

Организация хранилища, разработка ETL- процессов и визуализация данных

Нетология, курс DataOps-инженер

Выполнил: Манохин И.В.

Цель дипломной работы:

Дипломная работа является комплексным проектом, включающим в себя анализ и обработку исходных данных, разработку и документирование хранилищ данных (NDS/DDS), создание и описание процессов ETL (Extract, Transform, Load), а также визуализацию полученных данных. Основная цель заключается в обеспечении эффективного хранения и управления данными, что способствует принятию обоснованных бизнес-решений.

Кроме основных этапов, дипломная работа также представляет возможность погрузиться в анализ данных и создание информационных систем. В процессе выполнения проекта необходимо не только освоить технические аспекты работы с данными, но и научиться видеть в них ценность для бизнеса.

Этапы выполнения дипломной работы

01

Изучение контекста
датасета, его
атрибутов

02

Статистический
анализ, поиск
ошибок

03

Формирование схем
баз данных и
создание **ER**-
диаграмм

04

Разработка **ETL**-
процессов

05

Визуализация
данных. Построение
дашбордов

06

Формирование
выводов

Используемые инструменты и ПО

<u>Анализ данных</u>	Дистрибутив Anaconda (Jupyter Lab)
<u>Разработка ER-диаграмм</u>	Сервис dbdiagram.io
<u>Базы данных</u>	ClickHouse, PostgreSQL
<u>Выполнение ETL-процессов</u>	Apache Airflow
<u>Визуализация данных</u>	Tableau Desktop

Все окружение для работы над проектом было поднято локально с использованием контейнеризации (Docker).

The background features abstract blue wavy lines on the left and bottom, and a network diagram of connected nodes and lines in the top right corner.

01

Данные

Описание датасета

Датасет Supermarket Sales представляет из себя срез исторических данных о продажах товаров в 3 филиалах компании за 3 месяца (с января по март 2019 года).

Атрибуты датасета:

1. Invoice ID: программно-генерируемый идентификационный номер счета-фактуры
2. Branch: название филиала компании
3. City: местонахождение филиала (город)
4. Customer Type: тип покупателя (наличие клубной карты)
5. Gender: пол покупателя
6. Product Line: продуктовая линейка
7. Unit Price: цена единицы товара в долларах
8. Quantity: количество проданных товаров
9. Tax: сумма взимаемого налога с продажи (5%)
10. Total: общая стоимость продажи, включая налоги
11. Date: дата продажи
12. Time: время продажи
13. Payment: метод оплаты
14. COGS: себестоимость проданных товаров
15. Gross Profit Percentage: процент прибыли
16. Gross Revenue: прибыль с продажи
17. Rating: рейтинг покупки от покупателя (по шкале от 1 до 10)

The background features abstract blue wavy lines on the left and bottom, and a network diagram of connected dots on the top right.

02

Анализ

Получение статистики и поиск ошибок в данных

Для анализа исходных данных использовались библиотеки Pandas, Matplotlib и Seaborn.

Вывод количества строк, наличие Null-значений и типов данных:

```
[47]: # проверка типов данных
df.info()

<class 'pandas.core.frame.DataFrame'>
RangeIndex: 1000 entries, 0 to 999
Data columns (total 17 columns):
#   Column                Non-Null Count  Dtype
---  -
0   Invoice ID              1000 non-null   object
1   Branch                 1000 non-null   object
2   City                   1000 non-null   object
3   Customer type          1000 non-null   object
4   Gender                 1000 non-null   object
5   Product line           1000 non-null   object
6   Unit price             1000 non-null   float64
7   Quantity               1000 non-null   int64
8   Tax 5%                 1000 non-null   float64
9   Total                  1000 non-null   float64
10  Date                   1000 non-null   object
11  Time                   1000 non-null   object
12  Payment                1000 non-null   object
13  cogs                   1000 non-null   float64
14  gross margin percentage 1000 non-null   float64
15  gross income           1000 non-null   float64
16  Rating                 1000 non-null   float64
dtypes: float64(7), int64(1), object(9)
memory usage: 132.9+ KB
```

```
memory usage: 132.9+ KB
qslbq: t[0] 1000 non-null object
10 1000 non-null object
12 1000 non-null object
14 1000 non-null object
16 1000 non-null object
```

Проверка данных на дубликаты и вывод статистики по всем столбцам:

```
[48]: # проверка строк DataFrame на дубликаты
df.duplicated().any()
```

```
[48]: False
```

```
[49]: # воспользуемся методом Pandas .describe() для генерации статистики по всем столбцам
df.describe(include='all')
```

Вывод списков категориальных переменных:

```
[50]: print(df['Branch'].unique())
print(df['City'].unique())
print(df['Customer type'].unique())
print(df['Gender'].unique())
print(df['Product line'].unique())
print(df['Payment'].unique())
```

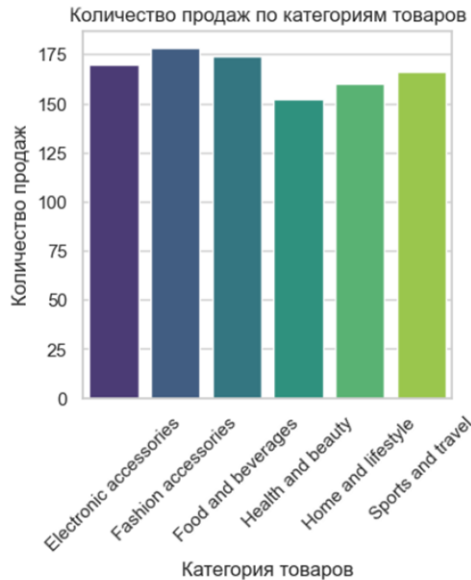
```
['A' 'C' 'B']
['Yangon' 'Naypyitaw' 'Mandalay']
['Member' 'Normal']
['Female' 'Male']
['Health and beauty' 'Electronic accessories' 'Home and lifestyle'
'Sports and travel' 'Food and beverages' 'Fashion accessories']
['Wallet' 'Cash' 'Credit card']
```

```
['Health and beauty', 'Electronic accessories', 'Home and lifestyle',
'Sports and travel', 'Food and beverages', 'Fashion accessories']
['Wallet', 'Cash', 'Credit card']
```

