

**Итоговая работа  
по курсу «DWH (Data Warehouse)»  
студента группы DWH-8 (DEG-7)  
Синельниковой Наталии Леонидовны**

## Запуск БД

В работе использовался удаленный тип подключения к базе данных, используя ПО Docker Desktop.

Пустая БД была запущена при помощи следующего скрипта в командной строке:

```
docker run --name netology_DWH -p "5433:5432" -e  
POSTGRES_PASSWORD=DWH123 -e POSTGRES_USER=DWH -e POSTGRES_DB=DWH-db-  
final -d postgres:12.8
```

Далее были созданы схема dim с таблицами измерений:

- dim.calendar - справочник дат
- dim.passengers - справочник пассажиров
- dim.aircrafts - справочник самолетов
- dim.airports - справочник аэропортов
- dim.tariff - справочник тарифов

схема fact с таблицей фактов

- fact.flights - содержит совершенные перелеты.

и схема rejected

-Поля данных rejected-таблицы идентичны полям таблиц справочников, за исключением:

1. Поля не содержат PRIMARY KEY и FOREIGN KEY
2. Поля не содержат какие либо ограничения
3. Тип автоматически сформированных полей serial изменен на int4
4. Добавлено поле reason\_for\_rejection, в которое будет помещена причина отсеивания данных в rejected-таблицы.

Диаграмма базы данных представлена на рисунке 1.

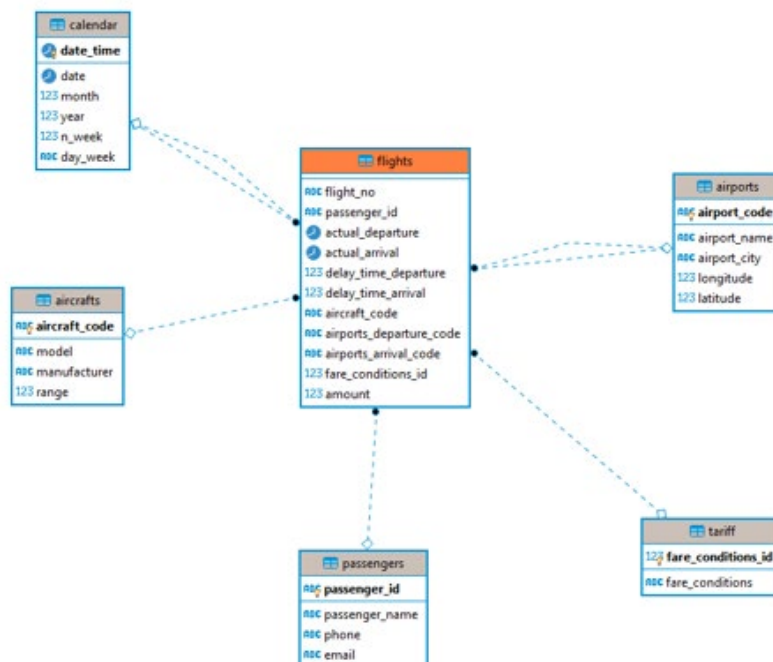


Рисунок 1 Диаграмма БД

Заполнение таблиц производилось с использованием локальной базы данных Bookings. Для написания ETL-трансформаций использовано ПО Pentaho Data Integration.

## Описание и наполнение таблиц базы данных

### Таблица dim.calendar

Таблица представляет собой справочник дат и времени совершения вылетов самолетов. Гранулярность – 1 мин. Тип SCD – 0, подразумеваем, что значения в таблице не изменяются.

Таблица создана с помощью SQL запроса

```
CREATE TABLE dim.calendar (  
    date_time timestamptz PRIMARY KEY NOT NULL, -- дата-время - ключ  
    "date" date NOT NULL, -- дата  
    "month" int4 NOT NULL, -- месяц  
    YEAR int4 NOT NULL, -- год  
    n_week int4 NOT NULL, -- неделя  
    day_week varchar(10) NOT NULL-- день недели  
);
```

Таблица содержит поля:

date\_time – дата и время в формате timestamptz, также являющаяся ключом;

date – дата в формате date;

month – номер месяца;

year – номер года;

n\_week – номер недели в году;

day\_week – день недели

Наполнение таблицы произведено при помощи SQL-запроса:

```
INSERT INTO dim.calendar(date_time, "date", "month", "year", n_week, day_week)  
SELECT gs AS date_time  
    , gs::date, date_part('month', gs)  
    , date_part('year', gs)  
    , date_part('week', gs)  
    , to_char(gs, 'day')  
FROM generate_series('2016-09-13', current_date, interval '1 minute') as gs;
```

### Таблица dim\_passengers

Таблица представляет собой справочник пассажиров. Гранулярность – один пассажир.

