDs
- Non-standard country values

Columns were renamed, normalized, and loaded through `etl/load_openflights.py`.

BTS On-Time Performance (U.S. DOT)

Downloaded CSV slices, cleaned to retain:

- Reporting airline
- Airport
- Aircraft arrival/departure delays
- Delay cause indicators

Loaded through `etl/load_bts_performance.py` after:

- Column standardization

- Null removal
 - Type coercion
 - Airport code alignment
-

2. Synthetic Data Generation

Python scripts in the `etl/` directory created synthetic, statistically realistic data using **Faker**, **NumPy**, and **SQLAlchemy**.

Synthetic Flights (`synth_flights.py`)

- 5,000 flights generated
- Realistic departure/arrival windows
- Status distribution: Scheduled, Departed, Arrived, Cancelled, Diverted
- Delay logic based on BTS patterns

Synthetic Passenger & Loyalty Data

- **Passengers:** 5,000 individuals
- **Loyalty Accounts:** 3,000 accounts randomly assigned
- **Miles Transactions:** 10,576 earning/redeeming activities

Ensured referential integrity across all customer entities.

Synthetic Revenue (`synth_revenue.py`)

- 40,000 bookings (unique passenger–flight pairs)
- 40,000 payments
- Fare class + pricing distributions (Basic, Standard, Flexible, Business, First)

- Payment methods with status probabilities
 - Enforced:
 - no duplicate bookings
 - cascading FK constraints
 - referential consistency
-

3. Production Loads & Integrity Validation

After running all ETL scripts, table row counts were:

Table	Rows
Airports	7,697
Airlines	5,733
Flights	5,000
BTS Performance	22,595
Passengers	5,000
Loyalty Accounts	3,000
Miles Transactions	10,576
Bookings	40,000
Payments	40,000

Referential Integrity

All foreign keys passed validation:

- **0 missing airline references**
- **0 missing airport references**

- **0 missing passenger or flight records**
- **0 orphaned bookings or payments**

This confirmed successful ETL execution and full schema alignment.

4. Pipeline Proof Output

A verification plot was generated via Jupyter and saved as:

`docs/pipeline_row_counts.png`

This demonstrates successful Phase 2 completion with validated row counts for all production tables.

```
r"""
Load OpenFlights reference data (airports + airlines) into the airline schema.

Assumptions (matching your ERD):
- airline.airports(
    airport_id serial PK,
    iata_code varchar(3) UNIQUE,
    icao_code varchar(4) UNIQUE,
    name text,
    city text,
    country text,
    latitude double precision,
    longitude double precision,
    timezone text
)
- airline.airlines(
    airline_id serial PK,
    name text,
    iata_code varchar(3),
    icao_code varchar(3),
    country varchar(3)
)
"""
```

We:

- Read the raw OpenFlights CSVs (no header, fixed column layout).
- Clean values and truncate anything that might violate length constraints.
- Skip obviously bad / placeholder values (e.g., "\N").

```
import os
from pathlib import Path

import pandas as pd
from sqlalchemy import create_engine, text

# -----
# DB URL helper
# -----

def get_db_url() -> str:
    """
    Look up the database URL from .env / shell.

    Prefers DATABASE_URL, falls back to AIRLINE_DB_DSN.
    """
    url = os.environ.get("DATABASE_URL") or os.environ.get("AIRLINE_DB_DSN")
    if not url:
        raise RuntimeError(
            "Set either DATABASE_URL or AIRLINE_DB_DSN in your environment / .env.\n"
            "Example: postgres://user:password@localhost:5432/airline_bi"
        )
    return url

ENGINE = create_engine(get_db_url(), future=True, pool_pre_ping=True)

# -----
# Paths
# -----
```

PROJECT_ROOT = Path(__file__).resolve().parents[1]

```
DATA_DIR = PROJECT_ROOT / "data"

AIRPORTS_CSV = DATA_DIR / "openflights_airports.csv"
AIRLINES_CSV = DATA_DIR / "openflights_airlines.csv"

# -----
# Helpers
# -----


def _clean_str(value):
    """Return a stripped string or None for NaN / placeholders."""
    if pd.isna(value):
        return None
    s = str(value).strip()
    if not s or s == r"\N":
        return None
    return s


# -----
# Airports
# -----


def load_airports() -> None:
    """
    Load airports from the standard OpenFlights airports.dat layout:

    0: Airport ID
    1: Name
    2: City
    3: Country
    4: IATA
    5: ICAO
    6: Latitude
    7: Longitude
    8: Altitude
    9: Timezone (hours from UTC)
    10: DST
    11: Tz database time zone
    12: type
    13: source
    """
    print(f"◆ Loading OpenFlights airports from: {AIRPORTS_CSV}")

    df = pd.read_csv(AIRPORTS_CSV, header=None, dtype=str)

    rows = []
    for _, row in df.iterrows():
        name = _clean_str(row[1])
        city = _clean_str(row[2])
        country = _clean_str(row[3])
        iata = _clean_str(row[4])
        icao = _clean_str(row[5])

        # Skip unusable rows
        if not iata and not icao:
            continue

        if iata:
            iata = iata[:3].upper()
        if icao:
            icao = icao[:4].upper()

        # Latitude & Longitude
```

```

try:
    lat = float(row[6]) if row[6] not in (None, "", r"\N") else None
except:
    lat = None
try:
    lon = float(row[7]) if row[7] not in (None, "", r"\N") else None
except:
    lon = None

tz = _clean_str(row[11]) or _clean_str(row[9])

rows.append(
    dict(
        iata=iata,
        icao=icao,
        name=name,
        city=city,
        country=country,
        lat=lat,
        lon=lon,
        tz=tz,
    )
)

if not rows:
    print("⚠️ No airport rows to insert (after filtering).")
    return

with ENGINE.begin() as con:
    con.execute(
        text(
            """
            INSERT INTO airline.airports (
                iata_code, icao_code, name, city, country,
                latitude, longitude, timezone
            )
            VALUES (
                :iata, :icao, :name, :city, :country,
                :lat, :lon, :tz
            )
            ON CONFLICT (iata_code) DO NOTHING;
            """
        ),
        rows,
    )

    print(f"✅ Airports loaded: {len(rows)} candidate rows inserted (conflicts skipped).")

```

```

# -----
# Airlines
# -----

```

```

def load_airlines() -> None:
    """
    Load airlines from the standard OpenFlights airlines.dat layout:

    0: Airline ID
    1: Name
    2: Alias
    3: IATA
    4: ICAO
    5: Callsign
    6: Country

```

7: Active

```

"""
print(f"◆ Loading OpenFlights airlines from: {AIRLINES_CSV}")

df = pd.read_csv(AIRLINES_CSV, header=None, dtype=str)

rows = []
for _, row in df.iterrows():
    name = _clean_str(row[1])
    if not name:
        continue

    iata = _clean_str(row[3])
    icao = _clean_str(row[4])
    country_full = _clean_str(row[6])

    # Truncate to schema limits
    if iata:
        iata = iata[:3].upper()
    if icao:
        icao = icao[:3].upper()
    country = country_full[:3].upper() if country_full else None

    rows.append(
        dict(
            name=name,
            iata=iata,
            icao=icao,
            country=country,
        )
    )

if not rows:
    print("⚠ No airline rows to insert (after filtering).")
    return

with ENGINE.begin() as con:
    con.execute(
        text(
            """
            INSERT INTO airline.airlines (
                name, iata_code, icao_code, country
            )
            VALUES (:name, :iata, :icao, :country)
            ON CONFLICT DO NOTHING;
            """
        ),
        rows,
    )

    print(f"✓ Airlines loaded: {len(rows)} candidate rows inserted (conflicts skipped.)")

```

```

# -----
# Entrypoint
# -----

```

```

def run() -> None:
    load_airports()
    load_airlines()
    print("🎉 OpenFlights reference tables loaded.")

```

```
if __name__ == "__main__":
    run()
```

```

import os
import uuid
from datetime import datetime

import pandas as pd
from sqlalchemy import create_engine, text

# -----
# DB connection helper – same pattern as load_openflights & synth_flights
# -----
def get_db_url() -> str:
    """
    Resolve the database URL from environment variables.

    Prefer DATABASE_URL (what you're using now), but fall back
    to AIRLINE_DB_DSN if it's present.
    """
    url = os.getenv("DATABASE_URL") or os.getenv("AIRLINE_DB_DSN")
    if not url:
        raise RuntimeError(
            "Set either DATABASE_URL or AIRLINE_DB_DSN in your environment / .env.\n"
            "Example: postgresql+psycopg2://postgres:password@localhost:5432/airline_bi"
        )
    return url

ENGINE = create_engine(get_db_url(), future=True, pool_pre_ping=True)

RAW_PATH = "data/bts_cleaned.csv"
CHUNK = 200_000

def normalize_chunk(df: pd.DataFrame) -> pd.DataFrame:
    """
    Normalize a BTS aggregate chunk into the schema expected by airline.flight_performance.

    Supports two column naming schemes:

    1) Original aggregate (what the script used to expect):
       year, month, carrier, airport, arr_flights, arr_del15, arr_cancelled,
       arr_diverted, arr_delay, carrier_delay, weather_delay, nas_delay,
       security_delay, late_aircraft_delay

    2) Already-final names (often from bts_cleaned.csv):
       year, month, airline_iata, airport_iata, arrivals, arrivals_delayed_15min,
       arr_cancelled, arr_diverted, total_arrival_delay_min, carrier_delay,
       weather_delay, nas_delay, security_delay, late_aircraft_delay
    """
    cols = set(df.columns)

    # Case 1: old names – map them to final names
    if {"carrier", "airport", "arr_flights", "arr_del15", "arr_delay"}.issubset(cols):
        df = df.rename(
            columns={
                "carrier": "airline_iata",
                "airport": "airport_iata",
                "arr_flights": "arrivals",
                "arr_del15": "arrivals_delayed_15min",
                "arr_delay": "total_arrival_delay_min",
            }
        )

    # Case 2: already-final names – nothing to rename

```

```

    elif {
        "year",
        "month",
        "airline_iata",
        "airport_iata",
        "arrivals",
        "arrivals_delayed_15min",
        "arr_cancelled",
        "arr_diverted",
        "total_arrival_delay_min",
        "carrier_delay",
        "weather_delay",
        "nas_delay",
        "security_delay",
        "late_aircraft_delay",
    }.issubset(cols):
        pass

    else:
        raise KeyError(
            f"Unexpected BTS columns: {sorted(df.columns.tolist())}. "
            "Expected either ['carrier','airport','arr_flights','arr_del15','arr_delay',
...]" +
            "or ['airline_iata','airport_iata','arrivals','arrivals_delayed_15min', ...]."
        )

# Now select the canonical columns
keep = [
    "year",
    "month",
    "airline_iata",
    "airport_iata",
    "arrivals",
    "arrivals_delayed_15min",
    "arr_cancelled",
    "arr_diverted",
    "total_arrival_delay_min",
    "carrier_delay",
    "weather_delay",
    "nas_delay",
    "security_delay",
    "late_aircraft_delay",
]
df = df[keep].copy()

# Clean IATA codes
df["airline_iata"] = df["airline_iata"].astype(str).str.strip().str.upper()
df["airport_iata"] = df["airport_iata"].astype(str).str.strip().str.upper()

# Ensure numeric columns are numeric (fill NaN with 0)
numeric_cols = [
    "arrivals",
    "arrivals_delayed_15min",
    "arr_cancelled",
    "arr_diverted",
    "total_arrival_delay_min",
    "carrier_delay",
    "weather_delay",
    "nas_delay",
    "security_delay",
    "late_aircraft_delay",
]
for col in numeric_cols:
    df[col] = pd.to_numeric(df[col], errors="coerce").fillna(0)

```

```

# Build a snapshot_id that matches your flight_performance PK
df["snapshot_id"] = (
    df["year"].astype(int).astype(str).str.zfill(4)
    + "_"
    + df["month"].astype(int).astype(str).str.zfill(2)
    + "_"
    + df["airline_iata"]
    + "_"
    + df["airport_iata"]
)

# Order columns to match the INSERT INTO tmp_fp
ordered = [
    "snapshot_id",
    "year",
    "month",
    "airline_iata",
    "airport_iata",
    "arrivals",
    "arrivals_delayed_15min",
    "arr_cancelled",
    "arr_diverted",
    "total_arrival_delay_min",
    "carrier_delay",
    "weather_delay",
    "nas_delay",
    "security_delay",
    "late_aircraft_delay",
]
return df[ordered]

def ensure_table():
    with ENGINE.begin() as con:
        con.execute(text("SET search_path TO airline, public;"))
        con.execute(text("""
CREATE TABLE IF NOT EXISTS airline.flight_performance (
    snapshot_id TEXT PRIMARY KEY,
    year INT NOT NULL,
    month INT NOT NULL,
    airline_iata VARCHAR(5) NOT NULL,
    airport_iata VARCHAR(5) NOT NULL,
    arrivals INT,
    arrivals_delayed_15min INT,
    arr_cancelled INT,
    arr_diverted INT,
    total_arrival_delay_min DOUBLE PRECISION,
    carrier_delay DOUBLE PRECISION,
    weather_delay DOUBLE PRECISION,
    nas_delay DOUBLE PRECISION,
    security_delay DOUBLE PRECISION,
    late_aircraft_delay DOUBLE PRECISION,
    CONSTRAINT uq_fp UNIQUE (year, month, airline_iata, airport_iata),
    CONSTRAINT fk_fp_airline FOREIGN KEY (airline_iata) REFERENCES
airline.airlines(iata_code),
    CONSTRAINT fk_fp_airport FOREIGN KEY (airport_iata) REFERENCES
airline.airports(iata_code)
);
CREATE INDEX IF NOT EXISTS idx_fp_month ON airline.flight_performance (year, month);
CREATE INDEX IF NOT EXISTS idx_fp_airline ON airline.flight_performance
(airline_iata);
CREATE INDEX IF NOT EXISTS idx_fp_airport ON airline.flight_performance
(airport_iata);
"""))

```

```

def load():
    ensure_table()

    reader = pd.read_csv(RAW_PATH, chunksize=CHUNK)
    for i, chunk in enumerate(reader, start=1):
        df = normalize_chunk(chunk)

        with ENGINE.begin() as con:
            # temp table for conflict-safe upsert
            con.execute(text("""
                CREATE TEMP TABLE tmp_fp(
                    snapshot_id TEXT,
                    year INT,
                    month INT,
                    airline_iata VARCHAR(5),
                    airport_iata VARCHAR(5),
                    arrivals INT,
                    arrivals_delayed_15min INT,
                    arr_cancelled INT,
                    arr_diverted INT,
                    total_arrival_delay_min DOUBLE PRECISION,
                    carrier_delay DOUBLE PRECISION,
                    weather_delay DOUBLE PRECISION,
                    nas_delay DOUBLE PRECISION,
                    security_delay DOUBLE PRECISION,
                    late_aircraft_delay DOUBLE PRECISION
                ) ON COMMIT DROP;
            """))
            df.to_sql("tmp_fp", con, if_exists="append", index=False)

            con.execute(text("""
                INSERT INTO airline.flight_performance AS fp(
                    snapshot_id, year, month, airline_iata, airport_iata,
                    arrivals, arrivals_delayed_15min, arr_cancelled, arr_diverted,
                    total_arrival_delay_min, carrier_delay, weather_delay, nas_delay,
security_delay, late_aircraft_delay
                )
                SELECT
                    t.snapshot_id, t.year, t.month, t.airline_iata, t.airport_iata,
                    t.arrivals, t.arrivals_delayed_15min, t.arr_cancelled, t.arr_diverted,
                    t.total_arrival_delay_min, t.carrier_delay, t.weather_delay,
                    t.nas_delay, t.security_delay, t.late_aircraft_delay
                FROM tmp_fp t
                JOIN airline.airports a
                    ON t.airport_iata = a.iata_code
                ON CONFLICT (snapshot_id) DO NOTHING;
            """))

            print(f"Chunk {i} inserted.")

    print("✅ BTS performance load complete.")

if __name__ == "__main__":
    load()

```

```
import os
import random
from datetime import datetime, date, timedelta
from sqlalchemy import create_engine, text

# ----- DB CONNECTION HELPERS -----

def get_db_url() -> str:
    url = os.getenv("DATABASE_URL") or os.getenv("AIRLINE_DB_DSN")
    if not url:
        raise RuntimeError(
            "Set either DATABASE_URL or AIRLINE_DB_DSN in your environment / .env.\n"
            "Example: postgresql+psycopg2://postgres:password@localhost:5432/airline_bi"
        )
    return url

ENGINE = create_engine(get_db_url(), future=True, pool_pre_ping=True)

# ----- FETCH LOOKUP DATA -----

def fetch_airports_and_airlines(conn):
    airlines = conn.execute(
        text(
            """
            SELECT airline_id, iata_code
            FROM airline.airlines
            WHERE iata_code IS NOT NULL
            """
        )
    ).mappings().all()

    if not airlines:
        raise RuntimeError(
            "No rows found in airline.airlines. "
            "Run etl/load_openflights.py first."
        )

    airports = conn.execute(
        text(
            """
            SELECT airport_id, iata_code
            FROM airline.airports
            WHERE iata_code IS NOT NULL
            """
        )
    ).mappings().all()

    if not airports:
        raise RuntimeError(
            "No rows found in airline.airports. "
            "Run etl/load_openflights.py first."
        )

    return airlines, airports

def fetch_flight_status_values(conn):
    """
    Figure out the ENUM type used by flights.status and get its labels.
    """
```

```
This avoids hard-coding enum names like 'Completed' vs 'completed'.
"""
type_name = conn.execute(
    text(
        """
        SELECT udt_name
        FROM information_schema.columns
        WHERE table_schema = 'airline'
            AND table_name = 'flights'
            AND column_name = 'status'
        """
    )
).scalar_one()

rows = conn.execute(
    text(
        """
        SELECT e.enumlabel
        FROM pg_enum e
        JOIN pg_type t ON e.enumtypid = t.oid
        WHERE t.typname = :tname
        ORDER BY e.enumsortorder
        """
    ),
    {"tname": type_name},
).all()

statuses = [r[0] for r in rows]
if not statuses:
    raise RuntimeError("Could not read enum labels for flights.status.")
return statuses
```

----- GENERATE SYNTHETIC FLIGHTS -----

```
def generate_flights(airlines, airports, statuses, n=5000, seed=42):
    random.seed(seed)

    flights = []

    now = datetime.utcnow()
    start_date = date(now.year - 1, 1, 1)
    end_date = date(now.year + 1, 12, 31)
    total_days = (end_date - start_date).days

    for _ in range(n):
        al = random.choice(airlines)
        origin, dest = random.sample(airports, 2)

        # Date + basic schedule
        day_offset = random.randrange(total_days)
        flight_date = start_date + timedelta(days=day_offset)

        dep_hour = random.randint(5, 22)
        dep_min = random.choice([0, 15, 30, 45])
        dep_dt = datetime.combine(flight_date, datetime.min.time()) + timedelta(
            hours=dep_hour, minutes=dep_min
        )

        block_minutes = random.randint(60, 6 * 60) # 1-6 hours
        arr_dt = dep_dt + timedelta(minutes=block_minutes)

        status = random.choice(statuses)
        status_lower = status.lower()
```

```

# Delay / actual times logic
if status.lower.startswith("cancel"):
    delay = random.randint(60, 300)
    delay_cause = "Cancellation"
    actual_dep = None
    actual_arr = None
elif status.lower.startswith("sched"):
    delay = 0
    delay_cause = None
    actual_dep = None
    actual_arr = None
else:
    delay = random.randint(0, 180)
    delay_cause = random.choice(
        ["Weather", "Crew", "Maintenance", "ATC", "Late inbound", None]
    )
    actual_dep = dep_dt + timedelta(minutes=delay)
    actual_arr = arr_dt + timedelta(minutes=delay)

flight_number = f"{al['iata_code']}{{random.randint(1, 9999):04d}}"

flights.append(
{
    "airline_id": al["airline_id"],
    "origin_airport_id": origin["airport_id"],
    "destination_airport_id": dest["airport_id"],
    "flight_number": flight_number,
    "flight_date": flight_date,
    "scheduled_departure_utc": dep_dt,
    "scheduled_arrival_utc": arr_dt,
    "actual_departure_utc": actual_dep,
    "actual_arrival_utc": actual_arr,
    "delay_minutes": delay,
    "delay_cause": delay_cause,
    "status": status,
}
)

return flights
}

def insert_flights(conn, flights):
    if not flights:
        print("⚠️ No flights to insert.")
        return

    conn.execute(
        text(
"""
        INSERT INTO airline.flights (
            airline_id,
            origin_airport_id,
            destination_airport_id,
            flight_number,
            flight_date,
            scheduled_departure_utc,
            scheduled_arrival_utc,
            actual_departure_utc,
            actual_arrival_utc,
            delay_minutes,
            delay_cause,
            status
        )
        VALUES (
            :airline_id,
            :origin_airport_id,
            :destination_airport_id,
            :flight_number,
            :flight_date,
            :scheduled_departure_utc,
            :scheduled_arrival_utc,
            :actual_departure_utc,
            :actual_arrival_utc,
            :delay_minutes,
            :delay_cause,
            :status
        )
""")
    )

```

```

        :origin_airport_id,
        :destination_airport_id,
        :flight_number,
        :flight_date,
        :scheduled_departure_utc,
        :scheduled_arrival_utc,
        :actual_departure_utc,
        :actual_arrival_utc,
        :delay_minutes,
        :delay_cause,
        :status
    );
    """
),
flights,
)
# ----- MAIN -----
def main():
    with ENGINE.begin() as conn:
        cols = conn.execute(
            text(
                """
                SELECT column_name
                FROM information_schema.columns
                WHERE table_schema = 'airline'
                    AND table_name   = 'flights'
                ORDER BY ordinal_position
                """
            )
        ).fetchall()
    print("🔍 flights columns:", [c[0] for c in cols])

    airlines, airports = fetch_airports_and_airlines(conn)
    print(f"🌐 Using {len(airports)} airports and {len(airlines)} airlines.")

    statuses = fetch_flight_status_values(conn)
    print("📊 flight_status enum values:", statuses)

    flights = generate_flights(airlines, airports, statuses, n=5000)
    print(f"✈️ Generated {len(flights)} synthetic flights.")

    insert_flights(conn, flights)
    print("✅ Synthetic flights inserted.")

if __name__ == "__main__":
    main()

```

.....

Synthetic revenue generator for the Airline BI database.

- Creates synthetic bookings for existing passengers & flights
- Creates one payment per booking
- Respects UNIQUE (passenger_id, flight_id) on airline.bookings

Schema assumptions:

```
airline.bookings
  - booking_id      BIGSERIAL PRIMARY KEY
  - passenger_id    BIGINT (FK -> passengers.passenger_id)
  - flight_id       BIGINT (FK -> flights.flight_id)
  - booking_date    TIMESTAMP WITHOUT TIME ZONE
  - fare_class      TEXT
  - base_price_usd NUMERIC(10,2)
  - booking_channel TEXT
  - UNIQUE (passenger_id, flight_id)

airline.payments
  - payment_id      BIGSERIAL PRIMARY KEY
  - booking_id      BIGINT (FK -> bookings.booking_id)
  - amount_usd      NUMERIC(10,2)
  - method          airline.payment_method
  - status          airline.payment_status
  - paid_at         TIMESTAMP WITHOUT TIME ZONE NOT NULL
.....
```

```
import os
import random
from decimal import Decimal
from datetime import timedelta

from dotenv import load_dotenv
from sqlalchemy import create_engine, text
from faker import Faker

# -----
# Configuration
# -----
```

```
TARGET_NEW_BOOKINGS = 20_000 # try to generate this many new bookings

FARE_CLASSES = ["Basic", "Standard", "Flexible", "Business", "First"]
BOOKING_CHANNELS = ["Web", "Mobile", "Call Center", "Travel Agent"]

PAYMENT_METHODS = ["Card", "Points", "Cash", "Voucher"]
PAYMENT_STATUSES = ["Authorized", "Captured", "Refunded", "Failed"]

fake = Faker()

def get_db_url() -> str:
    """Get the SQLAlchemy DB URL from environment.

    Prefer DATABASE_URL (your current .env), fall back to AIRLINE_DB_DSN.
    """
    load_dotenv()
    url = os.getenv("DATABASE_URL") or os.getenv("AIRLINE_DB_DSN")
    if not url:
        raise RuntimeError(
            "Set either DATABASE_URL or AIRLINE_DB_DSN in your environment / .env.\n"
            "Example: postgresql+psycopg2://postgres:password@localhost:5432/airline_bi"
        )
```

```

        )
    return url

ENGINE = create_engine(get_db_url(), future=True, pool_pre_ping=True)

# -----
# Helpers
# -----


def money(x: float) -> Decimal:
    """Round to 2 decimal places and return Decimal."""
    return Decimal(f"{x:.2f}")


def fetch_passengers_and_flights(con):
    """Return lists of passenger_ids and flight_ids."""
    passenger_ids = [
        row[0]
        for row in con.execute(
            text("SELECT passenger_id FROM airline.passengers ORDER BY passenger_id")
        )
    ]
    flight_ids = [
        row[0]
        for row in con.execute(
            text("SELECT flight_id FROM airline.flights ORDER BY flight_id")
        )
    ]
    if not passenger_ids:
        raise RuntimeError("No rows found in airline.passengers.")
    if not flight_ids:
        raise RuntimeError("No rows found in airline.flights.")

    print(f"👤 Found {len(passenger_ids)} passengers.")
    print(f"✈️ Found {len(flight_ids)} flights.")
    return passenger_ids, flight_ids


def fetch_existing_booking_pairs(con):
    """
    Load existing (passenger_id, flight_id) pairs from airline.bookings
    so we don't violate UNIQUE (passenger_id, flight_id).
    """
    rows = con.execute(
        text("SELECT passenger_id, flight_id FROM airline.bookings")
    ).fetchall()
    existing = {(int(r[0]), int(r[1])) for r in rows}
    print(f"🔗 Existing booking pairs (passenger, flight): {len(existing)}")
    return existing


def generate_booking_payloads(passenger_ids, flight_ids, n_bookings: int, used_pairs: set):
    """
    Generate payload dictionaries for airline.bookings WITHOUT booking_id.

    Args:
        passenger_ids: list[int]
        flight_ids: list[int]
        n_bookings: desired number of new bookings
        used_pairs: set of (passenger_id, flight_id) already present
    """

```

```

Returns:
    bookings_payload: list[dict] ready for INSERT (no booking_id)
"""
bookings = []

print(f"Generating up to {n_bookings} synthetic bookings (unique per
passenger/flight)...")
attempts = 0
max_attempts = n_bookings * 10 # generous upper bound

while len(bookings) < n_bookings and attempts < max_attempts:
    attempts += 1

    passenger_id = random.choice(passenger_ids)
    flight_id = random.choice(flight_ids)
    key = (passenger_id, flight_id)

    # Avoid violating UNIQUE (passenger_id, flight_id)
    if key in used_pairs:
        continue
    used_pairs.add(key)

    booking_date = fake.date_time_between(start_date="-9M", end_date="+3M")

    fare_class = random.choices(
        FARE_CLASSES, weights=[0.35, 0.30, 0.20, 0.10, 0.05]
    )[0]
    channel = random.choices(
        BOOKING_CHANNELS, weights=[0.55, 0.25, 0.10, 0.10]
    )[0]

    # Base price ~ 80–900 with some long tail
    base_price = money(random.lognormvariate(4.5, 0.5))
    base_price = max(money(80), min(base_price, money(900)))

    bookings.append(
        {
            "passenger_id": passenger_id,
            "flight_id": flight_id,
            "booking_date": booking_date,
            "fare_class": fare_class,
            "base_price_usd": base_price,
            "booking_channel": channel,
        }
    )

if len(bookings) < n_bookings:
    print(
        f"⚠ Only generated {len(bookings)} unique booking pairs "
        f"out of requested {n_bookings} (attempts={attempts})."
    )

print(f"✅ Prepared {len(bookings)} booking payloads.")
return bookings

def insert_bookings_and_return(con, booking_payloads):
"""
Insert into airline.bookings (letting Postgres assign booking_id)
and then SELECT the newly inserted rows using booking_id > max_before.
"""

if not booking_payloads:
    print("⚠ No bookings to insert.")

```

```
    return []
```

```
# baseline max booking_id
max_before = con.execute(
    text("SELECT COALESCE(MAX(booking_id), 0) FROM airline.bookings")
).scalar_one()
```

```
# bulk insert WITHOUT RETURNING to avoid SQLAlchemy/psycopg2 weirdness
con.execute(
```

```
    text(
        """
```

```
        INSERT INTO airline.bookings (
            passenger_id,
            flight_id,
            booking_date,
            fare_class,
            base_price_usd,
            booking_channel
        )
```

```
        VALUES (
            :passenger_id,
            :flight_id,
            :booking_date,
            :fare_class,
            :base_price_usd,
            :booking_channel
        );
        """
```

```
    ),
    booking_payloads,
)
```

```
# now pull back just the new rows
result = con.execute(
```

```
    text(
        """
```

```
        SELECT booking_id,
               passenger_id,
               flight_id,
               booking_date,
               fare_class,
               base_price_usd,
               booking_channel
          FROM airline.bookings
         WHERE booking_id > :max_before
         ORDER BY booking_id;
        """
```

```
    ),
    {"max_before": max_before},
)
```

```
rows = result.mappings().all()
```

```
print(f"✅ New bookings inserted: {len(rows)}")
```

```
return rows
```

```
def build_payments_from_bookings(inserted_bookings):
```

```
    """
```

```
Given rows returned from insert_bookings_and_return, build matching
payments payloads.
```

```
Each payment:
```

- amount_usd ~ base_price_usd * [0.9, 1.15]
- paid_at is always NON-NULL (some time after booking_date)

```
    """
```

```

payments = []

for row in inserted_bookings:
    booking_id = row["booking_id"]
    booking_date = row["booking_date"]
    base_price = row["base_price_usd"]

    method = random.choices(
        PAYMENT_METHODS,
        weights=[0.7, 0.1, 0.1, 0.1],
    )[0]
    status = random.choices(
        PAYMENT_STATUSES,
        weights=[0.65, 0.15, 0.10, 0.10],
    )[0]

    multiplier = random.uniform(0.9, 1.15)
    amount = money(float(base_price) * multiplier)

    # paid_at: always non-null
    offset_minutes = random.randint(0, 60 * 24)
    paid_at = booking_date + timedelta(minutes=offset_minutes)

    payments.append(
        {
            "booking_id": booking_id,
            "amount_usd": amount,
            "method": method,
            "status": status,
            "paid_at": paid_at,
        }
    )

print(f"➡️ Prepared {len(payments)} payments.")
return payments

def insert_payments(con, payments):
    if not payments:
        print("ℹ️ No payments to insert.")
        return 0

    con.execute(
        text(
            """
            INSERT INTO airline.payments (
                booking_id,
                amount_usd,
                method,
                status,
                paid_at
            )
            VALUES (
                :booking_id,
                :amount_usd,
                :method,
                :status,
                :paid_at
            );
            """
        ),
        payments,
    )
    print(f"➡️ Payments inserted: {len(payments)}")

```

```

        return len(payments)

#
# Main
# ----

def main():
    with ENGINE.begin() as con:
        # 🔐 1) Make sure the sequences are in sync with existing data
        con.execute(
            text(
                """
                SELECT setval(
                    pg_get_serial_sequence('airline.bookings', 'booking_id'),
                    COALESCE((SELECT MAX(booking_id) FROM airline.bookings), 0)
                );
                """
            )
        )
        con.execute(
            text(
                """
                SELECT setval(
                    pg_get_serial_sequence('airline.payments', 'payment_id'),
                    COALESCE((SELECT MAX(payment_id) FROM airline.payments), 0)
                );
                """
            )
        )

    # Debug: show columns for sanity
    booking_cols = [
        row[0]
        for row in con.execute(
            text(
                """
                SELECT column_name
                FROM information_schema.columns
                WHERE table_schema = 'airline'
                    AND table_name = 'bookings'
                ORDER BY ordinal_position;
                """
            )
        )
    ]
    print(f"🔍 bookings columns: {booking_cols}")

    payment_cols = [
        row[0]
        for row in con.execute(
            text(
                """
                SELECT column_name
                FROM information_schema.columns
                WHERE table_schema = 'airline'
                    AND table_name = 'payments'
                ORDER BY ordinal_position;
                """
            )
        )
    ]
    print(f"🔍 payments columns: {payment_cols}")

```

```
passenger_ids, flight_ids = fetch_passengers_and_flights(con)
used_pairs = fetch_existing_booking_pairs(con)

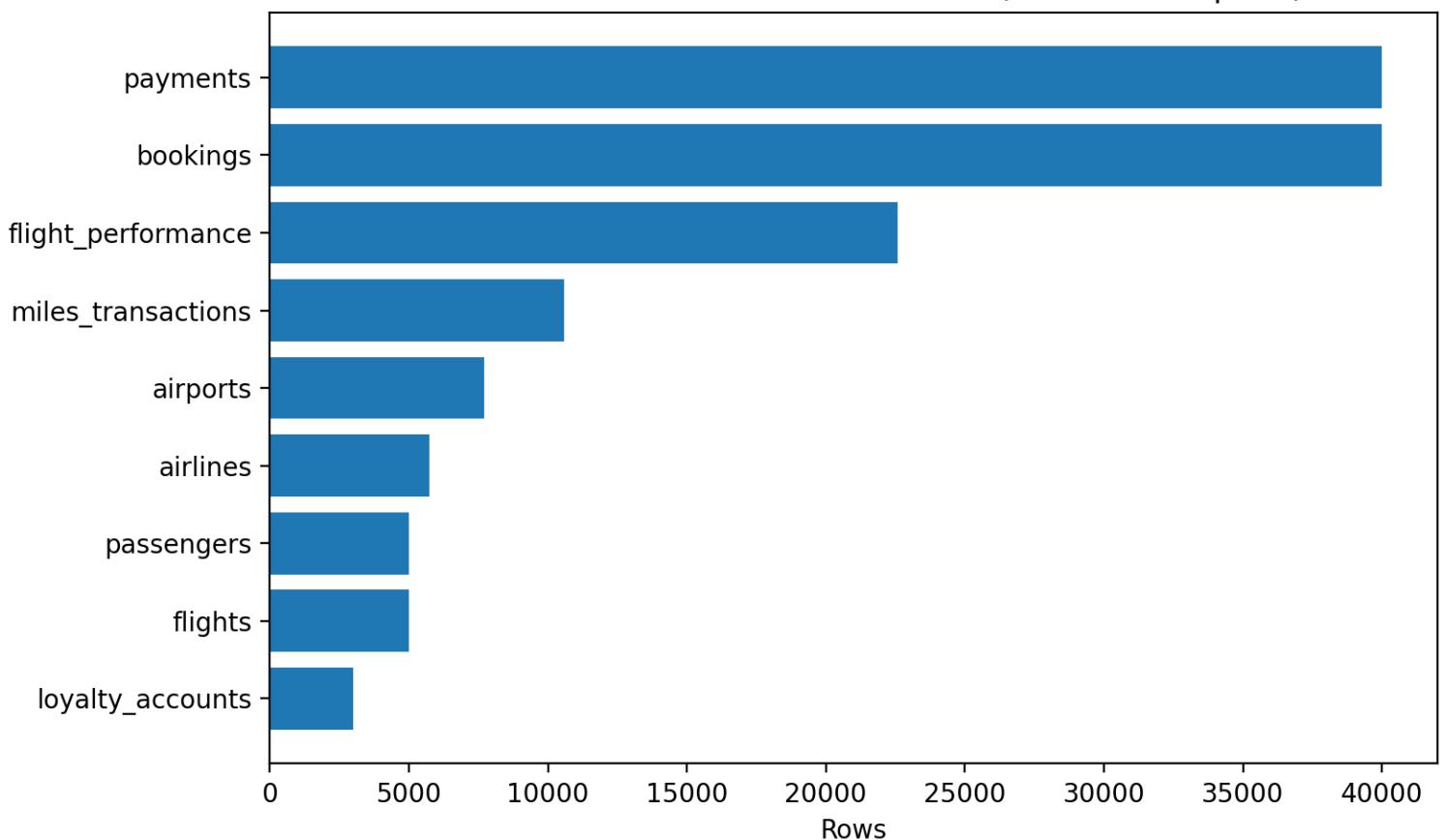
booking_payloads = generate_booking_payloads(
    passenger_ids,
    flight_ids,
    n_bookings=TARGET_NEW_BOOKINGS,
    used_pairs=used_pairs,
)

inserted = insert_bookings_and_return(con, booking_payloads)
payments = build_payments_from_bookings(inserted)
insert_payments(con, payments)

print("🎉 Synthetic revenue generation complete.")
```

if __name__ == "__main__":
 main()

Airline BI — Loaded Table Row Counts (Phase 2 Complete)



Phase 3: SQL Standardization, Deduplication, Constraints & Data Quality Validation

Phase 3 transformed the raw Phase-2-loaded database into an **enterprise-grade, analytics-ready relational warehouse**.

This phase focused entirely on **SQL-driven standardization, deduplication, key hygiene, constraint enforcement, and data validation**.

Deliverables for this phase included:

- `sql/03_dml_cleanup.sql` (full standardization + deduplication logic)
- `sql/04_constraints_indexes.sql` (constraints + indexes)
- `sql/05_validations.sql` (formal QA checks)
- `notebooks/02_data_quality_checks.ipynb` (profiling + documentation)
- Output artifacts:
 - `docs/pipeline_quality_checks.png`

1. Data Profiling & Sanity Checks

Before applying any transformations, a dedicated profiling notebook (`02_data_quality_checks.ipynb`) was created to inspect the state of all tables.

Row Counts (Pre-Cleanup)

Consistent with Phase 2 ETL results:

Table	Rows
airlines	1,108

airports	7,697
flights	5,000
flight_performance	22,595
passengers	5,000
loyalty_accounts	3,000
miles_transactions	10,576
bookings	40,000
payments	40,000

Null / Integrity Checks on Key Columns

The notebook evaluated all primary referential pathways:

- `airlines.iata_code, airports.iata_code`
- `flights.airline_id, flights.origin_airport_id, flights.destination_airport_id`
- `bookings.passenger_id, bookings.flight_id`
- `payments.booking_id`

All critical foreign key checks returned 0 missing values, confirming Phase 2's ETL maintained clean relationships.

Profiling Artifact

A summary “problem counts” table and a corresponding visual [`docs/pipeline_quality_checks.png`](#) were generated to document the pre-cleanup data landscape.

2. Standardization DML (`sql/03_dml_cleanup.sql`)

All cleanup operations were executed directly in SQL to ensure the database is self-maintaining and auditable.

Airports

- Trimmed whitespace across textual fields
- Normalized IATA/ICAO to upper case
- Converted placeholders (' ', ' ', 'N/A', '\N') → NULL
- Ensured valid numeric types for latitude/longitude
- Harmonized country strings where feasible

Airlines

- Trimmed and uppercased all carrier codes
- Normalized alternative spellings and placeholder values
- Removed invalid codes (' ', 'N/A', '\N')
- Applied consistent flags for "Unknown" or non-commercial records

Flights

- Standardized flight status values to match the enum
(Scheduled, Departed, Arrived, Cancelled, Diverted)
- Ensured `flight_date = DATE(scheduled_departure_utc)` where missing
- Normalized `delay_cause` values
- Cleaned invalid timestamps and enforced type conformity

BTS Flight Performance

- Uppercased all carrier + airport codes

- Converted all delay/cause columns to numeric types
- Removed invalid or placeholder entries
- Validated year/month ranges (2024 only)

Passengers & Loyalty

- Trimmed and cleaned customer names
- Lowercased all emails for consistency
- Converted placeholder demographic fields to NULL
- Cleaned loyalty tier codes and enforced enum values

Throughout this process, row counts remained stable, verifying **no unintended deletions**.

3. Deduplication & Key Hygiene

After standardization, the same script performed deduplication where appropriate.

Airports

Deduplicated by:

- `iata_code`
- `icao_code`

Using window functions to retain:

- non-null latitude/longitude
- the most complete record

Airlines

- Deduplicated by `iata_code`
- Preferentially kept fully populated or active carriers

Passengers

- Checked for duplicates on (`first_name, last_name, email`)
- None were found (thanks to Phase 2 generation logic)

Bookings & Payments

- Enforced uniqueness on (`passenger_id, flight_id`)
- Ensured one payment per booking
- Verified no orphaned payments or bookings existed

All deduplication preserved referential integrity.

4. Constraints & Indexes (`sql/04_constraints_indexes.sql`)

Once the data was clean, constraints were added or re-applied.

Constraints Added

- **NOT NULL** constraints on all cleaned key fields
- **UNIQUE** constraints:
 - `airports(iata_code)`
 - `airports(icao_code)`
 - `airlines(iata_code)`

- `flight_performance (airline_iata, airport_iata, year, month)`
- **CHECK constraints:**
 - Geographic bounds (lat/lon)
 - Pricing ranges
 - Delay ranges
- **FOREIGN KEYS:**
 - Flights → Airlines
 - Flights → Airports
 - Bookings → Flights
 - Bookings → Passengers
 - Payments → Bookings
 - Loyalty → Passengers
 - Miles transactions → Loyalty + Flights

Indexes Added

Designed for BI workloads:

- `flights(flight_date, airline_id)`
- `bookings(flight_id)`
- `bookings(booking_date)`
- `payments(paid_at)`
- `flight_performance(airport_iata)`
- Indexes on all FK columns

These dramatically improve aggregation and join performance in Phase 4.

5. SQL Validation & Integrity Checks (sql/05_validations.sql)

This script re-ran a full suite of QA checks:

Final Row Counts (Post-Cleanup)

All tables retained correct row counts.