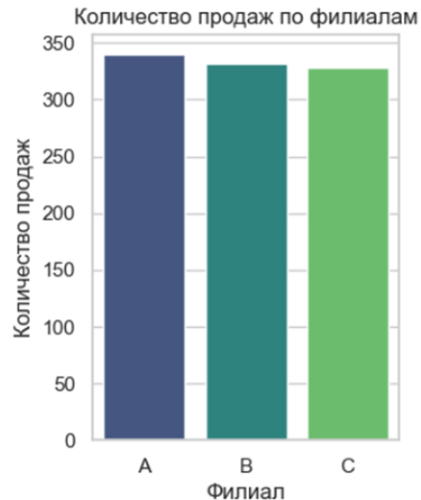

	Invoice ID	Branch	City	Customer type	Gender	Product line	Unit price	Quantity	Tax 5%	Total	Date	Time	Payment	cogs	gross margin percentage	gross income	Rating
count	1000	1000	1000	1000	1000	1000	1000.000000	1000.000000	1000.000000	1000.000000	1000	1000	1000	1000.000000	1.000000e+03	1000.000000	1000.000000
unique	1000	3	3	2	2	6	NaN	NaN	NaN	NaN	89	506	3	NaN	NaN	NaN	NaN
top	750-67-8428	A	Yangon	Member	Female	Fashion accessories	NaN	NaN	NaN	NaN	2/7/2019	19:48	Ewallet	NaN	NaN	NaN	NaN
freq	1	340	340	501	501	178	NaN	NaN	NaN	NaN	20	7	345	NaN	NaN	NaN	NaN
mean	NaN	NaN	NaN	NaN	NaN	NaN	55.672130	5.510000	15.379369	322.966749	NaN	NaN	NaN	307.58738	4.761905e+00	15.379369	6.97270
std	NaN	NaN	NaN	NaN	NaN	NaN	26.494628	2.923431	11.708825	245.885335	NaN	NaN	NaN	234.17651	6.131498e-14	11.708825	1.71858
min	NaN	NaN	NaN	NaN	NaN	NaN	10.080000	1.000000	0.508500	10.678500	NaN	NaN	NaN	10.17000	4.761905e+00	0.508500	4.00000
25%	NaN	NaN	NaN	NaN	NaN	NaN	32.875000	3.000000	5.924875	124.422375	NaN	NaN	NaN	118.49750	4.761905e+00	5.924875	5.50000
50%	NaN	NaN	NaN	NaN	NaN	NaN	55.230000	5.000000	12.088000	253.848000	NaN	NaN	NaN	241.76000	4.761905e+00	12.088000	7.00000
75%	NaN	NaN	NaN	NaN	NaN	NaN	77.935000	8.000000	22.445250	471.350250	NaN	NaN	NaN	448.90500	4.761905e+00	22.445250	8.50000
max	NaN	NaN	NaN	NaN	NaN	NaN	99.960000	10.000000	49.650000	1042.650000	NaN	NaN	NaN	993.00000	4.761905e+00	49.650000	10.00000
week	1791	1791	1791	1791	1791	1791	88.980000	10.000000	48.920000	1047.920000	1791	1791	1791	889.00000	4.761902e+00	48.920000	10.00000
1281	1791	1791	1791	1791	1791	1791	11.932000	8.000000	33.445200	41.730000	1791	1791	1791	448.90000	4.761902e+00	33.445200	8.00000
2501	1791	1791	1791	1791	1791	1791	20.730000	2.000000	15.088000	52.794000	1791	1791	1791	741.18000	4.761902e+00	15.088000	5.00000
1501	1791	1791	1791	1791	1791	1791	30.910000	3.000000	20.000000	104.400000	1791	1791	1791	118.49750	4.761902e+00	20.000000	3.00000

Информация полученная из describe:

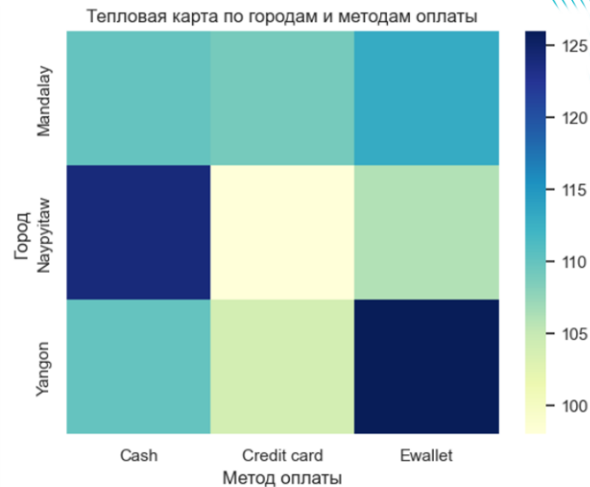
- В данных почти равное распределение как по типу клиента (Member/Normal), так и по полу (Female/Male). В обоих случаях это 501 на 499.
- Наибольшей популярностью пользуется категория Fashion accessories. На нее приходится 178 продаж.
- Средняя стоимость единицы товара \$55.67.
- Среднее количество товаров в покупке 5.51.
- Средний налог с продажи \$15.38.
- Средний доход с продажи \$322.97.
- Средняя прибыль с продажи \$15.38.
- Средний рейтинг покупки 6.97 баллов.