Таблица содержит поля:

passenger\_id – идентификатор пассажира в исходной таблице, является ключом;

passenger\_name – ФИО пассажира;

phone – контактный телефон;

email – адрес электронной почты.

Запрос SQL на создание таблицы

```
CREATE TABLE dim.passengers (  
    passenger_id varchar(20) PRIMARY KEY, -- ключ пассажира (bookings.tickets.passenger_id)  
    passenger_name text NOT NULL, -- ФИО пассажира (bookings.tickets.passenger)  
    phone varchar(20), -- телефон пассажира (bookings.tickets.contact_data ->> 'phone')  
    email varchar(150) -- email пассажира (bookings.tickets.contact_data ->> 'email')  
);
```

ETL-трансформация по наполнению и проверке данных в таблице dim.passengers представлена на рисунке 2.

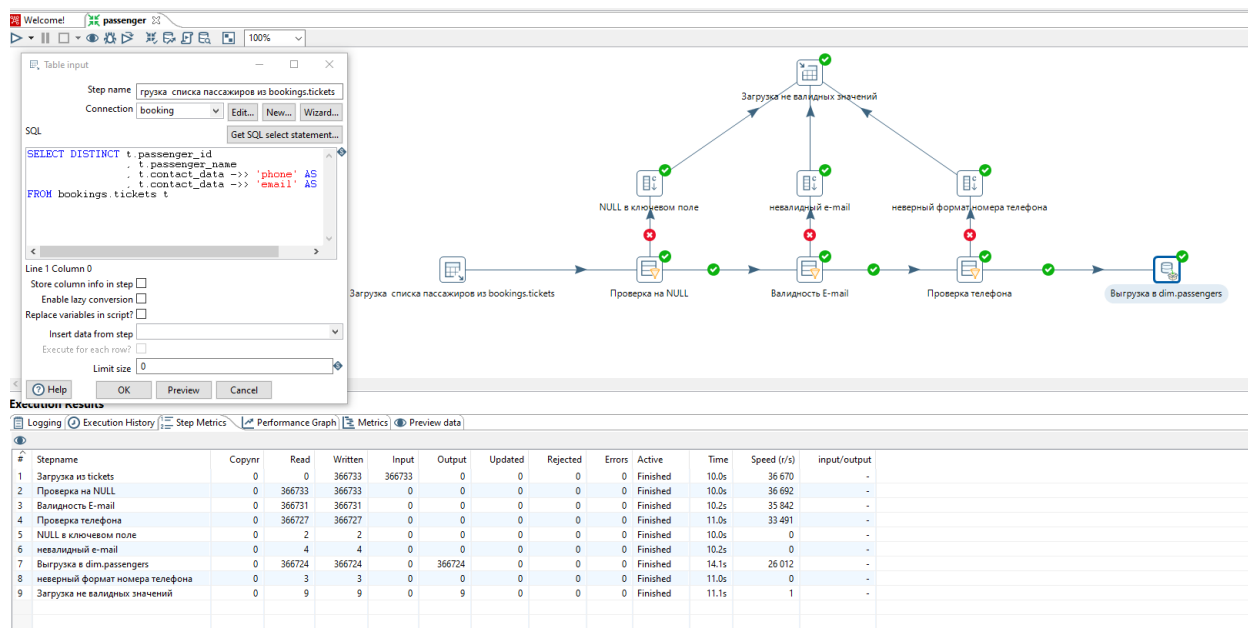


Рисунок 2 Заполнение dim.passengers

## Описание трансформации.

Первым шагом производится чтение данных из таблицы bookings.tickets.

Далее проверяем поле ключевые поля на NULL, отсутствующие значения записываем в таблицу rejected.passengers, созданную на этапе создания схемы rejected с помощью SQL скрипта

```
CREATE TABLE rejected.passengers (
  passenger_id varchar(20),
  passenger_name text,
  phone varchar(20),
  email varchar(150),
  reason_for_rejection TEXT -- поле с причиной отклонения
);
```

После проверяем формат адреса почты на корректность, а также номер телефона через регулярное выражение на отсутствие нечисловых символов и длины номера.