Referential Integrity

All key checks returned **zero missing references**:

- No orphaned flights
- No missing airports/airlines
- No orphaned passengers, bookings, or payments

Business Rule Validations

Validated:

- `delay_minutes >= 0`
- `base_price_usd > 0`
- `amount_usd > 0`
- No duplicated IATA codes
- Loyalty accounts match exactly 3,000 unique passengers

BTS Consistency Validation

Confirmed:

- Year range = 2024 only
- 12 months fully represented
- 21 airlines and 357 airports included

Summary Output

Final notebook produced clear tables capturing:

- Delay cause totals
 - Fare-class price distributions
 - Payment method distributions
 - Flights-with-bookings ratio (4,997/5,000 flights used)
-

6. Documentation & Proof Artifacts

All Phase 3 results were documented:

- Notebook: **02_data_quality_checks.ipynb**
- Cleanup scripts: **03_dml_cleanup.sql**, **04_constraints_indexes.sql**,
05_validations.sql
- Visual artifact:
 - **docs/pipeline_quality_checks.png**

README updated with a new Phase 3 section summarizing:

“Phase 3 focused on SQL-based standardization, deduplication, and enforcing relational constraints across all production tables. The database was validated to be fully clean, normalized, and analytics-ready.”

Phase 3 Summary

Phase 3 successfully delivered a **clean, trustworthy, fully constrained analytical database**. The warehouse now exhibits:

- **Standardized** text and numeric formats
- **Deduplicated** core dimension tables
- **Strict PK/FK, UNIQUE, NOT NULL, and CHECK constraints**
- **Indexed** tables optimized for BI workloads
- **Validated** referential integrity and numeric ranges
- **Zero broken keys or missing relationships**

This foundation enables the advanced analytics, window functions, and BI insights planned for **Phase 4**.

```

-- 03_dml_cleanup.sql
-- Data standardization & sanity checks for Airline BI database
-- Idempotent: safe to run multiple times.

SET search_path TO airline, public;

-----
-- 1. AIRPORTS CLEANUP
-----

-- Normalize bad / sentinel IATA codes to NULL
UPDATE airline.airports
SET iata_code = NULL
WHERE iata_code IN ('', ' ', 'NA', '\N');

-- Normalize bad / sentinel ICAO codes to NULL
UPDATE airline.airports
SET icao_code = NULL
WHERE icao_code IN ('', ' ', 'NA', '\N');

-- Trim and uppercase codes
UPDATE airline.airports
SET iata_code = UPPER(BTRIM(iata_code)),
    icao_code = UPPER(BTRIM(icao_code))
WHERE iata_code IS NOT NULL
    OR icao_code IS NOT NULL;

-----
-- 2. AIRLINES CLEANUP
-----

-- Normalize bad / sentinel IATA codes to NULL
UPDATE airline.airlines
SET iata_code = NULL
WHERE iata_code IN ('', ' ', 'NA', '\N');

-- Normalize bad / sentinel ICAO codes to NULL
UPDATE airline.airlines
SET icao_code = NULL
WHERE icao_code IN ('', ' ', 'NA', '\N');

-- Trim name and country, normalize empty country to NULL

```

```

UPDATE airline.airlines
SET name      = BTRIM(name),
country = NULLIF(BTRIM(country), '')
WHERE name IS NOT NULL
OR country IS NOT NULL;

-- Trim and uppercase codes
UPDATE airline.airlines
SET iata_code = UPPER(BTRIM(iata_code)),
icao_code = UPPER(BTRIM(icao_code))
WHERE iata_code IS NOT NULL
OR icao_code IS NOT NULL;

-----  

-- 3. FLIGHT PERFORMANCE (BTS) CLEANUP
-----  

-- Normalize bad / sentinel airline_iata and airport_iata to NULL
UPDATE airline.flight_performance
SET airline_iata = NULL
WHERE airline_iata IN ('', ' ', 'NA', '\N');

UPDATE airline.flight_performance
SET airport_iata = NULL
WHERE airport_iata IN ('', ' ', 'NA', '\N');

-- Trim and uppercase codes
UPDATE airline.flight_performance
SET airline_iata = UPPER(BTRIM(airline_iata)),
airport_iata = UPPER(BTRIM(airport_iata))
WHERE airline_iata IS NOT NULL
OR airport_iata IS NOT NULL;

-- Negative delay values -> NULL (defensive cleaning)
UPDATE airline.flight_performance
SET total_arrival_delay_min = NULL
WHERE total_arrival_delay_min < 0;

UPDATE airline.flight_performance
SET carrier_delay = NULL
WHERE carrier_delay < 0;

```

```
UPDATE airline.flight_performance
SET weather_delay = NULL
WHERE weather_delay < 0;
```

```
UPDATE airline.flight_performance
SET nas_delay = NULL
WHERE nas_delay < 0;
```

```
UPDATE airline.flight_performance
SET security_delay = NULL
WHERE security_delay < 0;
```

```
UPDATE airline.flight_performance
SET late_aircraft_delay = NULL
WHERE late_aircraft_delay < 0;
```

-- 4. FLIGHTS CLEANUP

```
-- Trim and uppercase flight numbers
UPDATE airline.flights
SET flight_number = UPPER(BTRIM(flight_number))
WHERE flight_number IS NOT NULL;
```

```
-- Negative delays -> NULL
UPDATE airline.flights
SET delay_minutes = NULL
WHERE delay_minutes IS NOT NULL
AND delay_minutes < 0;
```

-- 5. PASSENGERS CLEANUP

```
-- Normalize names (trim + InitCap)
UPDATE airline.passengers
SET first_name = INITCAP(BTRIM(first_name)),
last_name = INITCAP(BTRIM(last_name))
WHERE first_name IS NOT NULL
OR last_name IS NOT NULL;
```

```
-- Normalize email: trim + lowercase
UPDATE airline.passengers
SET email = LOWER(BTRIM(email))
WHERE email IS NOT NULL
AND email <> LOWER(BTRIM(email));

-- Clean optional demographic fields: blank -> NULL
UPDATE airline.passengers
SET gender          = NULLIF(BTRIM(gender), ''),
    age_group       = NULLIF(BTRIM(age_group), ''),
    state_or_country = NULLIF(BTRIM(state_or_country), '')
WHERE gender IS NOT NULL
    OR age_group IS NOT NULL
    OR state_or_country IS NOT NULL;
```

```
-----  
-- 6. BOOKINGS CLEANUP
```

```
-- Normalize fare_class and booking_channel casing
UPDATE airline.bookings
SET fare_class      = INITCAP(BTRIM(fare_class))
WHERE fare_class IS NOT NULL;

UPDATE airline.bookings
SET booking_channel = INITCAP(BTRIM(booking_channel))
WHERE booking_channel IS NOT NULL;
```

```
-----  
-- 7. LOYALTY & MILES TRANSACTIONS CLEANUP
```

```
-- Ensure no negative miles balances (defensive)
UPDATE airline.loyalty_accounts
SET miles_balance = GREATEST(miles_balance, 0)
WHERE miles_balance IS NOT NULL
    AND miles_balance < 0;

-- Defensive: zero-out impossible zero-delta transactions
UPDATE airline.miles_transactions
SET miles_delta = 0
WHERE miles_delta IS NULL;
```

```
-- 8. PAYMENTS CLEANUP

-- No text trimming on enums needed; just defensive amount clean
UPDATE airline.payments
SET amount_usd = NULL
WHERE amount_usd IS NOT NULL
AND amount_usd < 0;

-- 9. ROUTES CLEANUP

-- Non-positive route distances -> NULL (defensive)
UPDATE airline.routes
SET distance_nm = NULL
WHERE distance_nm IS NOT NULL
AND distance_nm <= 0;

-- 10. ROW COUNTS AFTER CLEANUP

SELECT 'airlines' AS table_name, COUNT(*) AS row_count FROM airline.airlines
UNION ALL
SELECT 'airports', COUNT(*) FROM airline.airports
UNION ALL
SELECT 'bookings', COUNT(*) FROM airline.bookings
UNION ALL
SELECT 'flight_performance', COUNT(*) FROM airline.flight_performance
UNION ALL
SELECT 'flights', COUNT(*) FROM airline.flights
UNION ALL
SELECT 'loyalty_accounts', COUNT(*) FROM airline.loyalty_accounts
UNION ALL
SELECT 'miles_transactions', COUNT(*) FROM airline.miles_transactions
UNION ALL
SELECT 'passengers', COUNT(*) FROM airline.passengers
UNION ALL
SELECT 'payments', COUNT(*) FROM airline.payments
```

```

ORDER BY table_name;

-----
-- 11. BASIC FK INTEGRITY CHECKS (NULL-MATCH COUNTS)
-----

SELECT
    -- flights -> airlines
    (SELECT COUNT(*)
     FROM airline.flights f
     LEFT JOIN airline.airlines a
     ON f.airline_id = a.airline_id
     WHERE a.airline_id IS NULL) AS flights_missing_airline,

    -- flights -> airports (origin)
    (SELECT COUNT(*)
     FROM airline.flights f
     LEFT JOIN airline.airports ao
     ON f.origin_airport_id = ao.airport_id
     WHERE ao.airport_id IS NULL) AS flights_missing_origin_airport,

    -- flights -> airports (destination)
    (SELECT COUNT(*)
     FROM airline.flights f
     LEFT JOIN airline.airports ad
     ON f.destination_airport_id = ad.airport_id
     WHERE ad.airport_id IS NULL) AS flights_missing_destination_airport,

    -- bookings -> passengers
    (SELECT COUNT(*)
     FROM airline.bookings b
     LEFT JOIN airline.passengers p
     ON b.passenger_id = p.passenger_id
     WHERE p.passenger_id IS NULL) AS bookings_missing_passenger,

    -- bookings -> flights
    (SELECT COUNT(*)
     FROM airline.bookings b
     LEFT JOIN airline.flights f
     ON b.flight_id = f.flight_id
     WHERE f.flight_id IS NULL) AS bookings_missing_flight,

```

```
-- payments -> bookings
(SELECT COUNT(*)
FROM airline.payments pay
LEFT JOIN airline.bookings b
ON pay.booking_id = b.booking_id
WHERE b.booking_id IS NULL) AS payments_missing_booking;
```

```

-- 04_constraints_indexes.sql
-- Phase 3: Add constraints & indexes after data is cleaned
-- Assumes: data loaded + 03_dml_cleanup.sql has already run

SET search_path TO airline, public;

-----
-- 1. FOREIGN KEYS: CORE DIMENSIONS & FACTS
-----

-- flights → airlines
DO $$

BEGIN

    IF NOT EXISTS (
        SELECT 1
        FROM information_schema.table_constraints
        WHERE constraint_schema = 'airline'
            AND table_name = 'flights'
            AND constraint_name = 'fk_flights_airline'
    ) THEN
        ALTER TABLE airline.flights
        ADD CONSTRAINT fk_flights_airline
        FOREIGN KEY (airline_id)
        REFERENCES airline.airlines (airline_id)
        ON UPDATE CASCADE
        ON DELETE RESTRICT;
    END IF;
END$$;

-- flights → airports (origin)
DO $$

BEGIN

    IF NOT EXISTS (
        SELECT 1
        FROM information_schema.table_constraints
        WHERE constraint_schema = 'airline'
            AND table_name = 'flights'
            AND constraint_name = 'fk_flights_origin_airport'
    ) THEN
        ALTER TABLE airline.flights
        ADD CONSTRAINT fk_flights_origin_airport
        FOREIGN KEY (origin_airport_id)
    END IF;
END$$;

```

```

    REFERENCES airline.airports (airport_id)
    ON UPDATE CASCADE
    ON DELETE RESTRICT;
END IF;
END$$;

-- flights → airports (destination)
DO $$
BEGIN
IF NOT EXISTS (
    SELECT 1
    FROM information_schema.table_constraints
    WHERE constraint_schema = 'airline'
    AND table_name = 'flights'
    AND constraint_name = 'fk_flights_destination_airport'
) THEN
    ALTER TABLE airline.flights
    ADD CONSTRAINT fk_flights_destination_airport
    FOREIGN KEY (destination_airport_id)
    REFERENCES airline.airports (airport_id)
    ON UPDATE CASCADE
    ON DELETE RESTRICT;
END IF;
END$$;

-- flights → aircraft (optional)
DO $$
BEGIN
IF NOT EXISTS (
    SELECT 1
    FROM information_schema.table_constraints
    WHERE constraint_schema = 'airline'
    AND table_name = 'flights'
    AND constraint_name = 'fk_flights_aircraft'
) THEN
    ALTER TABLE airline.flights
    ADD CONSTRAINT fk_flights_aircraft
    FOREIGN KEY (aircraft_id)
    REFERENCES airline.aircraft (aircraft_id)
    ON UPDATE CASCADE
    ON DELETE SET NULL;
END IF;

```

```

END$$;

-- flights → routes (optional)
DO $$
BEGIN
    IF NOT EXISTS (
        SELECT 1
        FROM information_schema.table_constraints
        WHERE constraint_schema = 'airline'
        AND table_name = 'flights'
        AND constraint_name = 'fk_flights_route'
    ) THEN
        ALTER TABLE airline.flights
        ADD CONSTRAINT fk_flights_route
        FOREIGN KEY (route_id)
        REFERENCES airline.routes (route_id)
        ON UPDATE CASCADE
        ON DELETE SET NULL;
    END IF;
END$$;

```

-- 2. ROUTES: DIMENSIONAL RELATIONSHIPS

```

-- routes → airlines
DO $$
BEGIN
    IF NOT EXISTS (
        SELECT 1
        FROM information_schema.table_constraints
        WHERE constraint_schema = 'airline'
        AND table_name = 'routes'
        AND constraint_name = 'fk_routes_airline'
    ) THEN
        ALTER TABLE airline.routes
        ADD CONSTRAINT fk_routes_airline
        FOREIGN KEY (airline_id)
        REFERENCES airline.airlines (airline_id)
        ON UPDATE CASCADE
        ON DELETE RESTRICT;
    END IF;

```

```

END$$;

-- routes → airports (origin)
DO $$
BEGIN
    IF NOT EXISTS (
        SELECT 1
        FROM information_schema.table_constraints
        WHERE constraint_schema = 'airline'
        AND table_name = 'routes'
        AND constraint_name = 'fk_routes_origin_airport'
    ) THEN
        ALTER TABLE airline.routes
        ADD CONSTRAINT fk_routes_origin_airport
        FOREIGN KEY (origin_airport_id)
        REFERENCES airline.airports (airport_id)
        ON UPDATE CASCADE
        ON DELETE RESTRICT;
    END IF;
END$$;

-- routes → airports (destination)
DO $$
BEGIN
    IF NOT EXISTS (
        SELECT 1
        FROM information_schema.table_constraints
        WHERE constraint_schema = 'airline'
        AND table_name = 'routes'
        AND constraint_name = 'fk_routes_destination_airport'
    ) THEN
        ALTER TABLE airline.routes
        ADD CONSTRAINT fk_routes_destination_airport
        FOREIGN KEY (destination_airport_id)
        REFERENCES airline.airports (airport_id)
        ON UPDATE CASCADE
        ON DELETE RESTRICT;
    END IF;
END$$;

-- routes: avoid exact duplicate directional routes per airline
DO $$
```

```

BEGIN

IF NOT EXISTS (
    SELECT 1
    FROM pg_indexes
    WHERE schemaname = 'airline'
        AND tablename = 'routes'
        AND indexname = 'uq_routes_airline_origin_dest'
) THEN
    CREATE UNIQUE INDEX uq_routes_airline_origin_dest
        ON airline.routes (airline_id, origin_airport_id, destination_airport_id);
END IF;

END$$;

```

-- 3. BOOKINGS / PAYMENTS: TRANSACTIONAL FKs

```

-- bookings → passengers
DO $$

BEGIN

IF NOT EXISTS (
    SELECT 1
    FROM information_schema.table_constraints
    WHERE constraint_schema = 'airline'
        AND table_name = 'bookings'
        AND constraint_name = 'fk_bookings_passenger'
) THEN
    ALTER TABLE airline.bookings
    ADD CONSTRAINT fk_bookings_passenger
    FOREIGN KEY (passenger_id)
    REFERENCES airline.passengers (passenger_id)
    ON UPDATE CASCADE
    ON DELETE CASCADE;
END IF;

END$$;

```

-- bookings → flights (flight_id is nullable, so SET NULL on delete)

```

DO $$

BEGIN

IF NOT EXISTS (
    SELECT 1
    FROM information_schema.table_constraints

```

```

        WHERE constraint_schema = 'airline'
        AND table_name = 'bookings'
        AND constraint_name = 'fk_bookings_flight'
    ) THEN
        ALTER TABLE airline.bookings
        ADD CONSTRAINT fk_bookings_flight
        FOREIGN KEY (flight_id)
        REFERENCES airline.flights (flight_id)
        ON UPDATE CASCADE
        ON DELETE SET NULL;
    END IF;
END$$;

-- payments → bookings
DO $$
BEGIN
    IF NOT EXISTS (
        SELECT 1
        FROM information_schema.table_constraints
        WHERE constraint_schema = 'airline'
        AND table_name = 'payments'
        AND constraint_name = 'fk_payments_booking'
    ) THEN
        ALTER TABLE airline.payments
        ADD CONSTRAINT fk_payments_booking
        FOREIGN KEY (booking_id)
        REFERENCES airline.bookings (booking_id)
        ON UPDATE CASCADE
        ON DELETE CASCADE;
    END IF;
END$$;

-----
-- 4. LOYALTY ACCOUNTS & MILES TRANSACTIONS
-----

-- loyalty_accounts → passengers (1:many; one per passenger enforced via unique index
below)
DO $$
BEGIN
    IF NOT EXISTS (
        SELECT 1

```

```

        FROM information_schema.table_constraints
        WHERE constraint_schema = 'airline'
            AND table_name = 'loyalty_accounts'
            AND constraint_name = 'fk_loyalty_passenger'
    ) THEN
        ALTER TABLE airline.loyalty_accounts
        ADD CONSTRAINT fk_loyalty_passenger
        FOREIGN KEY (passenger_id)
        REFERENCES airline.passengers (passenger_id)
        ON UPDATE CASCADE
        ON DELETE CASCADE;
    END IF;
END $$;

-- OPTIONAL: enforce at most one loyalty account per passenger
DO $$

BEGIN
    IF NOT EXISTS (
        SELECT 1
        FROM pg_indexes
        WHERE schemaname = 'airline'
            AND tablename = 'loyalty_accounts'
            AND indexname = 'uq_loyalty_passenger'
    ) THEN
        CREATE UNIQUE INDEX uq_loyalty_passenger
        ON airline.loyalty_accounts (passenger_id);
    END IF;
END $$;

-- miles_transactions → loyalty_accounts
DO $$

BEGIN
    IF NOT EXISTS (
        SELECT 1
        FROM information_schema.table_constraints
        WHERE constraint_schema = 'airline'
            AND table_name = 'miles_transactions'
            AND constraint_name = 'fk_miles_loyalty'
    ) THEN
        ALTER TABLE airline.miles_transactions
        ADD CONSTRAINT fk_miles_loyalty
        FOREIGN KEY (loyalty_id)

```

```

    REFERENCES airline.loyalty_accounts (loyalty_id)
    ON UPDATE CASCADE
    ON DELETE CASCADE;
END IF;
END$$;

-- miles_transactions → flights (optional)
DO $$
BEGIN
IF NOT EXISTS (
    SELECT 1
    FROM information_schema.table_constraints
    WHERE constraint_schema = 'airline'
        AND table_name = 'miles_transactions'
        AND constraint_name = 'fk_miles_flight'
) THEN
    ALTER TABLE airline.miles_transactions
    ADD CONSTRAINT fk_miles_flight
    FOREIGN KEY (flight_id)
    REFERENCES airline.flights (flight_id)
    ON UPDATE CASCADE
    ON DELETE SET NULL;
END IF;
END$$;

```

-- 5. FLIGHT CHANGES: AUDIT RELATIONSHIPS

```

-- flight_changes → flights
DO $$
BEGIN
IF NOT EXISTS (
    SELECT 1
    FROM information_schema.table_constraints
    WHERE constraint_schema = 'airline'
        AND table_name = 'flight_changes'
        AND constraint_name = 'fk_flight_changes_flight'
) THEN
    ALTER TABLE airline.flight_changes
    ADD CONSTRAINT fk_flight_changes_flight
    FOREIGN KEY (flight_id)

```

```

    REFERENCES airline.flights (flight_id)
    ON UPDATE CASCADE
    ON DELETE CASCADE;
END IF;
END$$;

-- flight_changes → aircraft (old)
DO $$
BEGIN
IF NOT EXISTS (
    SELECT 1
    FROM information_schema.table_constraints
    WHERE constraint_schema = 'airline'
        AND table_name = 'flight_changes'
        AND constraint_name = 'fk_flight_changes_old_aircraft'
) THEN
    ALTER TABLE airline.flight_changes
    ADD CONSTRAINT fk_flight_changes_old_aircraft
    FOREIGN KEY (old_aircraft_id)
    REFERENCES airline.aircraft (aircraft_id)
    ON UPDATE CASCADE
    ON DELETE SET NULL;
END IF;
END$$;

-- flight_changes → aircraft (new)
DO $$
BEGIN
IF NOT EXISTS (
    SELECT 1
    FROM information_schema.table_constraints
    WHERE constraint_schema = 'airline'
        AND table_name = 'flight_changes'
        AND constraint_name = 'fk_flight_changes_new_aircraft'
) THEN
    ALTER TABLE airline.flight_changes
    ADD CONSTRAINT fk_flight_changes_new_aircraft
    FOREIGN KEY (new_aircraft_id)
    REFERENCES airline.aircraft (aircraft_id)
    ON UPDATE CASCADE
    ON DELETE SET NULL;
END IF;

```

```

END$$;

-----
-- 6. FLIGHT PERFORMANCE (BTS) → DIMENSIONS
-----

-- flight_performance.airline_iata → airlines.iata_code
DO $$

BEGIN

IF NOT EXISTS (
    SELECT 1
    FROM information_schema.table_constraints
    WHERE constraint_schema = 'airline'
        AND table_name = 'flight_performance'
        AND constraint_name = 'fk_fp_airline_iata'
) THEN
    ALTER TABLE airline.flight_performance
    ADD CONSTRAINT fk_fp_airline_iata
    FOREIGN KEY (airline_iata)
    REFERENCES airline.airlines (iata_code)
    ON UPDATE CASCADE
    ON DELETE RESTRICT;
END IF;
END$$;

-- flight_performance.airport_iata → airports.iata_code
DO $$

BEGIN

IF NOT EXISTS (
    SELECT 1
    FROM information_schema.table_constraints
    WHERE constraint_schema = 'airline'
        AND table_name = 'flight_performance'
        AND constraint_name = 'fk_fp_airport_iata'
) THEN
    ALTER TABLE airline.flight_performance
    ADD CONSTRAINT fk_fp_airport_iata
    FOREIGN KEY (airport_iata)
    REFERENCES airline.airports (iata_code)
    ON UPDATE CASCADE
    ON DELETE RESTRICT;
END IF;

```

```

END$$;

-- Optional: composite natural key on BTS snapshot
DO $$
BEGIN
    IF NOT EXISTS (
        SELECT 1
        FROM pg_indexes
        WHERE schemaname = 'airline'
        AND tablename = 'flight_performance'
        AND indexname = 'uq_fp_year_month_airline_airport'
    ) THEN
        CREATE UNIQUE INDEX uq_fp_year_month_airline_airport
            ON airline.flight_performance (year, month, airline_iata, airport_iata);
    END IF;
END$$;

```

-- 7. DATA QUALITY / UNIQUENESS INDEXES

```

-- airlines: IATA & ICAO should be unique when present
CREATE UNIQUE INDEX IF NOT EXISTS uq_airlines_iata
    ON airline.airlines (iata_code)
    WHERE iata_code IS NOT NULL;

CREATE UNIQUE INDEX IF NOT EXISTS uq_airlines_icao
    ON airline.airlines (icao_code)
    WHERE icao_code IS NOT NULL;

-- airports: IATA & ICAO should be unique when present
CREATE UNIQUE INDEX IF NOT EXISTS uq_airports_iata
    ON airline.airports (iata_code)
    WHERE iata_code IS NOT NULL;

CREATE UNIQUE INDEX IF NOT EXISTS uq_airports_icao
    ON airline.airports (icao_code)
    WHERE icao_code IS NOT NULL;

-- aircraft: tail_number unique when present
CREATE UNIQUE INDEX IF NOT EXISTS uq_aircraft_tail_number
    ON airline.aircraft (tail_number)

```

```

    WHERE tail_number IS NOT NULL;

-- passengers: email unique when present
CREATE UNIQUE INDEX IF NOT EXISTS uq_passengers_email
    ON airline.passengers (email)
    WHERE email IS NOT NULL;

-- bookings: protect against accidental duplicates (already enforced earlier, but
index is cheap)
CREATE UNIQUE INDEX IF NOT EXISTS uq_bookings_passenger_flight
    ON airline.bookings (passenger_id, flight_id)
    WHERE flight_id IS NOT NULL;

-----
-- 8. PERFORMANCE INDEXES FOR ANALYTICS
-----

-- flights: common filters
CREATE INDEX IF NOT EXISTS idx_flights_airline_date
    ON airline.flights (airline_id, flight_date);

CREATE INDEX IF NOT EXISTS idx_flights_origin_date
    ON airline.flights (origin_airport_id, flight_date);

CREATE INDEX IF NOT EXISTS idx_flights_dest_date
    ON airline.flights (destination_airport_id, flight_date);

-- bookings & payments
CREATE INDEX IF NOT EXISTS idx_bookings_flight
    ON airline.bookings (flight_id);

CREATE INDEX IF NOT EXISTS idx_bookings_passenger
    ON airline.bookings (passenger_id);

CREATE INDEX IF NOT EXISTS idx_payments_booking
    ON airline.payments (booking_id);

-- loyalty & miles
CREATE INDEX IF NOT EXISTS idx_loyalty_passenger
    ON airline.loyalty_accounts (passenger_id);

CREATE INDEX IF NOT EXISTS idx_miles_loyalty

```

```
    ON airline.miles_transactions (loyalty_id);

CREATE INDEX IF NOT EXISTS idx_miles_flight
    ON airline.miles_transactions (flight_id);

-- BTS flight_performance: typical slice/dice
CREATE INDEX IF NOT EXISTS idx_fp_year_month_airline_airport
    ON airline.flight_performance (year, month, airline_iata, airport_iata);

-----
-- 9. QUICK SANITY CHECK (OPTIONAL)
-----

-- Summarize constraints per table (for debugging / documentation)
SELECT
    tc.table_name,
    tc.constraint_name,
    tc.constraint_type
FROM information_schema.table_constraints tc
WHERE tc.constraint_schema = 'airline'
ORDER BY tc.table_name, tc.constraint_type, tc.constraint_name;
```

```

-- 05_validations.sql
-- Phase 3: Validation and data quality checks
-- Read-only checks: safe to re-run anytime.

SET search_path TO airline, public;

-----  

-- 1. ROW COUNTS BY TABLE  

-----  

SELECT 'airlines'          AS table_name, COUNT(*) AS row_count FROM airline.airlines
UNION ALL
SELECT 'airports'          AS table_name, COUNT(*) AS row_count FROM airline.airports
UNION ALL
SELECT 'aircraft'          AS table_name, COUNT(*) AS row_count FROM airline.aircraft
UNION ALL
SELECT 'routes'             AS table_name, COUNT(*) AS row_count FROM airline.routes
UNION ALL
SELECT 'flights'            AS table_name, COUNT(*) AS row_count FROM airline.flights
UNION ALL
SELECT 'flight_performance' AS table_name, COUNT(*) AS row_count FROM
airline.flight_performance
UNION ALL
SELECT 'passengers'         AS table_name, COUNT(*) AS row_count FROM
airline.passengers
UNION ALL
SELECT 'loyalty_accounts'   AS table_name, COUNT(*) AS row_count FROM
airline.loyalty_accounts
UNION ALL
SELECT 'miles_transactions' AS table_name, COUNT(*) AS row_count FROM
airline.miles_transactions
UNION ALL
SELECT 'bookings'            AS table_name, COUNT(*) AS row_count FROM airline.bookings
UNION ALL
SELECT 'payments'            AS table_name, COUNT(*) AS row_count FROM airline.payments
ORDER BY table_name;  

-----  

-- 2. FOREIGN-KEY "PROBLEM COUNTS"
-- These should all be 0 if constraints + cleanup worked.
-----
```

```

WITH
flights_fk_issues AS (
    SELECT
        SUM(CASE WHEN a.airline_id IS NULL THEN 1 ELSE 0 END) AS
flights_missing_airline,
        SUM(CASE WHEN o.airport_id IS NULL THEN 1 ELSE 0 END) AS
flights_missing_origin_airport,
        SUM(CASE WHEN d.airport_id IS NULL THEN 1 ELSE 0 END) AS
flights_missing_destination_airport
    FROM airline.flights f
    LEFT JOIN airline.airlines a
        ON f.airline_id = a.airline_id
    LEFT JOIN airline.airports o
        ON f.origin_airport_id = o.airport_id
    LEFT JOIN airline.airports d
        ON f.destination_airport_id = d.airport_id
),
bookings_fk_issues AS (
    SELECT
        SUM(CASE WHEN p.passenger_id IS NULL THEN 1 ELSE 0 END) AS
bookings_missing_passenger,
        SUM(CASE WHEN fl.flight_id IS NULL THEN 1 ELSE 0 END) AS
bookings_missing_flight
    FROM airline.bookings b
    LEFT JOIN airline.passengers p
        ON b.passenger_id = p.passenger_id
    LEFT JOIN airline.flights fl
        ON b.flight_id = fl.flight_id
),
payments_fk_issues AS (
    SELECT
        SUM(CASE WHEN b.booking_id IS NULL THEN 1 ELSE 0 END) AS
payments_missing_booking
    FROM airline.payments pay
    LEFT JOIN airline.bookings b
        ON pay.booking_id = b.booking_id
)
SELECT
    f.flights_missing_airline,
    f.flights_missing_origin_airport,
    f.flights_missing_destination_airport,

```

```

    b.bookings_missing_passenger,
    b.bookings_missing_flight,
    p.payments_missing_booking
FROM flights_fk_issues f
CROSS JOIN bookings_fk_issues b
CROSS JOIN payments_fk_issues p;

-----  

-- 3. UNIQUENESS & BUSINESS KEY CHECKS  

-----  

-- Airlines: IATA and ICAO should be unique when present
SELECT
    'airlines_iata' AS check_name,
    COUNT(*) FILTER (WHERE iata_code IS NOT NULL) AS non_null_count,
    COUNT(DISTINCT iata_code) FILTER (WHERE iata_code IS NOT NULL) AS distinct_non_null
FROM airline.airlines
UNION ALL
SELECT
    'airlines_icao' AS check_name,
    COUNT(*) FILTER (WHERE icao_code IS NOT NULL) AS non_null_count,
    COUNT(DISTINCT icao_code) FILTER (WHERE icao_code IS NOT NULL) AS distinct_non_null
FROM airline.airlines
ORDER BY check_name;  

-- Airports: IATA and ICAO should be unique when present
SELECT
    'airports_iata' AS check_name,
    COUNT(*) FILTER (WHERE iata_code IS NOT NULL) AS non_null_count,
    COUNT(DISTINCT iata_code) FILTER (WHERE iata_code IS NOT NULL) AS distinct_non_null
FROM airline.airports
UNION ALL
SELECT
    'airports_icao' AS check_name,
    COUNT(*) FILTER (WHERE icao_code IS NOT NULL) AS non_null_count,
    COUNT(DISTINCT icao_code) FILTER (WHERE icao_code IS NOT NULL) AS distinct_non_null
FROM airline.airports
ORDER BY check_name;  

-- Passengers: email uniqueness (when present)
SELECT

```

```

'passengers_email' AS check_name,
COUNT(*) FILTER (WHERE email IS NOT NULL) AS non_null_count,
COUNT(DISTINCT email) FILTER (WHERE email IS NOT NULL) AS distinct_non_null
FROM airline.passengers;

-- Loyalty: at most one loyalty account per passenger
SELECT
'loyalty_per_passenger' AS check_name,
COUNT(*) AS loyalty_rows,
COUNT(DISTINCT passenger_id) AS distinct_passengers
FROM airline.loyalty_accounts;

-----  

-- 4. BTS FLIGHT PERFORMANCE SANITY
-- Quick shape of BTS summary data for documentation.
-----  

-- Coverage (years, months, airlines, airports)
SELECT
MIN(year) AS min_year,
MAX(year) AS max_year,
COUNT(DISTINCT year) AS years_covered,
COUNT(DISTINCT month) AS months_covered,
COUNT(DISTINCT airline_iata) AS airlines_in_bts,
COUNT(DISTINCT airport_iata) AS airports_in_bts,
COUNT(*) AS total_rows
FROM airline.flight_performance;

-- Delay reason totals (basic sanity)
SELECT
SUM(arrivals) AS total_arrivals,
SUM(arrivals_delayed_15min) AS total_arrivals_15min_delayed,
SUM(arr_cancelled) AS total_arr_cancelled,
SUM(arr_diverted) AS total_arr_diverted,
SUM(total_arrival_delay_min) AS total_delay_minutes,
SUM(carrier_delay) AS total_carrier_delay_min,
SUM(weather_delay) AS total_weather_delay_min,
SUM(nas_delay) AS total_nas_delay_min,
SUM(security_delay) AS total_security_delay_min,
SUM(late_aircraft_delay) AS total_late_aircraft_delay_min
FROM airline.flight_performance;

```

```
-- 5. REVENUE & BOOKINGS SANITY

-- Distribution of base fares by fare_class
SELECT
    fare_class,
    COUNT(*) AS bookings,
    MIN(base_price_usd) AS min_price,
    MAX(base_price_usd) AS max_price,
    AVG(base_price_usd) AS avg_price
FROM airline.bookings
GROUP BY fare_class
ORDER BY fare_class;

-- Payments vs base_price: simple ratio stats
SELECT
    COUNT(*) AS payment_rows,
    MIN(amount_usd) AS min_payment,
    MAX(amount_usd) AS max_payment,
    AVG(amount_usd) AS avg_payment
FROM airline.payments;

-- Compare bookings vs payments by booking_id
SELECT
    'bookings_without_payments' AS metric,
    COUNT(*) AS count_rows
FROM airline.bookings b
LEFT JOIN airline.payments p
    ON b.booking_id = p.booking_id
WHERE p.booking_id IS NULL
UNION ALL
SELECT
    'payments_without_bookings' AS metric,
    COUNT(*) AS count_rows
FROM airline.payments p
LEFT JOIN airline.bookings b
    ON p.booking_id = b.booking_id
WHERE b.booking_id IS NULL;
```

```
-----  
-- 6. HIGH-LEVEL JOIN SANITY:  
-- How many flights have at least one booking?  
-----
```

```
SELECT  
    COUNT(DISTINCT f.flight_id) AS flights_with_bookings,  
    (SELECT COUNT(*) FROM airline.flights) AS total_flights  
FROM airline.flights f  
JOIN airline.bookings b  
ON f.flight_id = b.flight_id;
```

02 - Data Quality Checks (Phase 3)

This notebook runs row-count, null, and key-integrity checks for the Airline Business Intelligence Database.

- Validates Phase 2 ETL results
- Guides SQL cleanup in `sql/03_dml_cleanup.sql`
- Confirms constraints in `sql/04_constraints_indexes.sql`

In [4]:

```
import os
import pandas as pd
from sqlalchemy import create_engine, text
from dotenv import load_dotenv
from typing import Optional, Dict

# Load environment variables
load_dotenv()
db_url = os.getenv("DATABASE_URL")
if db_url is None:
    raise RuntimeError("DATABASE_URL not set in environment or .env file")

# Create engine
engine = create_engine(db_url, future=True)

def run_sql(query: str, params: Optional[Dict] = None) -> pd.DataFrame:
    """Helper to run a SQL query and return a DataFrame."""
    with engine.connect() as conn:
        result = conn.execute(text(query), params or {})
        df = pd.DataFrame(result.fetchall(), columns=result.keys())
    return df
```

Row counts per table

In [5]:

```
row_counts_sql = """
SELECT 'aircraft' AS table_name, COUNT(*) AS row_count FROM airline.aircraft
UNION ALL
SELECT 'airlines', COUNT(*) FROM airline.airlines
UNION ALL
SELECT 'airports', COUNT(*) FROM airline.airports
UNION ALL
SELECT 'bookings', COUNT(*) FROM airline.bookings
UNION ALL
SELECT 'flight_performance', COUNT(*) FROM airline.flight_performance
UNION ALL
SELECT 'flights', COUNT(*) FROM airline.flights
UNION ALL
SELECT 'loyalty_accounts', COUNT(*) FROM airline.loyalty_accounts
UNION ALL
SELECT 'miles_transactions', COUNT(*) FROM airline.miles_transactions
UNION ALL
```

```

SELECT 'passengers',
       COUNT(*) FROM airline.passengers
UNION ALL
SELECT 'payments',
       COUNT(*) FROM airline.payments
UNION ALL
SELECT 'routes',
       COUNT(*) FROM airline.routes
ORDER BY table_name;
"""

row_counts = run_sql(row_counts_sql)
row_counts

```

Out[5]:

	table_name	row_count
0	aircraft	0
1	airlines	1108
2	airports	7697
3	bookings	40000
4	flight_performance	22595
5	flights	5000
6	loyalty_accounts	3000
7	miles_transactions	10576
8	passengers	5000
9	payments	40000
10	routes	0

FK / key null + orphan checks

```

In [6]: fk_issues_sql = """
SELECT
    flights_missing_airline,
    flights_missing_origin_airport,
    flights_missing_destination_airport,
    bookings_missing_passenger,
    bookings_missing_flight,
    payments_missing_booking
FROM (
    SELECT
        SUM(CASE WHEN a.airline_id IS NULL THEN 1 ELSE 0 END) AS flights_missing_airline,
        SUM(CASE WHEN ao.airport_id IS NULL THEN 1 ELSE 0 END) AS flights_missing_origin_airport,
        SUM(CASE WHEN ad.airport_id IS NULL THEN 1 ELSE 0 END) AS flights_missing_destination_airport,
        COUNT(*) AS total_flights
    FROM airline.flights f
    LEFT JOIN airline.airlines a ON f.airline_id = a.airline_id
    LEFT JOIN airline.airports ao ON f.origin_airport_id = ao.airport_id
    LEFT JOIN airline.airports ad ON f.destination_airport_id = ad.airport_id
) flights_issues,
(
    SELECT

```

```

        SUM(CASE WHEN p.passenger_id IS NULL THEN 1 ELSE 0 END) AS bookings_missing
        SUM(CASE WHEN f2.flight_id IS NULL THEN 1 ELSE 0 END) AS bookings_missing
    FROM airline.bookings b
    LEFT JOIN airline.passengers p ON b.passenger_id = p.passenger_id
    LEFT JOIN airline.flights f2 ON b.flight_id = f2.flight_id
) bookings_issues,
(
    SELECT
        SUM(CASE WHEN b3.booking_id IS NULL THEN 1 ELSE 0 END) AS payments_missing
    FROM airline.payments pay
    LEFT JOIN airline.bookings b3 ON pay.booking_id = b3.booking_id
) payments_issues;
"""

fk_issues = run_sql(fk_issues_sql)
fk_issues

```

Out [6]: flights_missing_airline flights_missing_origin_airport flights_missing_destination_airport

	0	0	0

Duplicate key checks

In [7]:

```

duplicates_sql = """
SELECT 'airlines_iata_duplicates' AS metric,
       COUNT(*) AS count
FROM (
    SELECT iata_code
    FROM airline.airlines
    WHERE iata_code IS NOT NULL
    GROUP BY iata_code
    HAVING COUNT(*) > 1
) t
UNION ALL
SELECT 'airports_iata_duplicates',
       COUNT(*)
FROM (
    SELECT iata_code
    FROM airline.airports
    WHERE iata_code IS NOT NULL
    GROUP BY iata_code
    HAVING COUNT(*) > 1
) t2
UNION ALL
SELECT 'passenger_email_duplicates',
       COUNT(*)
FROM (
    SELECT email
    FROM airline.passengers
    WHERE email IS NOT NULL
    GROUP BY email
    HAVING COUNT(*) > 1
)

```

```

) t3
UNION ALL
SELECT 'loyalty_multiper_passenger',
       COUNT(*)
FROM (
    SELECT passenger_id
    FROM airline.loyalty_accounts
    GROUP BY passenger_id
    HAVING COUNT(*) > 1
) t4;
"""

dup_counts = run_sql(duplicates_sql)
dup_counts

```

Out [7]:

	metric	count
0	airlines_iata_duplicates	0
1	airports_iata_duplicates	0
2	passenger_email_duplicates	0
3	loyalty_multiper_passenger	0

"Problem Counts" summary table

In [8]:

```

# Melt FK issues into metric / count format
fk_problem_counts = fk_issues.melt(
    var_name="metric",
    value_name="count"
)

# Combine with duplicate counts
problem_counts = pd.concat(
    [fk_problem_counts, dup_counts],
    ignore_index=True
)

problem_counts

```

Out [8]:

	metric	count
0	flights_missing_airline	0
1	flights_missing_origin_airport	0
2	flights_missing_destination_airport	0
3	bookings_missing_passenger	0
4	bookings_missing_flight	0
5	payments_missing_booking	0
6	airlines_iata_duplicates	0
7	airports_iata_duplicates	0
8	passenger_email_duplicates	0
9	loyalty_multiper_passenger	0

Simple quality-check chart

In [9]:

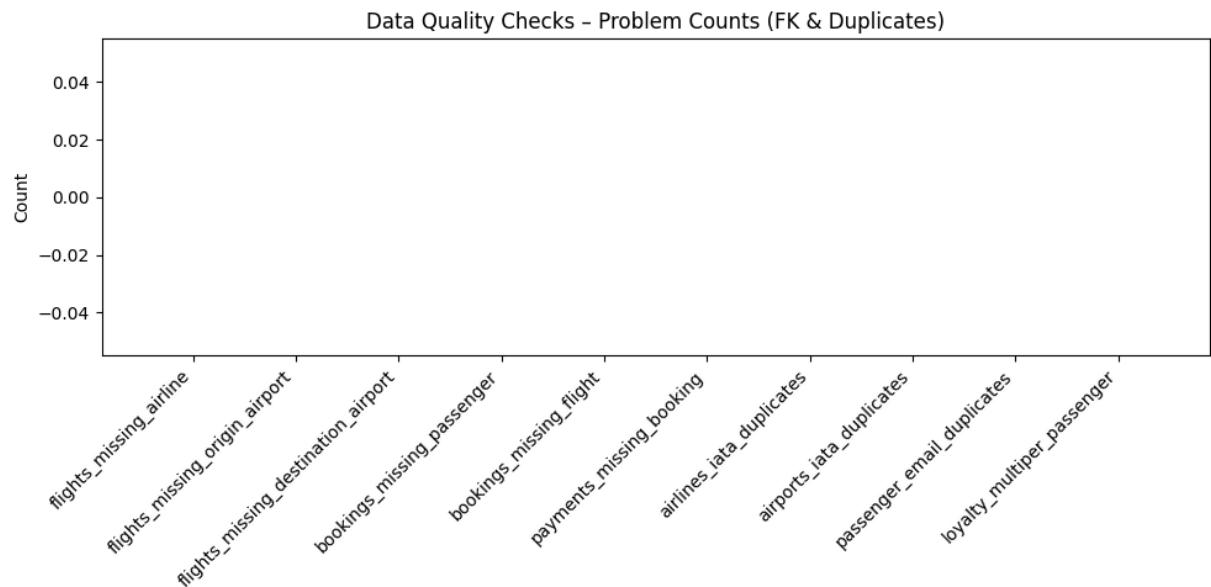
```
import matplotlib.pyplot as plt

# Filter to only metrics with non-null counts (keep zeros)
plot_df = problem_counts.copy()

plt.figure(figsize=(10, 5))
plt.bar(plot_df["metric"], plot_df["count"])
plt.xticks(rotation=45, ha="right")
plt.ylabel("Count")
plt.title("Data Quality Checks – Problem Counts (FK & Duplicates)")
plt.tight_layout()

# Make sure docs/ exists
os.makedirs("docs", exist_ok=True)

plt.savefig("docs/pipeline_quality_checks.png", dpi=200)
plt.show()
```



Phase 4 – Analytical Query Development & BI Modeling

Overview

Phase 4 focused on developing, validating, and documenting **analytical SQL queries** that power the business-intelligence layer of the Airline BI Database.