* Самая популярная категория Fashion Accessories, наименее популярная Health and beauty.



* Больше всего продаж приходится на филиал A, но в целом распределение почти равное.



* Метод оплаты Cash чаще всего используют в городе Naypyitaw

* Метод оплаты Credit Card наиболее часто используют в городе Mandalay

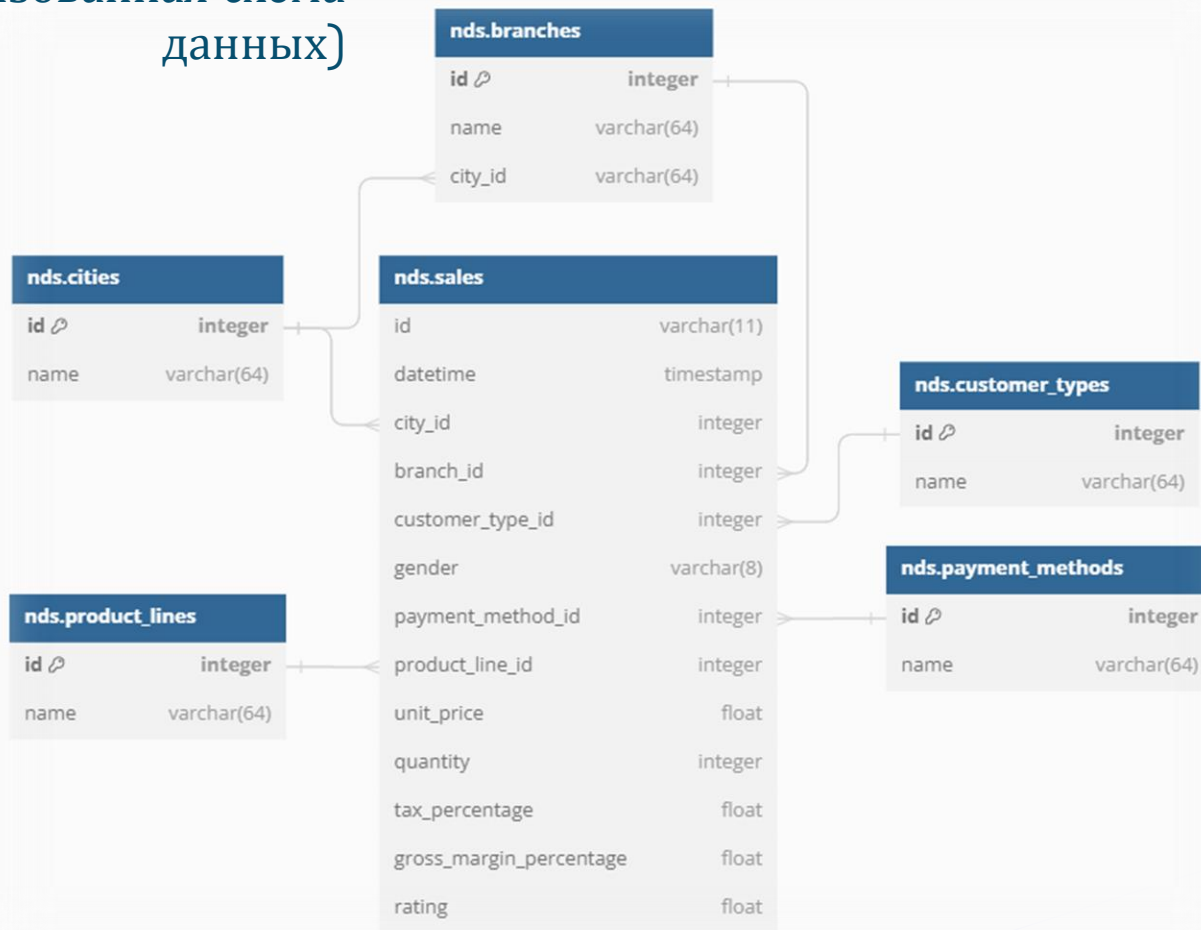
* Метод оплаты Ewallet наиболее часто используют в городе Yangon



03

Схемы данных

NDS (нормализованная схема данных)



Нормализованная схема включает в себя следующие таблицы:

1. Города (nds.cities):
 - Таблица для хранения названий городов.
2. Филиалы (nds.branches):
 - Таблица для хранения названий филиалов, зависит от таблицы nds.cities.
3. Продуктовые линейки (nds.product_lines):
 - Таблица для хранения категорий товаров.
4. Способы оплаты (nds.payment_methods):
 - Таблица для хранения доступных методов оплаты.
5. Типы клиентов (nds.customer_types):
 - Таблица для хранения различных типов клиентов.
6. Продажи (nds.sales):
 - Основная таблица для хранения информации о продажах. Отдельные столбцы для даты, времени, цены, количества и других характеристик продажи. Внешние ключи связывают записи в этой таблице с соответствующими записями в таблицах городов, филиалов, категорий товаров, способов оплаты и типов клиентов.