После всех проверок на последнем шаге отбираем нужные поля и производим запись в таблицу dim.passengers

## Таблица dim.aircrafts

Таблица представляет собой справочник самолетов

Таблица содержит поля:

- aircraft\_code – код самолета;
- model – наименование самолета;
- manufacturer - производитель
- range – дальность полета;

Запрос SQL на создание таблицы

```
CREATE TABLE dim.aircrafts (
  aircraft_code bpchar(3) PRIMARY KEY, -- ключ (bookings.aircrafts.aircraft_code)
  model varchar(20) NOT NULL, -- модель (bookings.aircrafts.model)
  manufacturer varchar(20) NOT NULL, -- производитель - первое слово в названии
  "range" int4 NOT NULL -- расстояние (bookings.aircrafts."range")
);
```

ETL-трансформация по наполнению таблицы представлена на рисунке 3.

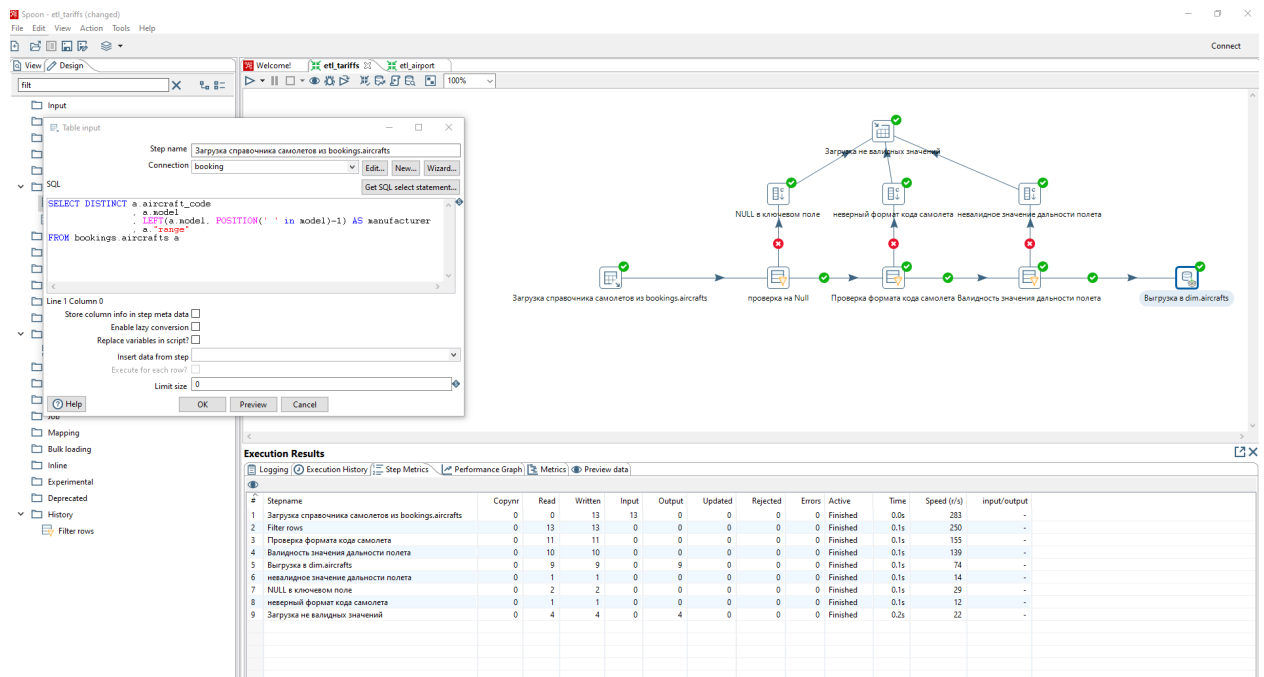


Рисунок 3 Заполнение dim.aircrafts

## Описание трансформации.

Первым шагом производится чтение данных из таблицы bookings.aircrafts. Далее проверяем поле ключевые поля на NULL, отсутствующие значения записываем в таблицу rejected.aircrafts, созданную на этапе создания схемы rejected с помощью SQL скрипта

```
CREATE TABLE rejected.aircrafts (
  aircraft_code bpchar(3),
  model varchar(20),
  manufacturer varchar(20),
  "range" int4,
  reason_for_rejection TEXT
);
```

После проверяем формат кода самолета на корректность, а также валидность дальности полета.