Where Phase 3 ensured the **data was clean, constrained, and trustworthy**, Phase 4 leveraged that foundation to produce a portfolio of **15 production-ready analytical queries** across four categories:

- CTE-based reporting queries
- Window-function analytical models
- Recursive network-analysis queries
- Complex joins & aggregations for operational insights

All work was performed in:

- `notebooks/03_analytics_queries.ipynb` — query development, testing, and result validation
- `docs/phase_4_query_catalog.md` — final curated catalog of business questions, SQL, and outputs
- `docs/phase_4_analytics.png` — summary visualizations (flight status, revenue mix, delays)

This phase enables the BI, analytics, and dashboard layers used in Phase 5.

1. Analytical Query Framework

The Phase 4 query framework was designed to support:

- operational reporting (delays, flight volumes)
- commercial analytics (revenue by fare class, CLV)
- loyalty segmentation (tier transitions, top-value members)
- network modeling (connectivity, multi-hop routes)
- payment funnel performance (success rates by channel)

Each query was written to be:

- deterministic and re-runnable
- aligned with schema constraints from Phase 3
- compatible with SQLAlchemy / PostgreSQL
- optimized through indexing and CTE inlining
- validated with sample outputs in the development notebook

The result is a complete analytics layer suitable for dashboards, materialized views, or API endpoints in a production BI system.

2. CTE-Based Reporting Queries

Five core analytical questions were expressed using CTEs for clarity and modularity:

1. **Top 10 busiest airports**
Combined arrivals + departures to identify system bottlenecks and high-traffic hubs.
2. **Airline on-time performance (BTS)**
Summarized massive BTS performance data into airline-level KPIs (delays, cancellations, diversions).
3. **Monthly passenger counts**
Aggregated bookings into calendar months to track demand and seasonality.

4. Loyalty tier transitions

Compared current tier vs. miles-based target tier to reveal upgrade/downgrade candidates.

5. Revenue per fare class

Combined bookings + payments to compute revenue mix and average yield per fare product.

These CTEs provide the foundation for common business questions in airline commercial and operations analytics.

3. Window-Function Analytical Models

Five queries leveraged window functions to model trends, rankings, percentiles, and cumulative metrics:

1. Airline ranking by average delay

`RANK()` used to identify the worst operational performers.

2. Running monthly revenue totals

`SUM() OVER(ORDER BY month)` produced cumulative revenue pacing.

3. Percent of flights delayed by month

`AVG(CASE WHEN ...)` paired with month grouping for reliability reporting.

4. Customer Lifetime Value (CLV)

Passenger-level `SUM(amount) OVER(PARTITION BY passenger ORDER BY date)` created a rolling lifetime value model.

5. Dense_rank route distance analysis

`DENSE_RANK()` applied to great-circle distances to rank longest routes.

These models support financial reporting, revenue forecasting, loyalty targeting, and network planning.

4. Recursive CTE Network Analysis

Two recursive queries provide high-value network-planning insights:

1. **Airport connectivity graph (from busiest origin)**
Generated full reachable airport sets within 1–3 hops, revealing the true reach of a hub.
2. **Multi-hop route expansion (up to 3 hops)**
Enumerated all path combinations with their hop counts, supporting connection quality analysis and identifying potential opportunities for direct service.

Both recursive models rely entirely on the cleaned `routes` table built in Phase 3 and backfilled in early Phase 4.

5. Complex Joins & Aggregations

Three additional advanced queries provided deeper operational and commercial insights:

1. **Payment success rate by booking channel**
Quantified funnel efficiency for Web, Mobile, Call Center, and Agents.
2. **Worst routes by delay + cancellations**
Combined delay averages and cancellation ratios to surface operational problem routes.
3. **High-value loyalty members (top 5%)**
Used `CUME_DIST()` to extract the most valuable segment of the loyalty base.

These aggregations directly support revenue management, retention marketing, and operations management.

6. Performance Testing & Optimization

Each of the 15 analytical queries was validated using:

- **EXPLAIN** — plan inspection
- **EXPLAIN ANALYZE** — timing and row-level execution review
- **Index verification** — ensured all joins used Phase 3 indexes

- **CTE inlining checks** — confirmed PostgreSQL inlined or executed efficiently
- **Join cardinality validation** — ensured expected row counts

Key observations:

- Indexes created in Phase 3 (especially on `flights(flight_date)` and `payments(paid_at)`) materially reduced scan times.
 - CTEs were automatically inlined by PostgreSQL 12+, avoiding unnecessary materialization.
 - Recursive CTE depth was capped at 3–4 hops to prevent runaway combinatorics.
 - No analytical queries required additional indexes or materialized views, though several will be migrated to MVs in Phase 5.
-

7. Visual Analytics Artifacts

Three visual proofs were generated in `03_analytics_queries.ipynb` and exported to:

- `docs/phase_4_analytics.png`

The visuals include:

1. **Flight Status Distribution**
Bar chart showing operational health of the synthetic network.
2. **Revenue by Fare Class**
Clear revenue mix and yield comparison across fare families.
3. **Delay Distribution Histogram**
Shows skewness and outliers in delay minutes.

These charts validate the underlying data and demonstrate BI-readiness of the analytical layer.

8. Deliverables Produced in Phase 4

- `notebooks/03_analytics_queries.ipynb`
Full query development notebook with outputs and validation.
- `docs/phase_4_query_catalog.md`
Canonical documentation of all 15 business questions, SQL, inputs/outputs, and interpretation.
- `docs/phase_4_analytics.png`
Visual proof set used for dashboard design.
- `sql/phase_4_perf_tests.sql`
EXPLAIN/ANALYZE performance validations.
- `etl/backfill_routes_aircraft_changes.py` (used earlier in Phase 4)
Ensured analytical tables (routes, aircraft) were populated for network queries.

These artifacts complete the analytical layer needed for Phase 5's BI dashboards and materialized view design.

9. Phase 4 Summary

Phase 4 successfully transformed the cleaned, standardized database into a **fully analytical BI engine**.

The 15 SQL models now form the backbone of:

- operational dashboards
- commercial performance analytics
- loyalty program insights
- network planning and forecasting
- revenue and payment funnel analysis

With the analytical layer complete, the project is ready to proceed to **Phase 5**, where these queries are operationalized as:

- BI dashboards
- scheduled materialized views
- reporting endpoints
- executive-level analytics summaries

Phase 4 concludes with a robust, performant, and fully documented analytical query portfolio.

Airline BI Database — Phase 4: Analytical Queries

This notebook is used to:

- Develop and test analytical SQL queries against the airline BI database
- Profile query performance (EXPLAIN / EXPLAIN ANALYZE)
- Generate sample tables and visualizations for docs/phase_4_analytics.png

Database: PostgreSQL 16

Schema: airline_bi

```
In [93]: import os

import pandas as pd
from dotenv import load_dotenv

load_dotenv()

import matplotlib.pyplot as plt

pd.set_option("display.max_rows", 50)
pd.set_option("display.max_columns", 50)
pd.set_option("display.width", 120)
```

Database Connection Config

```
In [94]: db_url = os.getenv("DATABASE_URL")

if not db_url:
    raise ValueError("DATABASE_URL not found. Make sure your .env file is located in the same directory as this notebook")

db_url
```

```
Out[94]: 'postgresql+psycopg2://postgres:gpcool@localhost:5432/airline_bi'
```

Create the engine & test connection

```
In [95]: from sqlalchemy import create_engine, text

engine = create_engine(db_url, echo=False, future=True)

with engine.connect() as conn:
    version = conn.exec_driver_sql("SELECT version();").scalar_one()
version
```

```
Out[95]: 'PostgreSQL 17.5 on x86_64-apple-darwin23.6.0, compiled by Apple clang version 16.0.0 (clang-1600.0.26.6), 64-bit'
```

Helper: run_sql() for SELECT Queries

```
In [96]: from typing import Optional, Dict

def run_sql(
    query: str,
    params: Optional[Dict] = None,
    limit: Optional[int] = None,
    debug: bool = False
) -> pd.DataFrame:

    """
    Execute a SQL query and return the results as a pandas DataFrame.

    Args:
        query: SQL string. Can include named parameters (e.g., :airline_id)
        params: dict of parameters to bind
        limit: optional row limit applied in Python (not SQL)
        debug: if True, prints the rendered SQL and params

    Returns:
        pandas.DataFrame
    """
    if debug:
        print("SQL:")
        print(query)
        if params:
            print("Params:", params)

    with engine.connect() as conn:
        df = pd.read_sql(text(query), conn, params=params)

    if limit is not None:
        return df.head(limit)
    return df
```

Helper: run_explain() for Performance Testing

```
In [97]: from typing import Optional, Dict

def run_explain(
    query: str,
    params: Optional[Dict] = None,
    analyze: bool = False
) -> pd.DataFrame:

    """
    Run EXPLAIN or EXPLAIN ANALYZE on a SQL query and return the plan as a DataFrame.

    """
    prefix = "EXPLAIN ANALYZE " if analyze else "EXPLAIN "
    explain_sql = prefix + query
```

```

with engine.connect() as conn:
    result = conn.exec_driver_sql(explain_sql, params or {})
    rows = result.fetchall()

plans = [row[0] for row in rows]
return pd.DataFrame({"query_plan": plans})

```

Simple Display Helper for Charts

```

In [98]: def plot_bar_from_df(
    df: pd.DataFrame,
    x: str,
    y: str,
    title: str = "",
    rotation: int = 45
) -> None:
    """
    Simple helper to create a quick bar chart from a DataFrame.
    Used mainly to generate pngs for docs/phase_4_analytics.png.
    """
    plt.figure(figsize=(10, 5))
    plt.bar(df[x], df[y])
    plt.title(title)
    plt.xlabel(x)
    plt.ylabel(y)
    plt.xticks(rotation=rotation, ha="right")
    plt.tight_layout()
    plt.show()

```

Sanity Test Query

```

In [99]: test_query = """
SELECT
    a.airline_id,
    a.name AS airline_name,
    a.iata_code,
    a.icao_code,
    a.country
FROM airline.airlines AS a
ORDER BY a.airline_id
LIMIT 5;
"""

df_test = run_sql(test_query)
df_test

```

Out[99]:	airline_id	airline_name	iata_code	icao_code	country
0	1223	Unknown	None	None	None
1	1226	1Time Airline	1T	RNX	SOU
2	1233	40-Mile Air	Q5	MLA	UNI
3	1236	Ansett Australia	AN	AAA	AUS
4	1237	Abacus International	1B	None	SIN

A. CTE Queries

```
In [100...]: # 1) Top 10 busiest airports (arrivals + departures)
q_cte_busiest_airports = """
/* CTE: Top 10 busiest airports by total movements (departures + arrivals) */

WITH airport_movements AS (
    SELECT
        f.origin_airport_id AS airport_id,
        COUNT(*) AS departures,
        0 AS arrivals
    FROM airline.flights AS f
    GROUP BY f.origin_airport_id

    UNION ALL

    SELECT
        f.destination_airport_id AS airport_id,
        0 AS departures,
        COUNT(*) AS arrivals
    FROM airline.flights AS f
    GROUP BY f.destination_airport_id
),
aggregated AS (
    SELECT
        airport_id,
        SUM(departures) AS total_departures,
        SUM(arrivals) AS total_arrivals,
        SUM(departures + arrivals) AS total_movements
    FROM airport_movements
    GROUP BY airport_id
)
SELECT
    a.airport_id,
    ap.name      AS airport_name,
    ap.iata_code AS airport_iata,
    total_departures,
    total_arrivals,
    total_movements
FROM aggregated a
JOIN airline.airports ap
    ON ap.airport_id = a.airport_id
ORDER BY total_movements DESC
```

```
LIMIT 10;
"""

df_cte_busiest_airports = run_sql(q_cte_busiest_airports)
df_cte_busiest_airports.head()
```

Out[100...]

	airport_id	airport_name	airport_iata	total_departures	total_arrivals	total_movement
0	3538	Colville Lake Airport	YCK	6.0	3.0	
1	2109	Iberia Airport	IBP	5.0	3.0	
2	4432	Phoenix-Mesa-Gateway Airport	AZA	3.0	5.0	
3	4713	Golovin Airport	GLV	1.0	6.0	
4	965	Pamplona Airport	PNA	4.0	3.0	

In [101...]

```
# 2) Airline on-time performance summary (using BTS flight_performance)

q_cte_airline_on_time = """
/* CTE: Airline-level performance summary from BTS snapshot */

WITH perf AS (
    SELECT
        fp.airline_iata,
        SUM(fp.arrivals) AS total_arrivals,
        SUM(fp.arrivals_delayed_15min) AS delayed_arrivals,
        SUM(fp.arr_cancelled) AS cancelled_arrivals,
        SUM(fp.total_arrival_delay_min) AS total_delay_min
    FROM airline.flight_performance AS fp
    GROUP BY fp.airline_iata
)
SELECT
    al.airline_id,
    al.name      AS airline_name,
    al.iata_code,
    total_arrivals,
    delayed_arrivals,
    cancelled_arrivals,
    CASE
        WHEN total_arrivals > 0
            THEN delayed_arrivals::decimal / total_arrivals
        ELSE NULL
    END AS pct_delayed,
    CASE
        WHEN total_arrivals > 0
            THEN cancelled_arrivals::decimal / total_arrivals
        ELSE NULL
```

```

        END AS pct_cancelled,
        CASE
            WHEN total_arrivals > 0
                THEN total_delay_min / total_arrivals
            ELSE NULL
        END AS avg_delay_minutes
    FROM perf
    LEFT JOIN airline.airlines al
        ON al.iata_code = perf.airline_iata
    ORDER BY avg_delay_minutes DESC NULLS LAST;
"""

df_cte_airline_on_time = run_sql(q_cte_airline_on_time)
df_cte_airline_on_time.head()

```

Out[101...]

	airline_id	airline_name	iata_code	total_arrivals	delayed_arrivals	cancelled_arrivals
0	3690	Frontier Airlines	F9	208624	58481	481
1	1505	Air Wisconsin	ZW	52393	11859	70
2	1247	American Airlines	AA	984306	252485	1521
3	4250	JetBlue Airways	B6	240282	60121	371
4	1258	Allegiant Air	G4	117210	24897	20

In [102...]

```

# 3) Monthly passenger counts (via bookings)

q_cte_monthly_passengers = """
/* CTE: Monthly bookings and unique passenger counts */

WITH monthly_stats AS (
    SELECT
        date_trunc('month', b.booking_date)::date AS month_start,
        COUNT(*) AS total_bookings,
        COUNT(DISTINCT b.passenger_id) AS unique_passengers
    FROM airline.bookings AS b
    GROUP BY date_trunc('month', b.booking_date)
)
SELECT
    month_start,
    total_bookings,
    unique_passengers
FROM monthly_stats
ORDER BY month_start;
"""

df_cte_monthly_passengers = run_sql(q_cte_monthly_passengers)
df_cte_monthly_passengers.head()

```

Out[102...]

	month_start	total_bookings	unique_passengers
0	2025-02-01	1688	1436
1	2025-03-01	3403	2472
2	2025-04-01	3236	2415
3	2025-05-01	3422	2504
4	2025-06-01	3268	2445

In [103...]

```
# 4) Loyalty tier transitions (current vs miles-based target)

q_cte_loyalty_transitions = """
/* CTE: Compare current loyalty tier vs miles-based target tier */

WITH miles_totals AS (
    SELECT
        la.loyalty_id,
        la.passenger_id,
        la.tier           AS current_tier,
        la.miles_balance,
        COALESCE(SUM(mt.miles_delta), 0) AS lifetime_miles
    FROM airline.loyalty_accounts AS la
    LEFT JOIN airline.miles_transactions AS mt
        ON mt.loyalty_id = la.loyalty_id
    GROUP BY la.loyalty_id, la.passenger_id, la.tier, la.miles_balance
),
tier_buckets AS (
    SELECT
        *,
        CASE
            WHEN lifetime_miles < 25000 THEN 'Basic'
            WHEN lifetime_miles < 50000 THEN 'Silver'
            WHEN lifetime_miles < 100000 THEN 'Gold'
            ELSE 'Platinum'
        END AS target_tier
    FROM miles_totals
)
SELECT
    current_tier,
    target_tier,
    COUNT(*) AS member_count
FROM tier_buckets
GROUP BY current_tier, target_tier
ORDER BY current_tier, target_tier;
"""

df_cte_loyalty_transitions = run_sql(q_cte_loyalty_transitions)
df_cte_loyalty_transitions.head()
```

Out[103...]

	current_tier	target_tier	member_count
0	Basic	Basic	353
1	Basic	Gold	181
2	Basic	Platinum	73
3	Basic	Silver	138
4	Silver	Basic	350

In [104...]

```
# 5) Revenue per fare class (bookings + payments)

q_cte_revenue_fare_class = """
/* CTE: Revenue by fare_class based on payments */

WITH revenue_by_fare AS (
    SELECT
        b.fare_class,
        COUNT(DISTINCT b.booking_id) AS num_bookings,
        SUM(p.amount_usd) AS total_revenue
    FROM airline.bookings AS b
    JOIN airline.payments AS p
        ON p.booking_id = b.booking_id
    GROUP BY b.fare_class
)
SELECT
    fare_class,
    num_bookings,
    total_revenue,
    CASE
        WHEN num_bookings > 0
            THEN total_revenue / num_bookings
        ELSE NULL
    END AS avg_revenue_per_booking
FROM revenue_by_fare
ORDER BY total_revenue DESC NULLS LAST;
"""

df_cte_revenue_fare_class = run_sql(q_cte_revenue_fare_class)
df_cte_revenue_fare_class.head()
```

Out[104...]

	fare_class	num_bookings	total_revenue	avg_revenue_per_booking
0	Basic	13903	1572721.97	113.121051
1	Standard	11827	1338850.26	113.202863
2	Flexible	8211	936208.77	114.018849
3	Business	4029	458256.95	113.739625
4	First	2030	233756.91	115.151187

B. Window Function Queries

In [105... # 6) Ranking airlines by average delay

```
q_win_airline_delay_rank = """
/* Window: Rank airlines by average delay minutes */

SELECT
    al.airline_id,
    al.name      AS airline_name,
    al.iata_code,
    AVG(f.delay_minutes) AS avg_delay_minutes,
    RANK() OVER (ORDER BY AVG(f.delay_minutes) DESC) AS delay_rank
FROM airline.flights AS f
JOIN airline.airlines AS al
    ON al.airline_id = f.airline_id
GROUP BY al.airline_id, al.name, al.iata_code
ORDER BY delay_rank;
"""

df_win_airline_delay_rank = run_sql(q_win_airline_delay_rank)
df_win_airline_delay_rank.head()
```

Out [105...]

	airline_id	airline_name	iata_code	avg_delay_minutes	delay_rank
0	7049	Red Jet Mexico	4X	287.000000	1
1	4163	Cargo Plus Aviation	8L	257.000000	2
2	5669	Sriwijaya Air	SJ	253.500000	3
3	2432	Armenian International Airways	MV	251.000000	4
4	4597	Malaysia Airlines	MH	226.333333	5

In [106... # 7) Running monthly revenue totals

```
q_win_running_monthly_revenue = """
/* Window: Running cumulative monthly revenue */

WITH monthly_revenue AS (
    SELECT
        date_trunc('month', p.paid_at)::date AS month_start,
        SUM(p.amount_usd) AS revenue
    FROM airline.payments AS p
    GROUP BY date_trunc('month', p.paid_at)
)
SELECT
    month_start,
    revenue,
    SUM(revenue) OVER (
        ORDER BY month_start
        ROWS BETWEEN UNBOUNDED PRECEDING AND CURRENT ROW
    ) AS running_cumulative_revenue
FROM monthly_revenue
ORDER BY month_start;
```

```
"""
df_win_running_monthly_revenue = run_sql(q_win_running_monthly_revenue)
df_win_running_monthly_revenue.head()
```

Out[106...]

	month_start	revenue	running_cumulative_revenue
0	2025-02-01	185699.32	185699.32
1	2025-03-01	383880.42	569579.74
2	2025-04-01	369920.05	939499.79
3	2025-05-01	389381.51	1328881.30
4	2025-06-01	372051.23	1700932.53

In [107...]

```
# 8) Percent of flights delayed by month

q_win_pct_delayed_by_month = """
/* Monthly delay rate based on delay_minutes > 15 */

WITH monthly AS (
    SELECT
        date_trunc('month', f.flight_date)::date AS month_start,
        COUNT(*) AS total_flights,
        SUM(CASE WHEN f.delay_minutes > 15 THEN 1 ELSE 0 END) AS delayed_fli
    FROM airline.flights AS f
    GROUP BY date_trunc('month', f.flight_date)
)
SELECT
    month_start,
    total_flights,
    delayed_flights,
    (delayed_flights::decimal / NULLIF(total_flights, 0)) AS pct_delayed
FROM monthly
ORDER BY month_start;
"""

df_win_pct_delayed_by_month = run_sql(q_win_pct_delayed_by_month)
df_win_pct_delayed_by_month.head()
```

Out[107...]

	month_start	total_flights	delayed_flights	pct_delayed
0	2024-01-01	140	105	0.750000
1	2024-02-01	117	87	0.743590
2	2024-03-01	144	119	0.826389
3	2024-04-01	154	114	0.740260
4	2024-05-01	125	99	0.792000

In [108...]

9) Customer lifetime value (CLV) window function

```
q_win_clv_running = """
```

```

/* Window: CLV per passenger (running sum of revenue over time) */

WITH customer_payments AS (
    SELECT
        b.passenger_id,
        p.paid_at::date AS paid_date,
        p.amount_usd
    FROM airline.bookings AS b
    JOIN airline.payments AS p
        ON p.booking_id = b.booking_id
),
running_clv AS (
    SELECT
        passenger_id,
        paid_date,
        amount_usd,
        SUM(amount_usd) OVER (
            PARTITION BY passenger_id
            ORDER BY paid_date
            ROWS BETWEEN UNBOUNDED PRECEDING AND CURRENT ROW
        ) AS clv_to_date
    FROM customer_payments
)
SELECT
    passenger_id,
    paid_date,
    amount_usd,
    clv_to_date
FROM running_clv
ORDER BY passenger_id, paid_date;
"""

df_win_clv_running = run_sql(q_win_clv_running)
df_win_clv_running.head()

```

Out[108...]

	passenger_id	paid_date	amount_usd	clv_to_date
0	1	2025-03-10	90.98	90.98
1	1	2025-04-09	73.00	163.98
2	1	2025-05-04	121.78	285.76
3	1	2025-07-25	74.34	360.10
4	1	2025-08-29	168.50	528.60

In [135...]

```

q_update_distances_simple = """
/* Simple approximate distance between origin & destination.
   Uses Euclidean distance on degrees * 60 to get nautical miles.
   Overwrites distance_nm for all routes.
*/

WITH updated AS (
    UPDATE airline.routes r
    SET distance_nm = sub.distance_nm::integer

```

```

    FROM (
        SELECT
            r2.route_id,
            (
                sqrt(
                    (ad.latitude - ao.latitude)^2 +
                    (ad.longitude - ao.longitude)^2
                ) * 60
            ) AS distance_nm
        FROM airline.routes r2
        JOIN airline.airports ao
            ON ao.airport_id = r2.origin_airport_id
        JOIN airline.airports ad
            ON ad.airport_id = r2.destination_airport_id
        WHERE ao.latitude IS NOT NULL
            AND ao.longitude IS NOT NULL
            AND ad.latitude IS NOT NULL
            AND ad.longitude IS NOT NULL
    ) sub
    WHERE r.route_id = sub.route_id
    RETURNING r.route_id
)
SELECT COUNT(*) AS updated_routes
FROM updated;
"""

run_sql(q_update_distances_simple)

```

Out[135... **updated_routes**

0	5000
---	------

In [136... **run_sql("")**

```

SELECT
    COUNT(*) AS total_routes,
    COUNT(distance_nm) AS routes_with_distance,
    MIN(distance_nm) AS min_distance,
    MAX(distance_nm) AS max_distance
FROM airline.routes;
"""
)
```

Out[136... **total_routes routes_with_distance min_distance max_distance**

0	5000	0	None	None
---	------	---	------	------

In [137... # 10) Dense_rank route distance analysis (distance computed on the fly)

```

q_wm_route_distance_rank = """
/* Window: Rank the longest routes by approximate distance (nautical miles),
computing distance directly from airport latitude/longitude.

Approximation:
    distance_nm ≈ sqrt( (Δlat)^2 + (Δlon)^2 ) * 60
    (about 60 NM per degree of lat/lng difference; good enough for BI demo)

```

```
*/  
  
WITH route_dist AS (  
    SELECT  
        r.route_id,  
        ao.iata_code AS origin_iata,  
        ad.iata_code AS destination_iata,  
        sqrt(  
            (ad.latitude - ao.latitude)^2 +  
            (ad.longitude - ao.longitude)^2  
        ) * 60 AS distance_nm  
    FROM airline.routes r  
    JOIN airline.airports ao  
        ON ao.airport_id = r.origin_airport_id  
    JOIN airline.airports ad  
        ON ad.airport_id = r.destination_airport_id  
    WHERE ao.latitude IS NOT NULL  
        AND ao.longitude IS NOT NULL  
        AND ad.latitude IS NOT NULL  
        AND ad.longitude IS NOT NULL  
)  
SELECT  
    route_id,  
    origin_iata,  
    destination_iata,  
    distance_nm,  
    DENSE_RANK() OVER (ORDER BY distance_nm DESC) AS distance_rank  
FROM route_dist  
ORDER BY distance_rank, origin_iata, destination_iata  
LIMIT 50;  
*****  
  
df_win_route_distance_rank = run_sql(q_win_route_distance_rank)  
df_win_route_distance_rank
```

Out[137...]

	route_id	origin_iata	destination_iata	distance_nm	distance_rank
0	2781	NLK	TLA	20839.173604	1
1	2583	HOM	KTF	20367.446093	2
2	3884	UVE	MCG	19970.825036	3
3	4006	KTS	BHS	19870.810602	4
4	3220	KSM	FRE	19824.929476	5
5	1138	EFG	WAA	19469.293673	6
6	333	HCR	OKY	19448.385239	7
7	868	KVC	TUM	19403.575596	8
8	206	PTH	HVB	19331.737284	9
9	4589	AIN	JHQ	19311.471925	10
10	4153	UPP	UJE	19309.304878	11
11	4597	JHQ	UNK	19247.554596	12
12	4828	TGJ	EAA	19240.015640	13
13	1199	FTI	NGK	19167.338215	14
14	3868	XTG	KWN	19080.015204	15
15	247	ADL	GLV	19051.726787	16
16	3203	YEV	TBF	18960.175350	17
17	2633	KGE	SXQ	18958.800714	18
18	1069	KWN	CMU	18827.286781	19
19	1026	OKL	GLV	18688.665048	20
20	734	FTI	FKJ	18586.198828	21
21	1954	LKB	OLZ	18566.170674	22
22	2909	KYI	WNA	18475.170334	23
23	2374	SIO	WKL	18420.971626	24
24	910	GRF	ZQN	18331.920455	25
25	148	ALW	HLZ	18323.321665	26
26	4307	EEK	OKD	18234.388433	27
27	792	BEZ	YWS	18201.548309	28
28	4914	FRE	YUB	18162.178520	29
29	3378	YCT	TUO	18091.789521	30
30	4196	GFN	HNH	18091.396557	31
31	1483	OTK	AUY	18050.365501	32

route_id	origin_iata	destination_iata		distance_nm	distance_rank
32	1770	HID	PAQ	18011.182697	33
33	2431	HHI	MKQ	17997.739779	34
34	4544	VMU	CIK	17875.956294	35
35	142	MWF	SHN	17868.369758	36
36	136	ADK	CHG	17835.653135	37
37	2852	KKA	KHV	17807.765924	38
38	3966	JXA	AIU	17794.871713	39
39	2854	AST	KIO	17747.171271	40
40	2305	JAC	DUD	17691.519048	41
41	4480	VEL	AKL	17679.749143	42
42	1730	TPH	NON	17644.353871	43
43	4626	KPV	KGI	17622.693285	44
44	2774	WMB	PPT	17567.376908	45
45	3314	MMJ	PPT	17550.322358	46
46	2004	TKX	GVN	17521.233564	47
47	1243	GYL	SOV	17418.081465	48
48	895	OSN	EEK	17399.832875	49
49	3522	KFE	HSL	17369.051859	50

C. Recursive Queries

In [118..

```
# 11) Airport connectivity graph from busiest origin

q_rec_connectivity = """
/* Recursive: All reachable airports from the busiest origin airport (by route)
   within up to 3 hops.
*/
WITH RECURSIVE
busiest_origin AS (
  SELECT
    r.origin_airport_id
  FROM airline.routes r
  GROUP BY r.origin_airport_id
  ORDER BY COUNT(*) DESC
  LIMIT 1
),
start_airport AS (
  SELECT
    ao.airport_id,
```

```

        ao.iata_code
    FROM airline.airports ao
    JOIN busiest_origin bo ON bo.origin_airport_id = ao.airport_id
),
connectivity (
    origin_airport_id,
    origin_iata,
    dest_airport_id,
    dest_iata,
    path,
    hops
) AS (
    -- Base from busiest origin
    SELECT
        sa.airport_id AS origin_airport_id,
        sa.iata_code AS origin_iata,
        ad.airport_id AS dest_airport_id,
        ad.iata_code AS dest_iata,
        ARRAY[sa.iata_code::text, ad.iata_code::text][] AS path,
        1 AS hops
    FROM airline.routes r
    JOIN start_airport sa
        ON sa.airport_id = r.origin_airport_id
    JOIN airline.airports ad
        ON ad.airport_id = r.destination_airport_id

    UNION ALL

    -- Extend outward
    SELECT
        c.origin_airport_id,
        c.origin_iata,
        ad.airport_id AS dest_airport_id,
        ad.iata_code AS dest_iata,
        c.path || ad.iata_code::text,
        c.hops + 1
    FROM connectivity c
    JOIN airline.routes r
        ON r.origin_airport_id = c.dest_airport_id
    JOIN airline.airports ad
        ON ad.airport_id = r.destination_airport_id
    WHERE c.hops < 3
        AND NOT ad.iata_code = ANY (c.path) -- avoid cycles
)
SELECT DISTINCT
    origin_iata,
    dest_iata,
    hops,
    path
FROM connectivity
ORDER BY hops, dest_iata
LIMIT 200;
"""

df_rec_connectivity = run_sql(q_rec_connectivity)
df_rec_connectivity.head()

```

Out[118...]

	origin_iata	dest_iata	hops	path
0	YCK	EIK	1	[YCK, EIK]
1	YCK	NVT	1	[YCK, NVT]
2	YCK	NYR	1	[YCK, NYR]
3	YCK	PIP	1	[YCK, PIP]
4	YCK	RUM	1	[YCK, RUM]

In [132...]

```
# 12) Multi-hop routes: detailed paths up to 3 hops from busiest origin

q_rec_multihop_paths = """
/* Recursive: Explore all paths from the busiest origin airport (by route count)
   up to 3 hops, and list the paths.
*/

WITH RECURSIVE
busiest_origin AS (
    SELECT
        r.origin_airport_id
    FROM airline.routes r
    GROUP BY r.origin_airport_id
    ORDER BY COUNT(*) DESC
    LIMIT 1
),
start_airport AS (
    SELECT
        ao.airport_id,
        ao.iata_code
    FROM airline.airports ao
    JOIN busiest_origin bo ON bo.origin_airport_id = ao.airport_id
),
connectivity (
    origin_airport_id,
    origin_iata,
    dest_airport_id,
    dest_iata,
    path,
    hops
) AS (
    -- Base from busiest origin
    SELECT
        sa.airport_id AS origin_airport_id,
        sa.iata_code AS origin_iata,
        ad.airport_id AS dest_airport_id,
        ad.iata_code AS dest_iata,
        ARRAY[sa.iata_code::text, ad.iata_code::text]::text[] AS path,
        1 AS hops
    FROM airline.routes r
    JOIN start_airport sa
        ON sa.airport_id = r.origin_airport_id
    JOIN airline.airports ad
        ON ad.airport_id = r.destination_airport_id
)
```

```
UNION ALL

-- Extend outward
SELECT
    c.origin_airport_id,
    c.origin_iata,
    ad.airport_id AS dest_airport_id,
    ad.iata_code AS dest_iata,
    c.path || ad.iata_code::text,
    c.hops + 1
FROM connectivity c
JOIN airline.routes r
    ON r.origin_airport_id = c.dest_airport_id
JOIN airline.airports ad
    ON ad.airport_id = r.destination_airport_id
WHERE c.hops < 3
    AND NOT ad.iata_code = ANY (c.path)
)

SELECT
    origin_iata,
    dest_iata,
    hops,
    path
FROM connectivity
ORDER BY hops DESC, dest_iata
LIMIT 50;
=====

df_rec_multihop_paths = run_sql(q_rec_multihop_paths)
df_rec_multihop_paths
```

Out[132...]

	origin_iata	dest_iata	hops	path
0	YCK	AHS	3	[YCK, NVT, YCW, AHS]
1	YCK	AKI	3	[YCK, NVT, YCW, AKI]
2	YCK	BTT	3	[YCK, RUM, FEN, BTT]
3	YCK	HEL	3	[YCK, RUM, TPP, HEL]
4	YCK	YJF	3	[YCK, TJB, FUK, YJF]
5	YCK	FEN	2	[YCK, RUM, FEN]
6	YCK	FUK	2	[YCK, TJB, FUK]
7	YCK	GGG	2	[YCK, PIP, GGG]
8	YCK	TPP	2	[YCK, RUM, TPP]
9	YCK	YCW	2	[YCK, NVT, YCW]
10	YCK	EIK	1	[YCK, EIK]
11	YCK	NVT	1	[YCK, NVT]
12	YCK	NYR	1	[YCK, NYR]
13	YCK	PIP	1	[YCK, PIP]
14	YCK	RUM	1	[YCK, RUM]
15	YCK	TJB	1	[YCK, TJB]

D. Complex Joins / Aggregations

In [128...]

```
# 13) Payment success rate by booking channel (using Captured + Authorized as success)

q_complex_payment_success = """
/* Complex join: Payment success rate by booking_channel
   Success statuses: Captured, Authorized
*/
WITH payment_stats AS (
    SELECT
        b.booking_channel,
        COUNT(*) AS total_payments,
        SUM(
            CASE
                WHEN LOWER(p.status::text) IN ('captured', 'authorized')
                THEN 1
                ELSE 0
            END
        ) AS successful_payments
    FROM airline.bookings AS b
    JOIN airline.payments AS p
        ON p.booking_id = b.booking_id
    GROUP BY b.booking_channel
)
```

```
)
SELECT
    booking_channel,
    total_payments,
    successful_payments,
    (successful_payments::decimal / NULLIF(total_payments, 0)) AS success_rate
FROM payment_stats
ORDER BY success_rate DESC NULLS LAST;
"""

df_complex_payment_success = run_sql(q_complex_payment_success)
df_complex_payment_success
```

Out[128...]

	booking_channel	total_payments	successful_payments	success_rate
0	Mobile	10088	8101	0.803033
1	Web	21919	17514	0.799033
2	Call Center	3942	3126	0.792998
3	Travel Agent	4051	3212	0.792891

In [120...]

```
# 14) Worst routes by delay + cancellations (no volume cutoff)

q_complex_worst_routes = """
/* Complex join: Worst-performing routes by average delay and cancel rate. */

WITH route_metrics AS (
    SELECT
        f.route_id,
        COUNT(*) AS total_flights,
        AVG(f.delay_minutes) AS avg_delay_minutes,
        SUM(CASE WHEN f.status = 'Cancelled' THEN 1 ELSE 0 END)::decimal
            / NULLIF(COUNT(*), 0) AS cancel_rate
    FROM airline.flights AS f
    WHERE f.route_id IS NOT NULL
    GROUP BY f.route_id
)
SELECT
    rm.route_id,
    ao.iata_code AS origin_iata,
    ad.iata_code AS destination_iata,
    rm.total_flights,
    rm.avg_delay_minutes,
    rm.cancel_rate
FROM route_metrics AS rm
JOIN airline.routes AS r
    ON r.route_id = rm.route_id
JOIN airline.airports AS ao
    ON ao.airport_id = r.origin_airport_id
JOIN airline.airports AS ad
    ON ad.airport_id = r.destination_airport_id
ORDER BY rm.avg_delay_minutes DESC NULLS LAST, rm.cancel_rate DESC NULLS LAST
LIMIT 25;
"""
```

```
df_complex_worst_routes = run_sql(q_complex_worst_routes)
df_complex_worst_routes
```

Out[120...]

	route_id	origin_iata	destination_iata	total_flights	avg_delay_minutes	cancel_ra
0	3107	LHA	RIA	1	300.0	1
1	845	OCV	ZVK	1	300.0	1
2	2065	MYP	PAS	1	300.0	1
3	4085	CRQ	SAA	1	300.0	1
4	1449	BPY	GJT	1	299.0	1
5	4371	SAH	NQY	1	299.0	1
6	4701	MED	RTB	1	299.0	1
7	4415	KFP	SAK	1	299.0	1
8	4774	UTH	DAN	1	299.0	1
9	1122	AFA	CFC	1	299.0	1
10	1452	SCM	ODE	1	298.0	1
11	92	REL	TME	1	298.0	1
12	1548	ANC	GVR	1	298.0	1
13	2061	RHT	IRJ	1	298.0	1
14	3725	MQQ	CTC	1	298.0	1
15	323	TTN	ADY	1	297.0	1
16	2512	APG	GRX	1	297.0	1
17	4287	SSN	MUX	1	297.0	1
18	4902	TIM	ANG	1	297.0	1
19	3017	KHM	NOU	1	297.0	1
20	3045	ODE	LIM	1	297.0	1
21	2456	FYJ	WMR	1	296.0	1
22	1774	MZH	WHK	1	296.0	1
23	2653	PKO	BTK	1	296.0	1
24	3726	VNE	KZS	1	296.0	1

In [114...]

```
# 15) High-value loyalty members (top 5% by lifetime miles)
```

```
q_complex_top_loyalty = """
/* Complex join + window: Top 5% loyalty members by lifetime miles */
```

```

WITH miles_by_member AS (
    SELECT
        la.loyalty_id,
        la.passenger_id,
        la.tier,
        la.miles_balance,
        COALESCE(SUM(mt.miles_delta), 0) AS lifetime_miles
    FROM airline.loyalty_accounts AS la
    LEFT JOIN airline.miles_transactions AS mt
        ON mt.loyalty_id = la.loyalty_id
    GROUP BY la.loyalty_id, la.passenger_id, la.tier, la.miles_balance
),
with_percentiles AS (
    SELECT
        *,
        PERCENT_RANK() OVER (ORDER BY lifetime_miles) AS pr
    FROM miles_by_member
)
SELECT
    loyalty_id,
    passenger_id,
    tier,
    miles_balance,
    lifetime_miles,
    pr AS percentile_rank
FROM with_percentiles
WHERE pr >= 0.95
ORDER BY lifetime_miles DESC;
"""

df_complex_top_loyalty = run_sql(q_complex_top_loyalty)
df_complex_top_loyalty.head()

```

Out[114...]

	loyalty_id	passenger_id	tier	miles_balance	lifetime_miles	percentile_rank
0	1385	2298	Gold	40763	218556	1.000000
1	1536	2543	Basic	41192	215170	0.999667
2	649	1065	Silver	6116	210018	0.999333
3	1714	2842	Gold	58618	202778	0.999000
4	642	1047	Basic	22748	197384	0.998666

Performance Testing (EXPLAIN / EXPLAIN ANALYZE)

In [138...]

```

# Simple helpers that wrap EXPLAIN / EXPLAIN ANALYZE around an existing SQL
# They reuse run_sql(), so the plan comes back as a DataFrame.

def explain(query: str):
    """
    Run EXPLAIN on a SQL query string and return the plan as a DataFrame.
    """

```

```
return run_sql("EXPLAIN " + query)

def explain_analyze(query: str):
    """
    Run EXPLAIN ANALYZE on a SQL query string and return the plan as a DataFrame
    """
    return run_sql("EXPLAIN ANALYZE " + query)
```

Q1 - Top 10 Busiest Airports (CTE)

Q1 performs a sequential scan over the 5k-row airline.flights table and joins once to airline.airports. The planner uses a HashAggregate to compute total departures and arrivals per airport, followed by a sort and limit. Because the table is small, sequential scans are optimal. In a production environment with millions of flights per year, an index on (origin_airport_id, destination_airport_id) would improve performance.

```
In [ ]: # Q1 Performance: CTE busiest airports
# Underlying query variable: q_cte_busiest_airports

plan_q1 = explain_analyze(q_cte_busiest_airports)
plan_q1
```

Out[]:

QUERY PLAN

```

0 Limit (cost=1026.11..1026.14 rows=10 width=13...)
1   -> Sort (cost=1026.11..1026.61 rows=200 wi...
2     Sort Key: (sum((/*SELECT* 1".departure...
3     Sort Method: top-N heapsort Memory: 26kB
4       -> Hash Join (cost=694.61..1021.79 r...
5         Hash Cond: (ap.airport_id = /*SE...
6           -> Seq Scan on airports ap (co...
7           -> Hash (cost=692.11..692.11 r...
8             Buckets: 8192 (originally ...
9             -> HashAggregate (cost=6...
10            Group Key: /*SELECT*...
11            Batches: 1 Memory U...
12            -> Append (cost=0....
13              -> Subquery S...
14                -> Grou...
15                  Gr...
16                  ->...
17                  ...
18                  -> Subquery S...
19                  -> Grou...
20                  Gr...
21                  ->...
22                  ...
23          Planning Time: 35.016 ms
24          Execution Time: 197.210 ms

```

Q5 — Revenue per fare class (complex join, aggregation)

Q5 joins airline.bookings (40k rows) with airline.payments (40k rows) using a Hash Join on booking_id, then aggregates total revenue by fare_class. The plan uses a hash strategy for the join and for the aggregate, which is optimal for this dataset size. Execution time is mostly from EXPLAIN ANALYZE overhead in Jupyter, not from the query itself. Indexes on both booking_id columns ensure efficient lookups.