DDS

(таблицы фактов и измерений по схеме звезда)

dds.sales_dates	
id	integer
date	date
year	integer
month	integer
month_name	varchar(9)
day	integer
day_of_the_week	varchar(9)

dds.cities	
id	integer
name	varchar(64)

dds.customer_types	
id	integer
name	varchar(64)

dds.product_lines	
id	integer
name	varchar(64)

dds.sales	
id	varchar(11)
date_id	integer
time_id	integer
city_id	integer
branch_id	integer
customer_type_id	integer
gender	varchar(8)
payment_method_id	integer
product_line_id	integer
rating_id	integer
unit_price	float
quantity	integer
tax_amount	float
total	float
cogs	float
gross_income	float

dds.sales_time	
id	integer
time	time
hour	integer
minute	integer
second	integer

dds.branches	
id	integer
name	varchar(64)
city	varchar(64)

dds.payment_methods	
id	integer
name	varchar(64)

dds.sales_rating	
id	integer
rating	float

Схема типа «звезда» включает в себя следующие таблицы:

1. Города (dds.cities):
 - Таблица для хранения названий городов.
2. Филиалы (dds.branches):
 - Таблица для хранения информации о филиалах, включающая информацию о городе.
3. Продуктовые линейки (dds.product_lines):
 - Таблица для хранения категорий товаров.
4. Способы оплаты (dds.payment_methods):
 - Таблица для хранения доступных методов оплаты.
5. Типы клиентов (dds.customer_types):
 - Таблица для хранения различных типов клиентов.
6. Дата продажи (dds.sales_dates):
 - Таблица для хранения информации о датах продажи, включая элементы дат.
7. Время продажи (dds.sales_time):
 - Таблица для хранения информации о времени продажи, включая элементы времени.
8. Рейтинг продаж (dds.sales_rating):
 - Таблица для хранения рейтингов продаж.
9. Продажи (dds.sales):
 - Основная таблица для хранения фактов о продажах. Содержит информацию о продажах, включая дату, время, место, тип клиента, продукт, стоимость, налог и другие характеристики. Связана с таблицами измерений через внешние ключи.

Временная и историческая таблицы

unprocessed_data и processed_data

unprocessed_data	
id	varchar
datetime	datetime
city	varchar
branch	varchar
customer_type	varchar
gender	varchar
product_line	varchar
unit_price	float
quantity	integer
tax	float
total	float
payment_method	varchar
cogs	float
gross_margin_percentage	float
gross_income	float
rating	float
insert_time	datetime



dbdiagram.io



dbdiagram.io

unprocessed_data

processed_data

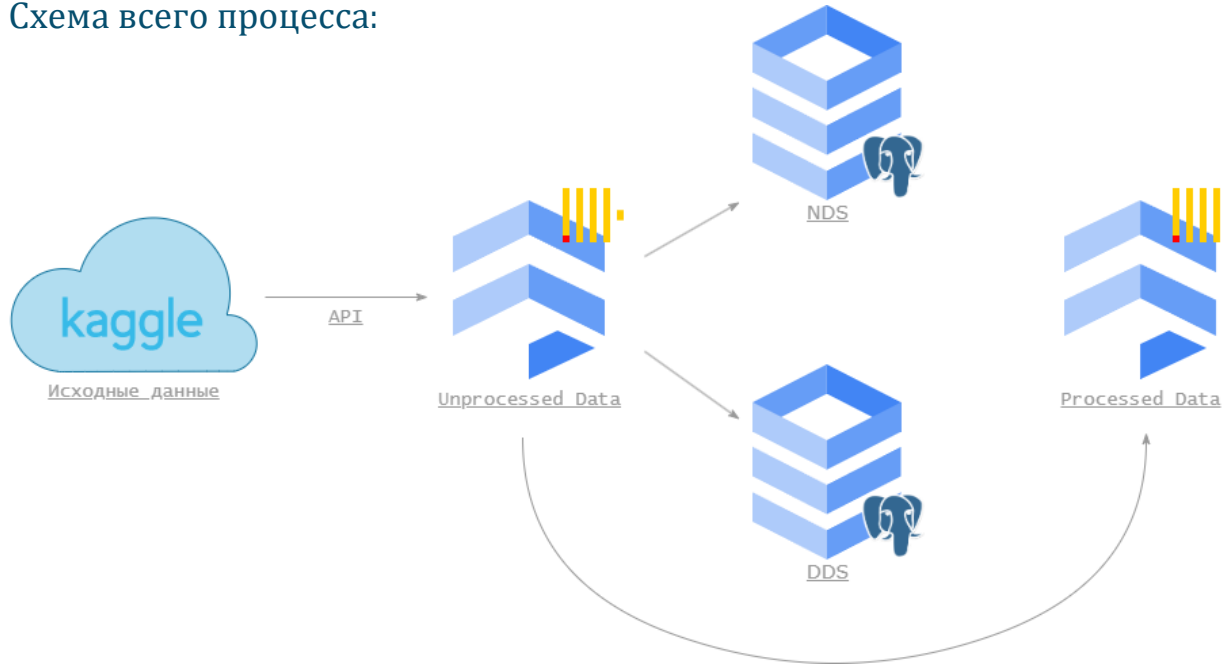


04

ETL-процессы

ETL-процессы были реализованы в 2 этапа. Первый DAG отвечает за выгрузку данных из Kaggle (через Kaggle API), второй DAG отвечает за загрузку данных в хранилище.