После всех проверок, на последнем шаге отбираем нужные поля и производим запись в таблицу dim.aircrafts

## Таблица dim.airports

Таблица представляет собой справочник самолетов

Запрос SQL на создание таблицы

```
CREATE TABLE dim.airports (
  airport_code bpchar(3) PRIMARY KEY, -- ключ (bookings.airports.airport_code)
  airport_name varchar(50) NOT NULL, -- название аэропорта (bookings.airports.airport_name)
  airport_city varchar(50) NOT NULL, -- город аэропорта (bookings.airports.city)
  longitude float(8) NOT NULL, -- долгота города (bookings.airports.longitude)
  latitude float(8) NOT NULL -- широта города (bookings.airports.latitude)
);
```

ETL-трансформация по наполнению таблицы представлена на рисунке 4.

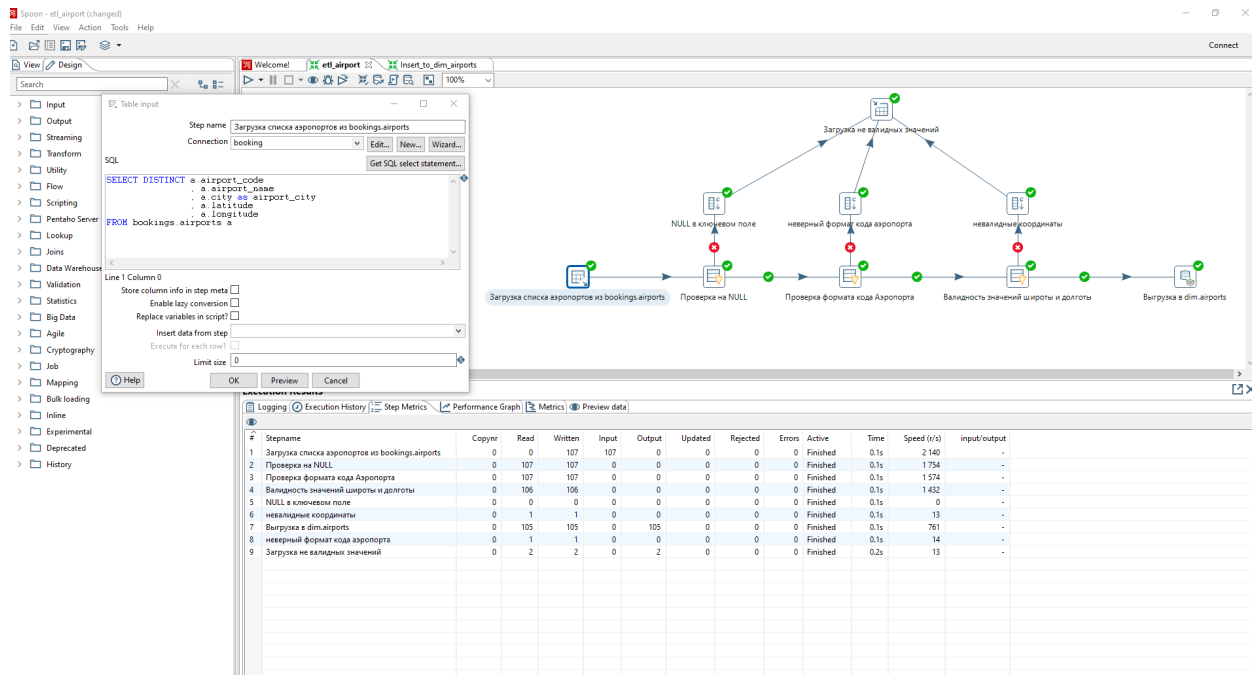


Рисунок 4. Заполнение dim.airports

## Описание трансформации.

Первым шагом производится чтение данных из таблицы bookings.airports.

Далее проверяем поле ключевые поля на NULL, отсутствующие значения записываем в таблицу rejected.airports, созданную на этапе создания схемы rejected с помощью SQL скрипта

```
CREATE TABLE rejected.airports (
  airport_code bpchar(3),
  airport_name varchar(50),
  airport_city varchar(50),
  longitude float(8),
  latitude float(8),
  reason_for_rejection TEXT -- поле с причиной отклонения
);
```

После проверяем формат кода аэропорта на корректность, а также валидность значений широты и долготы.