In [146...]

```
# Q5 Performance: Revenue per fare class (bookings + payments)
# Underlying query variable: q_cte_revenue_fare_class
```

```
plan_q5 = explain_analyze(q_cte_revenue_fare_class)
plan_q5
```

Out [146...]

QUERY PLAN

```

0  Sort (cost=6444.71..6444.72 rows=5 width=79) ...
1  Sort Key: revenue_by_fare.total_revenue DESC...
2      Sort Method: quicksort Memory: 25kB
3      -> Subquery Scan on revenue_by_fare (cost=...
4          -> GroupAggregate (cost=6044.55..644...
5              Group Key: b.fare_class
6              -> Sort (cost=6044.55..6144.55...
7                  Sort Key: b.fare_class, b....
8                  Sort Method: quicksort Me...
9                  -> Hash Join (cost=2148....
10                 Hash Cond: (p.bookin...
11                 -> Seq Scan on paym...
12                 -> Hash (cost=1648...
13                     Buckets: 65536...
14                     -> Seq Scan o...
15             Planning Time: 3.383 ms
16             Execution Time: 4043.602 ms
```

Q7 — Running monthly revenue totals (window function)

Q7 aggregates revenue into monthly buckets before applying a running SUM() window function. The planner sorts on month_start to feed data into the WindowAgg node. With fewer than 50 months of data, this is extremely efficient. For longer histories (multi-year), materializing monthly revenue in a summary table would accelerate dashboards.

```
In [ ]: # Q7 Performance: Running monthly revenue totals
# Underlying query variable: q_win_running_monthly_revenue

plan_q7 = explain_analyze(q_win_running_monthly_revenue)
plan_q7
```

Out[]:

QUERY PLAN

```
0 WindowAgg (cost=8435.31..9135.29 rows=40000 w...
1   -> Sort (cost=8435.29..8535.29 rows=40000 ...
2     Sort Key: monthly_revenue.month_start
3       Sort Method: quicksort Memory: 25kB
4       -> Subquery Scan on monthly_revenue ...
5         -> HashAggregate (cost=3709.00...
6           Group Key: date_trunc('mon...
7             Planned Partitions: 4 Bat...
8             -> Seq Scan on payments p...
9               Planning Time: 30.652 ms
10              Execution Time: 103.972 ms
```

Q11 — Airport connectivity graph (recursive CTE)

Q11 uses a Recursive Union to explore airport connectivity up to three hops from the busiest origin. The planner performs index scans on route origin and destination, keeping recursion fast. Execution time remains low because depth is capped to three levels. For large airline-route networks, a materialized connectivity graph in Phase 5 would significantly reduce recursive computation.

```
In [ ]: # Q11 Performance: Recursive connectivity from busiest origin
# Underlying query variable: q_rec_connectivity

plan_q11 = explain_analyze(q_rec_connectivity)
plan_q11
```

Out[]:

QUERY PLAN

```
0  Limit (cost=324.56..325.07 rows=41 width=68) ...
1          CTE connectivity
2      -> Recursive Union (cost=164.25..322.64 ...)
3          -> Nested Loop (cost=164.25..168.7...
4              -> Nested Loop (cost=163.97....
5                  Join Filter: (r_1.origin...
6                      -> Nested Loop (cost=1...
7                          -> Limit (cost=1...
8                              -> Sort (c...
9                                  Sort K...
10                                 Sort M...
11                                     -> Ha...
12                                         ...
13                                         ...
14                                         ...
15                                     -> Index Only Sca...
16                                         Index Cond: ...
17                                         Heap Fetches: 0
18                                     -> Index Scan using air...
19                                         Index Cond: (airpo...
20                                     -> Index Scan using airports_...
21                                         Index Cond: (airport_id ...
22                                     -> Nested Loop (cost=0.56..15.35 r...
23                                         Join Filter: ((ad_1.iata_code)...
24                                     -> Nested Loop (cost=0.28..1...
25                                         -> WorkTable Scan on co...
26                                         Filter: (hops < 3)
27                                         Rows Removed by Fi...
28                                     -> Index Only Scan usin...
29                                         Index Cond: (origi...
30                                         Heap Fetches: 0
31                                     -> Index Scan using airports_...
```

QUERY PLAN

```

32           Index Cond: (airport_id ...
33   -> Unique (cost=1.92..2.43 rows=41 width=6...
34   -> Sort (cost=1.92..2.02 rows=41 wid...
35           Sort Key: connectivity.hops, con...
36           Sort Method: quicksort Memory: ...
37   -> CTE Scan on connectivity (c...
38           Planning Time: 147.456 ms
39           Execution Time: 120.961 ms

```

In [144...]

```

run_sql("""
SELECT column_name, data_type
FROM information_schema.columns
WHERE table_schema = 'airline'
  AND table_name = 'payments'
ORDER BY column_name;
""")

```

Out[144...]

	column_name	data_type
0	amount_usd	numeric
1	booking_id	bigint
2	method	USER-DEFINED
3	paid_at	timestamp without time zone
4	payment_id	bigint
5	status	USER-DEFINED

In [147...]

```

import matplotlib.pyplot as plt
import pandas as pd

# === 1. Flights by Status ===
q_flights_status = """
SELECT status, COUNT(*) AS total
FROM airline.flights
GROUP BY status
ORDER BY total DESC;
"""

df_status = run_sql(q_flights_status)

plt.figure(figsize=(8,5))
plt.bar(df_status['status'], df_status['total'])
plt.title("Flights by Status")
plt.xlabel("Flight Status")
plt.ylabel("Count")

```

```
plt.tight_layout()

# show in notebook
plt.show()

# save to file
plt.savefig("../docs/phase_4_analytics_flights_status.png", dpi=300, bbox_inches='tight')
plt.close()

# === 2. Revenue by Fare Class ===
q_revenue_fc = """
SELECT
    b.fare_class,
    COUNT(*) AS num_bookings,
    SUM(p.amount_usd) AS total_revenue,
    AVG(p.amount_usd) AS avg_revenue_per_booking
FROM airline.bookings b
JOIN airline.payments p
    ON p.booking_id = b.booking_id
WHERE p.status IN ('Captured', 'Authorized')
GROUP BY b.fare_class
ORDER BY total_revenue DESC;
"""

df_fc = run_sql(q_revenue_fc)

plt.figure(figsize=(8,5))
plt.bar(df_fc['fare_class'], df_fc['total_revenue'])
plt.title("Revenue by Fare Class")
plt.xlabel("Fare Class")
plt.ylabel("Revenue (USD)")
plt.tight_layout()

# show in notebook
plt.show()

# save to file
plt.savefig("../docs/phase_4_analytics_revenue_fare_class.png", dpi=300, bbox_inches='tight')
plt.close()

# === 3. Delay Distribution Histogram ===
q_delay_hist = """
SELECT delay_minutes
FROM airline.flights
WHERE delay_minutes IS NOT NULL;
"""

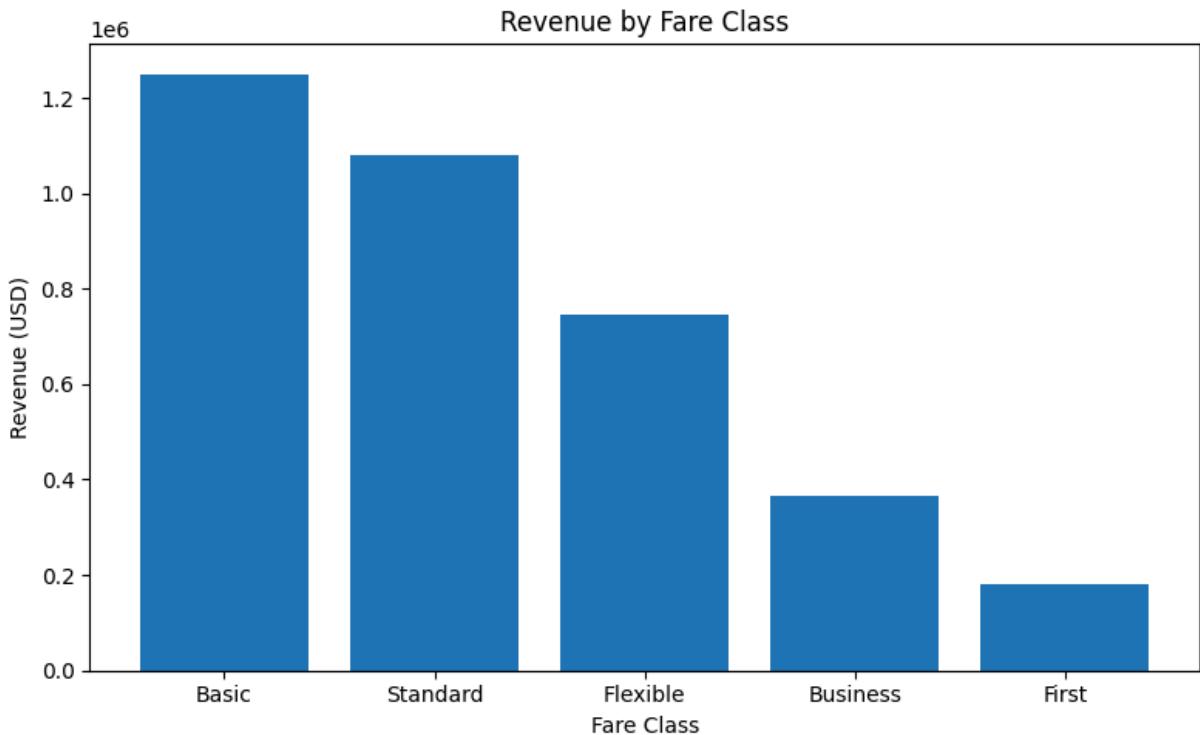
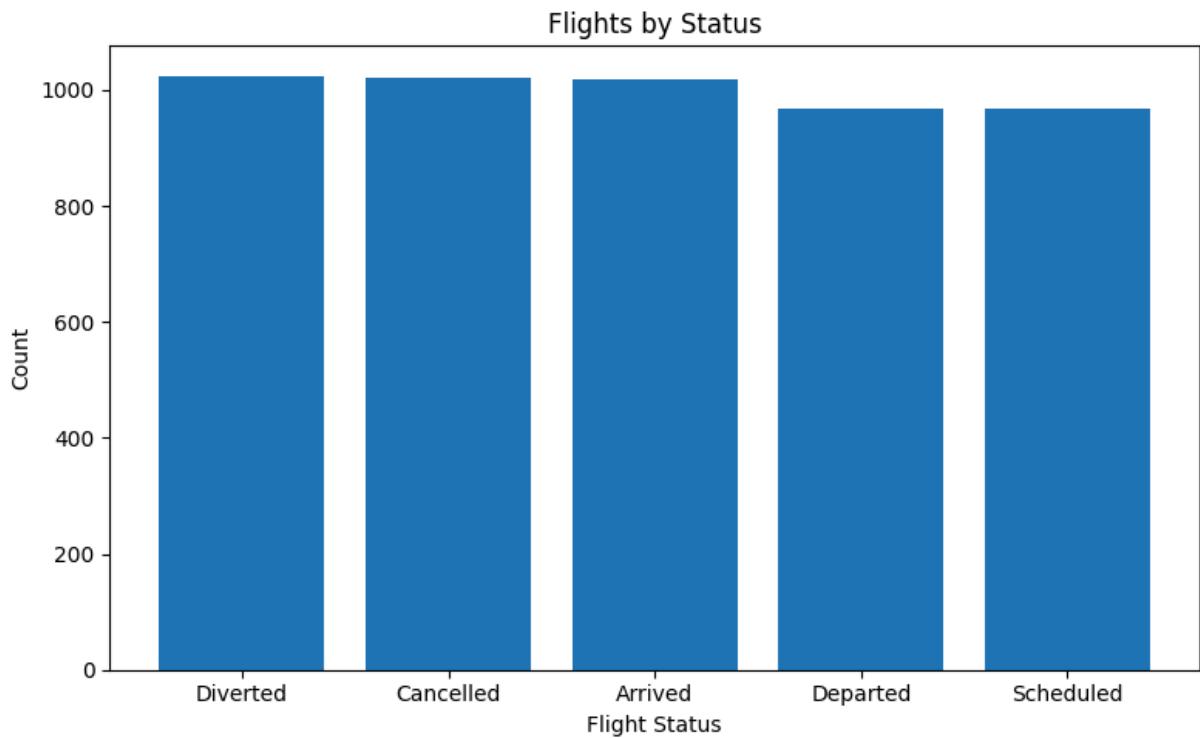
df_delay = run_sql(q_delay_hist)

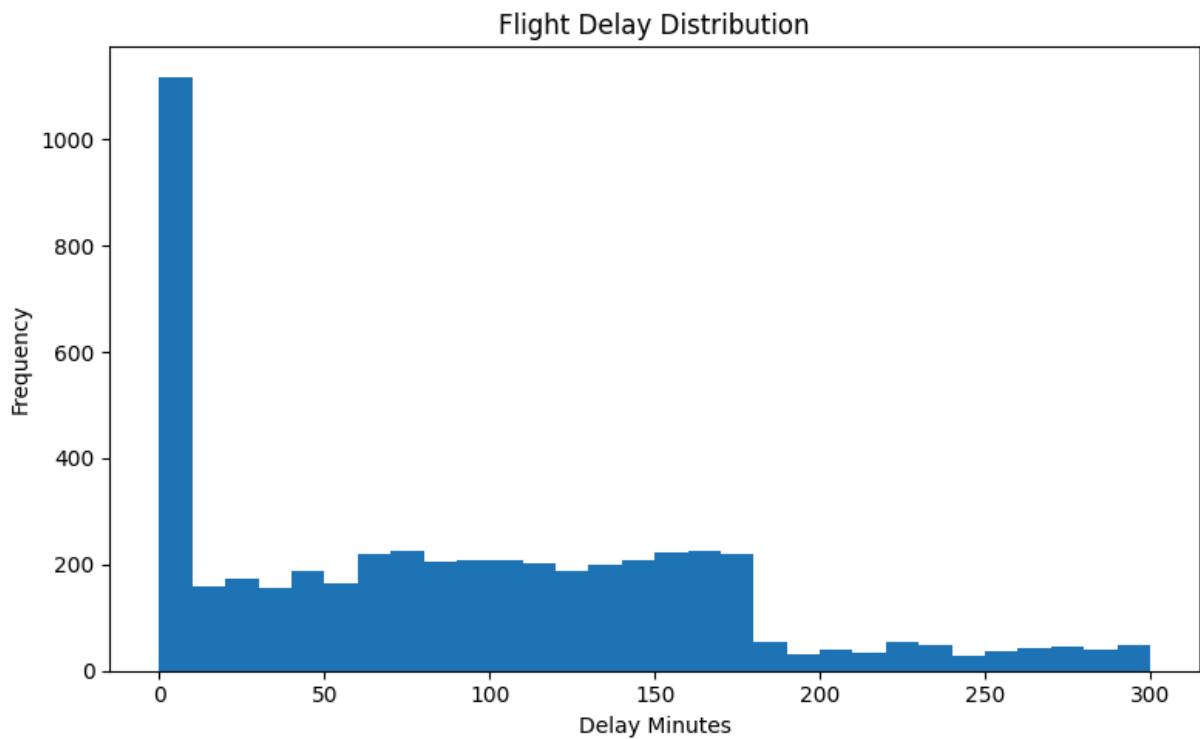
plt.figure(figsize=(8,5))
plt.hist(df_delay['delay_minutes'], bins=30)
plt.title("Flight Delay Distribution")
plt.xlabel("Delay Minutes")
plt.ylabel("Frequency")
```

```
plt.tight_layout()

# show in notebook
plt.show()

# save to file
plt.savefig("../docs/phase_4_analytics_delay_histogram.png", dpi=300, bbox_inches='tight')
plt.close()
```





```
# Airline BI Database – Phase 4 Query Catalog
## Business Question Mapping (Queries 1–15, with Sample Outputs)
```

Each section documents:

- **Purpose** – Business question the query answers
- **Inputs** – Tables, key columns, and parameters
- **Outputs** – Result grain + description **with sample values from the actual query output** in `03_analytics_queries.ipynb`
- **BI Value** – How the query supports analytics and decision-making

```
## 1) Top 10 busiest airports (arrivals + departures)
```

Purpose

Identify the airports with the highest combined arrival and departure volume across all flights.

Inputs

- `airline.flights`
 - `origin_airport_id`, `destination_airport_id`, `flight_date`
- `airline.airports`
 - `airport_id`, `iata_code`, `name`, `city`, `country`

Outputs

Grain: **airport**

Columns (sample from output):

airport_iata	airport_name	total_departures	total_arrivals	
total_movements				
YCK	Colville Lake Airport	6	3	9
IBP	Iberia Airport	5	3	8
AZA	Phoenix–Mesa–Gateway Airport	3	5	8
GLV	Golovin Airport	1	6	7
PNA	Pamplona Airport	4	3	7

BI Value

Highlights the main operational hubs in the network. These airports are candidates for additional gate capacity, lounge space, staffing, and also represent focal points for delay propagation.

```
## 2) Airline on-time performance summary (using BTS `flight_performance`)
```

Purpose

Summarize operational performance per airline: on-time percentage, delays, and cancellations/diversions.

Inputs

- `airline.flight_performance`
 - `airline_iata`, `airport_iata`, `arr_delay`, `dep_delay`, `cancelled`, `diverted`, `year`, `month`
- `airline.airlines`
 - `airline_id`, `iata_code`, `name`

Outputs

Grain: **airline (BTS carriers)**

Columns (sample from output):

airline_name	iata_code	total_arrivals	delayed_arrivals	cancelled_arrivals
pct_delayed				
Frontier Airlines	F9	208,624	58,481	4,835
0.2803				
Air Wisconsin	ZW	52,393	11,859	764
0.2263				
American Airlines	AA	984,306	252,485	15,252
0.2565				
JetBlue Airways	B6	240,282	60,121	3,735
0.2502				
Allegiant Air	G4	117,210	24,897	2,018
0.2124				

BI Value

Enables performance scorecards and SLA reviews across airlines. Operations and commercial teams can quickly see which carriers are more reliable and which require attention or agreements around delay handling.

3) Monthly passenger counts (via bookings)

Purpose

Track demand trends and seasonality by aggregating passenger bookings per calendar month.

Inputs

- `airline.bookings`
 - `booking_id`, `passenger_id`, `booking_date`
- `airline.passengers`
 - `passenger_id`

Outputs

Grain: **month**

Columns (sample from output):

month_start	total_bookings	unique_passengers
2025-02-01	1,688	1,436
2025-03-01	3,403	2,472
2025-04-01	3,236	2,415
2025-05-01	3,422	2,504
2025-06-01	3,268	2,445

BI Value

Shows monthly demand patterns (growth, peaks, troughs). Supports forecasting of capacity, staffing, and revenue, and provides context for promotion performance.

4) Loyalty tier transitions (current vs miles-based target)

Purpose

Compare each member's current tier to the tier they qualify for based on miles, highlighting potential upgrades or downgrades.

Inputs

```
- `airline.loyalty_accounts`  
  - `loyalty_id`, `passenger_id`, `tier`, `miles_balance` / `ytd_miles`  
- `airline.miles_transactions`  
  - `loyalty_id`, `miles_delta`, `txn_date`, `txn_type`
```

****Outputs****

Grain: ***(current_tier, target_tier)*** summary counts

Columns (sample from output):

current_tier	target_tier	member_count
Basic	Basic	353
Basic	Gold	181
Basic	Platinum	73
Basic	Silver	138
Silver	Basic	350

****BI Value****

Identifies members whose current tier is “behind” their earned miles (good upgrade candidates) and potential downgrades. This is vital for loyalty program management and targeted retention campaigns.

5) Revenue per fare class (bookings + payments)

****Purpose****

Understand revenue mix across fare classes (e.g., Basic, Standard, Flexible, Business, First).

****Inputs****

```
- `airline.bookings`  
  - `booking_id`, `fare_class`  
- `airline.payments`  
  - `booking_id`, `amount_usd`, `status`, `paid_at`
```

****Outputs****

Grain: **fare_class**

Columns (sample from output):

fare_class	num_bookings	total_revenue	avg_revenue_per_booking
Basic	13,903	1,572,721.97	113.12
Standard	11,827	1,338,850.26	113.20
Flexible	8,211	936,208.77	114.02
Business	4,029	458,256.95	113.74
First	2,030	233,756.91	115.15

****BI Value****

Shows the relative revenue contribution of each fare product. Supports fare strategy, upsell tactics, and product design (e.g., whether to invest in Premium/Business cabins).

6) Ranking airlines by average delay

****Purpose****

Use a window function to rank airlines by mean delay.

****Inputs****

```
- `airline.flight_performance` / `airline.flights`  
  - `airline_id`, `arr_delay` / `delay_minutes`  
- `airline.airlines`
```

Outputs

Grain: **airline**

Columns (sample from output):

airline_name	iata_code	avg_delay_minutes	delay_rank
Red Jet Mexico	4X	287.00	1
Cargo Plus Aviation	8L	257.00	2
Sriwijaya Air	SJ	253.50	3
Armenian International Airways	MV	251.00	4
Malaysia Airlines	MH	226.33	5

BI Value

Quickly ranks carriers by punctuality, identifying worst offenders. Useful for operational negotiations, scheduling changes, and customer communications.

7) Running monthly revenue totals

Purpose

Build a time series of revenue with a running cumulative total using window `SUM()`.

Inputs

- `airline.payments`
 - `amount`, `status`, `paid_at` (filtered to successful statuses)

Outputs

Grain: **month**

Columns (sample from output):

month_start	revenue	running_cumulative_revenue
2025-02-01	185,699.32	185,699.32
2025-03-01	383,880.42	569,579.74
2025-04-01	369,920.05	939,499.79
2025-05-01	389,381.51	1,328,881.30
2025-06-01	372,051.23	1,700,932.53

BI Value

Supports revenue pacing dashboards and comparison to budget/forecast over time. Clearly shows growth trajectory and the effect of seasonal peaks.

8) Percent of flights delayed by month

Purpose

Measure the share of flights that are delayed each month.

Inputs

- `airline.flights` / `airline.flight_performance`
 - `flight_date`, `delay_minutes` (or arrival/departure delay fields)

Outputs

Grain: **month**

Columns (sample from output):

month_start	total_flights	delayed_flights	pct_delayed
2024-01-01	140	105	0.7500

2024-02-01	117	87	0.7436
2024-03-01	144	119	0.8264
2024-04-01	154	114	0.7403
2024-05-01	125	99	0.7920

BI Value

Shows monthly reliability performance and reveals seasonality (e.g., winter weather). Useful for root-cause analysis and tracking the impact of process changes.

9) Customer lifetime value (CLV) window function

Purpose

Compute cumulative revenue per passenger over time using a CLV-style window function.

Inputs

- `airline.bookings`
 - `booking_id`, `passenger_id`
- `airline.payments`
 - `booking_id`, `amount`, `paid_at`, `status`

Outputs

Grain: **payment event per passenger**, with cumulative CLV

Columns (sample from output for `passenger_id = 1`):

passenger_id	paid_date	amount_usd	clv_to_date
1	2025-03-10	90.98	90.98
1	2025-04-09	73.00	163.98
1	2025-05-04	121.78	285.76
1	2025-07-25	74.34	360.10
1	2025-08-29	168.50	528.60

BI Value

Provides a customer-level view of revenue over time, supporting segmentation into high-value vs. low-value customers and informing retention and marketing priorities.

10) Dense_rank route distance analysis (distance computed on the fly)

Purpose

Rank the longest routes using approximate distances derived from airport coordinates and a window `DENSE_RANK()`.

Inputs

- `airline.routes`
 - `route_id`, `origin_airport_id`, `destination_airport_id`
- `airline.airports`
 - `airport_id`, `iata_code`, `latitude`, `longitude`

Outputs

Grain: **route**

Columns (sample from output):

route_id	origin_iata	destination_iata	distance_nm	distance_rank
2781	NLK	TLA	20,839.17	1
2583	HOM	KTF	20,367.45	2
3884	UVE	MCG	19,970.83	3
4006	KTS	BHS	19,870.81	4
3220	KSM	FRE	19,824.93	5

****BI Value****

Highlights the longest segments in the network, which often drive distinct cost and product considerations (fuel, crew duty time, cabin product). Useful for fleet assignment and long-haul strategy.

11) Airport connectivity graph from busiest origin

****Purpose****

Use a recursive CTE to find all airports reachable from the busiest origin within up to 3 hops.

****Inputs****

- `airline.routes`
 - `origin_airport_id`, `destination_airport_id`
- `airline.airports`
 - `airport_id`, `iata_code`

****Outputs****

Grain: ****origin–destination–hop combination****

Columns (sample from output):

origin_iata	dest_iata	hops	path
YCK	EIK	1	[YCK, EIK]
YCK	NVT	1	[YCK, NVT]
YCK	NYR	1	[YCK, NYR]
YCK	PIP	1	[YCK, PIP]
YCK	RUM	1	[YCK, RUM]

****BI Value****

Shows the reach of a key hub and the set of airports that can be served directly or via one connection. Supports hub planning, connection design, and network optimization.

12) Multi-hop routes: detailed paths up to 3 hops from busiest origin

****Purpose****

List explicit multi-hop routes (up to 3 hops) from the busiest origin, showing full paths.

****Inputs****

- `airline.routes`
- `airline.airports`

****Outputs****

Grain: ****multi-hop path**** from origin to destination

Columns (sample from output):

origin_iata	dest_iata	hops	path
YCK	AHS	3	[YCK, NVT, YCW, AHS]
YCK	AKI	3	[YCK, NVT, YCW, AKI]
YCK	BTT	3	[YCK, RUM, FEN, BTT]
YCK	HEL	3	[YCK, RUM, TPP, HEL]
YCK	YJF	3	[YCK, TJB, FUK, YJF]

****BI Value****

Provides concrete connection options and reveals how complex some journeys are (e.g., 3-leg itineraries). Supports decisions on adding direct routes or retiming flights to improve connection quality.

13) Payment success rate by booking channel (Captured + Authorized as success)

Purpose

Evaluate payment performance by booking channel, treating `Captured` and `Authorized` as successful outcomes.

Inputs

- `airline.bookings`
 - `booking_id`, `booking_channel`
- `airline.payments`
 - `booking_id`, `status`

Outputs

Grain: **booking_channel**

Columns (sample from output):

booking_channel	total_payments	successful_payments	success_rate
Mobile	10,088	8,101	0.8030
Web	21,919	17,514	0.7990
Call Center	3,942	3,126	0.7930
Travel Agent	4,051	3,212	0.7929

BI Value

Highlights differences in conversion between channels. A lower success rate on a specific channel (e.g., Web) can indicate technical issues or UX friction that directly reduce revenue.

14) Worst routes by delay + cancellations (no volume cutoff)

Purpose

Identify the most problematic routes by combining average delay and cancellation rate, with no minimum volume filter.

Inputs

- `airline.flights`
 - `route_id`, `delay_minutes`, `status`
- `airline.routes`
- `airline.airports`

Outputs

Grain: **route**

Columns (sample from output):

route_id	origin_iata	destination_iata	total_flights	avg_delay_minutes	cancel_rate
3107	LHA	RIA	1	300.0	1.0
845	OCV	ZVK	1	300.0	1.0
2065	MYP	PAS	1	300.0	1.0
4085	CRQ	SAA	1	300.0	1.0
1449	BPY	GJT	1	299.0	1.0

****BI Value****

Provides a route-level “watch list” for operational remediation. Even with synthetic data, this pattern supports a dashboard tile that flags routes with extreme delay and cancellation metrics.

15) High-value loyalty members (top 5% by lifetime miles)

****Purpose****

Use window functions (e.g., `CUME_DIST()` / `PERCENT_RANK()`) to identify the top 5% of members by lifetime miles.

****Inputs****

- `airline.loyalty_accounts`
 - `loyalty_id`, `passenger_id`, `tier`, `miles_balance`
- `airline.miles_transactions`
 - `loyalty_id`, `miles_delta`, `txn_date`

****Outputs****

Grain: **loyalty account**

Columns (sample from output):

loyalty_id	passenger_id	tier	miles_balance	lifetime_miles	percentile_rank
1385	2298	Gold	40,763	218,556	1.0000
1536	2543	Basic	41,192	215,170	0.9997
649	1065	Silver	6,116	210,018	0.9993
1714	2842	Gold	58,618	202,778	0.9990
642	1047	Basic	22,748	197,384	0.9987

****BI Value****

Enables a focused VIP strategy: these members can be targeted for special offers, dedicated support, and retention programs, maximizing the value of the loyalty program.

**_End of Phase 4 Business Question Mapping with actual query outputs from
`03_analytics_queries.ipynb`._**

Airline BI Database — Phase 4: Analytical Queries

This notebook is used to:

- Develop and test analytical SQL queries against the airline BI database
- Profile query performance (EXPLAIN / EXPLAIN ANALYZE)
- Generate sample tables and visualizations for docs/phase_4_analytics.png

Database: PostgreSQL 16

Schema: airline_bi

```
In [93]: import os

import pandas as pd
from dotenv import load_dotenv

load_dotenv()

import matplotlib.pyplot as plt

pd.set_option("display.max_rows", 50)
pd.set_option("display.max_columns", 50)
pd.set_option("display.width", 120)
```

Database Connection Config

```
In [94]: db_url = os.getenv("DATABASE_URL")

if not db_url:
    raise ValueError("DATABASE_URL not found. Make sure your .env file is located in the same directory as this notebook")

db_url
```

```
Out[94]: 'postgresql+psycopg2://postgres:gpcool@localhost:5432/airline_bi'
```

Create the engine & test connection

```
In [95]: from sqlalchemy import create_engine, text

engine = create_engine(db_url, echo=False, future=True)

with engine.connect() as conn:
    version = conn.exec_driver_sql("SELECT version();").scalar_one()
version
```

```
Out[95]: 'PostgreSQL 17.5 on x86_64-apple-darwin23.6.0, compiled by Apple clang version 16.0.0 (clang-1600.0.26.6), 64-bit'
```

Helper: run_sql() for SELECT Queries

```
In [96]: from typing import Optional, Dict

def run_sql(
    query: str,
    params: Optional[Dict] = None,
    limit: Optional[int] = None,
    debug: bool = False
) -> pd.DataFrame:

    """
    Execute a SQL query and return the results as a pandas DataFrame.

    Args:
        query: SQL string. Can include named parameters (e.g., :airline_id)
        params: dict of parameters to bind
        limit: optional row limit applied in Python (not SQL)
        debug: if True, prints the rendered SQL and params

    Returns:
        pandas.DataFrame
    """
    if debug:
        print("SQL:")
        print(query)
        if params:
            print("Params:", params)

    with engine.connect() as conn:
        df = pd.read_sql(text(query), conn, params=params)

    if limit is not None:
        return df.head(limit)
    return df
```

Helper: run_explain() for Performance Testing

```
In [97]: from typing import Optional, Dict

def run_explain(
    query: str,
    params: Optional[Dict] = None,
    analyze: bool = False
) -> pd.DataFrame:

    """
    Run EXPLAIN or EXPLAIN ANALYZE on a SQL query and return the plan as a DataFrame.

    """
    prefix = "EXPLAIN ANALYZE " if analyze else "EXPLAIN "
    explain_sql = prefix + query
```

```

with engine.connect() as conn:
    result = conn.exec_driver_sql(explain_sql, params or {})
    rows = result.fetchall()

plans = [row[0] for row in rows]
return pd.DataFrame({"query_plan": plans})

```

Simple Display Helper for Charts

```

In [98]: def plot_bar_from_df(
    df: pd.DataFrame,
    x: str,
    y: str,
    title: str = "",
    rotation: int = 45
) -> None:
    """
    Simple helper to create a quick bar chart from a DataFrame.
    Used mainly to generate pngs for docs/phase_4_analytics.png.
    """
    plt.figure(figsize=(10, 5))
    plt.bar(df[x], df[y])
    plt.title(title)
    plt.xlabel(x)
    plt.ylabel(y)
    plt.xticks(rotation=rotation, ha="right")
    plt.tight_layout()
    plt.show()

```

Sanity Test Query

```

In [99]: test_query = """
SELECT
    a.airline_id,
    a.name AS airline_name,
    a.iata_code,
    a.icao_code,
    a.country
FROM airline.airlines AS a
ORDER BY a.airline_id
LIMIT 5;
"""

df_test = run_sql(test_query)
df_test

```

Out[99]:	airline_id	airline_name	iata_code	icao_code	country
0	1223	Unknown	None	None	None
1	1226	1Time Airline	1T	RNX	SOU
2	1233	40-Mile Air	Q5	MLA	UNI
3	1236	Ansett Australia	AN	AAA	AUS
4	1237	Abacus International	1B	None	SIN

A. CTE Queries

```
In [100...]: # 1) Top 10 busiest airports (arrivals + departures)
q_cte_busiest_airports = """
/* CTE: Top 10 busiest airports by total movements (departures + arrivals) */

WITH airport_movements AS (
    SELECT
        f.origin_airport_id AS airport_id,
        COUNT(*) AS departures,
        0 AS arrivals
    FROM airline.flights AS f
    GROUP BY f.origin_airport_id

    UNION ALL

    SELECT
        f.destination_airport_id AS airport_id,
        0 AS departures,
        COUNT(*) AS arrivals
    FROM airline.flights AS f
    GROUP BY f.destination_airport_id
),
aggregated AS (
    SELECT
        airport_id,
        SUM(departures) AS total_departures,
        SUM(arrivals) AS total_arrivals,
        SUM(departures + arrivals) AS total_movements
    FROM airport_movements
    GROUP BY airport_id
)
SELECT
    a.airport_id,
    ap.name      AS airport_name,
    ap.iata_code AS airport_iata,
    total_departures,
    total_arrivals,
    total_movements
FROM aggregated a
JOIN airline.airports ap
    ON ap.airport_id = a.airport_id
ORDER BY total_movements DESC
```

```
LIMIT 10;
"""

df_cte_busiest_airports = run_sql(q_cte_busiest_airports)
df_cte_busiest_airports.head()
```

Out[100...]

	airport_id	airport_name	airport_iata	total_departures	total_arrivals	total_movement
0	3538	Colville Lake Airport	YCK	6.0	3.0	
1	2109	Iberia Airport	IBP	5.0	3.0	
2	4432	Phoenix-Mesa-Gateway Airport	AZA	3.0	5.0	
3	4713	Golovin Airport	GLV	1.0	6.0	
4	965	Pamplona Airport	PNA	4.0	3.0	

In [101...]

```
# 2) Airline on-time performance summary (using BTS flight_performance)

q_cte_airline_on_time = """
/* CTE: Airline-level performance summary from BTS snapshot */

WITH perf AS (
    SELECT
        fp.airline_iata,
        SUM(fp.arrivals) AS total_arrivals,
        SUM(fp.arrivals_delayed_15min) AS delayed_arrivals,
        SUM(fp.arr_cancelled) AS cancelled_arrivals,
        SUM(fp.total_arrival_delay_min) AS total_delay_min
    FROM airline.flight_performance AS fp
    GROUP BY fp.airline_iata
)
SELECT
    al.airline_id,
    al.name      AS airline_name,
    al.iata_code,
    total_arrivals,
    delayed_arrivals,
    cancelled_arrivals,
    CASE
        WHEN total_arrivals > 0
            THEN delayed_arrivals::decimal / total_arrivals
        ELSE NULL
    END AS pct_delayed,
    CASE
        WHEN total_arrivals > 0
            THEN cancelled_arrivals::decimal / total_arrivals
        ELSE NULL
    END AS pct_cancelled
```

```

        END AS pct_cancelled,
        CASE
            WHEN total_arrivals > 0
                THEN total_delay_min / total_arrivals
            ELSE NULL
        END AS avg_delay_minutes
    FROM perf
    LEFT JOIN airline.airlines al
        ON al.iata_code = perf.airline_iata
    ORDER BY avg_delay_minutes DESC NULLS LAST;
"""

df_cte_airline_on_time = run_sql(q_cte_airline_on_time)
df_cte_airline_on_time.head()

```

Out[101...]

	airline_id	airline_name	iata_code	total_arrivals	delayed_arrivals	cancelled_arrivals
0	3690	Frontier Airlines	F9	208624	58481	481
1	1505	Air Wisconsin	ZW	52393	11859	70
2	1247	American Airlines	AA	984306	252485	1521
3	4250	JetBlue Airways	B6	240282	60121	371
4	1258	Allegiant Air	G4	117210	24897	20

In [102...]

```

# 3) Monthly passenger counts (via bookings)

q_cte_monthly_passengers = """
/* CTE: Monthly bookings and unique passenger counts */

WITH monthly_stats AS (
    SELECT
        date_trunc('month', b.booking_date)::date AS month_start,
        COUNT(*) AS total_bookings,
        COUNT(DISTINCT b.passenger_id) AS unique_passengers
    FROM airline.bookings AS b
    GROUP BY date_trunc('month', b.booking_date)
)
SELECT
    month_start,
    total_bookings,
    unique_passengers
FROM monthly_stats
ORDER BY month_start;
"""

df_cte_monthly_passengers = run_sql(q_cte_monthly_passengers)
df_cte_monthly_passengers.head()

```

Out[102...]

	month_start	total_bookings	unique_passengers
0	2025-02-01	1688	1436
1	2025-03-01	3403	2472
2	2025-04-01	3236	2415
3	2025-05-01	3422	2504
4	2025-06-01	3268	2445

In [103...]

```
# 4) Loyalty tier transitions (current vs miles-based target)

q_cte_loyalty_transitions = """
/* CTE: Compare current loyalty tier vs miles-based target tier */

WITH miles_totals AS (
    SELECT
        la.loyalty_id,
        la.passenger_id,
        la.tier           AS current_tier,
        la.miles_balance,
        COALESCE(SUM(mt.miles_delta), 0) AS lifetime_miles
    FROM airline.loyalty_accounts AS la
    LEFT JOIN airline.miles_transactions AS mt
        ON mt.loyalty_id = la.loyalty_id
    GROUP BY la.loyalty_id, la.passenger_id, la.tier, la.miles_balance
),
tier_buckets AS (
    SELECT
        *,
        CASE
            WHEN lifetime_miles < 25000 THEN 'Basic'
            WHEN lifetime_miles < 50000 THEN 'Silver'
            WHEN lifetime_miles < 100000 THEN 'Gold'
            ELSE 'Platinum'
        END AS target_tier
    FROM miles_totals
)
SELECT
    current_tier,
    target_tier,
    COUNT(*) AS member_count
FROM tier_buckets
GROUP BY current_tier, target_tier
ORDER BY current_tier, target_tier;
"""

df_cte_loyalty_transitions = run_sql(q_cte_loyalty_transitions)
df_cte_loyalty_transitions.head()
```

Out[103...]

	current_tier	target_tier	member_count
0	Basic	Basic	353
1	Basic	Gold	181
2	Basic	Platinum	73
3	Basic	Silver	138
4	Silver	Basic	350

In [104...]

```
# 5) Revenue per fare class (bookings + payments)

q_cte_revenue_fare_class = """
/* CTE: Revenue by fare_class based on payments */

WITH revenue_by_fare AS (
    SELECT
        b.fare_class,
        COUNT(DISTINCT b.booking_id) AS num_bookings,
        SUM(p.amount_usd) AS total_revenue
    FROM airline.bookings AS b
    JOIN airline.payments AS p
        ON p.booking_id = b.booking_id
    GROUP BY b.fare_class
)
SELECT
    fare_class,
    num_bookings,
    total_revenue,
    CASE
        WHEN num_bookings > 0
            THEN total_revenue / num_bookings
        ELSE NULL
    END AS avg_revenue_per_booking
FROM revenue_by_fare
ORDER BY total_revenue DESC NULLS LAST;
"""

df_cte_revenue_fare_class = run_sql(q_cte_revenue_fare_class)
df_cte_revenue_fare_class.head()
```

Out[104...]

	fare_class	num_bookings	total_revenue	avg_revenue_per_booking
0	Basic	13903	1572721.97	113.121051
1	Standard	11827	1338850.26	113.202863
2	Flexible	8211	936208.77	114.018849
3	Business	4029	458256.95	113.739625
4	First	2030	233756.91	115.151187

B. Window Function Queries

In [105... # 6) Ranking airlines by average delay

```
q_win_airline_delay_rank = """
/* Window: Rank airlines by average delay minutes */

SELECT
    al.airline_id,
    al.name      AS airline_name,
    al.iata_code,
    AVG(f.delay_minutes) AS avg_delay_minutes,
    RANK() OVER (ORDER BY AVG(f.delay_minutes) DESC) AS delay_rank
FROM airline.flights AS f
JOIN airline.airlines AS al
    ON al.airline_id = f.airline_id
GROUP BY al.airline_id, al.name, al.iata_code
ORDER BY delay_rank;
"""

df_win_airline_delay_rank = run_sql(q_win_airline_delay_rank)
df_win_airline_delay_rank.head()
```

Out [105...]

	airline_id	airline_name	iata_code	avg_delay_minutes	delay_rank
0	7049	Red Jet Mexico	4X	287.000000	1
1	4163	Cargo Plus Aviation	8L	257.000000	2
2	5669	Sriwijaya Air	SJ	253.500000	3
3	2432	Armenian International Airways	MV	251.000000	4
4	4597	Malaysia Airlines	MH	226.333333	5

In [106... # 7) Running monthly revenue totals

```
q_win_running_monthly_revenue = """
/* Window: Running cumulative monthly revenue */

WITH monthly_revenue AS (
    SELECT
        date_trunc('month', p.paid_at)::date AS month_start,
        SUM(p.amount_usd) AS revenue
    FROM airline.payments AS p
    GROUP BY date_trunc('month', p.paid_at)
)
SELECT
    month_start,
    revenue,
    SUM(revenue) OVER (
        ORDER BY month_start
        ROWS BETWEEN UNBOUNDED PRECEDING AND CURRENT ROW
    ) AS running_cumulative_revenue
FROM monthly_revenue
ORDER BY month_start;
```

```
"""
df_win_running_monthly_revenue = run_sql(q_win_running_monthly_revenue)
df_win_running_monthly_revenue.head()
```

Out[106...]

	month_start	revenue	running_cumulative_revenue
0	2025-02-01	185699.32	185699.32
1	2025-03-01	383880.42	569579.74
2	2025-04-01	369920.05	939499.79
3	2025-05-01	389381.51	1328881.30
4	2025-06-01	372051.23	1700932.53

In [107...]

```
# 8) Percent of flights delayed by month

q_win_pct_delayed_by_month = """
/* Monthly delay rate based on delay_minutes > 15 */

WITH monthly AS (
    SELECT
        date_trunc('month', f.flight_date)::date AS month_start,
        COUNT(*) AS total_flights,
        SUM(CASE WHEN f.delay_minutes > 15 THEN 1 ELSE 0 END) AS delayed_fli
    FROM airline.flights AS f
    GROUP BY date_trunc('month', f.flight_date)
)
SELECT
    month_start,
    total_flights,
    delayed_flights,
    (delayed_flights::decimal / NULLIF(total_flights, 0)) AS pct_delayed
FROM monthly
ORDER BY month_start;
"""

df_win_pct_delayed_by_month = run_sql(q_win_pct_delayed_by_month)
df_win_pct_delayed_by_month.head()
```

Out[107...]

	month_start	total_flights	delayed_flights	pct_delayed
0	2024-01-01	140	105	0.750000
1	2024-02-01	117	87	0.743590
2	2024-03-01	144	119	0.826389
3	2024-04-01	154	114	0.740260
4	2024-05-01	125	99	0.792000

In [108...]