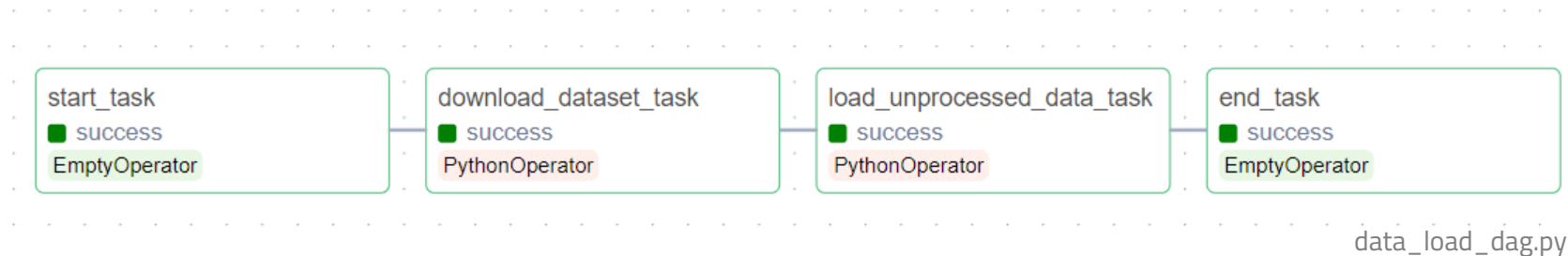
Схема всего процесса:



Выгрузка данных из источника (Kaggle) и загрузка в ClickHouse

Краткое описание: DAG выполняет выгрузку данных из источника (через Kaggle API), валидацию (проверка на Null значения, дубликаты строк, наличие валидных значений атрибутов, а также соответствие формата Invoice ID регулярному выражению), преобразование типов (раздельные атрибуты даты и времени преобразовываются в один) и загрузку данных в ClickHouse.

Граф задач:



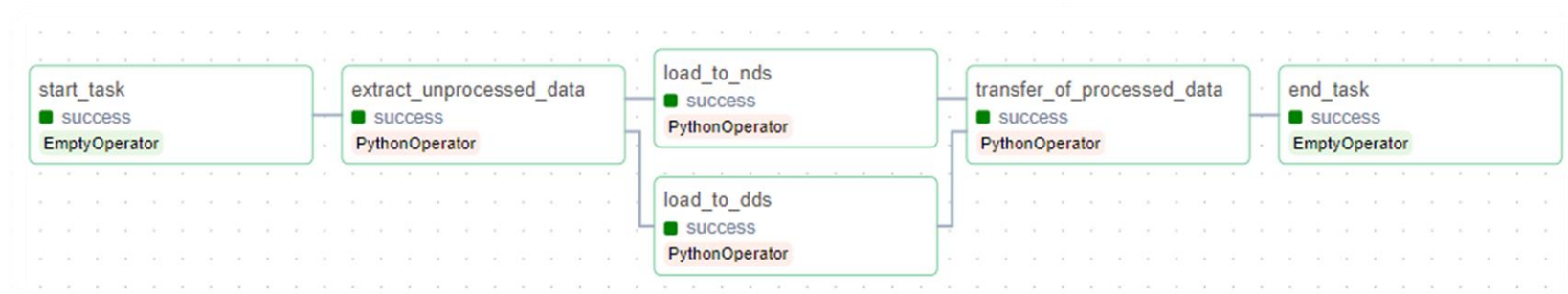
Подробное описание задач DAG:

1. Начальная задача (start_task):
 - Тип задачи: DummyOperator
 - Роль: фиктивная задача, обозначающая начало DAG.
2. Задача загрузки датасета из Kaggle (download_dataset_task):
 - Тип задачи: PythonOperator
 - Python-функция: download_dataset_from_kaggle
 - Параметры:
 - dataset_id: идентификатор датасета в Kaggle ('aungpyaeap/supermarket-sales').
 - dataset_path: путь для сохранения загруженного датасета ('./supermarket-sales').
 - Роль: загрузка датасета из Kaggle и распаковка в указанную директорию.
3. Задача загрузки необработанных данных в ClickHouse (load_unprocessed_data_task):
 - Тип задачи: PythonOperator
 - Python-функция: load_into_clickhouse
 - Параметры:
 - dataset_path: путь к загруженному датасету ('./supermarket-sales').
 - clickhouse_params: параметры подключения к ClickHouse (clickhouse_params).
 - table_name: имя таблицы в ClickHouse, куда будут загружены данные ('unprocessed_data').
 - batch_size: размер батча при загрузке данных (100).
 - Роль: обработка и валидация данных из датасета, а затем загрузка их в ClickHouse.
4. Завершающая задача (end_task):
 - Тип задачи: DummyOperator
 - Роль: фиктивная задача, обозначающая конец DAG.

Загрузка данных в хранилище

Краткое описание: DAG выполняет выгрузку данных из источника (таблица исходных данных в ClickHouse), маппинг значений (преобразование в id), загрузку данных в NDS и DDS, а также перемещение исходных данных из таблицы `unprocessed_data` в таблицу `processed_data`, очистку таблицы `unprocessed_data`.

Граф задач:



data_processing_dag.py

Подробное описание задач DAG:

1. Начальная задача (start_task):
 - Тип задачи: DummyOperator
 - Роль: фиктивная задача, обозначающая начало DAG.
2. Задача извлечения данных из ClickHouse (extract_task):
 - Тип задачи: PythonOperator
 - Python-функция: extract_unprocessed_data
 - Параметры:
 - clickhouse_params: параметры подключения к ClickHouse (clickhouse_params).
 - table_name: имя таблицы с необработанными данными ('unprocessed_data').
 - Роль: извлечение данных из ClickHouse и передача их для дальнейшей обработки.
3. Задача загрузки данных в PostgreSQL (NDS) (load_nds_task):
 - Тип задачи: PythonOperator
 - Python-функция: load_to_nds
 - Параметры:
 - postgres_params: параметры подключения к PostgreSQL (postgres_params).
 - column_map: словарь для маппинга значений (column_map).
 - schema: имя схемы в PostgreSQL ('nds').
 - Роль: загрузка обработанных данных в таблицу PostgreSQL для схемы NDS.