После всех проверок, на последнем шаге отбираем нужные поля и производим запись в таблицу dim.airports.

## Таблица dim.tariff

Таблица представляет собой справочник тарифов

Запрос SQL на создание таблицы

```
CREATE TABLE dim.tariff (  
    fare_conditions_id serial PRIMARY KEY, -- ключ тарифа - не задан, задаем сами  
    fare_conditions varchar(10) NOT NULL -- тариф (bookings.seats.fare_conditions)  
);
```

ETL-трансформация по наполнению таблицы представлена на рисунке 4.

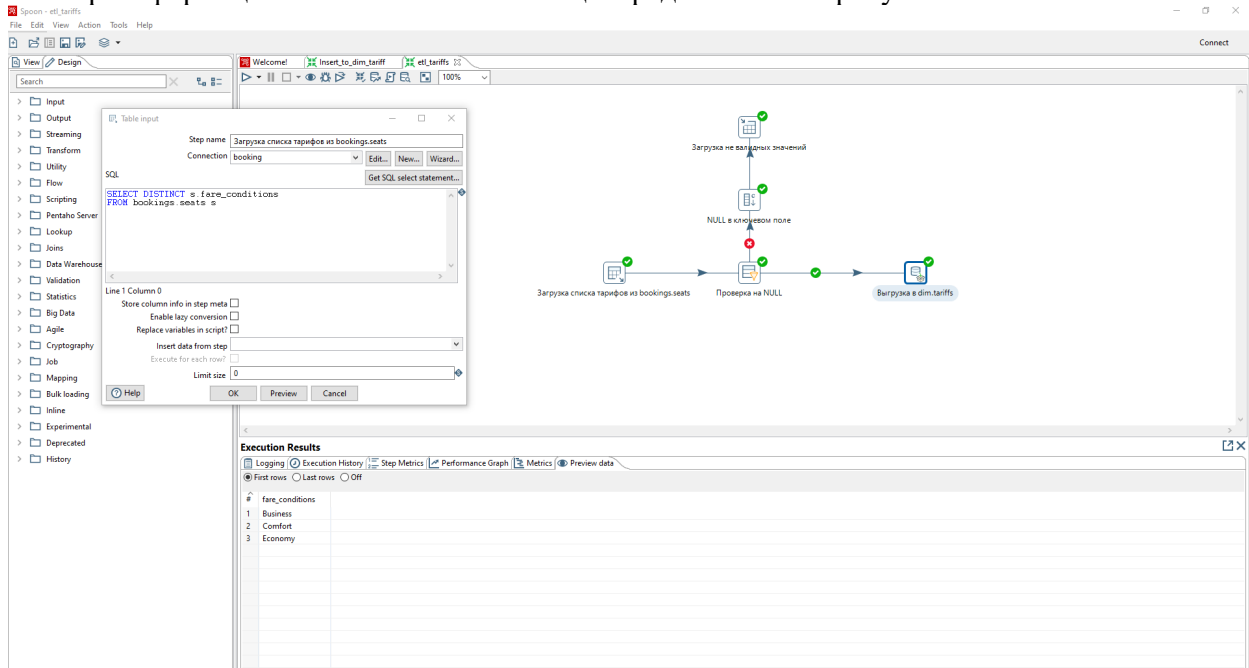


Рисунок 4. Заполнение dim.tariff

## Описание трансформации.

Первым шагом производится чтение данных из таблицы bookings.seats.

Далее проверяем поле ключевые поля на NULL, отсутствующие значения записываем в таблицу rejected.tariff, созданную на этапе создания схемы rejected с помощью SQL скрипта

```
CREATE TABLE rejected.tariff (  
    fare_conditions_id int4,  
    fare_conditions varchar(10),  
    reason_for_rejection TEXT -- поле с причиной отклонения  
);
```