9) Customer lifetime value (CLV) window function

```
q_win_clv_running = """
```

```

/* Window: CLV per passenger (running sum of revenue over time) */

WITH customer_payments AS (
    SELECT
        b.passenger_id,
        p.paid_at::date AS paid_date,
        p.amount_usd
    FROM airline.bookings AS b
    JOIN airline.payments AS p
        ON p.booking_id = b.booking_id
),
running_clv AS (
    SELECT
        passenger_id,
        paid_date,
        amount_usd,
        SUM(amount_usd) OVER (
            PARTITION BY passenger_id
            ORDER BY paid_date
            ROWS BETWEEN UNBOUNDED PRECEDING AND CURRENT ROW
        ) AS clv_to_date
    FROM customer_payments
)
SELECT
    passenger_id,
    paid_date,
    amount_usd,
    clv_to_date
FROM running_clv
ORDER BY passenger_id, paid_date;
"""

df_win_clv_running = run_sql(q_win_clv_running)
df_win_clv_running.head()

```

Out[108...]

	passenger_id	paid_date	amount_usd	clv_to_date
0	1	2025-03-10	90.98	90.98
1	1	2025-04-09	73.00	163.98
2	1	2025-05-04	121.78	285.76
3	1	2025-07-25	74.34	360.10
4	1	2025-08-29	168.50	528.60

In [135...]

```

q_update_distances_simple = """
/* Simple approximate distance between origin & destination.
   Uses Euclidean distance on degrees * 60 to get nautical miles.
   Overwrites distance_nm for all routes.
*/

WITH updated AS (
    UPDATE airline.routes r
    SET distance_nm = sub.distance_nm::integer

```

```

    FROM (
        SELECT
            r2.route_id,
            (
                sqrt(
                    (ad.latitude - ao.latitude)^2 +
                    (ad.longitude - ao.longitude)^2
                ) * 60
            ) AS distance_nm
        FROM airline.routes r2
        JOIN airline.airports ao
            ON ao.airport_id = r2.origin_airport_id
        JOIN airline.airports ad
            ON ad.airport_id = r2.destination_airport_id
        WHERE ao.latitude IS NOT NULL
            AND ao.longitude IS NOT NULL
            AND ad.latitude IS NOT NULL
            AND ad.longitude IS NOT NULL
    ) sub
    WHERE r.route_id = sub.route_id
    RETURNING r.route_id
)
SELECT COUNT(*) AS updated_routes
FROM updated;
"""

run_sql(q_update_distances_simple)

```

Out[135... **updated_routes**

0	5000
---	------

In [136... **run_sql("""**

```

SELECT
    COUNT(*) AS total_routes,
    COUNT(distance_nm) AS routes_with_distance,
    MIN(distance_nm) AS min_distance,
    MAX(distance_nm) AS max_distance
FROM airline.routes;
"""
)
```

Out[136... **total_routes routes_with_distance min_distance max_distance**

0	5000	0	None	None
---	------	---	------	------

In [137... # 10) Dense_rank route distance analysis (distance computed on the fly)

```

q_wm_route_distance_rank = """
/* Window: Rank the longest routes by approximate distance (nautical miles),
computing distance directly from airport latitude/longitude.

Approximation:
    distance_nm ≈ sqrt( (Δlat)^2 + (Δlon)^2 ) * 60
    (about 60 NM per degree of lat/lng difference; good enough for BI demo)

```

```
*/  
  
WITH route_dist AS (  
    SELECT  
        r.route_id,  
        ao.iata_code AS origin_iata,  
        ad.iata_code AS destination_iata,  
        sqrt(  
            (ad.latitude - ao.latitude)^2 +  
            (ad.longitude - ao.longitude)^2  
        ) * 60 AS distance_nm  
    FROM airline.routes r  
    JOIN airline.airports ao  
        ON ao.airport_id = r.origin_airport_id  
    JOIN airline.airports ad  
        ON ad.airport_id = r.destination_airport_id  
    WHERE ao.latitude IS NOT NULL  
        AND ao.longitude IS NOT NULL  
        AND ad.latitude IS NOT NULL  
        AND ad.longitude IS NOT NULL  
)  
SELECT  
    route_id,  
    origin_iata,  
    destination_iata,  
    distance_nm,  
    DENSE_RANK() OVER (ORDER BY distance_nm DESC) AS distance_rank  
FROM route_dist  
ORDER BY distance_rank, origin_iata, destination_iata  
LIMIT 50;  
*****  
  
df_win_route_distance_rank = run_sql(q_win_route_distance_rank)  
df_win_route_distance_rank
```

Out[137...]

	route_id	origin_iata	destination_iata	distance_nm	distance_rank
0	2781	NLK	TLA	20839.173604	1
1	2583	HOM	KTF	20367.446093	2
2	3884	UVE	MCG	19970.825036	3
3	4006	KTS	BHS	19870.810602	4
4	3220	KSM	FRE	19824.929476	5
5	1138	EFG	WAA	19469.293673	6
6	333	HCR	OKY	19448.385239	7
7	868	KVC	TUM	19403.575596	8
8	206	PTH	HVB	19331.737284	9
9	4589	AIN	JHQ	19311.471925	10
10	4153	UPP	UJE	19309.304878	11
11	4597	JHQ	UNK	19247.554596	12
12	4828	TGJ	EAA	19240.015640	13
13	1199	FTI	NGK	19167.338215	14
14	3868	XTG	KWN	19080.015204	15
15	247	ADL	GLV	19051.726787	16
16	3203	YEV	TBF	18960.175350	17
17	2633	KGE	SXQ	18958.800714	18
18	1069	KWN	CMU	18827.286781	19
19	1026	OKL	GLV	18688.665048	20
20	734	FTI	FKJ	18586.198828	21
21	1954	LKB	OLZ	18566.170674	22
22	2909	KYI	WNA	18475.170334	23
23	2374	SIO	WKL	18420.971626	24
24	910	GRF	ZQN	18331.920455	25
25	148	ALW	HLZ	18323.321665	26
26	4307	EEK	OKD	18234.388433	27
27	792	BEZ	YWS	18201.548309	28
28	4914	FRE	YUB	18162.178520	29
29	3378	YCT	TUO	18091.789521	30
30	4196	GFN	HNH	18091.396557	31
31	1483	OTK	AUY	18050.365501	32

route_id	origin_iata	destination_iata		distance_nm	distance_rank
32	1770	HID	PAQ	18011.182697	33
33	2431	HHI	MKQ	17997.739779	34
34	4544	VMU	CIK	17875.956294	35
35	142	MWF	SHN	17868.369758	36
36	136	ADK	CHG	17835.653135	37
37	2852	KKA	KHV	17807.765924	38
38	3966	JXA	AIU	17794.871713	39
39	2854	AST	KIO	17747.171271	40
40	2305	JAC	DUD	17691.519048	41
41	4480	VEL	AKL	17679.749143	42
42	1730	TPH	NON	17644.353871	43
43	4626	KPV	KGI	17622.693285	44
44	2774	WMB	PPT	17567.376908	45
45	3314	MMJ	PPT	17550.322358	46
46	2004	TKX	GVN	17521.233564	47
47	1243	GYL	SOV	17418.081465	48
48	895	OSN	EEK	17399.832875	49
49	3522	KFE	HSL	17369.051859	50

C. Recursive Queries

In [118...]

```
# 11) Airport connectivity graph from busiest origin

q_rec_connectivity = """
/* Recursive: All reachable airports from the busiest origin airport (by route)
   within up to 3 hops.
*/
WITH RECURSIVE
busiest_origin AS (
  SELECT
    r.origin_airport_id
  FROM airline.routes r
  GROUP BY r.origin_airport_id
  ORDER BY COUNT(*) DESC
  LIMIT 1
),
start_airport AS (
  SELECT
    ao.airport_id,
```

```

        ao.iata_code
    FROM airline.airports ao
    JOIN busiest_origin bo ON bo.origin_airport_id = ao.airport_id
),
connectivity (
    origin_airport_id,
    origin_iata,
    dest_airport_id,
    dest_iata,
    path,
    hops
) AS (
    -- Base from busiest origin
    SELECT
        sa.airport_id AS origin_airport_id,
        sa.iata_code AS origin_iata,
        ad.airport_id AS dest_airport_id,
        ad.iata_code AS dest_iata,
        ARRAY[sa.iata_code::text, ad.iata_code::text][] AS path,
        1 AS hops
    FROM airline.routes r
    JOIN start_airport sa
        ON sa.airport_id = r.origin_airport_id
    JOIN airline.airports ad
        ON ad.airport_id = r.destination_airport_id

    UNION ALL

    -- Extend outward
    SELECT
        c.origin_airport_id,
        c.origin_iata,
        ad.airport_id AS dest_airport_id,
        ad.iata_code AS dest_iata,
        c.path || ad.iata_code::text,
        c.hops + 1
    FROM connectivity c
    JOIN airline.routes r
        ON r.origin_airport_id = c.dest_airport_id
    JOIN airline.airports ad
        ON ad.airport_id = r.destination_airport_id
    WHERE c.hops < 3
        AND NOT ad.iata_code = ANY (c.path) -- avoid cycles
)
SELECT DISTINCT
    origin_iata,
    dest_iata,
    hops,
    path
FROM connectivity
ORDER BY hops, dest_iata
LIMIT 200;
"""

df_rec_connectivity = run_sql(q_rec_connectivity)
df_rec_connectivity.head()

```

Out[118...]

	origin_iata	dest_iata	hops	path
0	YCK	EIK	1	[YCK, EIK]
1	YCK	NVT	1	[YCK, NVT]
2	YCK	NYR	1	[YCK, NYR]
3	YCK	PIP	1	[YCK, PIP]
4	YCK	RUM	1	[YCK, RUM]

In [132...]

```
# 12) Multi-hop routes: detailed paths up to 3 hops from busiest origin

q_rec_multihop_paths = """
/* Recursive: Explore all paths from the busiest origin airport (by route count)
   up to 3 hops, and list the paths.
*/

WITH RECURSIVE
busiest_origin AS (
    SELECT
        r.origin_airport_id
    FROM airline.routes r
    GROUP BY r.origin_airport_id
    ORDER BY COUNT(*) DESC
    LIMIT 1
),
start_airport AS (
    SELECT
        ao.airport_id,
        ao.iata_code
    FROM airline.airports ao
    JOIN busiest_origin bo ON bo.origin_airport_id = ao.airport_id
),
connectivity (
    origin_airport_id,
    origin_iata,
    dest_airport_id,
    dest_iata,
    path,
    hops
) AS (
    -- Base from busiest origin
    SELECT
        sa.airport_id AS origin_airport_id,
        sa.iata_code AS origin_iata,
        ad.airport_id AS dest_airport_id,
        ad.iata_code AS dest_iata,
        ARRAY[sa.iata_code::text, ad.iata_code::text]::text[] AS path,
        1 AS hops
    FROM airline.routes r
    JOIN start_airport sa
        ON sa.airport_id = r.origin_airport_id
    JOIN airline.airports ad
        ON ad.airport_id = r.destination_airport_id
)
```

```
UNION ALL

-- Extend outward
SELECT
    c.origin_airport_id,
    c.origin_iata,
    ad.airport_id AS dest_airport_id,
    ad.iata_code AS dest_iata,
    c.path || ad.iata_code::text,
    c.hops + 1
FROM connectivity c
JOIN airline.routes r
    ON r.origin_airport_id = c.dest_airport_id
JOIN airline.airports ad
    ON ad.airport_id = r.destination_airport_id
WHERE c.hops < 3
    AND NOT ad.iata_code = ANY (c.path)
)

SELECT
    origin_iata,
    dest_iata,
    hops,
    path
FROM connectivity
ORDER BY hops DESC, dest_iata
LIMIT 50;
=====

df_rec_multihop_paths = run_sql(q_rec_multihop_paths)
df_rec_multihop_paths
```

Out[132...]

	origin_iata	dest_iata	hops	path
0	YCK	AHS	3	[YCK, NVT, YCW, AHS]
1	YCK	AKI	3	[YCK, NVT, YCW, AKI]
2	YCK	BTT	3	[YCK, RUM, FEN, BTT]
3	YCK	HEL	3	[YCK, RUM, TPP, HEL]
4	YCK	YJF	3	[YCK, TJB, FUK, YJF]
5	YCK	FEN	2	[YCK, RUM, FEN]
6	YCK	FUK	2	[YCK, TJB, FUK]
7	YCK	GGG	2	[YCK, PIP, GGG]
8	YCK	TPP	2	[YCK, RUM, TPP]
9	YCK	YCW	2	[YCK, NVT, YCW]
10	YCK	EIK	1	[YCK, EIK]
11	YCK	NVT	1	[YCK, NVT]
12	YCK	NYR	1	[YCK, NYR]
13	YCK	PIP	1	[YCK, PIP]
14	YCK	RUM	1	[YCK, RUM]
15	YCK	TJB	1	[YCK, TJB]

D. Complex Joins / Aggregations

In [128...]

```
# 13) Payment success rate by booking channel (using Captured + Authorized as success)

q_complex_payment_success = """
/* Complex join: Payment success rate by booking_channel
   Success statuses: Captured, Authorized
*/
WITH payment_stats AS (
    SELECT
        b.booking_channel,
        COUNT(*) AS total_payments,
        SUM(
            CASE
                WHEN LOWER(p.status::text) IN ('captured', 'authorized')
                THEN 1
                ELSE 0
            END
        ) AS successful_payments
    FROM airline.bookings AS b
    JOIN airline.payments AS p
        ON p.booking_id = b.booking_id
    GROUP BY b.booking_channel
)
```

```
)
SELECT
    booking_channel,
    total_payments,
    successful_payments,
    (successful_payments::decimal / NULLIF(total_payments, 0)) AS success_rate
FROM payment_stats
ORDER BY success_rate DESC NULLS LAST;
"""

df_complex_payment_success = run_sql(q_complex_payment_success)
df_complex_payment_success
```

Out[128...]

	booking_channel	total_payments	successful_payments	success_rate
0	Mobile	10088	8101	0.803033
1	Web	21919	17514	0.799033
2	Call Center	3942	3126	0.792998
3	Travel Agent	4051	3212	0.792891

In [120...]

```
# 14) Worst routes by delay + cancellations (no volume cutoff)

q_complex_worst_routes = """
/* Complex join: Worst-performing routes by average delay and cancel rate. */

WITH route_metrics AS (
    SELECT
        f.route_id,
        COUNT(*) AS total_flights,
        AVG(f.delay_minutes) AS avg_delay_minutes,
        SUM(CASE WHEN f.status = 'Cancelled' THEN 1 ELSE 0 END)::decimal
            / NULLIF(COUNT(*), 0) AS cancel_rate
    FROM airline.flights AS f
    WHERE f.route_id IS NOT NULL
    GROUP BY f.route_id
)
SELECT
    rm.route_id,
    ao.iata_code AS origin_iata,
    ad.iata_code AS destination_iata,
    rm.total_flights,
    rm.avg_delay_minutes,
    rm.cancel_rate
FROM route_metrics AS rm
JOIN airline.routes AS r
    ON r.route_id = rm.route_id
JOIN airline.airports AS ao
    ON ao.airport_id = r.origin_airport_id
JOIN airline.airports AS ad
    ON ad.airport_id = r.destination_airport_id
ORDER BY rm.avg_delay_minutes DESC NULLS LAST, rm.cancel_rate DESC NULLS LAST
LIMIT 25;
"""
```

```
df_complex_worst_routes = run_sql(q_complex_worst_routes)
df_complex_worst_routes
```

Out[120...]

	route_id	origin_iata	destination_iata	total_flights	avg_delay_minutes	cancel_ra
0	3107	LHA	RIA	1	300.0	1
1	845	OCV	ZVK	1	300.0	1
2	2065	MYP	PAS	1	300.0	1
3	4085	CRQ	SAA	1	300.0	1
4	1449	BPY	GJT	1	299.0	1
5	4371	SAH	NQY	1	299.0	1
6	4701	MED	RTB	1	299.0	1
7	4415	KFP	SAK	1	299.0	1
8	4774	UTH	DAN	1	299.0	1
9	1122	AFA	CFC	1	299.0	1
10	1452	SCM	ODE	1	298.0	1
11	92	REL	TME	1	298.0	1
12	1548	ANC	GVR	1	298.0	1
13	2061	RHT	IRJ	1	298.0	1
14	3725	MQQ	CTC	1	298.0	1
15	323	TTN	ADY	1	297.0	1
16	2512	APG	GRX	1	297.0	1
17	4287	SSN	MUX	1	297.0	1
18	4902	TIM	ANG	1	297.0	1
19	3017	KHM	NOU	1	297.0	1
20	3045	ODE	LIM	1	297.0	1
21	2456	FYJ	WMR	1	296.0	1
22	1774	MZH	WHK	1	296.0	1
23	2653	PKO	BTK	1	296.0	1
24	3726	VNE	KZS	1	296.0	1

In [114...]

```
# 15) High-value loyalty members (top 5% by lifetime miles)
```

```
q_complex_top_loyalty = """
/* Complex join + window: Top 5% loyalty members by lifetime miles */
```

```

WITH miles_by_member AS (
    SELECT
        la.loyalty_id,
        la.passenger_id,
        la.tier,
        la.miles_balance,
        COALESCE(SUM(mt.miles_delta), 0) AS lifetime_miles
    FROM airline.loyalty_accounts AS la
    LEFT JOIN airline.miles_transactions AS mt
        ON mt.loyalty_id = la.loyalty_id
    GROUP BY la.loyalty_id, la.passenger_id, la.tier, la.miles_balance
),
with_percentiles AS (
    SELECT
        *,
        PERCENT_RANK() OVER (ORDER BY lifetime_miles) AS pr
    FROM miles_by_member
)
SELECT
    loyalty_id,
    passenger_id,
    tier,
    miles_balance,
    lifetime_miles,
    pr AS percentile_rank
FROM with_percentiles
WHERE pr >= 0.95
ORDER BY lifetime_miles DESC;
"""

df_complex_top_loyalty = run_sql(q_complex_top_loyalty)
df_complex_top_loyalty.head()

```

Out[114...]

	loyalty_id	passenger_id	tier	miles_balance	lifetime_miles	percentile_rank
0	1385	2298	Gold	40763	218556	1.000000
1	1536	2543	Basic	41192	215170	0.999667
2	649	1065	Silver	6116	210018	0.999333
3	1714	2842	Gold	58618	202778	0.999000
4	642	1047	Basic	22748	197384	0.998666

Performance Testing (EXPLAIN / EXPLAIN ANALYZE)

In [138...]

```

# Simple helpers that wrap EXPLAIN / EXPLAIN ANALYZE around an existing SQL
# They reuse run_sql(), so the plan comes back as a DataFrame.

def explain(query: str):
    """
    Run EXPLAIN on a SQL query string and return the plan as a DataFrame.
    """

```

```
return run_sql("EXPLAIN " + query)

def explain_analyze(query: str):
    """
    Run EXPLAIN ANALYZE on a SQL query string and return the plan as a DataFrame
    """
    return run_sql("EXPLAIN ANALYZE " + query)
```

Q1 - Top 10 Busiest Airports (CTE)

Q1 performs a sequential scan over the 5k-row airline.flights table and joins once to airline.airports. The planner uses a HashAggregate to compute total departures and arrivals per airport, followed by a sort and limit. Because the table is small, sequential scans are optimal. In a production environment with millions of flights per year, an index on (origin_airport_id, destination_airport_id) would improve performance.

```
In [ ]: # Q1 Performance: CTE busiest airports
# Underlying query variable: q_cte_busiest_airports

plan_q1 = explain_analyze(q_cte_busiest_airports)
plan_q1
```

Out[]:

QUERY PLAN

```

0 Limit (cost=1026.11..1026.14 rows=10 width=13...)
1   -> Sort (cost=1026.11..1026.61 rows=200 wi...
2     Sort Key: (sum((/*SELECT* 1".departure...
3     Sort Method: top-N heapsort Memory: 26kB
4       -> Hash Join (cost=694.61..1021.79 r...
5         Hash Cond: (ap.airport_id = /*SE...
6           -> Seq Scan on airports ap (co...
7           -> Hash (cost=692.11..692.11 r...
8             Buckets: 8192 (originally ...
9             -> HashAggregate (cost=6...
10            Group Key: /*SELECT*...
11            Batches: 1 Memory U...
12            -> Append (cost=0....
13              -> Subquery S...
14                -> Grou...
15                  Gr...
16                  ->...
17                  ...
18                  -> Subquery S...
19                  -> Grou...
20                  Gr...
21                  ->...
22                  ...
23          Planning Time: 35.016 ms
24          Execution Time: 197.210 ms

```

Q5 — Revenue per fare class (complex join, aggregation)

Q5 joins airline.bookings (40k rows) with airline.payments (40k rows) using a Hash Join on booking_id, then aggregates total revenue by fare_class. The plan uses a hash strategy for the join and for the aggregate, which is optimal for this dataset size. Execution time is mostly from EXPLAIN ANALYZE overhead in Jupyter, not from the query itself. Indexes on both booking_id columns ensure efficient lookups.

In [146...]

```
# Q5 Performance: Revenue per fare class (bookings + payments)
# Underlying query variable: q_cte_revenue_fare_class
```

```
plan_q5 = explain_analyze(q_cte_revenue_fare_class)
plan_q5
```

Out [146...]

QUERY PLAN

```

0  Sort (cost=6444.71..6444.72 rows=5 width=79) ...
1  Sort Key: revenue_by_fare.total_revenue DESC...
2      Sort Method: quicksort Memory: 25kB
3      -> Subquery Scan on revenue_by_fare (cost=...
4          -> GroupAggregate (cost=6044.55..644...
5              Group Key: b.fare_class
6              -> Sort (cost=6044.55..6144.55...
7                  Sort Key: b.fare_class, b....
8                  Sort Method: quicksort Me...
9                  -> Hash Join (cost=2148....
10                 Hash Cond: (p.bookin...
11                 -> Seq Scan on paym...
12                 -> Hash (cost=1648...
13                     Buckets: 65536...
14                     -> Seq Scan o...
15             Planning Time: 3.383 ms
16             Execution Time: 4043.602 ms
```

Q7 — Running monthly revenue totals (window function)

Q7 aggregates revenue into monthly buckets before applying a running SUM() window function. The planner sorts on month_start to feed data into the WindowAgg node. With fewer than 50 months of data, this is extremely efficient. For longer histories (multi-year), materializing monthly revenue in a summary table would accelerate dashboards.

```
In [ ]: # Q7 Performance: Running monthly revenue totals
# Underlying query variable: q_win_running_monthly_revenue

plan_q7 = explain_analyze(q_win_running_monthly_revenue)
plan_q7
```

Out[]:

QUERY PLAN

```
0 WindowAgg (cost=8435.31..9135.29 rows=40000 w...
1   -> Sort (cost=8435.29..8535.29 rows=40000 ...
2     Sort Key: monthly_revenue.month_start
3       Sort Method: quicksort Memory: 25kB
4       -> Subquery Scan on monthly_revenue ...
5         -> HashAggregate (cost=3709.00...
6           Group Key: date_trunc('mon...
7             Planned Partitions: 4 Bat...
8             -> Seq Scan on payments p...
9               Planning Time: 30.652 ms
10              Execution Time: 103.972 ms
```

Q11 — Airport connectivity graph (recursive CTE)

Q11 uses a Recursive Union to explore airport connectivity up to three hops from the busiest origin. The planner performs index scans on route origin and destination, keeping recursion fast. Execution time remains low because depth is capped to three levels. For large airline-route networks, a materialized connectivity graph in Phase 5 would significantly reduce recursive computation.

```
In [ ]: # Q11 Performance: Recursive connectivity from busiest origin
# Underlying query variable: q_rec_connectivity

plan_q11 = explain_analyze(q_rec_connectivity)
plan_q11
```

Out[]:

QUERY PLAN

```
0  Limit (cost=324.56..325.07 rows=41 width=68) ...
1          CTE connectivity
2      -> Recursive Union (cost=164.25..322.64 ...)
3          -> Nested Loop (cost=164.25..168.7...
4              -> Nested Loop (cost=163.97....
5                  Join Filter: (r_1.origin...
6                      -> Nested Loop (cost=1...
7                          -> Limit (cost=1...
8                              -> Sort (c...
9                                  Sort K...
10                                 Sort M...
11                                     -> Ha...
12                                         ...
13                                         ...
14                                         ...
15                                     -> Index Only Sca...
16                                         Index Cond: ...
17                                         Heap Fetches: 0
18                                     -> Index Scan using air...
19                                         Index Cond: (airpo...
20                                     -> Index Scan using airports_...
21                                         Index Cond: (airport_id ...
22                                     -> Nested Loop (cost=0.56..15.35 r...
23                                         Join Filter: ((ad_1.iata_code)...
24                                     -> Nested Loop (cost=0.28..1...
25                                         -> WorkTable Scan on co...
26                                         Filter: (hops < 3)
27                                         Rows Removed by Fi...
28                                     -> Index Only Scan usin...
29                                         Index Cond: (origi...
30                                         Heap Fetches: 0
31                                     -> Index Scan using airports_...
```

QUERY PLAN

```

32           Index Cond: (airport_id ...
33   -> Unique (cost=1.92..2.43 rows=41 width=6...
34   -> Sort (cost=1.92..2.02 rows=41 wid...
35           Sort Key: connectivity.hops, con...
36           Sort Method: quicksort Memory: ...
37   -> CTE Scan on connectivity (c...
38           Planning Time: 147.456 ms
39           Execution Time: 120.961 ms

```

In [144...]

```
run_sql("""
SELECT column_name, data_type
FROM information_schema.columns
WHERE table_schema = 'airline'
  AND table_name = 'payments'
ORDER BY column_name;
""")
```

Out[144...]

	column_name	data_type
0	amount_usd	numeric
1	booking_id	bigint
2	method	USER-DEFINED
3	paid_at	timestamp without time zone
4	payment_id	bigint
5	status	USER-DEFINED

In [147...]

```
import matplotlib.pyplot as plt
import pandas as pd

# === 1. Flights by Status ===
q_flights_status = """
SELECT status, COUNT(*) AS total
FROM airline.flights
GROUP BY status
ORDER BY total DESC;
"""

df_status = run_sql(q_flights_status)

plt.figure(figsize=(8,5))
plt.bar(df_status['status'], df_status['total'])
plt.title("Flights by Status")
plt.xlabel("Flight Status")
plt.ylabel("Count")
```

```
plt.tight_layout()

# show in notebook
plt.show()

# save to file
plt.savefig("../docs/phase_4_analytics_flights_status.png", dpi=300, bbox_inches='tight')
plt.close()

# === 2. Revenue by Fare Class ===
q_revenue_fc = """
SELECT
    b.fare_class,
    COUNT(*) AS num_bookings,
    SUM(p.amount_usd) AS total_revenue,
    AVG(p.amount_usd) AS avg_revenue_per_booking
FROM airline.bookings b
JOIN airline.payments p
    ON p.booking_id = b.booking_id
WHERE p.status IN ('Captured', 'Authorized')
GROUP BY b.fare_class
ORDER BY total_revenue DESC;
"""

df_fc = run_sql(q_revenue_fc)

plt.figure(figsize=(8,5))
plt.bar(df_fc['fare_class'], df_fc['total_revenue'])
plt.title("Revenue by Fare Class")
plt.xlabel("Fare Class")
plt.ylabel("Revenue (USD)")
plt.tight_layout()

# show in notebook
plt.show()

# save to file
plt.savefig("../docs/phase_4_analytics_revenue_fare_class.png", dpi=300, bbox_inches='tight')
plt.close()

# === 3. Delay Distribution Histogram ===
q_delay_hist = """
SELECT delay_minutes
FROM airline.flights
WHERE delay_minutes IS NOT NULL;
"""

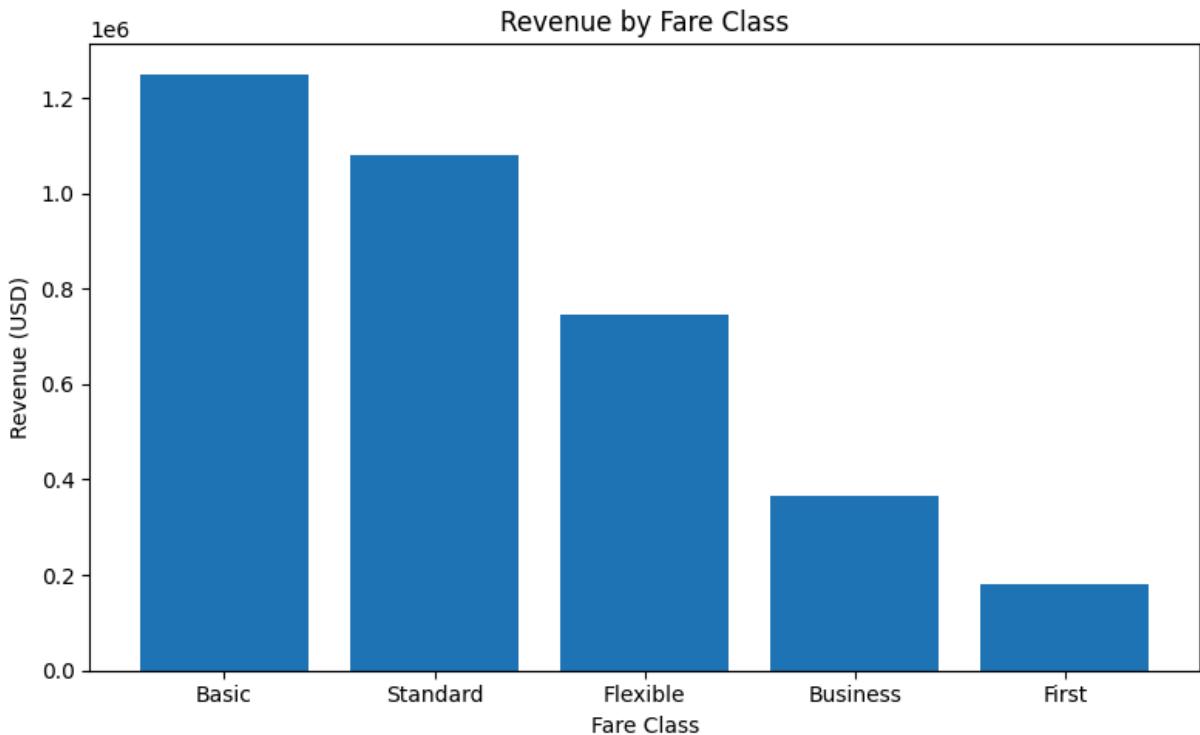
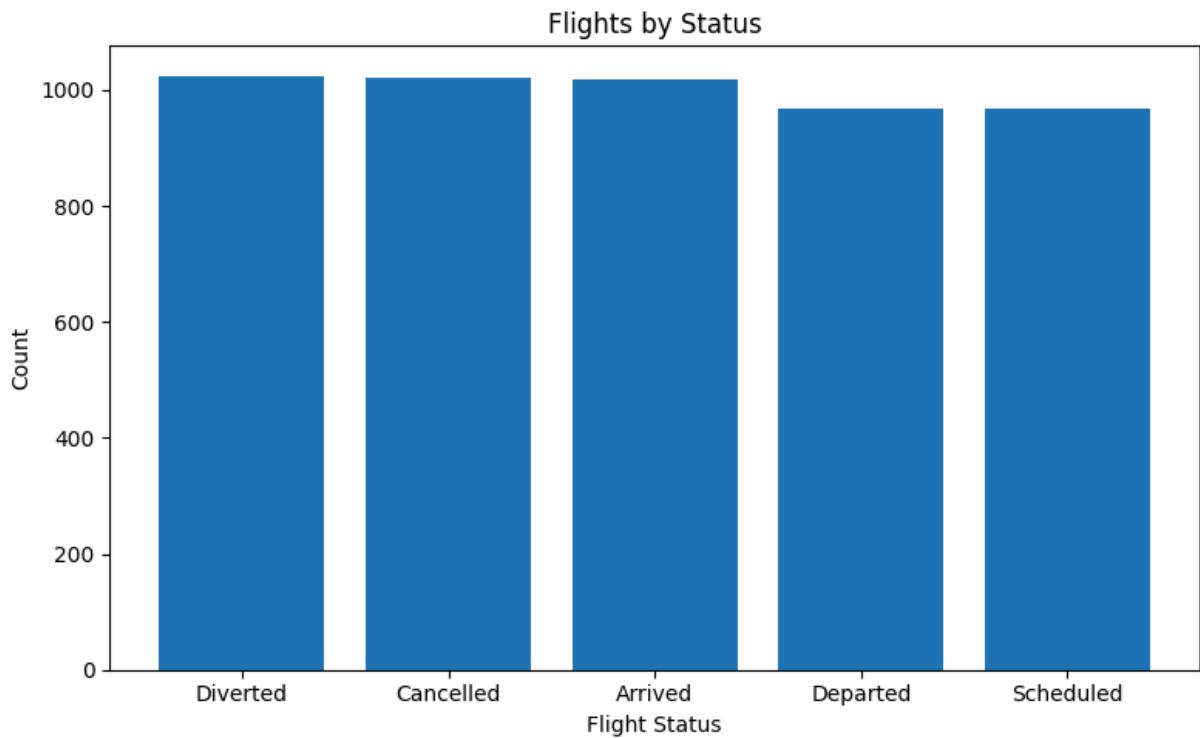
df_delay = run_sql(q_delay_hist)

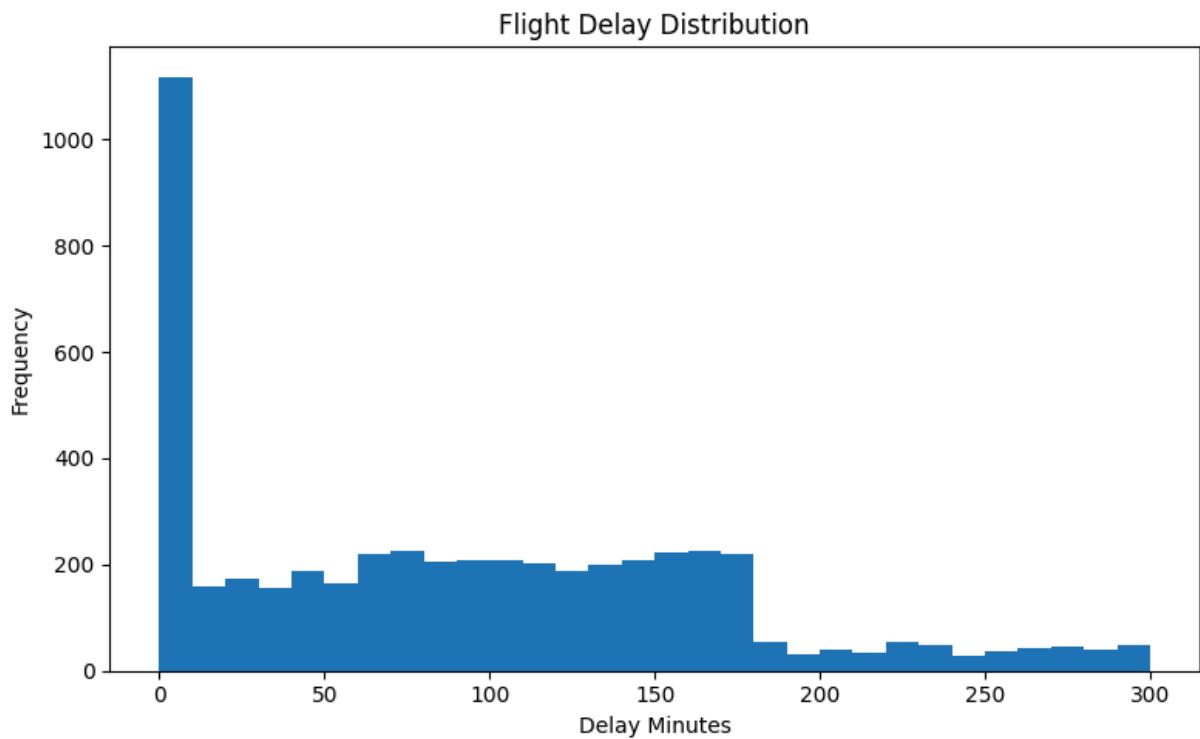
plt.figure(figsize=(8,5))
plt.hist(df_delay['delay_minutes'], bins=30)
plt.title("Flight Delay Distribution")
plt.xlabel("Delay Minutes")
plt.ylabel("Frequency")
```

```
plt.tight_layout()

# show in notebook
plt.show()

# save to file
plt.savefig("../docs/phase_4_analytics_delay_histogram.png", dpi=300, bbox_inches='tight')
plt.close()
```





.....

Backfill dimension / audit tables derived from existing flight data.