4. Задача загрузки данных в PostgreSQL (DDS) (load_dds_task):

- Тип задачи: PythonOperator
- Python-функция: load_to_dds
- Параметры:
 - postgres_params: параметры подключения к PostgreSQL (postgres_params).
 - column_map: словарь для маппинга значений (column_map).
 - schema: Имя схемы в PostgreSQL ('dds').
- Роль: загрузка обработанных данных в таблицы PostgreSQL для системы DDS.

5. Задача передачи обработанных данных в ClickHouse (transfer_task):

- Тип задачи: PythonOperator
- Python-функция: transfer_of_processed_data
- Параметры:
 - clickhouse_params: параметры подключения к ClickHouse (clickhouse_params).
 - table_name: имя таблицы с обработанными данными ('processed_data').
 - batch_size: размер батча при загрузке данных (100).
- Роль: загрузка обработанных данных в ClickHouse и очистка таблицы с необработанными данными.

6. Завершающая задача (end_task):

- Тип задачи: DummyOperator
- Роль: фиктивная задача, обозначающая конец DAG.

Проверка данных в хранилище

Нормализованная схема данных

```
SELECT p.id AS invoice_id,  
       p.datetime,  
       c."name" AS city,  
       b."name" AS branch,  
       ct."name" AS customer_type,  
       pm."name" AS payment_method,  
       pl."name" AS product_line,  
       p.gender,  
       p.unit_price,  
       p.quantity,  
       p.unit_price * p.quantity AS cogs,  
       p.tax_percentage,  
       p.unit_price * p.quantity / 100 * p.tax_percentage AS tax_amount,  
       (p.unit_price * p.quantity) + (p.unit_price * p.quantity / 100 *  
p.tax_percentage) AS total,  
       p.gross_margin_percentage,  
       ((p.unit_price * p.quantity) + (p.unit_price * p.quantity / 100 *  
p.tax_percentage)) / 100 * p.gross_margin_percentage AS gross_income,  
       p.rating  
FROM nds.sales p  
LEFT JOIN nds.cities c ON p.city_id = c.id  
LEFT JOIN nds.branches b ON p.branch_id = b.id  
LEFT JOIN nds.customer_types ct ON p.customer_type_id = ct.id  
LEFT JOIN nds.payment_methods pm ON p.payment_method_id = pm.id  
LEFT JOIN nds.product_lines pl ON p.product_line_id = pl.id  
WHERE p.id = '761-49-0439';
```



Name	Value
invoice_id	761-49-0439
datetime	2019-01-19 10:17:00.000
city	Mandalay
branch	B
customer_type	Member
payment_method	Ewallet
product_line	Electronic accessories
gender	Female
unit_price	12.1
quantity	8
cogs	96.8
tax_percentage	5.0
tax_amount	4.84
total	101.64
gross_margin_percentage	4.761905
gross_income	4.8400002419999995
rating	8.6

rating	8.6
gross_income	4.8400002419999995
gross_margin_percentage	4.761905
total	101.64
tax_amount	4.84
tax_percentage	5.0

Схема данных типа «звезда»

```
SELECT p.id AS invoice_id,  
       sd."date",  
       sd."year",  
       sd."month",  
       sd."day",  
       sd."day_of_the_week",  
       st.datetime::time AS "time",  
       c."name" AS city,  
       b."name" AS branch,  
       ct."name" AS customer_type,  
       pm."name" AS payment_method,  
       pl."name" AS product_line,  
       sr.rating,  
       p.gender,  
       p.unit_price,  
       p.quantity,  
       p.tax_amount,  
       p.total,  
       p.cogs,  
       p.gross_income  
FROM dds.sales p  
LEFT JOIN dds.sales_dates sd ON p.date_id = sd.id  
LEFT JOIN dds.sales_time st ON p.time_id = st.id  
LEFT JOIN dds.cities c ON p.city_id = c.id  
LEFT JOIN dds.branches b ON p.branch_id = b.id  
LEFT JOIN dds.customer_types ct ON p.customer_type_id = ct.id  
LEFT JOIN dds.payment_methods pm ON p.payment_method_id = pm.id  
LEFT JOIN dds.product_lines pl ON p.product_line_id = pl.id  
LEFT JOIN dds.sales_rating sr ON p.rating_id = sr.id  
WHERE p.id = '253-12-6086';
```



Name	Value
invoice_id	253-12-6086
date	2019-03-12
year	2019
month	3
day	12
day_of_the_week	Tuesday
time	12:43:00
city	Yangon
branch	A
customer_type	Member
payment_method	Credit card
product_line	Sports and travel
rating	8.7
gender	Female
unit_price	98.4
quantity	7
tax_amount	34.44
total	723.24
cogs	688.8
gross_income	34.44

gross_income	34.44
cogs	688.8
total	723.24
tax_amount	34.44
quantity	7
unit_price	98.4