## Таблица fact.flights

Таблица содержит информацию о фактически выполненных перелетах. Гранулярность – один пассажир. Таблица связана с таблицами измерений по соответствующим id по схеме «звезда»

```
CREATE TABLE fact.flights (  
    flight_no bpchar(6) NOT NULL, -- совокупность номера рейса (flight_no) и даты отправления  
    (scheduled_departure или actual_departure) являются естественным ключом, поэтому добавление  
    данного поля упростит работу конечного пользователя с таблицей фактов  
    passenger_id text REFERENCES dim.passengers(passenger_id),  
    -- Пассажир (bookings.tickets.passenger_name)  
    actual_departure timestamp NOT NULL REFERENCES dim.calendar(date_time),  
    -- Дата и время вылета (факт) (bookings.flights.actual_departure)  
    actual_arrival timestamp NOT NULL REFERENCES dim.calendar(date_time),  
    -- Дата и время прилета (факт) (bookings.flights.actual_arrival)
```

```

delay_time_departure int4 NOT NULL,
-- Задержка вылета (разница между фактической и запланированной датой в секундах)
(bookings.flights.actual_departure - bookings.flights.scheduled_departure)
delay_time_arrival int4 NOT NULL,
-- Задержка прилета (разница между фактической и запланированной датой в секундах)
(bookings.flights.actual_arrival - bookings.flights.scheduled_arrival)
aircraft_code varchar(30) NOT NULL REFERENCES dim.aircrafts(aircraft_code),
-- код самолета (bookings.flights.aircraft_code)
airports_departure_code varchar(30) NOT NULL REFERENCES dim.airports(airport_code),
-- Аэропорт вылета (bookings.flights.airports_departure)
airports_arrival_code varchar(30) NOT NULL REFERENCES dim.airports(airport_code),
-- Аэропорт прилета (bookings.flights.airports_arrival)
fare_conditions_id int REFERENCES dim.tariff(fare_conditions_id),
-- ключ класса обслуживания (dim.tariff)
amount numeric(10,2)
-- стоимость (bookings.ticket_flights.amount)

```

ETL-трансформация по наполнению таблицы представлена на рисунке 5.