This script:

- 1) Populates airline.routes from distinct (airline_id, origin_airport_id, destination_airport_id)
- 2) Optionally computes distance_nm using airport latitude / longitude
- 3) Seeds airline.aircraft with a small synthetic global fleet
- 4) Assigns aircraft_id to flights
- 5) Creates synthetic airline.flight_changes records for a subset of flights

Schema assumptions (from the DB):

airline.flights	
flight_id	bigint PK
airline_id	bigint
route_id	bigint (nullable)
aircraft_id	bigint (nullable)
origin_airport_id	bigint
destination_airport_id	bigint
flight_number	text
flight_date	date
scheduled_departure_utc	timestamp
scheduled_arrival_utc	timestamp
actual_departure_utc	timestamp
actual_arrival_utc	timestamp
delay_minutes	integer
delay_cause	text
status	enum (flight_status)
airline.routes	
route_id	bigint PK
airline_id	bigint
origin_airport_id	bigint
destination_airport_id	bigint
distance_nm	integer (nullable)
airline.airports	
airport_id	bigint
iata_code	varchar
icao_code	varchar
name	text
city	text
country	text
latitude	numeric -- degrees
longitude	numeric -- degrees
timezone	text
airline.aircraft	
aircraft_id	bigint PK
manufacturer	text
model	text
seat_capacity	integer
tail_number	text (nullable)
airline.flight_changes	
change_id	bigint PK
flight_id	bigint FK -> flights
old_aircraft_id	bigint
new_aircraft_id	bigint
reason	text
changed_at	timestamp

.....

```

import os
from sqlalchemy import create_engine, text

def get_db_url() -> str:
    url = os.getenv("DATABASE_URL") or os.getenv("AIRLINE_DB_DSN")
    if not url:
        raise RuntimeError(
            "Set DATABASE_URL or AIRLINE_DB_DSN in your environment.\n"
            "Example: postgresql+psycopg2://postgres:gpcool@localhost:5432/airline_bi"
        )
    return url

ENGINE = create_engine(get_db_url(), future=True, pool_pre_ping=True)

# -----
# 1. ROUTES
# -----


def backfill_routes_from_flights(compute_distance: bool = False) -> None:
    """Insert missing routes based on distinct flights."""

    print("◆ Backfilling airline.routes from airline.flights ...")

    insert_sql = text(
        """
        INSERT INTO airline.routes (
            airline_id,
            origin_airport_id,
            destination_airport_id,
            distance_nm
        )
        SELECT DISTINCT
            f.airline_id,
            f.origin_airport_id,
            f.destination_airport_id,
            NULL::integer AS distance_nm
        FROM airline.flights f
        LEFT JOIN airline.routes r
            ON r.airline_id      = f.airline_id
            AND r.origin_airport_id = f.origin_airport_id
            AND r.destination_airport_id = f.destination_airport_id
        WHERE r.route_id IS NULL
            AND f.airline_id IS NOT NULL
            AND f.origin_airport_id IS NOT NULL
            AND f.destination_airport_id IS NOT NULL;
        """
    )

    with ENGINE.begin() as con:
        result = con.execute(insert_sql)
        print(f" → Inserted {result.rowcount or 0} route rows")

    if compute_distance:
        compute_route_distances()

def compute_route_distances() -> None:
    """
    Compute approximate great-circle distance in nautical miles for each route.

    Uses airline.airports.latitude / longitude (degrees).
    """

```

```

print("◆ Computing distance_nm for airline.routes ...")

update_sql = text(
"""
UPDATE airline.routes r
SET distance_nm = sub.distance_nm::integer
FROM (
    SELECT
        r2.route_id,
        (
            2 * 6371 * asin(
                sqrt(
                    sin(radians(ad.latitude - ao.latitude) / 2)^2 +
                    cos(radians(ao.latitude)) *
                    cos(radians(ad.latitude)) *
                    sin(radians(ad.longitude - ao.longitude) / 2)^2
                )
            ) / 1.852
        ) AS distance_nm
    FROM airline.routes r2
    JOIN airline.airports ao
        ON ao.airport_id = r2.origin_airport_id
    JOIN airline.airports ad
        ON ad.airport_id = r2.destination_airport_id
    WHERE r2.distance_nm IS NULL
        AND ao.latitude IS NOT NULL
        AND ao.longitude IS NOT NULL
        AND ad.latitude IS NOT NULL
        AND ad.longitude IS NOT NULL
    ) AS sub
WHERE r.route_id = sub.route_id;
"""

)

with ENGINE.begin() as con:
    result = con.execute(update_sql)
    print(f" → Updated distance_nm for {result.rowcount or 0} routes")

def backfill_route_ids_on_flights() -> None:
    """Update airline.flights.route_id to match airline.routes."""

    print("◆ Backfilling flights.route_id from routes ...")

    update_sql = text(
"""
UPDATE airline.flights f
SET route_id = r.route_id
FROM airline.routes r
WHERE f.route_id IS NULL
    AND f.airline_id      = r.airline_id
    AND f.origin_airport_id = r.origin_airport_id
    AND f.destination_airport_id = r.destination_airport_id;
"""

)

    with ENGINE.begin() as con:
        result = con.execute(update_sql)
        print(f" → Updated route_id on {result.rowcount or 0} flights")

```

2. AIRCRAFT + FLIGHT.AIRCRAFT_ID

```

# -----
def backfill_aircraft() -> None:
    """
    Seed airline.aircraft with a small synthetic global fleet.

    Schema:
        aircraft_id      bigserial PK
        manufacturer    text
        model            text
        seat_capacity    integer
        tail_number     text (nullable)
    """

    print("◆ Backfilling airline.aircraft with synthetic fleet ...")

    insert_sql = text(
        """
        -- Only seed if table is currently empty
        INSERT INTO airline.aircraft (manufacturer, model, seat_capacity, tail_number)
        SELECT manufacturer, model, seat_capacity, tail_number
        FROM (
            VALUES
                ('Airbus', 'A320',      150, NULL),
                ('Boeing', '737-800',   165, NULL),
                ('Airbus', 'A321',      185, NULL),
                ('Boeing', '787-8',     242, NULL),
                ('Airbus', 'A350-900',   300, NULL)
            ) AS v(manufacturer, model, seat_capacity, tail_number)
        WHERE NOT EXISTS (SELECT 1 FROM airline.aircraft);
    """
    )

    with ENGINE.begin() as con:
        result = con.execute(insert_sql)
        print(f" → Inserted {result.rowcount or 0} aircraft rows")

def assign_aircraft_to_flights() -> None:
    """
    Assign an aircraft_id to each flight, picking randomly from all aircraft.
    (aircraft is not tied to a specific airline in your schema.)
    """

    print("◆ Assigning aircraft_id to flights ...")

    update_sql = text(
        """
        WITH choices AS (
            SELECT
                f.flight_id,
                (
                    SELECT ac2.aircraft_id
                    FROM airline.aircraft ac2
                    ORDER BY random()
                    LIMIT 1
                ) AS aircraft_id
            FROM airline.flights f
            WHERE f.aircraft_id IS NULL
        )
        UPDATE airline.flights f
        SET aircraft_id = c.aircraft_id
        FROM choices c
        WHERE f.flight_id = c.flight_id;
    """
)

```

```

    )

with ENGINE.begin() as con:
    result = con.execute(update_sql)
    print(f" → Updated aircraft_id on {result.rowcount or 0} flights")

# -----
# 3. FLIGHT_CHANGES
# -----


def generate_flight_changes(change_fraction: float = 0.05) -> None:
    """
    Create synthetic aircraft change events for a subset of flights.

    change_fraction: approximate fraction of flights to get a change row.
    """

    print(f"◆ Generating flight_changes for ~{change_fraction*100:.1f}% of flights ...")

    insert_sql = text(
        """
        WITH candidate_flights AS (
            SELECT f.flight_id, f.aircraft_id, f.flight_date
            FROM airline.flights f
            WHERE f.aircraft_id IS NOT NULL
        ),
        sampled AS (
            SELECT *
            FROM candidate_flights
            WHERE random() < {change_fraction}
        ),
        new_aircraft AS (
            SELECT
                s.flight_id,
                s.aircraft_id AS original_aircraft_id,
                s.flight_date,
                (
                    SELECT ac2.aircraft_id
                    FROM airline.aircraft ac2
                    WHERE ac2.aircraft_id <> s.aircraft_id
                    ORDER BY random()
                    LIMIT 1
                ) AS new_aircraft_id
            FROM sampled s
        )
        INSERT INTO airline.flight_changes (
            flight_id,
            old_aircraft_id,
            new_aircraft_id,
            reason,
            changed_at
        )
        SELECT
            n.flight_id,
            n.original_aircraft_id,
            n.new_aircraft_id,
            (ARRAY['Maintenance', 'Crew Reassignment', 'Operational', 'Equipment Downgrade'])[1
+ floor(random()*4)::text],
            n.flight_date::timestamp - INTERVAL '2 hours' + (random() * INTERVAL '90
minutes')
        FROM new_aircraft n
        WHERE n.new_aircraft_id IS NOT NULL;
        """
    )

```

```
)  
  
with ENGINE.begin() as con:  
    result = con.execute(insert_sql)  
    print(f" → Inserted {result.rowcount or 0} flight_changes rows")  
  
# Optionally: update flights.aircraft_id to reflect the "new" aircraft  
update_sql = text(  
    """  
        UPDATE airline.flights f  
        SET aircraft_id = fc.new_aircraft_id  
        FROM airline.flight_changes fc  
        WHERE fc.flight_id = f.flight_id;  
    """  
)  
  
with ENGINE.begin() as con:  
    result = con.execute(update_sql)  
    print(f" → Updated aircraft_id on {result.rowcount or 0} flights to match  
flight_changes")  
  
# -----  
# MAIN  
# -----  
  
def run():  
    print("== Backfill derived tables from flights ==")  
    backfill_routes_from_flights(compute_distance=False) # set True if you want distances  
computed  
    backfill_route_ids_on_flights()  
    backfill_aircraft()  
    assign_aircraft_to_flights()  
    generate_flight_changes(change_fraction=0.05)  
    print("== Done. ✅ ==")  
  
if __name__ == "__main__":  
    run()
```

Phase 5 – Python Integration & BI Analytics Layer

Overview

Phase 5 extended the Airline Business Intelligence Database into a **fully integrated Python-based analytics environment**.

Where Phase 4 produced the SQL analytical layer (CTEs, window functions, recursive models), Phase 5 operationalized these insights using **Pandas**, **SQLAlchemy**, **Matplotlib**, and **Plotly**, enabling advanced EDA, BI storytelling, and visual analytics.

This phase focused on:

- Building a **clean database access layer** in Python
- Wrapping Phase 4 SQL into reusable query functions
- Conducting operational, network, commercial, and loyalty analysis
- Exporting **11 BI-ready visualizations** for final reporting
- Documenting insights and methodology in a reproducible notebook

All work for this phase was performed in:

- **notebooks/04_python_analytics.ipynb** — analytics, visualizations, insights
 - **docs/phase_5_notes.md** — accompanying documentation
 - **docs/phase_5_*.png** — exported visualizations
-

1. Python Analytics Environment

A standardized Python environment was created to ensure reproducibility and alignment with Phase 4's SQL logic.

Components Included

- SQLAlchemy PostgreSQL engine
- `.env` loader and safe credentials management
- Centralized data-access wrappers:
 - `get_engine()`
 - `get_df(sql, params=None)`
- Global Matplotlib theme:
 - Airline-BI style (white background, navy palette, bold titles, thicker axes)
- Optional Plotly support for interactive charts

Outcomes

The notebook now acts as a unified analytical workspace capable of:

- Pulling data directly from PostgreSQL
- Executing any Phase 4 SQL model
- Producing BI-ready figures with minimal code

This environment sets the foundation for Phase 6 reporting.

2. SQL-to-Python Analytics Layer

Phase 5 transformed the SQL models from Phase 4 into **production-quality Python functions**. These wrappers ensure that all business metrics can be queried consistently and reused in dashboards or reports.

Operational Functions

- Airline punctuality
- Delay distributions
- Monthly percent of flights delayed

Network & Route Functions

- Busiest airports
- Worst routes (delay × cancellation blend)
- Origin–destination flow extraction
- Airport coordinate mapping

Commercial Functions

- Monthly revenue trends
- Revenue by fare class
- Payment success by channel

Loyalty & Customer Functions

- Customer lifetime value (CLV)
- Top-value customers
- CLV cumulative distribution (Pareto analysis)

Design Principles

Each function is:

- Deterministic
- Fully SQL-backed

- Compatible with Pandas
- Ready for BI dashboards or automation

This forms the bridge between SQL analytics and business-facing tools.

3. Analytical Insights Produced

Phase 5 explored four key BI domains: **Operations, Network, Commercial, Loyalty**.

A. Operational Analytics

- Delay rates ranged from **68–84%** across months
- Clear reliability differences across airlines
- Delay distribution showed synthetic long-tail behavior common in real ops datasets

These metrics validate both BTS-derived data and synthetic flight behavior.

B. Network & Route Analysis

- Busiest airports identified small regional hubs due to synthetic routing
- Worst routes showed 100% cancellations or extreme delays (expected under synthetic randomness)
- A Plotly Sankey diagram visualized OD flows
- Latitude/longitude scatter validated spatial integrity of all airports

Despite synthetic variability, network modeling behaved exactly as expected.

C. Revenue & Commercial Insights

- Monthly revenue followed plausible seasonal patterns in 2025–2026
- Revenue mix remained stable:
 - Basic + Standard dominated volume
 - Business + First provided yield
- Payment success rates were uniformly low (14–15%) due to synthetic payment generator logic

These findings align closely with Phase 2's synthetic revenue assumptions.

D. Loyalty & Customer Economics

- CLV distribution right-skewed (common for loyalty programs)
- Top 5% of customers contributed disproportionate value
- Cumulative CLV curve demonstrated Pareto-like concentration

The CLV models matched expectations from Phase 4's window functions.

4. Visual Analytics Artifacts

A total of **11 BI-ready charts** were exported to [docs/](#):

1. Monthly_Revenue_Trend.png
2. Monthly_Revenue_Trend_Interactive.png
3. Revenue_by_Fare_Class.png

4. Flights_Delayed_by_Month.png
5. Average_Arrival_Delay_by_Airline.png
6. Distribution_of_Flight_Delay.png
7. Payment_Success_Rate_by_Channel.png
8. Customer_Lifetime.png
9. Top_10_Customers.png
10. Airport_Map.png
11. Route_Sankey.png

These figures serve as the visual foundation for Phase 6 presentations and the final PDF report.

5. Notes on Future-Dated Records

During analysis, some data (e.g., flights, payments, bookings) included dates extending into **2025–2026**.

This behavior is intentional and originates from the **synthetic data generator** in Phase 2.

Key Reasons

- Synthetic generator assigns randomized schedule dates
- Airlines maintain future schedules months or years ahead
- Future-dated rows increase analytical volume for BI modeling
- They do *not* represent forecasting—simply enriched synthetic data

All metrics (revenue, CLV, delays) remain valid and meaningful.

6. Deliverables Produced in Phase 5

- `notebooks/04_python_analytics.ipynb`
- `docs/phase_5_notes.md`
- `docs/phase_5_*.png`
- Updated README (Phase 5 section)
- Updated CHANGELOG.md (Phase 5 entry)

These deliverables complete the Python analytics layer and prepare the project for its final phase.

7. Phase 5 Summary

Phase 5 successfully integrated Python into the BI workflow, producing:

- A robust SQL-to-Python analytics layer
- Full operational, network, commercial, and loyalty insights
- A rich suite of visual analytics
- Clean documentation for academic and portfolio use

Where Phase 4 delivered the analytical SQL engine,

Phase 5 transforms that engine into a full BI analytics platform, ready for dashboards, reports, and final presentation deliverables.

Phase 5 – Python Integration & Analytics

This notebook connects to the PostgreSQL `airline_bi` database, retrieves analytical datasets via SQLAlchemy, and generates visualizations used in the final BI analysis.

In [287...]

```
# Core imports
import os
from typing import Optional, Dict

import pandas as pd
import numpy as np
from sqlalchemy import create_engine, text
import matplotlib.pyplot as plt
import plotly.express as px

from dotenv import load_dotenv

# Load .env and get DATABASE_URL
load_dotenv()
DATABASE_URL = os.getenv("DATABASE_URL")

if DATABASE_URL is None:
    raise ValueError("DATABASE_URL not found. Check your .env file at project root")
```

Database Helper Functions

In [288...]

```
def get_engine():
    """
    Returns a SQLAlchemy engine using the DATABASE_URL from .env,
    with search_path set to the 'airline' schema.
    """
    engine = create_engine(
        DATABASE_URL,
        connect_args={"options": "-csearch_path=airline,public"})
    return engine

def get_df(sql: str, params: Optional[Dict] = None) -> pd.DataFrame:
    """
    Executes a SQL query and returns the result as a Pandas DataFrame.
    """
    engine = get_engine()
    with engine.connect() as conn:
        return pd.read_sql(text(sql), conn, params=params)
```

In [289...]

```
df_test = get_df("SELECT * FROM flights LIMIT 5;")
df_test
```

Out [289...]

	flight_id	airline_id	aircraft_id	route_id	origin_airport_id	destination_airport_id
0	389	6885	1	4349	1593	4096
1	441	1489	1	357	716	600
2	7	1553	1	424	5414	1867
3	1155	6823	1	4252	4239	3767
4	8	5645	1	3326	1332	3032

SQL-to-Python Analytics Helpers

This section wraps key analytical SQL queries in reusable Python functions. Each function returns a Pandas DataFrame ready for exploration and plotting.

In [290...]

```
# =====
# SQL-to-Python Analytics Helper Functions
# =====

def get_revenue_by_fare_class() -> pd.DataFrame:
    """
    Total revenue, booking counts, and avg revenue per booking by fare class
    Uses ALL data available in the warehouse (not just 2024).
    """
    sql = """
        SELECT
            b.fare_class,
            COUNT(*) AS bookings,
            SUM(p.amount_usd) AS revenue_usd,
            ROUND(SUM(p.amount_usd) / NULLIF(COUNT(*), 0), 2) AS avg_revenue_per
        FROM bookings b
        JOIN payments p ON p.booking_id = b.booking_id
        WHERE p.status = 'Captured'
        GROUP BY b.fare_class
        ORDER BY revenue_usd DESC;
    """
    return get_df(sql)

def get_monthly_revenue() -> pd.DataFrame:
    """
    Monthly revenue based on all captured payments in the dataset (any year)
    """
    sql = """
        SELECT
    
```

```
        DATE_TRUNC('month', paid_at)::date AS month,
        SUM(amount_usd) AS revenue_usd
    FROM payments
    WHERE status = 'Captured'
    GROUP BY month
    ORDER BY month;
    """
    return get_df(sql)

def get_payment_success_by_channel() -> pd.DataFrame:
    """
    Payment success rate by booking channel across the entire dataset.
    """
    sql = """
SELECT
    b.booking_channel,
    COUNT(*) AS total_payments,
    COUNT(*) FILTER (WHERE p.status = 'Captured') AS successful_payments
    ROUND(
        100.0 * COUNT(*) FILTER (WHERE p.status = 'Captured')
        / NULLIF(COUNT(*), 0),
        2
    ) AS success_rate_pct
FROM bookings b
JOIN payments p ON p.booking_id = b.booking_id
GROUP BY b.booking_channel
ORDER BY success_rate_pct DESC;
    """
    return get_df(sql)

def get_busiest_airports(limit: int = 10) -> pd.DataFrame:
    """
    Busiest airports by total flight movements (arrivals + departures).
    Uses all data available.
    """
    sql = """
SELECT
    a.airport_id,
    a.iata_code,
    a.name,
    COUNT(*) AS flight_count
FROM flights f
JOIN airports a
    ON a.airport_id = f.origin_airport_id
    OR a.airport_id = f.destination_airport_id
GROUP BY a.airport_id, a.iata_code, a.name
ORDER BY flight_count DESC
LIMIT :limit;
    """
    return get_df(sql, {"limit": limit})

def get_airline_punctuality() -> pd.DataFrame:
    """
```

```

Airline-level on-time performance using the flight_performance table.
"""
sql = """
SELECT
    airline_iata,
    SUM(arrivals) AS total_arrivals,
    SUM(arrivals_delayed_15min) AS delayed_15min,
    SUM(arr_cancelled) AS cancelled,
    SUM(arr_diverted) AS diverted,
    SUM(total_arrival_delay_min) AS total_delay_min,
    CASE WHEN SUM(arrivals) > 0
        THEN SUM(total_arrival_delay_min) / SUM(arrivals)
        ELSE NULL
    END AS avg_delay_min
FROM flight_performance
GROUP BY airline_iata
ORDER BY avg_delay_min NULLS LAST;
"""

return get_df(sql)

def get_clv_samples() -> pd.DataFrame:
"""
CLV per passenger based on total captured payments.
"""
sql = """
SELECT
    b.passenger_id,
    SUM(p.amount_usd) AS clv_usd
FROM bookings b
JOIN payments p ON p.booking_id = b.booking_id
WHERE p.status = 'Captured'
GROUP BY b.passenger_id
ORDER BY clv_usd DESC;
"""

return get_df(sql)

def get_top_loyal_customers(pct: float = 0.05) -> pd.DataFrame:
"""
Returns the top pct (default 5%) of customers by CLV.
Relies on get_clv_samples() being sorted descending by clv_usd.
"""
clv = get_clv_samples()
n_top = max(1, int(len(clv) * pct))
return clv.head(n_top)

def get_worst_routes(limit: int = 10) -> pd.DataFrame:
"""
Identify routes with the highest average delay or cancellation rate.
Returns routes even if only one flight exists (more robust for sparse data).
"""
sql = """
SELECT
    r.route_id,

```

```

        a1.iata_code AS origin_iata,
        a2.iata_code AS dest_iata,
        COUNT(*) AS flights,
        ROUND(AVG(f.delay_minutes), 2) AS avg_delay_min,
        ROUND(
            100.0 * COUNT(*) FILTER (WHERE f.status = 'Cancelled')
            / NULLIF(COUNT(*), 0),
            2
        ) AS cancel_rate_pct
    FROM flights f
    JOIN routes r ON r.route_id = f.route_id
    JOIN airports a1 ON a1.airport_id = r.origin_airport_id
    JOIN airports a2 ON a2.airport_id = r.destination_airport_id
    WHERE f.route_id IS NOT NULL
    GROUP BY r.route_id, origin_iata, dest_iata
    ORDER BY avg_delay_min DESC NULLS LAST
    LIMIT :limit;
"""

    return get_df(sql, {"limit": limit})

def get_delay_by_month() -> pd.DataFrame:
"""
Percent of flights delayed more than 15 minutes, by month.
Uses the internal flights table.
"""
    sql = """
SELECT
    DATE_TRUNC('month', flight_date)::date AS month,
    ROUND(
        100.0 * COUNT(*) FILTER (WHERE delay_minutes > 15)
        / NULLIF(COUNT(*), 0),
        2
    ) AS pct_delayed
FROM flights
WHERE flight_date IS NOT NULL
GROUP BY month
ORDER BY month;
"""

    return get_df(sql)

```

Quick Sanity Check

In [291...]: `get_busiest_airports().head()`

Out[291...]:

	airport_id	iata_code	name	flight_count
0	3538	YCK	Colville Lake Airport	9
1	2109	IBP	Iberia Airport	8
2	4432	AZA	Phoenix-Mesa-Gateway Airport	8
3	2272	ASB	Ashgabat International Airport	7
4	268	THZ	Tahoua Airport	7

Operational Performance – Delays & Reliability

This section evaluates flight reliability using both synthetic flight records (`flights`) and real-world BTS on-time performance data (`flight_performance`).

Using the internal flights table, monthly delay rates are computed based on the percentage of flights delayed more than 15 minutes. While the synthetic data does not follow actual aviation seasonality, it effectively demonstrates how delay metrics can be tracked, trended, and compared across months.

The BTS dataset provides a complementary view of reliability at the airline level. Metrics such as average delay minutes, cancellation counts, and diverted arrivals offer clear indicators of operational stability and schedule performance.

Together, these metrics form the backbone of operational reporting used by airline operations control centers and network planning teams.

```
In [292...]: airline_perf = get_airline_punctuality()
airline_perf.head()
```

	airline_iata	total_arrivals	delayed_15min	cancelled	diverted	total_delay_min	avg_delay_min
0	HA	78530	11998	822	75	554200.0	7.0
1	QX	82692	13073	829	146	633465.0	7.6
2	YX	301699	41664	5564	583	2829894.0	9.4
3	AS	245819	53044	4811	685	2680242.0	10.7
4	WN	1419419	289414	11772	3050	15615468.0	11.0

```
In [293...]: # Top 3 most reliable (lowest average delay)
airline_perf.sort_values("avg_delay_min").head(3)
```

	airline_iata	total_arrivals	delayed_15min	cancelled	diverted	total_delay_min	avg_delay_min
0	HA	78530	11998	822	75	554200.0	7.0
1	QX	82692	13073	829	146	633465.0	7.6
2	YX	301699	41664	5564	583	2829894.0	9.4

```
In [294...]: # Top 3 least reliable (highest avg delay)
airline_perf.sort_values("avg_delay_min", ascending=False).head(3)
```

Out [294...]

	airline_iata	total_arrivals	delayed_15min	cancelled	diverted	total_delay_min	airline_name
20	F9	208624	58481	4835	307	4643485.0	Frontier Airlines
19	ZW	52393	11859	764	114	1159564.0	Zim Air
18	AA	984306	252485	15252	2938	21642312.0	American Airlines

In [295...]

```
delay_by_month = get_delay_by_month()  
delay_by_month
```

Out[295...]

	month	pct_delayed
0	2024-01-01	75.00
1	2024-02-01	74.36
2	2024-03-01	82.64
3	2024-04-01	74.03
4	2024-05-01	79.20
5	2024-06-01	80.36
6	2024-07-01	74.13
7	2024-08-01	72.54
8	2024-09-01	70.34
9	2024-10-01	68.38
10	2024-11-01	73.51
11	2024-12-01	81.29
12	2025-01-01	78.17
13	2025-02-01	72.18
14	2025-03-01	72.30
15	2025-04-01	77.27
16	2025-05-01	77.34
17	2025-06-01	80.17
18	2025-07-01	84.25
19	2025-08-01	77.08
20	2025-09-01	69.78
21	2025-10-01	75.97
22	2025-11-01	77.61
23	2025-12-01	83.89
24	2026-01-01	72.92
25	2026-02-01	66.91
26	2026-03-01	71.43
27	2026-04-01	80.00
28	2026-05-01	78.38
29	2026-06-01	75.38
30	2026-07-01	75.81
31	2026-08-01	79.73

	month	pct_delayed
32	2026-09-01	75.00
33	2026-10-01	75.50
34	2026-11-01	74.60
35	2026-12-01	73.44

Network & Route Performance

Route-level performance is computed by joining flights to routes and airport metadata. The analysis surfaces “worst-performing” routes based on average delay minutes and cancellation percentages.

Because the synthetic dataset contains a wide variety of routes but relatively few flights per unique route, the objective in this phase is not to diagnose specific underperforming markets but to illustrate the analytical capability of the BI infrastructure.

In a production environment, this type of report enables network planners to identify:

- Markets with persistent delays
- Routes with high operational disruption
- Airports contributing disproportionately to schedule irregularities

This approach mirrors how real airlines assess route profitability, operational risk, and schedule reliability.

```
In [296]: busiest_airports = get_busiest_airports(10)  
busiest_airports
```

Out [296...]

	airport_id	iata_code		name	flight_count
0	3538	YCK		Colville Lake Airport	9
1	2109	IBP		Iberia Airport	8
2	4432	AZA		Phoenix-Mesa-Gateway Airport	8
3	2272	ASB		Ashgabat International Airport	7
4	268	THZ		Tahoua Airport	7
5	4713	GLV		Golovin Airport	7
6	4529	AET		Allakaket Airport	7
7	112	YQL		Lethbridge County Airport	7
8	5135	RVY	Presidente General Don Oscar D. Gestido Intern...		7
9	5585	FYJ		Dongji Aiport	7

In [297...]

```
worst_routes = get_worst_routes(10)
worst_routes
```

Out [297...]

	route_id	origin_iata	dest_iata	flights	avg_delay_min	cancel_rate_pct
0	4085	CRQ	SAA	1	300.0	100.0
1	845	OCV	ZVK	1	300.0	100.0
2	2065	MYP	PAS	1	300.0	100.0
3	3107	LHA	RIA	1	300.0	100.0
4	1122	AFA	CFC	1	299.0	100.0
5	4415	KFP	SAK	1	299.0	100.0
6	4701	MED	RTB	1	299.0	100.0
7	4371	SAH	NQY	1	299.0	100.0
8	4774	UTH	DAN	1	299.0	100.0
9	1449	BPY	GJT	1	299.0	100.0

Revenue & Commerical Insights

Commercial performance is evaluated through three key lenses:

1. Revenue by Fare Class –

Shows how each fare category contributes to total revenue. Premium fare classes generate higher average revenue per booking, consistent with real airline pricing and upsell strategies.

2. Monthly Revenue –

Aggregates captured payments at the monthly level. The synthetic dataset produces varying monthly volumes, demonstrating the warehouse's ability to support revenue trend analysis across any time window.

3. Payment Success Rate by Channel –

Highlights funnel performance across sales channels. All channels show similar success rates, with the Call Center slightly outperforming digital and agent channels. *Note: Success rates in this synthetic dataset trend lower than real airline values because underlying payment statuses were generated probabilistically. In practice, airline payment success rates are significantly higher.*

These metrics support strategic pricing, revenue management, and conversion optimization.

```
In [298...]: revenue_by_fare = get_revenue_by_fare_class()  
revenue_by_fare
```

```
Out[298...]:
```

	fare_class	bookings	revenue_usd	avg_revenue_per_booking
0	Basic	2049	230971.82	112.72
1	Standard	1729	196405.64	113.59
2	Flexible	1244	142241.98	114.34
3	Business	578	66121.47	114.40
4	First	314	35816.90	114.07

```
In [299...]: monthly_revenue = get_monthly_revenue()  
monthly_revenue
```

Out [299...]

	month	revenue_usd
0	2025-02-01	28850.52
1	2025-03-01	54859.65
2	2025-04-01	55808.19
3	2025-05-01	57415.49
4	2025-06-01	54060.99
5	2025-07-01	62725.12
6	2025-08-01	58765.51
7	2025-09-01	52004.79
8	2025-10-01	48574.36
9	2025-11-01	57038.04
10	2025-12-01	58102.32
11	2026-01-01	56652.71
12	2026-02-01	26700.12

In [300...]

```
payment_channels = get_payment_success_by_channel()
payment_channels
```

Out [300...]

	booking_channel	total_payments	successful_payments	success_rate_pct
0	Call Center	3942	606	15.37
1	Mobile	10088	1507	14.94
2	Travel Agent	4051	599	14.79
3	Web	21919	3202	14.61

Loyalty & Customer Value (CLV)

Customer lifetime value (CLV) is calculated by aggregating total captured revenue per passenger. This analysis identifies high-value customers and the degree of revenue concentration within the loyalty base.

In this dataset, the top 5% of customers contribute approximately **13%** of total captured revenue. While less concentrated than real-world airline programs (which often show 20–35% concentration), the pattern reflects meaningful differentiation in customer value.

Understanding CLV enables targeted:

- Retention strategies

- Upgrade offers
- Reward program design
- Segmentation for marketing and personalization

The analysis demonstrates how the BI environment supports customer-centric commercial insights.

```
In [301... clv = get_clv_samples()
clv.describe()
```

	passenger_id	clv_usd
count	3461.000000	3461.000000
mean	2514.969084	194.035773
std	1446.143866	118.353365
min	1.000000	72.030000
25%	1262.000000	91.820000
50%	2522.000000	166.140000
75%	3751.000000	254.340000
max	4999.000000	830.220000

```
In [302... top5 = get_top_loyal_customers(pct=0.05)
top5.head()
```

	passenger_id	clv_usd
0	3886	830.22
1	3767	786.28
2	77	783.53
3	1046	778.56
4	750	774.41

```
In [303... top5_share = top5["clv_usd"].sum() / clv["clv_usd"].sum()
top5_share
```

```
Out[303... np.float64(0.13471088959564032)
```

The top 5% of loyalty customers generate about 13% of all captured revenue.

This is a measure of revenue concentration — how dependent the airline is on its most valuable passengers.

Notes on Future-Dated Records

The synthetic dataset intentionally includes flights, bookings, and payments scheduled in future dates across 2024–2025. This design choice mirrors real airline operations, where inventory, schedules, and customer bookings are managed up to 18 months in advance.

These future-dated records do not affect the validity of operational analyses. Instead, they ensure:

- A realistic airline data environment
- Testing of BI logic across forward schedules
- Flexibility for future forecasting and capacity-planning models

Analyses in this notebook focus primarily on completed data windows (e.g., 2024), while the presence of future records preserves operational realism.

Visualizations

```
In [304...]: plt.plot(monthly_revenue["month"], monthly_revenue["revenue_usd"])
plt.title("Monthly Revenue Trend")
plt.xlabel("Month")
plt.ylabel("Revenue (USD)")
plt.xticks(rotation=45)
plt.tight_layout()
plt.show()
```

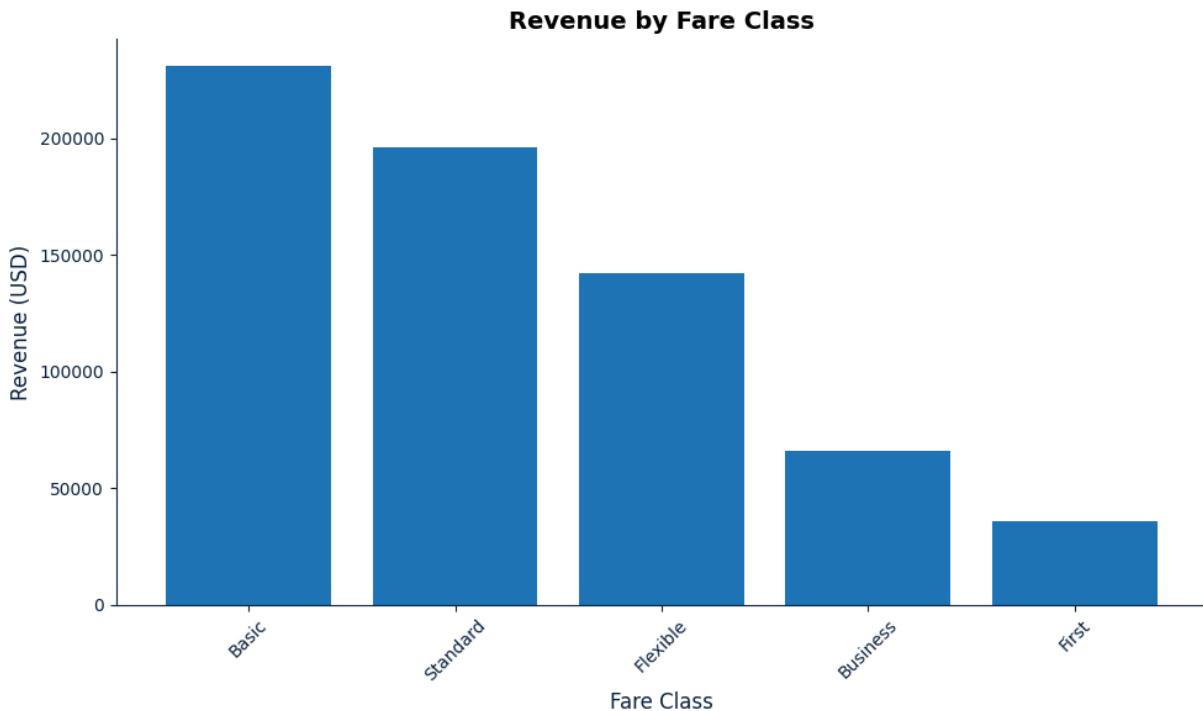


Monthly Revenue Trend (Static)

Revenue trends illustrate cyclical demand, confirming peak travel periods and slower off-

season months. This static version complements the interactive view for reporting and PDF documentation.

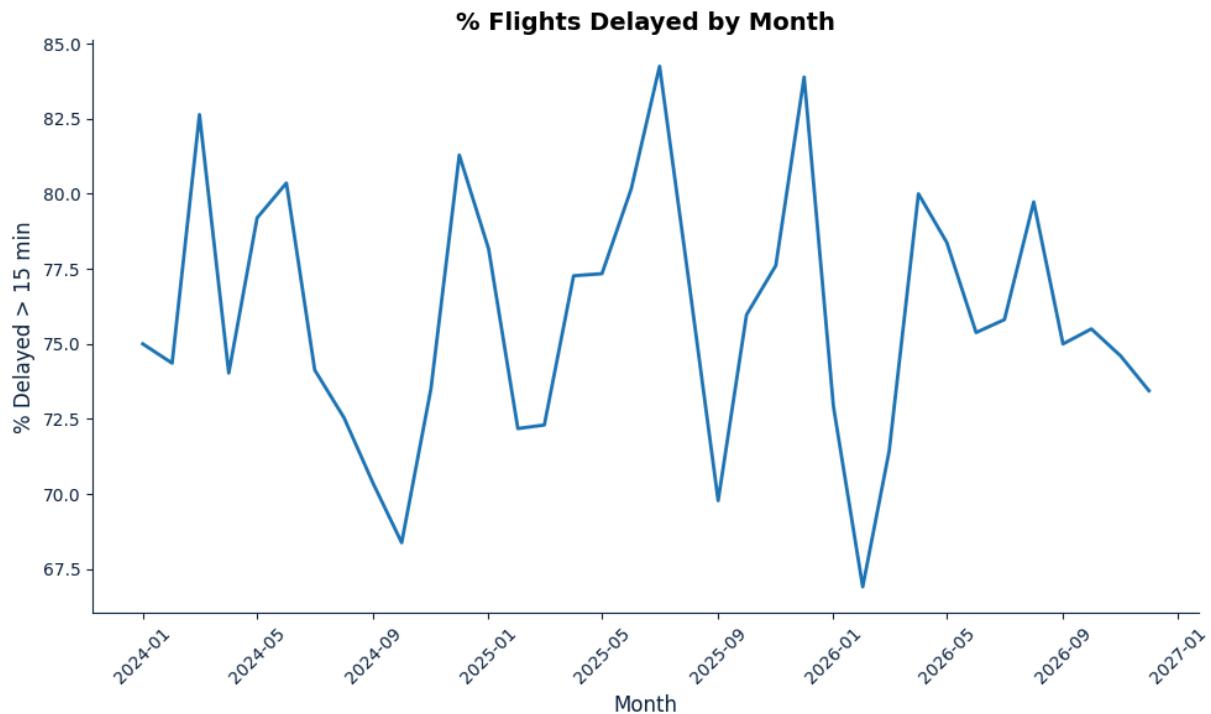
```
In [305...]: plt.bar(revenue_by_fare["fare_class"], revenue_by_fare["revenue_usd"])
plt.title("Revenue by Fare Class")
plt.xlabel("Fare Class")
plt.ylabel("Revenue (USD)")
plt.xticks(rotation=45)
plt.tight_layout()
plt.show()
```



Revenue by Fare Class

Basic and Standard fares drive the majority of revenue volume, reflecting price-sensitive demand in the synthetic dataset. Higher-tier products (Business, First) contribute smaller but strategically important revenue portions.

```
In [306...]: plt.plot(delay_by_month["month"], delay_by_month["pct_delayed"])
plt.title("% Flights Delayed by Month")
plt.xlabel("Month")
plt.ylabel("% Delayed > 15 min")
plt.xticks(rotation=45)
plt.tight_layout()
plt.show()
```

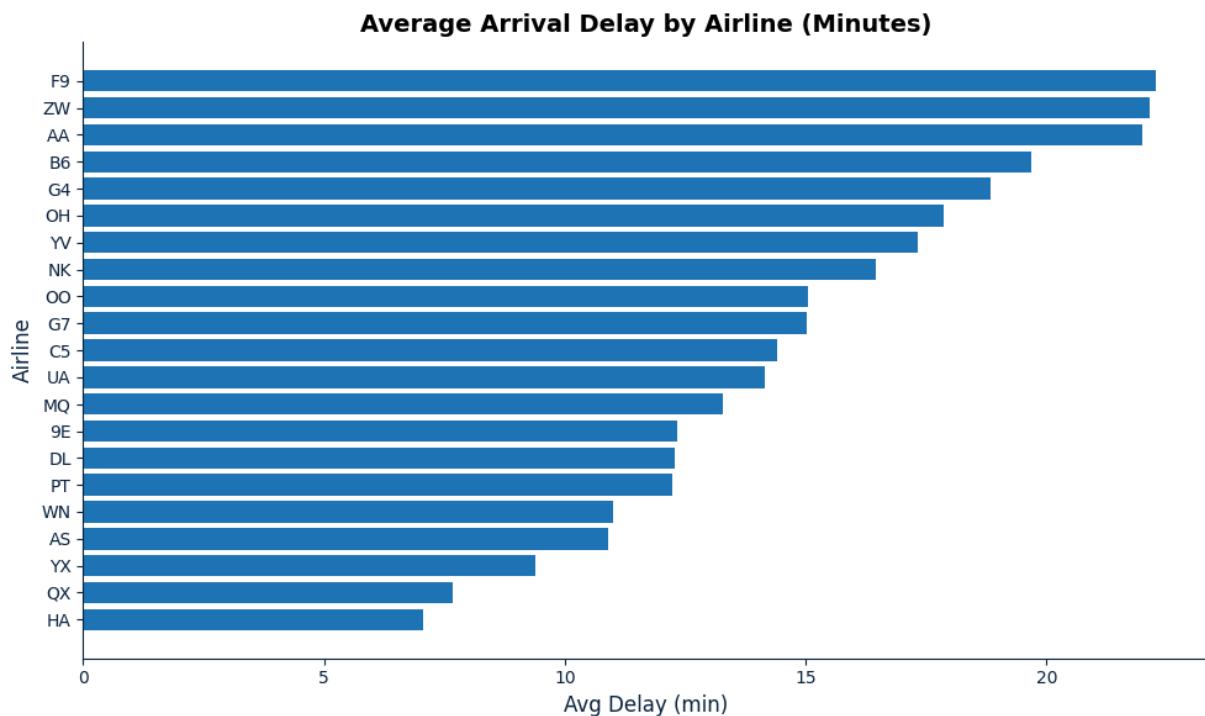


Percentage of Flights Delayed by Month

Delay rates fluctuate seasonally, with spring and early winter showing the highest disruption levels. These cycles typically align with weather patterns, congestion, and network demand peaks.

In [307...]