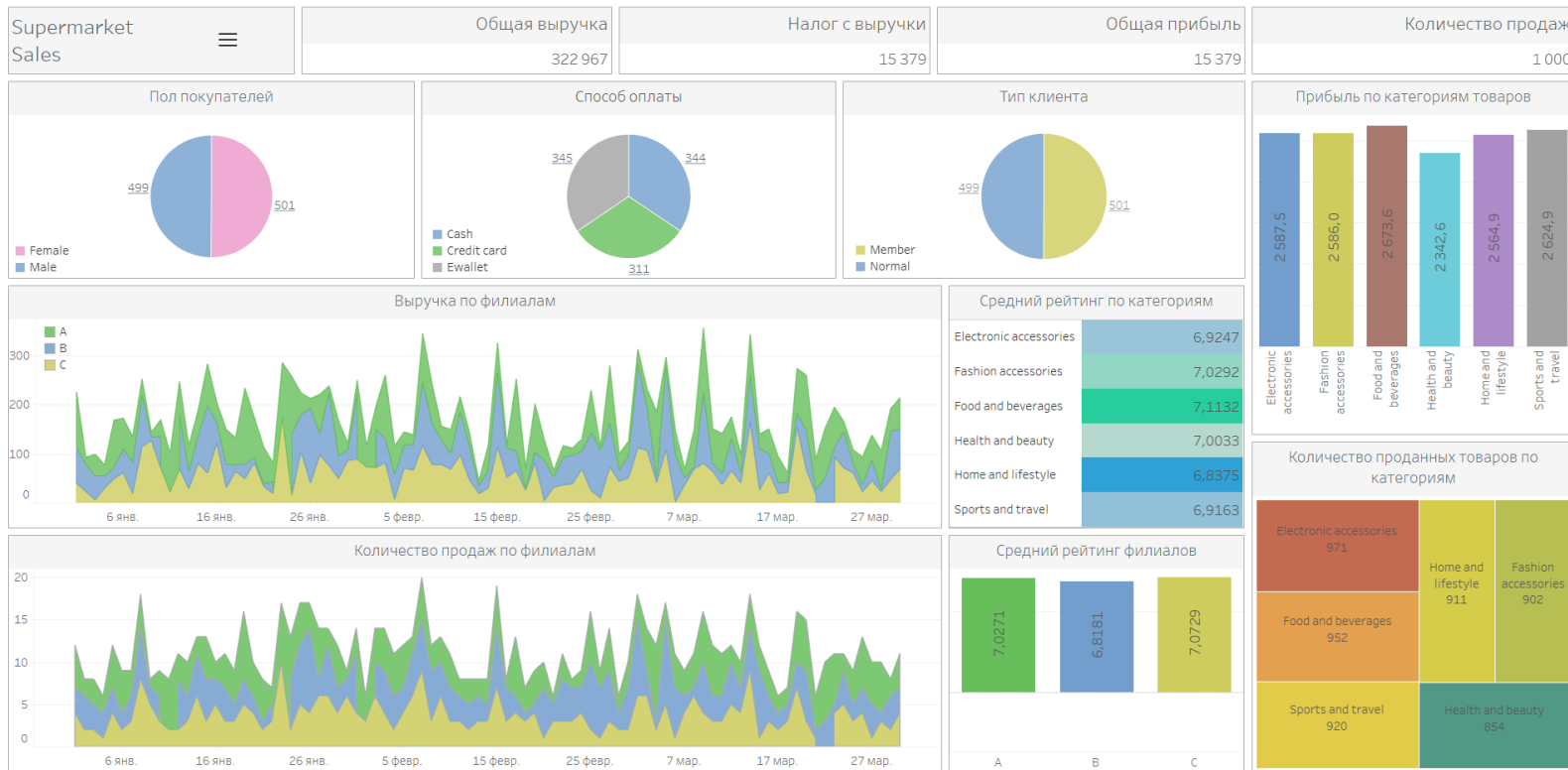
05

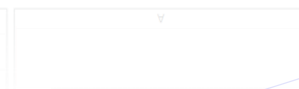
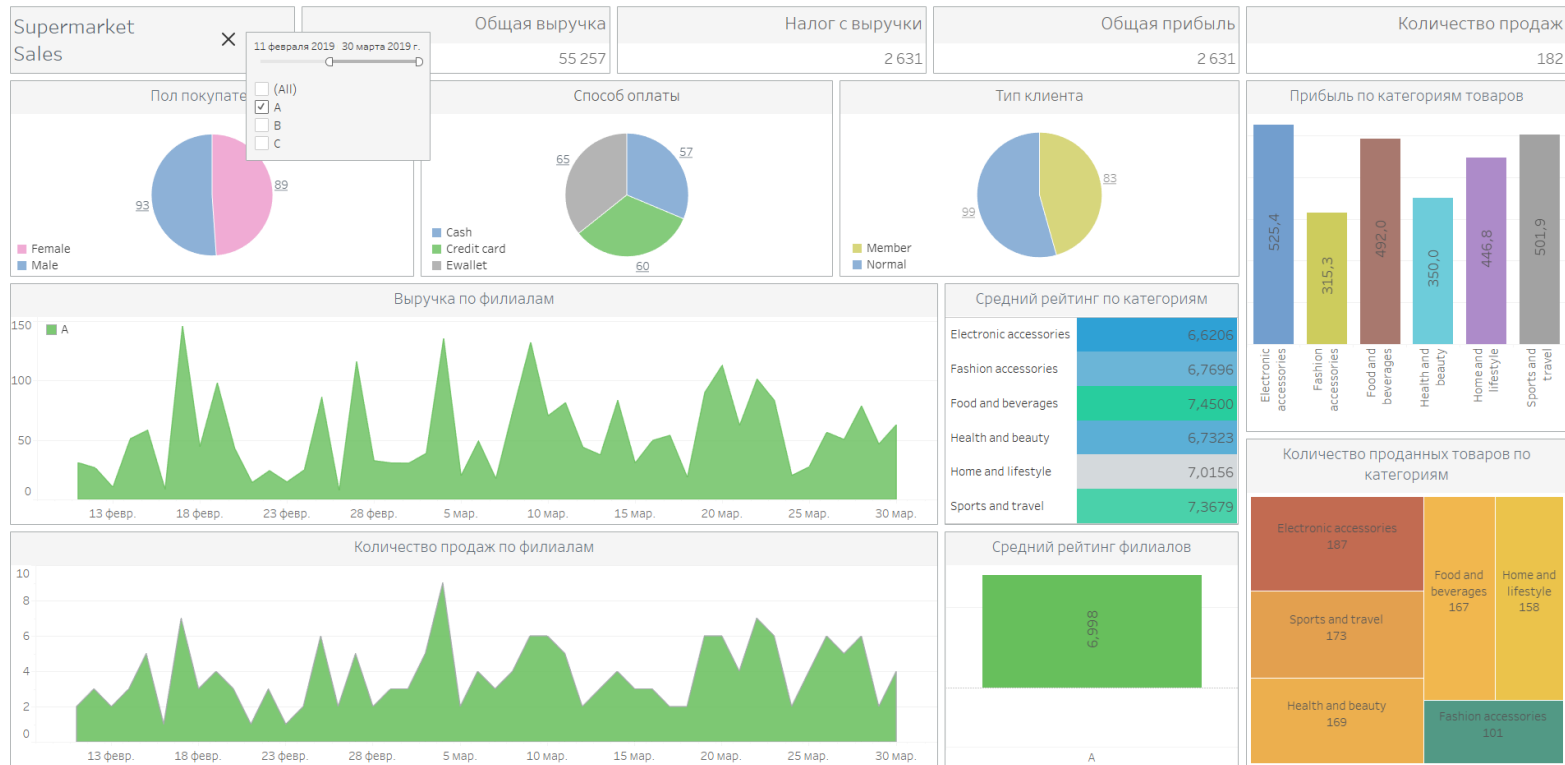
Визуализация данных