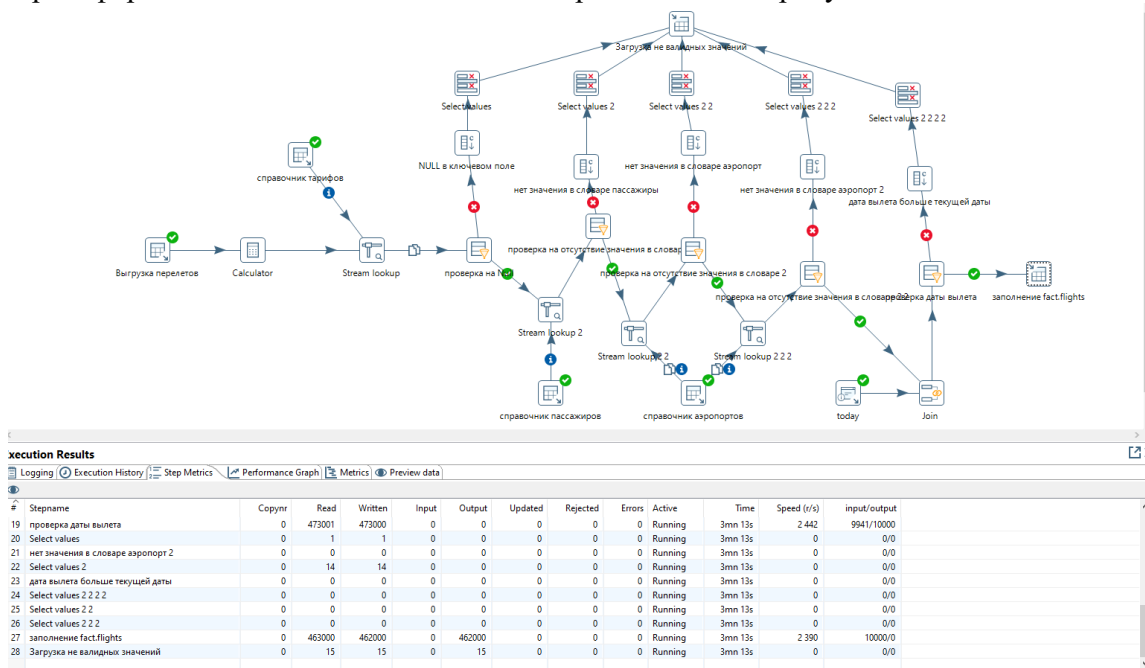


Рисунок 5. Заполнение fact.flights

## Описание трансформации.

Сперва выгружаем всю доступную информацию из БД bookings. рис 6

Table input

Step name: Выгрузка перелетов

Connection: booking

SQL

```
SELECT f.flight_no
      ,f.scheduled_departure
      ,f.scheduled_arrival
      ,f.departure_airport
      ,f.arrival_airport
      ,f.aircraft_code
      ,f.actual_departure
      ,f.actual_arrival
      ,tf.fare_conditions
      ,tf.amount
      ,t.passenger_id
FROM bookings.flights f
JOIN bookings.ticket_flights tf ON tf.flight_id = f.flight_id
JOIN bookings.tickets t ON t.ticket_no = tf.ticket_no
WHERE f.status = 'Arrived'
```

Line 1 Column 0

Store column info in step meta data ☐

Enable lazy conversion ☐

Replace variables in script? ☐

Insert data from step

Execute for each row? ☐

Limit size: 0

Help OK Preview Cancel

Рисунок 6. SQL запрос к БД источника

Далее рассчитываем время задержки вылета и прилета. Так как в БД bookings нет естественного ключа тарифов, то к значениям подтягиваем данные из справочника.

Далее проводим проверки на Null в ключевых полях и наличие значений словарей в данных таблицы, также проводится проверка даты вылета относительно текущей даты.

Отсутствующие и невалидные значения записываем в таблицу `rejected.flights`, созданную на этапе создания схемы `rejected` с помощью SQL скрипта

```
CREATE TABLE rejected.flights (
    flight_no bpchar(6),
    passenger_id text,
    actual_departure timestamp ,
    actual_arrival timestamp,
    delay_time_departure int4,
    delay_time_arrival int4,
    aircraft_code varchar(30),
    airports_departure_code varchar(30),
    airports_arrival_code varchar(30),
    fare_conditions_id int,
    amount numeric(10,2),
    reason_for_rejection TEXT -- поле с причиной отклонения
);
```

После всех проверок, на последнем шаге отбираем нужные поля, производим маппинг и производим запись в таблицу `fact.flights`

Фрагмент таблицы `fact_flights` приведён на рисунке 7.



CREATE TABLE rejected.tariff (
 fare\_conditions\_id int4,
 fare\_conditions varchar(10),
 reason\_for\_rejection TEXT -- поле с причиной отклонения
);