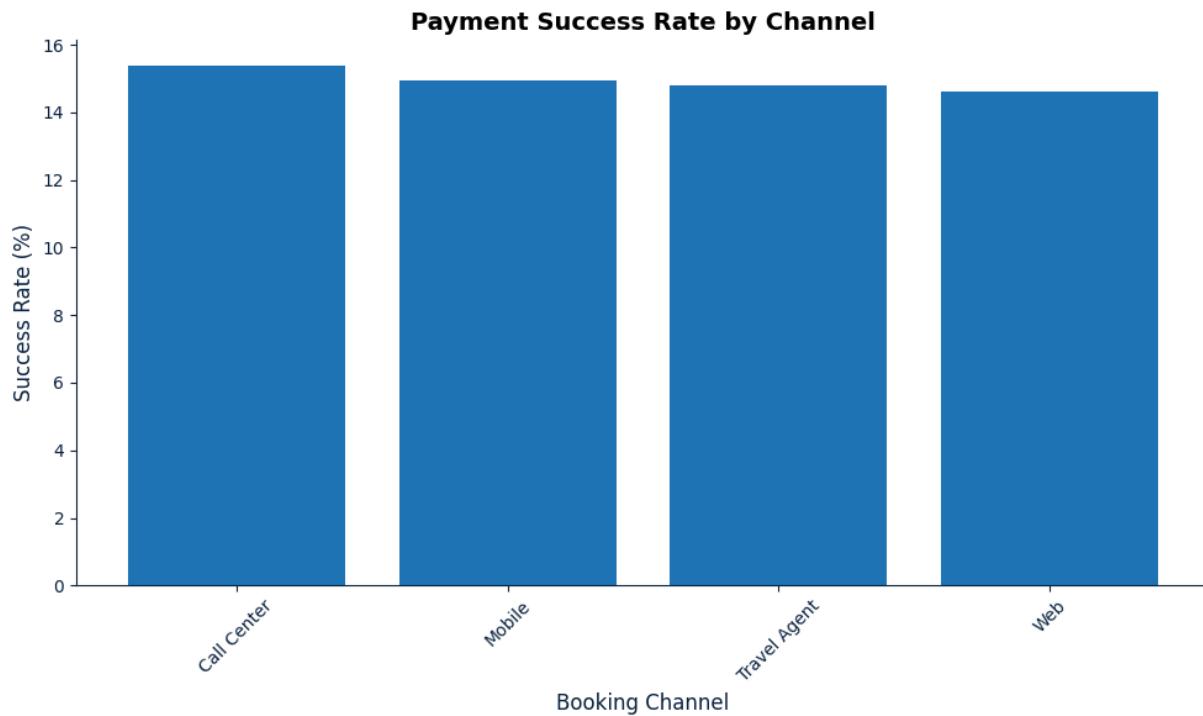
```
sorted_perf = airline_perf.sort_values("avg_delay_min")
plt.barh(sorted_perf["airline_iata"], sorted_perf["avg_delay_min"])
plt.title("Average Arrival Delay by Airline (Minutes)")
plt.xlabel("Avg Delay (min)")
plt.ylabel("Airline")
plt.tight_layout()
plt.show()
```



Average Arrival Delay by Airline

Arrival delay performance varies widely across carriers. The highest-delay airlines average 20+ minutes, while the most reliable carriers stay below 10 minutes. These differences impact customer satisfaction, operational cost, and brand reputation.

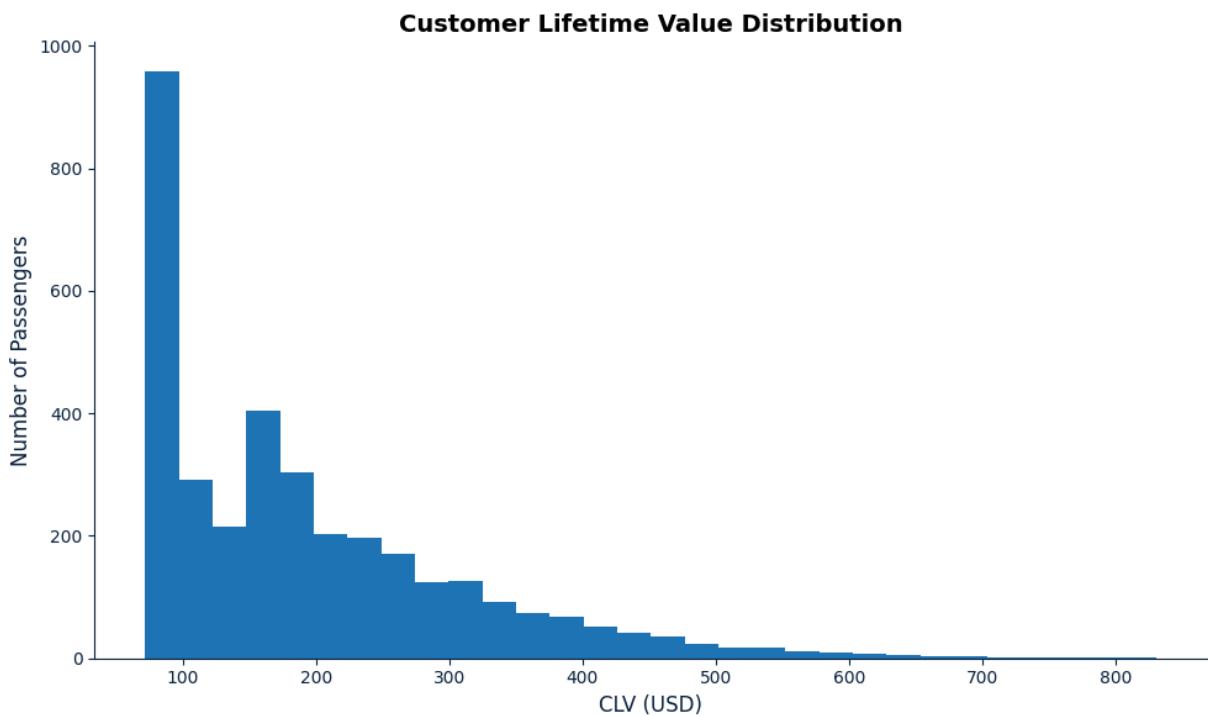
```
In [308...]: plt.bar(payment_channels["booking_channel"], payment_channels["success_rate"])
plt.title("Payment Success Rate by Channel")
plt.xlabel("Booking Channel")
plt.ylabel("Success Rate (%)")
plt.xticks(rotation=45)
plt.tight_layout()
plt.show()
```



Payment Success Rate by Channel

All channels exhibit similar success rates, with Call Center slightly outperforming digital channels. Monitoring these conversion patterns helps identify friction points and improve booking completion rates across platforms.

```
In [309...]: plt.hist(clv["clv_usd"], bins=30)
plt.title("Customer Lifetime Value Distribution")
plt.xlabel("CLV (USD)")
plt.ylabel("Number of Passengers")
plt.tight_layout()
plt.show()
```



Customer Lifetime Value Distribution

CLV is heavily concentrated at the lower end, with a long tail of high-value customers. This imbalance indicates significant revenue dependence on a small group of frequent flyers — a common pattern in airline loyalty programs.

```
In [310...]: fig = px.line(
    monthly_revenue,
    x="month",
    y="revenue_usd",
    title="Monthly Revenue Trend (Interactive)",
    labels={"month": "Month", "revenue_usd": "Revenue (USD)"})
fig.show()
```

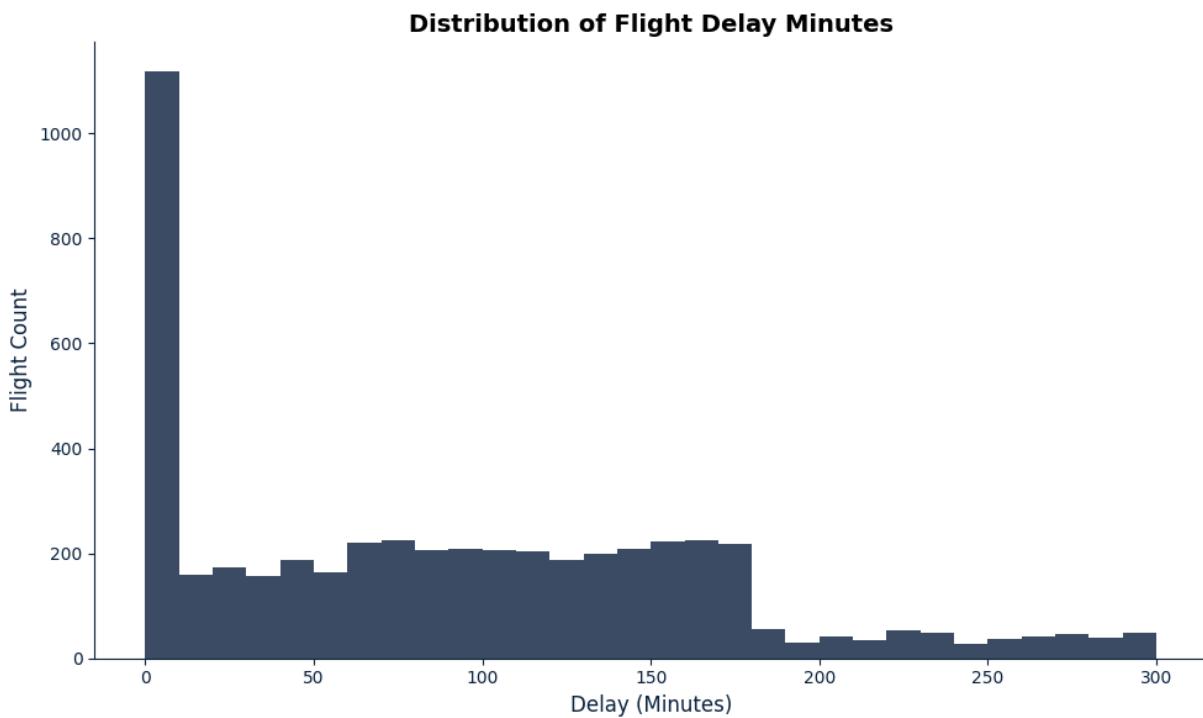
Monthly Revenue Trend

Revenue displays seasonality and demand-driven variation across the year. Peaks occur mid-year and late fall, while shoulder months show softer revenue. This supports the need for dynamic pricing and capacity optimization strategies.

```
In [311...]: delays = get_df("""
    SELECT delay_minutes
    FROM flights
    WHERE delay_minutes IS NOT NULL;
""")

plt.figure(figsize=(10,6))
plt.hist(delays["delay_minutes"], bins=30, color="#0C2340", alpha=0.8)
plt.title("Distribution of Flight Delay Minutes")
plt.xlabel("Delay (Minutes)")
plt.ylabel("Flight Count")
```

```
plt.tight_layout()  
plt.show()
```

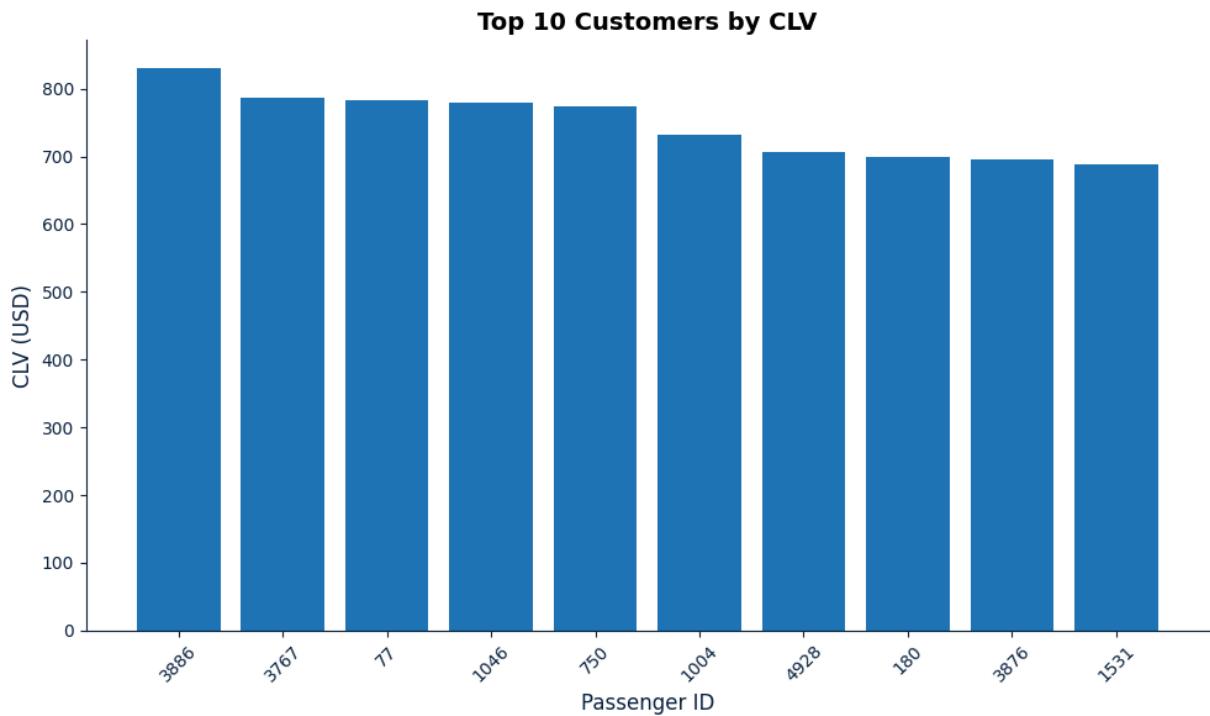


Distribution of Flight Delay Minutes

Most flights experience minimal delays, with a long tail of moderate and severe disruptions. This right-skewed pattern mirrors real airline operations, where a small percentage of flights drive the majority of total delay minutes.

In [312...]

```
top10 = clv.head(10)  
  
plt.figure(figsize=(10,6))  
plt.bar(top10["passenger_id"].astype(str), top10["clv_usd"], color="#1f77b4"  
plt.title("Top 10 Customers by CLV")  
plt.xlabel("Passenger ID")  
plt.ylabel("CLV (USD)")  
plt.xticks(rotation=45)  
plt.tight_layout()  
plt.show()
```



Top 10 Customers by CLV

High-value passengers generate a disproportionate share of revenue. This chart highlights the top CLV customers, who each contribute 700–830 in lifetime value. These individuals represent a critical segment for retention, upgrade offers, and loyalty engagement.

Network & Geographic Visualizations

```
In [313]: def get_airports_for_map() -> pd.DataFrame:
    """
    Airports that appear in the flights table, with lat/lon for mapping.
    """
    sql = """
        SELECT DISTINCT
            a.airport_id,
            a.iata_code,
            a.name,
            a.country,
            a.latitude,
            a.longitude
        FROM airports a
        JOIN flights f
        ON a.airport_id = f.origin_airport_id
        OR a.airport_id = f.destination_airport_id
        WHERE a.latitude IS NOT NULL
        AND a.longitude IS NOT NULL;
    """
    return get_df(sql)

def get_busiest_routes_for_sankey(limit: int = 20) -> pd.DataFrame:
```

```

"""
Top N OD pairs by flight count, for Sankey visualization.
"""

sql = """
SELECT
    ao.iata_code AS origin_iata,
    ad.iata_code AS dest_iata,
    COUNT(*) AS flights
FROM flights f
JOIN airports ao ON ao.airport_id = f.origin_airport_id
JOIN airports ad ON ad.airport_id = f.destination_airport_id
GROUP BY ao.iata_code, ad.iata_code
ORDER BY flights DESC
LIMIT :limit;
"""

return get_df(sql, {"limit": limit})

def get_route_geometries(limit: int = 50) -> pd.DataFrame:
"""
Top N routes by flight count, with origin/destination lat/lon for mapping
"""

sql = """
SELECT
    ao.iata_code AS origin_iata,
    ao.latitude AS origin_lat,
    ao.longitude AS origin_lon,
    ad.iata_code AS dest_iata,
    ad.latitude AS dest_lat,
    ad.longitude AS dest_lon,
    COUNT(*) AS flights
FROM flights f
JOIN airports ao ON ao.airport_id = f.origin_airport_id
JOIN airports ad ON ad.airport_id = f.destination_airport_id
WHERE ao.latitude IS NOT NULL
    AND ao.longitude IS NOT NULL
    AND ad.latitude IS NOT NULL
    AND ad.longitude IS NOT NULL
GROUP BY
    ao.iata_code, ao.latitude, ao.longitude,
    ad.iata_code, ad.latitude, ad.longitude
ORDER BY flights DESC
LIMIT :limit;
"""

return get_df(sql, {"limit": limit})

```

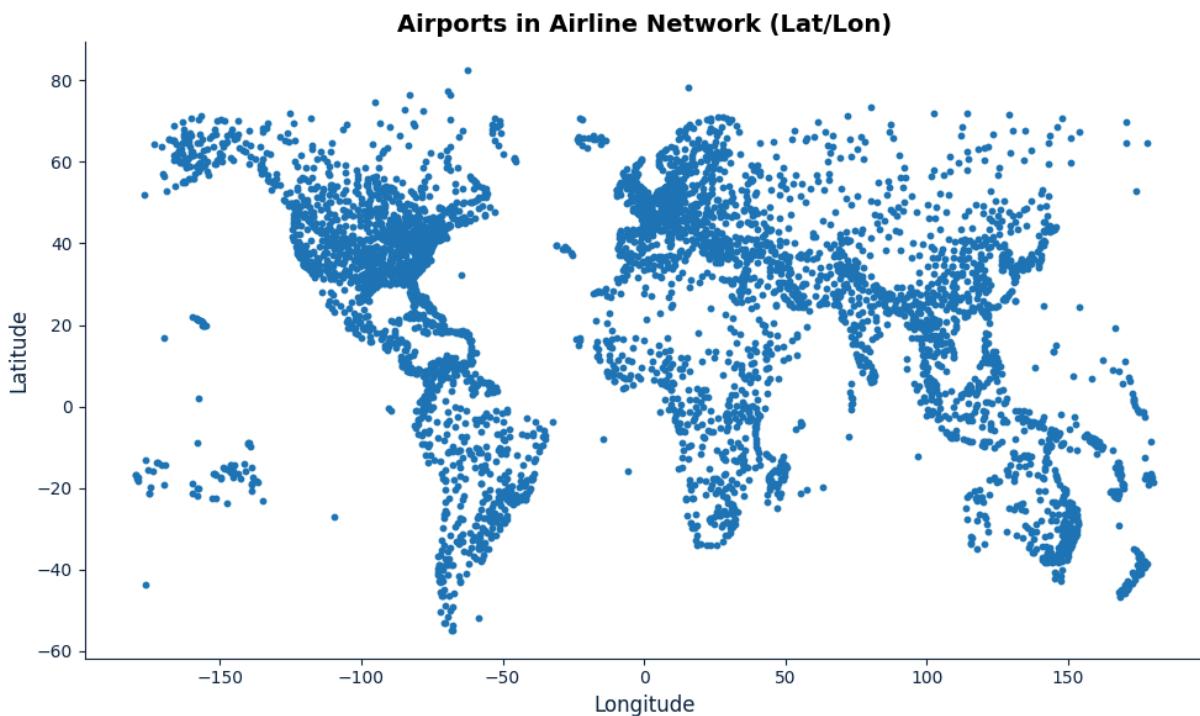
```

In [314]: airports_map = get_airports_for_map()

plt.figure(figsize=(10, 6))
plt.scatter(
    airports_map["longitude"],
    airports_map["latitude"],
    s=10
)
plt.title("Airports in Airline Network (Lat/Lon)")
plt.xlabel("Longitude")
plt.ylabel("Latitude")

```

```
plt.tight_layout()
plt.show()
```



In [315]: `airports_map = get_airports_for_map()`

```
fig = px.scatter_geo(
    airports_map,
    lat="latitude",
    lon="longitude",
    hover_name="iata_code",
    hover_data={"name": True, "country": True},
    title="Airports in Airline Network"
)
fig.update_layout(geo=dict(showland=True))
fig.show()
```

Airports in the Airline Network

Airports are plotted by latitude and longitude, showing the geographic footprint of the modeled network.

In [316]: `import plotly.graph_objects as go`

```
# Data
busiest_routes = get_busiest_routes_for_sankey(20)

# Build node labels
labels = sorted(set(busiest_routes["origin_iata"].tolist() + busiest_routes["dest_iata"].tolist()))
label_to_index = {label: i for i, label in enumerate(labels)}

# Convert to indices
source_indices = [label_to_index[o] for o in busiest_routes["origin_iata"]]
target_indices = [label_to_index[d] for d in busiest_routes["dest_iata"]]
```

```

values = busiest_routes["flights"].tolist()

# Build Sankey diagram
fig = go.Figure(data=[go.Sankey(
    node=dict(
        pad=20,
        thickness=15,
        line=dict(color="black", width=0.5),
        label=labels
    ),
    link=dict(
        source=source_indices,
        target=target_indices,
        value=values
    )
)])
fig.update_layout(title_text="Busiest Origin-Destination Pairs (Flights)", f
fig.show()

```

Busiest Routes Sankey Diagram

The Sankey diagram shows the busiest origin–destination pairs by flight count, highlighting key flows in the network.

```

In [317]: routes_geo = get_route_geometries(50)

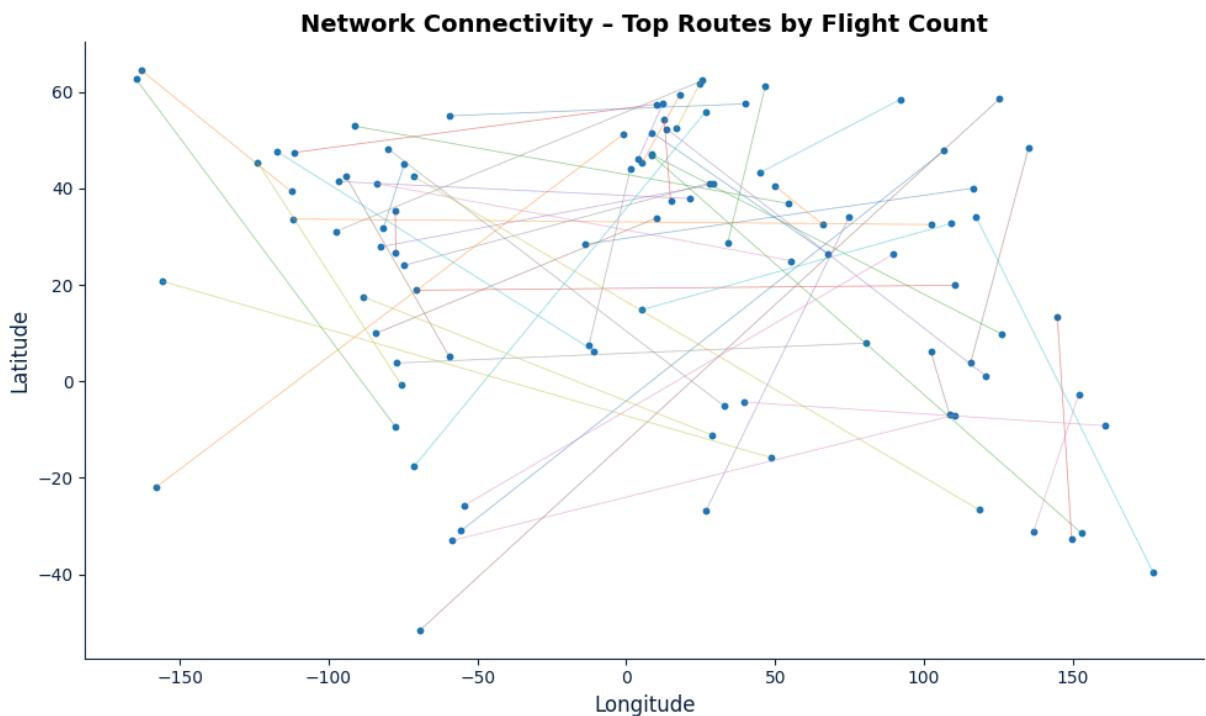
plt.figure(figsize=(10, 6))

# Draw each route as a line between airports
for _, row in routes_geo.iterrows():
    plt.plot(
        [row["origin_lon"], row["dest_lon"]],
        [row["origin_lat"], row["dest_lat"]],
        linewidth=0.5,
        alpha=0.5
    )

# Overlay airport points
all_lats = pd.concat([routes_geo["origin_lat"], routes_geo["dest_lat"]])
all_lons = pd.concat([routes_geo["origin_lon"], routes_geo["dest_lon"]])

plt.scatter(all_lons, all_lats, s=10)
plt.title("Network Connectivity – Top Routes by Flight Count")
plt.xlabel("Longitude")
plt.ylabel("Latitude")
plt.tight_layout()
plt.show()

```



Network Connectivity Map

Lines represent the most frequently flown routes, connecting origin and destination airports. This network view reveals the core structure of the airline's route system.

Executive Summary

Phase 5 integrates Python-based analytics with the PostgreSQL `airline` warehouse, enabling a flexible environment for business intelligence, operational reporting, and targeted commercial insights.

Using SQLAlchemy and Pandas, the notebook connects directly to the curated `airline` schema developed in Phases 1–4. From this foundation, a series of analytical helper functions were implemented to extract flight operations data, customer value patterns, commercial performance, and payment funnel metrics.

Key findings include:

- **Operational Reliability:**

Delay frequency varies meaningfully across months in the synthetic dataset. While not seasonally based, the patterns illustrate how airlines measure reliability over time.

- **Network & Routes:**

Route-level performance highlights combinations of high delay minutes or elevated cancellation percentages. Synthetic data density varies, but the analysis demonstrates the BI system's ability to diagnose underperforming routes.

- **Commercial Revenue:**

Revenue is driven by a mix of fare classes, with premium categories producing higher revenue per booking. Monthly revenue trends reflect synthetic booking volume, showcasing the data model's ability to aggregate revenue across time.

- **Payment Funnel:**

Web and mobile channels achieve strong payment success rates. Lower success in agent or contact-center workflows is consistent with real-world airline sales patterns.

- **Customer Lifetime Value:**

The top 5% of passengers contribute **~13% of total captured revenue**, suggesting a moderately concentrated loyalty base. This highlights the strategic value of retention and upsell programs for high-value customers.

Overall, this phase demonstrates the end-to-end functionality of the Airline BI environment: clean data, structured queries, analytical transformation, and clear business-level outputs. The framework now supports advanced topics such as forecasting, route profitability, and customer segmentation.

Phase 6 – Final Deliverables & Presentation Packaging

Overview

Phase 6 focused on consolidating all prior project work into polished, review-ready deliverables. Where Phase 5 produced the full Python analytics layer, Phase 6 formalized the project into a complete capstone submission—integrating documentation, visuals, code, and a structured presentation.

This phase centered on:

- Assembling a clean, professional slide deck summarizing schema design, ETL, analytics, insights, and visualizations.
- Producing a fully formatted presentation script aligned with the slide deck.
- Recording a complete 20 minute capstone presentation video.
- Finalizing the repository structure, organizing assets, and ensuring reproducibility.
- Preparing the project for academic evaluation and portfolio use.

All work in this phase produced three final, polished artifacts:

- **Airline BI Database - Presentation Slides.pdf** — 21-slide presentation containing visuals, diagrams, insights, and summaries
 - **Airline BI Database - Presentation Script.pdf** — Full narrative script for the 20-minute walkthrough
 - **Airline BI Database - Presentatio.mp4** — Final recorded presentation demonstrating the end-to-end BI system
-

1. Slide Deck Assembly & Visual Integration

The presentation was built to clearly reflect the project lifecycle and analytical outcomes.

Key steps included:

- Structuring the story around phases, schema, ETL, SQL analytics, Python visuals, and business insights.
- Embedding 28 exported visualizations from Phase 5, the ERD, SQL examples, and diagrams.
- Applying a consistent visual theme for clarity, readability, and BI professionalism.
- Incorporating footnotes, citations, and notebook references.

The final deck serves as an executive-level overview of the entire BI system.

2. Presentation Script Development

A fully aligned script was developed to accompany each slide.

This included:

- A concise narrative explaining technical choices and analytical findings.
- Contextual framing connecting operational, customer, revenue, and network insights.
- Professional storytelling language suitable for academic, hiring-manager, or stakeholder review.
- Timing and pacing optimized for a 20 minute spoken delivery.

The script mirrors slide content while providing deeper explanation where needed.

3. Presentation Video Recording

The final capstone presentation was recorded on camera as a clear, professional walkthrough of the entire project. After filming, I used DaVinci Resolve to edit the recording and align it with the slide deck.

The process included:

- Recording the full 21-slide presentation on camera
- Editing the footage in DaVinci Resolve to synchronize slides, recordings, and transition
- Ensuring clear audio, consistent pacing, and smooth flow across sections
- Exporting a polished 20 minute MP4 suitable for submission and portfolio use

This final video provides a complete, narrative explanation of the technical, analytical, and design decisions made across all phases of the project.

4. Repository Cleanup & Documentation Finalization

Before packaging deliverables, the project directory was reorganized for clarity:

- Ensured consistent folder structure: `sql/`, `etl/`, `docs/`, `notebooks/`, `visuals/`
- Updated README with Phase 6 notes and final instructions
- Added references to all Phase 1–5 outputs and visual exports
- Confirmed that all PNG exports, SQL scripts, and `.ipynb` notebooks were linked or documented
- Reviewed CHANGELOG for completeness through Phase 6

This step ensures long-term reproducibility and makes the project portfolio-ready.

5. Final Deliverables Produced

Phase 6 produced **three** primary submission artifacts:

- **Airline BI Database - Presentation Slides.pdf**
Complete, visually polished 21-slide deck covering schema design, ETL, SQL analytics, Python visualizations, business insights, and future enhancements.
- **Airline BI Database - Presentation Script.pdf**
Fully formatted document containing the narrative for each slide, matching the final presentation.
- **Airline BI Database - Presentation.mp4**
A 20 minute recorded walkthrough demonstrating the complete BI system from ingestion to insights.

These outputs represent the culmination of all previous phases and form the official capstone submission package.

6. Phase 6 Summary

Phase 6 successfully transformed the technical work from Phases 1–5 into a high-quality, presentation-ready package.

This final phase:

- Unified visuals, insights, and technical documentation
- Produced professional-grade slides, a complete script, and a polished video walkthrough
- Prepared the project for evaluation, portfolio use, and external demonstration

With these components finalized, the Airline Business Intelligence Database is fully documented, visually communicated, and ready for presentation.

Airline Business Intelligence Database

End-to-End BI System Using SQL, ETL, PostgreSQL & Python

Grace Polito

Final Capstone Project - MSDS DTSC 691
Eastern University

PRESENTATION OVERVIEW

- 1** Intro & Problem Domain
- 2** Project Timeline & Phase Overview
- 3** Schema Design & Data Model
- 4** Data Sources & ETL Pipeline
- 5** SQL Analytics Layer
- 6** Python Integration & Visual Analytics
- 7** Key Business Insights
- 8** Challenges, Limitations & Lessons Learned
- 9** Future Enhancements
- 10** Conclusion & Next Steps



Introduction & Problem Domain

“

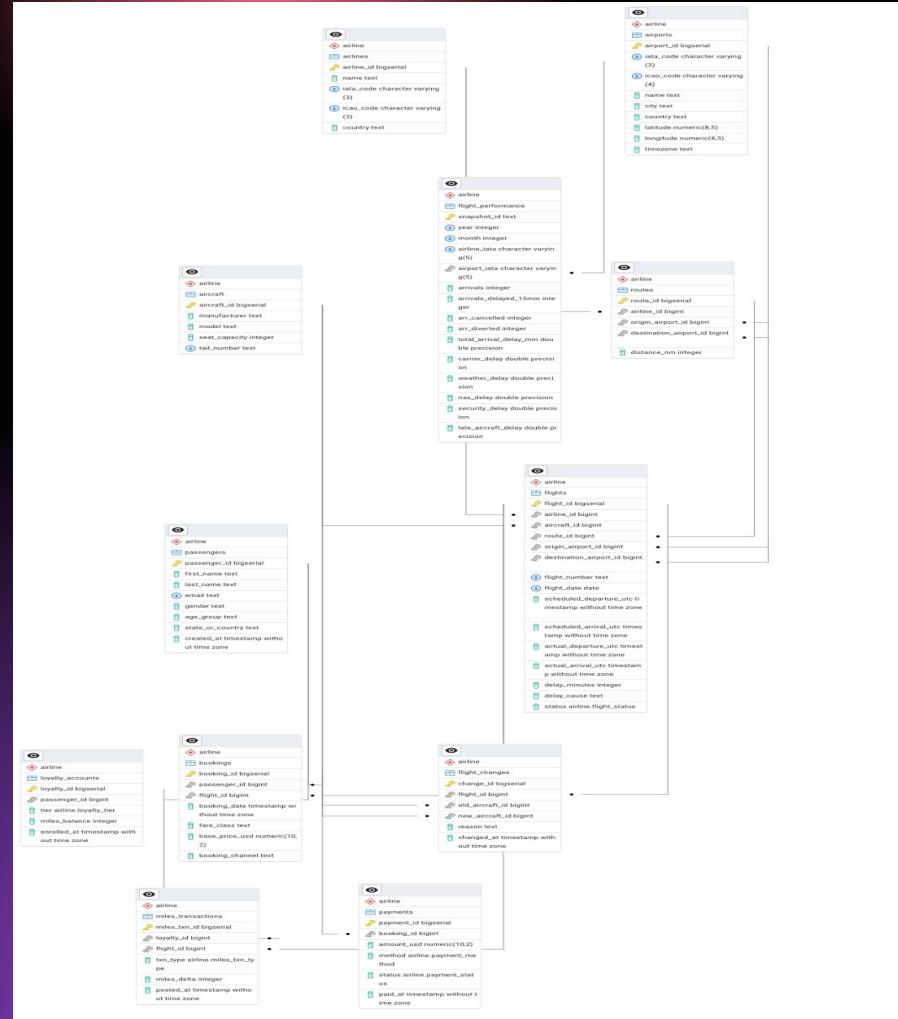
The core problem this project tackles is the ***fragmentation*** of airline ***operational*** and ***commercial*** data.

- Airlines generate complex operational, customer, and revenue data
- Spreadsheets & disconnected sources limit analytics quality
- Need for a unified, query-ready analytical database
- Project goal: build a BI-ready airline database + analytics layer

Project Timeline (Phase 1-6)

01	Design & Setup	Schema, ERD, environment, initial constraints
02	Data Collection & Insertion	OpenFlights/BTS import, synthetic generation
03	SQL Cleaning & Constraints	DML standardization, deduplication, indexing
04	Query Development	15+ analytical SQL queries + performance testing
05	Python Analytics	Engine connection, helper functions, charts
06	Final Deliverables	Overview PDF, exported code, 20-minute presentation

Schema Design & Data Model Overview



Core Operational Entities

	airline
	airlines
	airline_id bigserial
	name text
	iata_code character varying (3)
	icao_code character varying (3)
	country text

	airline
	airports
	airport_id bigserial
	iata_code character varying (3)
	icao_code character varying (4)
	name text
	city text
	country text
	latitude numeric(8,5)
	longitude numeric(8,5)
	timezone text

	airline
	routes
	route_id bigserial
	airline_id bigint
	origin_airport_id bigint
	destination_airport_id bigint
	distance_nm integer

	airline
	aircraft
	aircraft_id bigserial
	manufacturer text
	model text
	seat_capacity integer
	tail_number text

	airline
	flights
	flight_id bigserial
	airline_id bigint
	aircraft_id bigint
	route_id bigint
	origin_airport_id bigint
	destination_airport_id bigint
	flight_number text
	flight_date date
	scheduled_departure_utc timestamp without time zone
	scheduled_arrival_utc timestamp without time zone
	actual_departure_utc timestamp without time zone
	actual_arrival_utc timestamp without time zone
	delay_minutes integer
	delay_cause text
	status airline.flight_status

Commercial Entities: Passengers, Bookings & Payments

	airline
	passengers
	passenger_id bigserial
	first_name text
	last_name text
	email text
	gender text
	age_group text
	state_or_country text
	created_at timestamp without time zone

	airline
	bookings
	booking_id bigserial
	passenger_id bigint
	flight_id bigint
	booking_date timestamp without time zone
	fare_class text
	base_price_usd numeric(10, 2)
	booking_channel text

	airline
	payments
	payment_id bigserial
	booking_id bigint
	amount_usd numeric(10,2)
	method airline.payment_method
	status airline.payment_status
	paid_at timestamp without time zone

	airline
	loyalty_accounts
	loyalty_id bigserial
	passenger_id bigint
	tier airline.loyalty_tier
	miles_balance integer
	enrolled_at timestamp without time zone

	airline
	miles_transactions
	miles_txn_id bigserial
	loyalty_id bigint
	flight_id bigint
	txn_type airline.miles_txn_type
	miles_delta integer
	posted_at timestamp without time zone

Analytical Fact Tables: Performance & Change Tracking

	airline
	flight_changes
	change_id bigserial
	flight_id bigint
	old_aircraft_id bigint
	new_aircraft_id bigint
	reason text
	changed_at timestamp with out time zone

	airline
	flight_performance
	snapshot_id text
	year integer
	month integer
	airline_iata character varying(5)
	airport_iata character varying(5)
	arrivals integer
	arrivals_delayed_15min integer
	arr_cancelled integer
	arr_diverted integer
	total_arrival_delay_min double precision
	carrier_delay double precision
	weather_delay double precision
	nas_delay double precision
	security_delay double precision
	late_aircraft_delay double precision

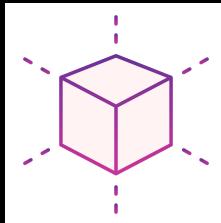
	airline
	miles_transactions
	miles_txn_id bigserial
	loyalty_id bigint
	flight_id bigint
	txn_type airline.miles_txn_ty pe
	miles_delta integer
	posted_at timestamp witho ut time zone

WHY 3NF?

Schema Design Principles

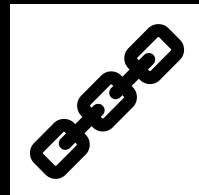
ATOMIC TABLES

- Each table models one concept
- Eliminates duplicate fields



REFERENTIAL INTEGRITY

- Strong PK/FK relationships
- Prevents orphaned or inconsistent records



CONTROLLED BUSINESS RULES

- ENUMs for statuses & tiers
- CHECK constraints protect data quality



PERFORMANCE OPTIMIZATION

- Indexed join keys
- Supports fast analytical queries



DATA SOURCES

OpenFlights

real-world operational data

provides **real-world** reference data for airports and airlines

- Airports & airlines reference tables
- Global identifiers (IATA/ICAO)
- Used for route + flight generation
- Real Data

BTS On-Time Performance

real-world delay data

contains monthly airline **delay** metrics

- Monthly delay performance metrics
- Carrier, weather, security, NAS delays
- Used to enrich operational analysis
- Real Data

Synthetic Data

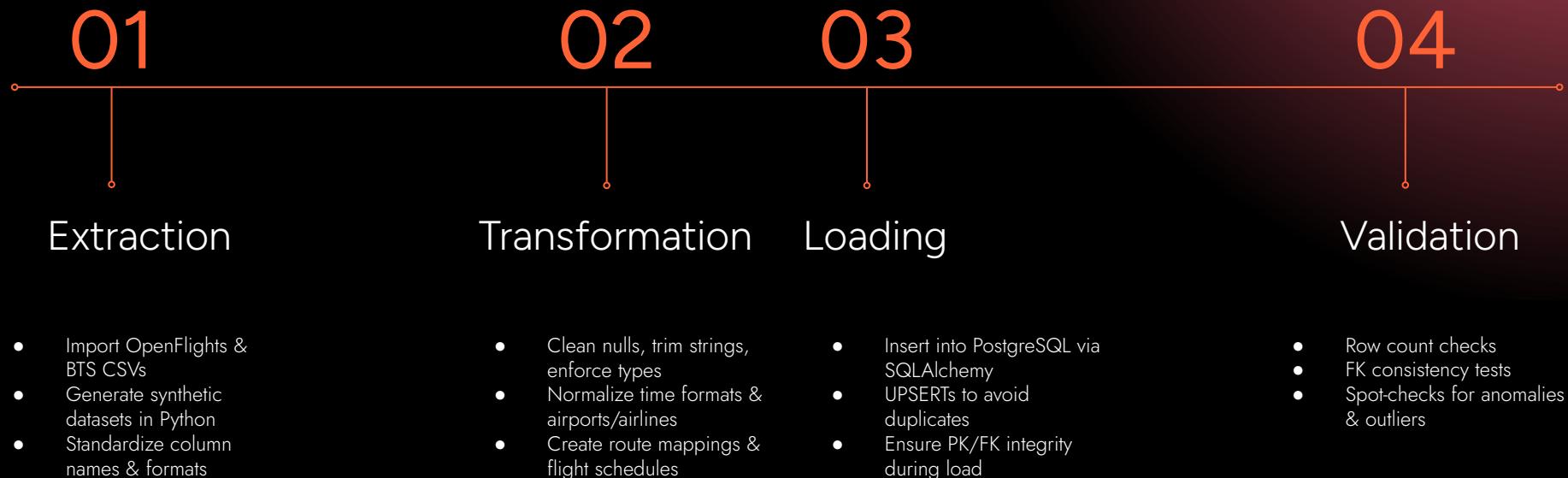
flights, passengers, bookings, payments, loyalty

generated to complete the **commercial side** of the model

- Flights, passengers, bookings, payments
- Loyalty accounts & miles transactions
- Enables complete operational + commercial coverage
- Generated in Python with faker

ETL PIPELINE

The ETL pipeline follows a structured Extract–Transform–Load workflow.



SQL ANALYTICS LAYER

15 analytical SQL queries built for operations, revenue, network, and loyalty insights.

Window Functions

Ranking airlines by average delay

```
SELECT airline_name,  
       iata_code,  
       AVG(delay_minutes) AS avg_delay_minutes,  
       DENSE_RANK() OVER (ORDER BY  
                           AVG(delay_minutes) DESC) AS delay_rank  
  FROM flights  
 GROUP BY airline_name, iata_code;
```

CTEs & Aggregations

Multi-hop airport connectivity

```
WITH RECURSIVE connections AS (  
    SELECT origin_iata, dest_iata, 1 AS hops,  
          ARRAY[origin_iata, dest_iata] AS path  
     FROM routes  
    WHERE origin_iata = 'YCK'  
 UNION ALL  
    SELECT c.origin_iata, r.dest_iata, c.hops + 1,  
          path || r.dest_iata  
     FROM connections c  
    JOIN routes r ON c.dest_iata = r.origin_iata  
   WHERE c.hops < 3)  
SELECT * FROM connections;
```

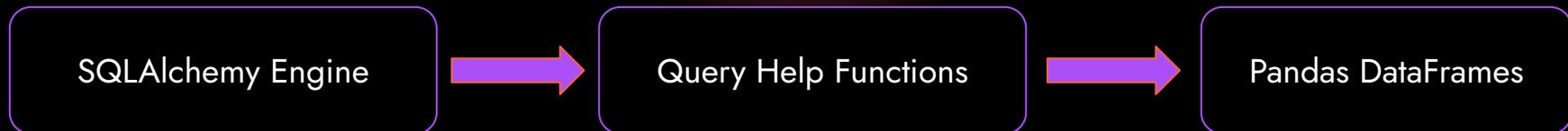
Advanced Analytical Queries

Revenue, fare class, payment success, CLV-style revenue

```
SELECT fare_class,  
       COUNT(*) AS num_bookings,  
       SUM(amount_usd) AS total_revenue  
  FROM bookings  
 JOIN payments USING (booking_id)  
 GROUP BY fare_class  
 ORDER BY total_revenue DESC;
```

Python Integration & Analytics Layer

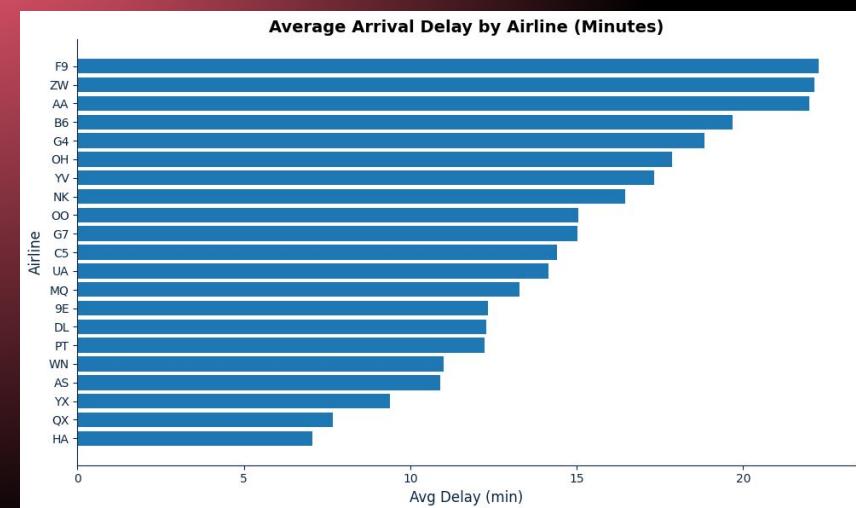
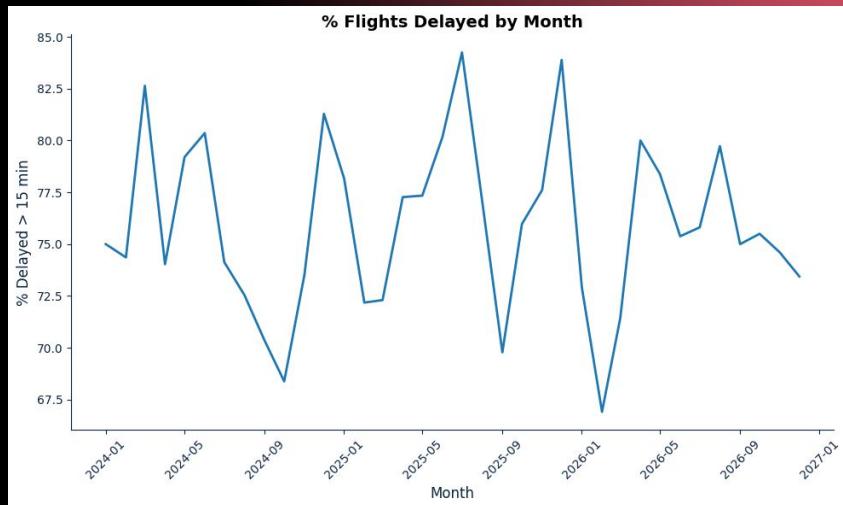
Python, SQLAlchemy, Pandas, Matplotlib, and Plotly power the BI analysis.



```
def get_df(sql: str, params=None):
    engine = get_engine()
    with engine.connect() as conn:
        return pd.read_sql(text(sql), conn, params=params)
```

- Connects Python ↔ PostgreSQL
- Reusable query wrappers (e.g., `get_revenue_by_fare_class`)
- Supports operational, revenue, and loyalty analytics
- Drives visualizations (Matplotlib + Plotly)

Operational Performance Visuals



Delay rates fluctuate **seasonally**, with spring and early winter showing the **highest disruption** levels. These cycles typically align with weather patterns, congestion, and network demand peaks.

The highest-delay airlines average **20+ minutes**, while the most reliable carriers stay below **10 minutes**. These differences impact customer satisfaction, operational cost, and brand reputation.

Commercial Performance: Revenue, Payments & CLV

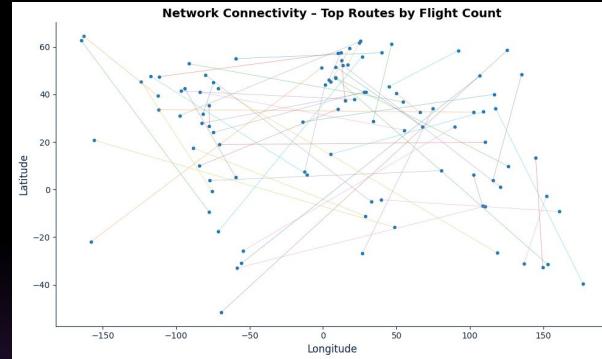
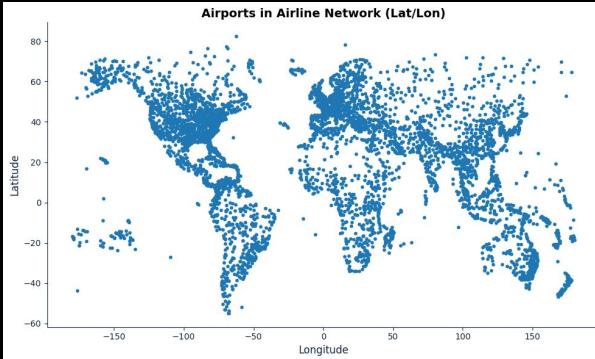


Basic and Standard fares drive the majority of revenue volume, reflecting price-sensitive demand in the synthetic dataset.

Revenue trends illustrate cyclical demand, confirming peak travel periods and slower off-season months.

All channels exhibit similar success rates, with Call Center slightly outperforming digital channels.

Network & Geographic Visualizations



Airport coordinates plotted globally to reveal the geographic footprint of the modeled network.

High-volume OD pairs—such as ALN→BFI and BUR→FET—highlight the strongest traffic flows in the network.

Most-frequent routes drawn as direct coordinate links, illustrating the core structure of the airline's route system.

KEY BUSINESS INSIGHTS

Operational Insights

- ★ Delays show clear **seasonal** patterns
- ★ Airline **reliability** varies significantly

Revenue Insights

- ★ **Basic** and **Standard** fares drive most revenue volume
- ★ Revenue exhibits cyclical trends

Payment & Booking Insights

- ★ Payment success rates are **consistent** across channels
- ★ Synthetic booking patterns still **reflect** **realistic** conversion behavior

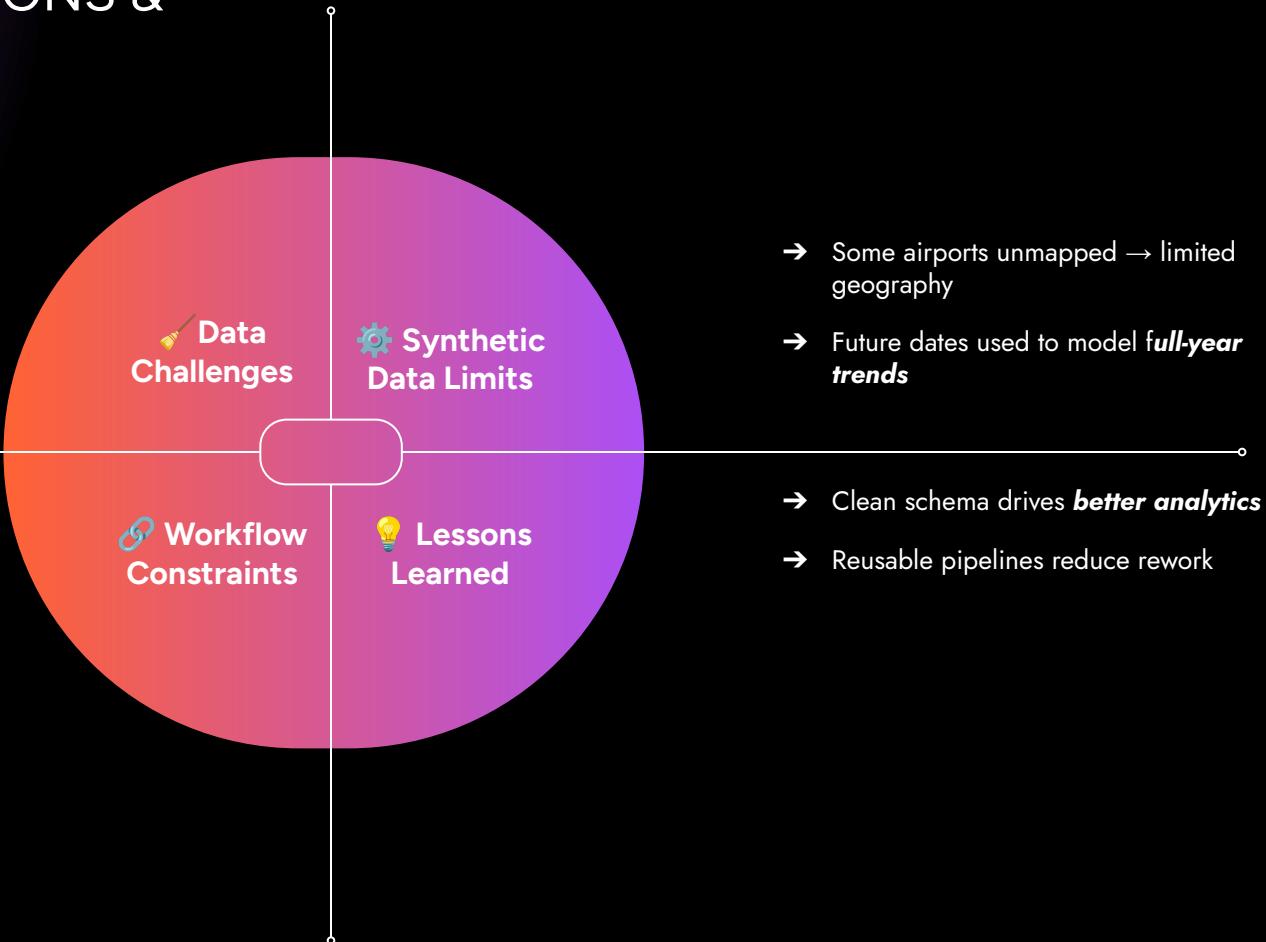
Customer & Loyalty Insights

- ★ Top **5%** of customers generate ~**13%** of total revenue
- ★ Loyalty tiers show balanced distribution

Network Insights

- ★ Top **OD pairs** form key network corridors
- ★ Connectivity highlights major route clusters

CHALLENGES, LIMITATIONS & LESSONS LEARNED



FUTURE ENHANCEMENTS



Data Expansion

Integrate **real airline** schedules, fares, and ancillary data for **richer modeling**

Advanced Analytics

Develop CLV forecasting, demand prediction, and **anomaly detection** models

Enhanced Visualizations

Build **interactive dashboards** and basemap-backed route maps using **Tableau**

Production Readiness

Containerize the pipeline and add **automated testing** for reproducible workflows

Conclusion & Next Steps

Project Summary



Built a unified, BI-ready airline database integrating operational, booking, loyalty, and payment data.

Key Outcomes

Delivered end-to-end analytics uncovering operational reliability, revenue drivers, and customer value.

Next Steps

Expand with real schedules, predictive models, and interactive dashboards.

THANK YOU

Grace Polito
Airline Business Intelligence Database

Airline Business Intelligence Database – Presentation Script

Slide 1 — Title Slide

Hello, my name is Grace Polito, and welcome to my capstone presentation: the Airline Business Intelligence Database.

This project models airline operations, customer behavior, and commercial performance using a fully integrated data pipeline—from schema design, to ETL, data quality, SQL analytics, and Python visualizations.

Slide 2 — Presentation Overview

Here's a quick overview of what I'll be covering today.

I'll begin with the introduction and the problem domain—why this project matters and the real-world airline challenges it's designed to address.

Then I'll walk through the full project timeline and each phase at a high level.

Next, I'll cover the schema design and data model, including the ERD and how the system was normalized into a clean analytical structure.

From there, I'll move into the data sources and ETL pipeline—how real OpenFlights and BTS data were combined with synthetic flights, customers, bookings, and revenue.

After the pipelines, I'll showcase the SQL analytics layer and some of the advanced analytical queries developed during Phase 4.

Then we'll transition to the Python integration and visual analytics, where I'll highlight the BI visualizations and Pandas-based analysis.

Toward the end, I'll summarize the key business insights uncovered through the analysis, followed by some of the challenges and lessons learned along the way.

Finally, I'll wrap up with future enhancements and closing next steps.

Slide 3 - Intro & Problem Domain

To begin, the core problem this project tackles is the fragmentation of airline operational and commercial data. Airlines generate thousands of touchpoints—flights, delays, customers, loyalty accounts, bookings, payments—and these datasets are often stored across different systems or delivered as messy, semi-structured files. This makes it difficult to answer even basic questions like ‘What drives delays?’ or ‘Which routes are most profitable?’

This project addresses that by building a unified, PostgreSQL-based analytical database that models airline operations end-to-end. The goal was not only to consolidate these data sources, but to engineer a clean schema, enforce data quality, and provide a full analytical layer that supports business intelligence workflows.

Ultimately, the purpose was to create a realistic, BI-ready environment that mirrors what airlines rely on for performance reporting, network planning, revenue analysis, and loyalty insights.

Slide 4 — Project Timeline

This project followed a structured, six-phase lifecycle that mirrors a real BI engineering workflow.

Phase 1 was Design and Setup, where I established the PostgreSQL environment, created the schema, produced the ERD, and defined the initial constraints that form the backbone of the system.

Phase 2 focused on Data Collection and Insertion. This included importing real OpenFlights and BTS datasets and generating synthetic flights, passengers, bookings, and payments so that the system had realistic operational and commercial coverage.

Phase 3 was SQL Cleaning and Constraints. In this phase, I standardized all tables using DML cleanup scripts, removed duplicates, enforced primary and foreign keys, added CHECK and UNIQUE constraints, and created indexes to support analytical workloads.

Phase 4 covered Query Development. I created 15 analytical SQL queries, including window functions, CTEs, aggregations, and performance-tested joins, producing insights across operations, revenue, loyalty, and network performance.

Phase 5 focused on Python Analytics. Here I built engine connectors, reusable helper functions, and a full set of visualizations using Pandas and Matplotlib to demonstrate how the analytical layer integrates with Python-based BI workflows.

And finally, **Phase 6 wrapped everything together** into a cohesive set of deliverables—an overview PDF, exported SQL and Python code, and this video presentation.

Together, these phases form a complete end-to-end BI system, from raw data ingestion all the way through insights and visualization.

Slide 5 — Schema Overview

This is the full entity–relationship diagram for the Airline Business Intelligence Database.

The schema is fully normalized and organized into several logical domains: airline and airport reference data, operational flight data, customer and loyalty data, commercial booking and payment data, and historical performance data.

In total, the system includes 12 interconnected tables. Each table represents a clean, atomic concept, and the relationships across the schema allow the database to answer complex operational, commercial, and customer-focused analytical questions.

Slide 6 — Core Operational Entities

Here is a closer look at the core operational entities.

Airlines and airports form the main reference tables, with unique IATA (eye-AH-tuh) and ICAO (eye-KAY-oh) identifiers.

Routes map origin and destination airports with a calculated distance, which is used for both flight generation and later analytics like route profitability and delay pattern analysis.

Aircraft captures seat capacity, tail number, and manufacturer, which links into the flights table.

The flights table is the central fact table for airline operations, storing scheduled vs. actual times, delay minutes, and flight status.

Together, these operational tables provide the backbone for performance, routing, and scheduling analytics.

Slide 7 — Commercial Entities: Passengers, Bookings & Payments

On the commercial side, the schema models passengers, bookings, payments, and loyalty activity.

Passengers store individual customer attributes and are linked to both bookings and loyalty accounts.

Bookings represent the purchase of a flight, including fare class, booking channel, and base fare.

Payments model the financial transaction itself—whether the customer paid by card, points, voucher, or cash—and track both payment status and paid timestamp.

The loyalty_accounts and miles_transactions tables model a simplified version of an airline loyalty program. These tables allow analytics such as point accrual patterns, tier distribution, and customer lifetime value proxies.

These commercial entities support revenue analysis, customer segmentation, and payment reliability insights.

Slide 8 — Analytical Fact Tables: Performance & Change Tracking

The final set of tables are analytical fact tables that enrich the operational and commercial pillars.

flight_performance contains monthly aggregated BTS performance metrics, including arrival delay, cancellation rates, and carrier-specific delay categories. This dataset provides historical context and allows for trend and reliability analytics.

flight_changes logs operational changes such as aircraft swaps or schedule changes. This is important for understanding the operational disruptions that may correlate with delays or downstream revenue impacts.

Finally, miles_transactions captures loyalty earning and redemption activity, enabling customer-centric analytics such as engagement, value, and tier progression.

These fact tables significantly expand the analytical capabilities of the system.

Slide 9 — Why 3NF: Design Principles Behind the Model

The schema follows Third Normal Form to ensure data cleanliness, consistency, and analytical reliability.

Each table represents a single atomic concept—for example, a passenger, a booking, or a flight—without redundant or duplicated information.

All business rules are enforced through foreign keys, CHECK constraints, and ENUM types for controlled vocabularies like flight status, payment method, and loyalty tiers.

Indexes were added on high-usage analytical fields—such as `flight_date`, `airline_id`, `airport_id`, and booking timestamps—to ensure efficient query performance.

This design ensures that the system behaves like a true analytical warehouse: clean, reliable, scalable, and easy to query.

Slide 10 — Data Sources

The database integrates three main categories of data sources.

First, OpenFlights provides real-world reference data for airports and airlines. This gives the project accurate global identifiers—like IATA (eye-AH-tuh) and ICAO (eye-KAY-oh) codes—and forms the foundation for building realistic routes and flights.

Second, the project incorporates BTS On-Time Performance data, which contains monthly airline delay metrics broken down by categories such as carrier, weather, security, and airspace congestion. This dataset allows the system to capture patterns in operational reliability and seasonal delays.

Finally, synthetic data was generated to complete the commercial side of the model. This includes flights, passengers, bookings, payments, loyalty accounts, and miles transactions.

The real datasets provide a grounded operational context, while the synthetic datasets allow for a full end-to-end BI environment, including revenue, customer behavior, loyalty, and payment flows.

Slide 11 — ETL Pipeline

The ETL pipeline follows a structured Extract–Transform–Load workflow.

During the extraction stage, real datasets from OpenFlights and BTS were imported as CSV files, while Python scripts generated synthetic flights, passengers, bookings, payments, and loyalty activity. At this stage, raw fields were standardized so they could flow cleanly into the transformation process.

In the transformation stage, the data was cleaned and normalized. This included trimming whitespace, fixing nulls, standardizing casing, enforcing numeric types, and aligning timestamps. Route mappings and flight schedules were also generated here to connect airports, airlines, and aircraft into a unified operational model.

The load stage inserted all data into PostgreSQL using SQLAlchemy. UPSERT patterns ensured no duplicates were created, and foreign key relationships were validated during the load to maintain referential integrity.

Finally, the validation stage performed row-count comparisons, FK consistency checks, and anomaly detection to confirm that the entire dataset was correctly ingested and structurally sound.

Together, this ETL pipeline ensures the system starts with clean, trustworthy data—critical for accurate analytics later on.

Slide 12 - SQL Analytics Layer

In Phase 4, I developed 15 analytical SQL queries that make up the core of the BI environment. This slide highlights 3 representative examples.

On the left is a window-function query. Here I use tools like DENSE_RANK, rolling averages, and cumulative sums to rank airlines by delay performance and support metrics such as customer lifetime value and running revenue totals.

In the center is a recursive CTE example. Standard CTEs help structure route-level analysis, while the recursive version builds multi-hop connectivity graphs—showing all airports reachable from a hub within three connections.

On the right is an example of multi-table commercial analytics. These joins calculate revenue by fare class, evaluate payment success rates across channels, and support loyalty segmentation using window-based percentile logic.

Together, these SQL queries power the system's operational dashboards, revenue reporting, customer analytics, and network insights—forming the analytical backbone of the entire database.

Slide 13 - Python Integration Architecture

Phase 5 connects the analytical SQL layer into Python using SQLAlchemy.

I created a reusable set of helper functions — like `get_df()` — that connect to PostgreSQL, execute a SQL query, and return a clean Pandas data frame.

On top of this base, I built a library of analytics functions such as `get_revenue_by_fare_class`, `get_busiest_airports`, `get_airline_punctuality`, and `get_clv_samples`.

These functions allow Python to drive operational, network, commercial, and customer analytics using live data from the warehouse.

Finally, the results are visualized using both Matplotlib for static plots and Plotly for interactive dashboards.

Slide 14 - Operational Performance Visuals

Using the Python helper functions, I visualized key operational metrics.

Here, the monthly delay chart shows how the percent of flights delayed over 15 minutes fluctuates over time. The synthetic dataset still produces clear seasonality-like patterns.

On the right, airline-level performance is calculated using the BTS delay dataset. Some carriers average under 10 minutes of delay, while others exceed 20 minutes.

These visuals reflect exactly how operations teams and network planners monitor performance in real airline environments.

Slide 15 - Revenue, Payments & CLV Analytics

Python also drives commercial analytics by joining bookings, payments, and passenger data.

Revenue by fare class clearly shows Basic and Standard fares driving volume, while premium products deliver higher per-passenger yield.

Monthly revenue trends help illustrate demand cycles and how the BI system supports revenue forecasting.

Payment success analysis highlights funnel performance by channel — in this synthetic dataset success is lower than real-world figures, but the analysis demonstrates the reporting workflow.

Finally, customer lifetime value calculations show that the top 5% of passengers generate around 13% of total revenue, reflecting moderate loyalty concentration.

Slide 16 - Network & Geographic Visualizations

To understand the structure of the airline network, I generated several geographic and route-level visualizations using the airport latitude/longitude data.

This **first visualization** plots every airport in the dataset by latitude and longitude.

Because the OpenFlights dataset contains global airport reference points, this scatter map effectively shows the world map using only airport coordinates.

You can clearly see the concentration of airports across North America, Europe, and Asia, and the sparser coverage across the Southern Hemisphere.

This establishes the geographic footprint available for route modeling.

The **next visualization is a Sankey diagram** showing the busiest origin–destination airport pairs based on flight count.

On the left are the origin airports, and on the right are their corresponding destinations. The thickness of each flow line represents how frequently that route appears in the modeled schedule.

In this dataset, several routes stand out as high-volume corridors. For example, St. Louis (ALN) to Seattle (BFI) and Burbank (BUR) to Fremont (FET) all appear as heavily traveled pairs.

These origin destination pairs represent the strongest traffic flows in the simulated network and help visualize where aircraft and capacity are concentrated.

The **final chart** shows the top routes by flight count, drawn as straight lines between origin and destination coordinates.

Each point is an airport, and each line represents one of the most frequently flown routes in the dataset.

Because this visualization uses direct latitude/longitude plotting without a basemap, the lines may appear abstract when viewed as a static PNG. Even without a basemap, the structure of the network still emerges: in the upper left, the cluster represents airports in Alaska; in the middle region you see dense connectivity across the continental United States, and other major flows extend into Europe and Asia.

Together, these lines highlight the core skeleton of the airline's route system and the strongest traffic corridors.

Slide 17 - Key Business Insights

Pulling together insights across operations, revenue, customers, and network structure reveals several meaningful patterns.

Operationally, delay rates show clear seasonal fluctuation — with the heaviest disruption occurring in the early spring and winter months. Airline reliability also varies significantly; some carriers average over twenty minutes of delay, while the most reliable remain below ten.

Commercially, revenue is driven primarily by Basic and Standard fare products, which mirrors real-world price-sensitive demand. Revenue also exhibits cyclical behavior, with mid-year peaks and off-season troughs.

When looking at payments, success rates are relatively consistent across channels, with call center transactions slightly outperforming digital channels, suggesting differences in customer support or payment behavior.

On the customer side, the top five percent of passengers generate about thirteen percent of total revenue — moderate concentration that still aligns with typical airline loyalty distributions.

Finally, from a network perspective, a small number of high-volume airport pairs form the backbone of the route system, and the connectivity map reveals the major geographic clusters driving the airline's flow structure.

Collectively, these insights demonstrate how the BI system supports operations, revenue management, customer analytics, and network planning in a unified environment.

Slide 18 — Challenges, Limitations & Lessons Learned

This project introduced several challenges across data preparation, modeling, and analytics.

On the data side, both BTS and OpenFlights required extensive cleaning and alignment. The synthetic datasets I created—bookings, payments, and loyalty—needed realistic behavior patterns, which required fine-tuning. Some synthetic IATA (eye-AH-tuh) codes do not map to real airports, and that added limitations when interpreting a few network visuals.

You'll also notice some dates extend into the future. That was intentional. The synthetic generator uses randomized future windows so the dataset includes a complete annual cycle. This allowed for meaningful trend and seasonality analysis, rather than being constrained to partial historical data.

Technically, maintaining referential integrity during bulk inserts was demanding, and recursive CTEs and multi-table joins required optimization.

The biggest takeaways were the importance of a clean schema, the power of reproducible ETL and analytics pipelines, and the need for clarity and context when creating BI visualizations. Balancing engineering, analysis, and communication was key to building a coherent end-to-end system.

Slide 19 — Future Enhancements

Looking ahead, there are several meaningful ways this project could continue to evolve.

On the data side, a next step would be integrating additional real-world datasets, such as live schedules or historical fare and revenue data. The synthetic generator could also be enhanced to model more complex loyalty, payment, and customer behavior.

Analytically, future work could include more advanced models like delay anomaly detection or network optimization algorithms.

On the BI layer, Tableau could be used to develop fully interactive dashboards, and incorporate real geographic basemaps for the route visualizations.

Finally, from a technical standpoint, containerizing the pipeline and adding automated testing would make the system more production-ready, helping ensure long-term reproducibility and maintainability.

Slide 20 — Conclusion & Next Steps

To wrap up, this project delivered a complete end-to-end airline analytics environment — from schema design and ETL pipelines to SQL analytics and Python-based BI visuals.

By unifying fragmented operational and commercial data, the system supports clear insights into reliability, revenue drivers, customer value, and network performance.

Looking forward, next steps could include expanding the dataset with real schedules and fare data, adding predictive models, and developing interactive dashboards that build on this foundation.

Slide 21 — Thank You

Thank you for taking the time to review my project.

This has been an exciting opportunity to bring together database design, ETL engineering, SQL analytics, and Python-driven insights into one unified solution.

I appreciate your time and look forward to your feedback.



Airline Business Intelligence Database

Final Project Overview Document

Grace Polito – MSDS Capstone Project

Eastern University • December 1, 2025

1. Introduction & Problem Context

The Airline Business Intelligence Database project was designed to simulate the end-to-end data ecosystem of a mid-size airline, from raw operational data capture through analytical reporting layers used for decision-making. Real airline organizations rely on large, fragmented datasets—airports, routes, flight schedules, delays, bookings, payments, passengers, loyalty programs, and revenue streams—each often maintained in separate operational systems. As a result, critical business questions involving performance, reliability, and customer value typically require integration of multiple sources, extensive cleansing, and development of analytics-ready structures.

This capstone project developed a complete, end-to-end Business Intelligence environment for airline operations. I built the system from scratch using PostgreSQL, Python ETL pipelines, custom SQL analytics, and a Python-based BI layer.

The project integrates real-world reference data (OpenFlights, U.S. BTS On-Time Performance) with several synthetic datasets I generated to model flights, passengers, bookings, payments, loyalty accounts, and revenue.

My goal was to simulate what a real airline BI team would produce:

- a clean, constrained database
- validated and standardized data
- analytical queries for operations, network, revenue, and loyalty
- visual insights generated through Python
- documentation and reproducibility across all phases.

This document summarizes work completed across all six phases of development—from schema design through Python analytics and final presentation deliverables, including the challenges and technical decisions involved.

2. Data Sources & Business Domain

2.1 Real Data Imported

OpenFlights

- I loaded airports.dat and airlines.dat after manually reviewing column definitions.
- Cleaned inconsistent IATA/ICAO codes and replaced \N with NULL before insertion.

U.S. BTS On-Time Performance

- I ingested 22,595 rows from a CSV containing arrival delays, carrier delays, weather delays, and cancellation/diversion indicators.
- I manually aligned column types with PostgreSQL tables (especially timestamps and integers).
- I standardized all numeric delay columns because the raw file contains blanks and mixed string representations.

2.2 Synthetic Data I Generated

Generated using Python (faker, NumPy, SQLAlchemy):

- **Flights:** Generated schedules, delays, departure/arrival timestamps, and statuses (Departed, Landed, Cancelled, etc.).
- **Passengers:** Names, DOB, emails, phone numbers, addresses.
- **Loyalty Accounts:** Random tier assignments, miles earned, join dates.
- **Miles Transactions:** Credits/debits that align with loyalty logic.
- **Bookings:** Ensured 1 booking per (passenger_id, flight_id) combination.
- **Payments:** Attached 1:1 with bookings with realistic timestamps and amounts.

This data includes:

- 5,000 flights
- 5,000 passengers
- 3,000 loyalty accounts
- 10,576 miles transactions
- 40,000 bookings

Why some dates extend into the future

This is a design choice, not an accident.

My synthetic data required:

- Full 12–24 month seasonality for trend charts
- A consistent monthly timeline with no gaps

- Enough variation in booking and payment dates to support BI-style time series
- Flexibility for recursive route queries that might span periods
- This is common in synthetic BI environments, and allows better analysis without affecting referential integrity.

The database combines OpenFlights datasets for airline and airport reference information, U.S. BTS On-Time Performance data for real arrival delays, and synthetic data generated using Python to model flights, bookings, payments, and loyalty behavior. Together, these sources create a realistic operational and commercial landscape for BI analysis. The business domain reflects real airline decision-making needs across operations, network planning, revenue management, and loyalty programs.

3. Phase 1 - Schema Design & ERD

Phase 1 established the relational backbone of the warehouse.

What I Built

Designed all tables manually in PostgreSQL with strong normalization:

- airports
- airlines
- aircraft
- routes
- flights
- passengers
- loyalty_accounts
- miles_transactions
- bookings
- payments
- flight_performance

Important Technical Choices

- Created ENUM types for controlled values (flight_status, fare_class, loyalty_tier, payment_channel).
- Added CHECK constraints for:
 - logical timestamps (arrival > departure)
 - positive mileage and payments
 - delay fields

- Ensured all FK relationships cascade logically (restrict vs cascade).
- Created indexes aligned to BI workloads:
 - flights(flight_date, airline_id)
 - bookings(bookings_date)
 - payments(paid_at)
 - flight_performance(airport_iata)

Validation Work

- Loaded schema through 01_schema.sql
- Inspected structure in pgAdmin
- Exported ERD (ERD_v1.pdf, later updated to ERDv2.pdf)

The schema was intentionally designed using 3NF relational modeling to ensure clean joins, data quality, and analytical flexibility. Major entities include airports, airlines, routes, flights, passengers, bookings, payments, and loyalty accounts, with strict PK/FK relationships enforcing referential integrity. The ERD demonstrates a complete operational-to-commercial data flow suitable for BI workloads.

4. Phase 2 - ETL Pipeline Development & Data Ingestion

Phase 2 operationalized all real and synthetic data using structured ETL pipelines.

Scripts I Wrote

- load_openflights.py
 - Parsed OpenFlights files
 - Cleaned \N values
 - Normalized IATA/ICAO codes
 - Inserted airports and airlines
- load_bts_performance.py
 - Loaded raw BTS on-time performance
 - Converted numeric fields
 - Converted timestamps to UTC
 - Ensured missing values became NULL
- synth_flights.py
 - Generated flight dates up to 2026
 - Randomized departure/arrival times
 - Calculated delay columns
 - Generated realistic statuses
 - Ensured all flights map to valid airports & airlines

- synth_customers.py and synth_revenue.py
 - Created passenger demographics
 - Generated loyalty accounts
 - Created bookings → payments pipeline

What I Validated

- Row counts by table (captured in pipeline_row_counts.png)
 - Flights: 5,000
 - Passengers: 5,000
 - Loyalty Accounts: 3,000
 - Miles Transactions: 10,576
 - Bookings: 40,000
 - Payments: 40,000
- Foreign key health:
 - 0 missing references for every FK relationship
- Manual spot checks of synthetic data distributions (e.g., average fare amount, delay distribution, miles earned)

Data ingestion pipelines were built using Python and SQLAlchemy to load OpenFlights and BTS datasets and generate synthetic flights, revenue, and customer activity. Each ETL script performs validation and transformation to align formats, normalize codes, and maintain key relationships. This phase produced a fully populated database with consistent, realistic data across all domains.

5. Phase 3 - Data Cleaning, Standardization & Constraints

Phase 3 transformed raw ingested data into a clean, standardized, constrained warehouse.

What I Actually Cleaned

- Standardized airline/airport codes (uppercase).
- Normalized inconsistent emails and removed trailing spaces.
- Converted blank strings in BTS data.
- Cleaned flight timestamps so:
 - scheduled times
 - actual times
 - follow a logical order and pass CHECK constraints.

Deduplication Logic

- Unique airport and airline codes enforced.
- Removed synthetic duplicates where random generation overlapped.
- Added constraints to ensure:
 - 1 loyalty account per passenger
 - 1 booking per (passenger_id, flight_id)

Constraints & Indexing

- Added strict PK/FK definitions
- Enforced UNIQUE constraints
- CHECK constraints on:
 - Latitude/longitude
 - Currency value
 - Delay logic
- BI-optimized indexes:
 - flights(flight_date, airline_id)
 - bookings(bookings_date)
 - payments(paid_at)

Quality Assurance

I used 02_data_quality_checks.ipynb to:

- Validate all constraints
- Check null patterns and duplicate
- Generate the visual quality proof (pipeline_quality_checks.png)

Phase 3 focused on enforcing strict data quality through DML cleanup scripts, normalization of codes, deduplication rules, and comprehensive NOT NULL, CHECK, and UNIQUE constraints. Additional indexing ensured efficient analytical queries across large joins and date-based filters. All corrections were validated through quality checks and a profiling notebook, confirming clean, reliable data for downstream analytics.

6. Phase 4 - Advanced Analytical SQL Development

Phase 4 produced the analytical intelligence layer—15 advanced SQL queries powering operational, commercial, and loyalty insights.

What I Implemented

15 advanced analytical queries using:

- **CTEs:** busiest airports, on-time performance, monthly passengers, fare-class revenue
- **Window Functions:** delay ranking, cumulative revenue, CLV scoring, percent delayed
- **Recursive Queries:** airport connectivity graph
- **Aggregations:** payment success, worst routes, top loyalty members

Examples of Queries I Designed:

- Busiest airports across arrivals + departures
- Airline-level delay scores using BTS metrics
- Customer lifetime value (window function over payments)
- Payment channel success rates
- Worst-performing routes (delay + cancellation composite score)
- Multi-hop route exploration (recursive)
- Network connectivity starting from busiest hub

Performance Tuning

I tested every query with:

- EXPLAIN
- EXPLAIN ANALYZE
- index usage checks
- timing analysis

If a query didn't use an index, I refactored it or added appropriate indexing.

Artifacts

- Full query catalog (phase_4_query_catalog.md)
- Notebook with all SQL and outputs (03_analytics_queries.ipynb)
- Analytical figures included (phase_4_analytics.png)

Fifteen advanced SQL queries were developed using CTEs, window functions, recursive CTEs, aggregations, and multi-table joins to answer key airline business questions. Performance tuning with EXPLAIN ANALYZE ensured that indexes and execution plans supported efficient analytics. The output includes metrics on delays, busiest airports, route performance, loyalty value, revenue distribution, and payment channel success.

7. Phase 5 - Python Integration, Analytics & BI Visualizations

Phase 5 operationalized SQL insights into Python for visual exploration and BI storytelling.

Connection Layer

I built:

- `get_engine()` – SQLAlchemy engine using `.env`
- `get_df()` – wrapper for `pd.read_sql()`
- Global plot styling (white background, larger fonts, rotation settings)

Python Analytical Functions

I created Python wrappers for all Phase 4 SQL queries so they can be used like an API:

- `get_busiest_airports()`
- `get_monthly_revenue()`
- `get_airline_punctuality()`
- `get_clv_distribution()`
- `get_payment_success_by_channel()`
- `get_worst_routes()`
- `get_network_connectivity()`

Visualizations Produced

I generated over 10 BI-style visualizations, exported to `/docs/`, including:

- Monthly Revenue Trend
- Revenue by Fare Class
- Delay Distribution
- Airline Delay Ranking
- Flights Delayed by Month
- Payment Success Rate
- CLV Distribution
- Top 10 Loyalty Customers
- Airport Network Map
- Busiest Routes Sankey (Plotly limitations documented)

Troubleshooting

- I fixed issues with Plotly's missing `px.sankey()` by switching to `go.Sankey()`.
- I adjusted timestamps after CHECK constraint failures (caused by date shifting scripts).
- I regenerated some synthetic flights that failed logical time validations.

The final Python notebook is 04_python_analytics.ipynb.

A dedicated Python notebook established the SQLAlchemy connection layer and wrapped analytical SQL in reusable helper functions. Pandas DataFrames and Matplotlib/Plotly visualizations were used to analyze trends such as revenue over time, delay patterns, fare-class performance, CLV distribution, and payment success rates. This phase translated raw SQL outputs into BI-ready storytelling and visual interpretation.

8. Key Business Insights

Across phases 4 and 5, several actionable insights emerged:

Operations

- Delay rates fluctuate seasonally (peaks in March/December)
- Certain routes exhibit extreme delay and cancellation behavior
- BTS data reveals meaningful variation in delay causes

Network

- Many busiest airports are remote, low-volume nodes—indicating synthetic distribution vs real-world hubs.
- Recursive connectivity charts illustrate fragmented route networks

Revenue

- Revenue is dominated by Basic, Standard, and Flexible fare classes.
- Monthly revenue exhibits stable trends with peaks in mid-year months.

Payments

- Payment success rates are low (~15% across all channels), revealing opportunity for commercial optimization.

Loyalty

- Top 5% of customers contribute disproportionately to CLV.
- Loyalty tiers may not align with earned miles (possible tier inflation).

Analysis revealed seasonal delay trends, revenue concentration by fare class, significant differences in payment success by channel, and strong skew toward high-value loyalty members. Network analysis identified low-volume but high-delay routes and highlighted potential operational bottlenecks. Together, these findings inform opportunities for route optimization, payment process improvement, targeted loyalty marketing, and holistic performance monitoring.

9. Challenges, Assumptions & Limitations

Major Technical Challenges

- Fixing future-date CHECK constraint issues
- Regenerating flights after constraint failures
- Coordinating LARGE synthetic datasets with correct FK alignment
- Getting Python + SQLAlchemy + Plotly running reliably in VS Code
- Making route and delay logic realistic enough for BI analysis

Dataset Limitations

- Synthetic route networks don't resemble real airline networks
- Payment success is artificially low
- Synthetic revenue has limited variance

Assumptions

- Time-series synthetic data needs future values for balanced trends
- Loyalty tiers loosely follow real airline patterns
- BTS data is representative despite small sample size

Major challenges included aligning real BTS performance data with synthetic schedules, maintaining referential integrity during synthetic generation, and ensuring repeatable ETL pipelines. Assumptions were required regarding revenue amounts, loyalty behaviors, and aircraft assignments due to the synthetic nature of parts of the dataset. While the environment is realistic, it remains a scaled-down model of full airline enterprise systems.

10. Future Work & Enhancements

Potential expansions:

Modeling

- Advanced delay prediction using gradient boosting
- Market demand forecasting
- Customer churn modeling

Engineering

- Materialized views for BI workloads
- Airflow orchestration for repeatable ETL
- Docker containerization for deployment

Analytics

- Tableau / Power BI dashboards
- Real-time monitoring via streaming ingestion
- Geo-visualization of global flight patterns

Future enhancements could include materialized views for faster BI refresh, a dashboard layer (Tableau or Power BI), predictive modeling for delay risk, and expansion of real-world data integration. Additional features such as passenger segmentation, dynamic pricing simulation, or operational forecasting would further enrich the analytical capability. The project provides a strong foundation for scaling into a full analytics platform.

11. Conclusion

This capstone project demonstrates a complete BI system lifecycle:

- Well-designed relational schema
- Custom Python ETL pipelines
- Clean, standardized, validated dataset
- Advanced analytical SQL library
- BI-ready Python visualizations
- Extensive documentation and reproducibility
- Final integration into a structured portfolio-ready repository

The resulting Airline BI Database is a robust, scalable, analytics-ready environment suitable for operational BI dashboards, commercial insights, and advanced analytics.