SELECT \* FROM fact.flights f

flight_no	passenger_id	actual_departure	actual_arrival	123 delay_time	123 delay_time	123 aircraft	123 airport	123 airport	123 amount
PG0592	3467 656288	2016-10-04 07:04:00	2016-10-04 11:00:00	7 440 000	7 500 000	SUS	TOF	SVO	29 000
PG0592	0397 273845	2016-10-04 07:04:00	2016-10-04 11:00:00	7 440 000	7 500 000	SUS	TOF	SVO	29 000
PG0592	8532 667662	2016-10-04 07:04:00	2016-10-04 11:00:00	7 440 000	7 500 000	SUS	TOF	SVO	29 000
PG0599	7562 427011	2016-10-04 14:58:00	2016-10-04 15:54:00	180 000	240 000	SUS	SVX	HMA	20 200
PG0599	5846 846245	2016-10-04 14:58:00	2016-10-04 15:54:00	180 000	240 000	SUS	SVX	HMA	6 700
PG0599	8853 133461	2016-10-04 14:58:00	2016-10-04 15:54:00	180 000	240 000	SUS	SVX	HMA	6 700
PG0600	7891 563261	2016-10-04 11:53:00	2016-10-04 12:48:00	180 000	180 000	SUS	SVX	HMA	7 400
PG0600	0128 671835	2016-10-04 11:53:00	2016-10-04 12:48:00	180 000	180 000	SUS	HMA	SVX	6 700
PG0600	0659 149538	2016-10-04 11:53:00	2016-10-04 12:48:00	180 000	180 000	SUS	HMA	SVX	6 700
PG0600	3351 517099	2016-10-04 11:53:00	2016-10-04 12:48:00	180 000	180 000	SUS	HMA	SVX	6 700
PG0600	8354 106746	2016-10-04 11:53:00	2016-10-04 12:48:00	180 000	180 000	SUS	HMA	SVX	6 700
PG0600	2160 936542	2016-10-04 11:53:00	2016-10-04 12:48:00	180 000	180 000	SUS	HMA	SVX	6 700
PG0605	5214 826036	2016-10-04 17:41:00	2016-10-04 18:36:00	60 000	60 000	SUS	DME	ULV	6 700
PG0605	5613 661636	2016-10-04 17:41:00	2016-10-04 18:36:00	60 000	60 000	SUS	DME	ULV	6 700
PG0605	4800 469873	2016-10-04 17:41:00	2016-10-04 18:36:00	60 000	60 000	SUS	DME	ULV	6 700
PG0605	7030 227429	2016-10-04 17:41:00	2016-10-04 18:36:00	60 000	60 000	SUS	DME	ULV	6 700
PG0605	3489 609017	2016-10-04 17:41:00	2016-10-04 18:36:00	60 000	60 000	SUS	DME	ULV	6 700
PG0605	2989 717492	2016-10-04 17:41:00	2016-10-04 18:36:00	60 000	60 000	SUS	DME	ULV	6 700
PG0605	7571 460873	2016-10-04 17:41:00	2016-10-04 18:36:00	60 000	60 000	SUS	DME	ULV	6 700
PG0605	4690 435740	2016-10-04 17:41:00	2016-10-04 18:36:00	60 000	60 000	SUS	DME	ULV	6 700
PG0605	7772 769960	2016-10-04 17:41:00	2016-10-04 18:36:00	60 000	60 000	SUS	DME	ULV	6 700
PG0605	3989 791770	2016-10-04 17:41:00	2016-10-04 18:36:00	60 000	60 000	SUS	DME	ULV	6 700
PG0605	5002 973113	2016-10-04 17:41:00	2016-10-04 18:36:00	60 000	60 000	SUS	DME	ULV	6 700
PG0605	9756 199022	2016-10-04 17:41:00	2016-10-04 18:36:00	60 000	60 000	SUS	DME	ULV	6 700
PG0606	7285 326407	2016-10-04 14:58:00	2016-10-04 15:53:00	180 000	180 000	SUS	ULV	DME	20 200
PG0606	6579 852994	2016-10-04 14:58:00	2016-10-04 15:53:00	180 000	180 000	SUS	ULV	DME	1 000
PG0606	9622 318535	2016-10-04 14:58:00	2016-10-04 15:53:00	180 000	180 000	SUS	ULV	DME	6 700
PG0606	3515 141793	2016-10-04 14:58:00	2016-10-04 15:53:00	180 000	180 000	SUS	ULV	DME	7 400
PG0606	0191 920575	2016-10-04 14:58:00	2016-10-04 15:53:00	180 000	180 000	SUS	ULV	DME	6 700
PG0606	8103 540710	2016-10-04 14:58:00	2016-10-04 15:53:00	180 000	180 000	SUS	ULV	DME	6 700
PG0606	8264 718739	2016-10-04 14:58:00	2016-10-04 15:53:00	180 000	180 000	SUS	ULV	DME	6 700
PG0606	2668 757001	2016-10-04 14:58:00	2016-10-04 15:53:00	180 000	180 000	SUS	ULV	DME	6 700
PG0606	0894 052639	2016-10-04 14:58:00	2016-10-04 15:53:00	180 000	180 000	SUS	ULV	DME	6 700
PG0606	0062 050016	2016-10-04 14:58:00	2016-10-04 15:53:00	180 000	180 000	SUS	ULV	DME	6 700
PG0606	1811 723159	2016-10-04 14:58:00	2016-10-04 15:53:00	180 000	180 000	SUS	ULV	DME	6 700
PG0606	6701 138915	2016-10-04 14:58:00	2016-10-04 15:53:00	180 000	180 000	SUS	ULV	DME	6 700
PG0606	3830 189040	2016-10-04 14:58:00	2016-10-04 15:53:00	180 000	180 000	SUS	ULV	DME	6 700

Рисунок 7 Фрагмент таблицы fact\_flights

## Приложение:

### SQL-скрипты создания таблиц

- dwh\_final.sql

### ETL-трансформации.

- etl\_aircrafts.ktr
- etl\_airport.ktr
- etl\_passenger.ktr
- etl\_tariffs.ktr
- etl\_flight.ktr

Ссылка на github: [https://github.com/finesun16/netology\\_dwh\\_final](https://github.com/finesun16/netology_dwh_final)