Выбранные метрики:

- **Общая выручка компании:**
 - Цель: оценить общую финансовую производительность компании.
- **Налог с выручки:**
 - Цель: обеспечить контроль за налоговой обязанностью компании и ее финансовой устойчивостью.
- **Общая прибыль:**
 - Цель: измерить общую прибыль и оценить финансовую эффективность бизнеса.
- **Количество выполненных продаж:**
 - Цель: оценить активность продаж и популярность товаров или услуг.
- **Пол покупателей:**
 - Цель: анализ гендерного распределения клиентов.
- **Способ оплаты покупки:**
 - Цель: изучение предпочтений клиентов по методам оплаты.
- **Тип клиента (наличие клубной карты):**
 - Цель: анализ клиентской базы и их участие в программе лояльности.
- **Прибыль по категориям товаров:**
 - Цель: оценка прибыльности различных категорий товаров.
- **Выручка по филиалам:**
 - Цель: измерение доли выручки, генерируемой каждым филиалом.
- **Средний рейтинг от клиентов по категориям товаров:**
 - Цель: оценка удовлетворенности клиентов различными категориями товаров.
- **Количество продаж по филиалам:**
 - Цель: измерение активности клиентов в различных филиалах.
- **Средний рейтинг филиалов:**
 - Цель: оценка удовлетворенности клиентов работой различных филиалов.
- **Количество проданных товаров по категориям:**
 - Цель: измерение популярности конкретных категорий товаров.

Итоговый дашборд





—

19 331

920,5

920,5

64

A pie chart showing the distribution of responses for 'How often do you use the Internet?'. The chart is divided into two segments: a pink segment labeled '25' and a blue segment labeled '49'.



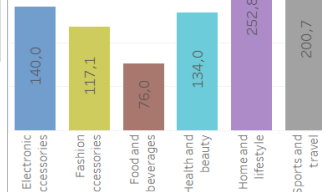
Number of Children	Percentage
1 child	64%
2 children	30%
3 children	6%



A pie chart with a blue sector and a yellow sector. The blue sector is labeled with the number 64.



Category	Value
Electronic accessories	140,0
Fashion accessories	117,1
Food and beverages	75,0
Health and beauty	134,0
Home and lifestyle	252,8
Sports and travel	200,7



График, иллюстрирующий динамику индекса В (ВВП) в России с 8 января по 29 марта 2020 года. Показывает значительный пик в конце января и колебания в феврале и марте.

Дата	Индекс В (ВВП)
8 янв.	18
15 янв.	80
22 янв.	18
29 янв.	18
5 февр.	40
12 февр.	38
19 февр.	0
26 февр.	42
3 мар.	62
10 мар.	15
17 мар.	18
24 мар.	28
29 мар.	42



Electronic accessories	7,650
Fashion accessories	6,973
Food and beverages	6,880
Health and beauty	6,282
Home and lifestyle	6,167
Sports and travel	7,275

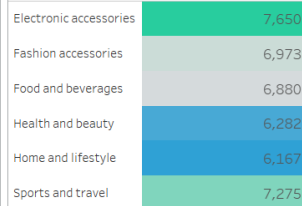
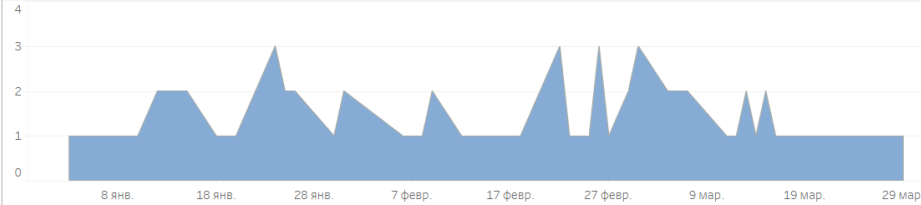


График показывает динамику индекса заболеваемости COVID-19 в Москве. По оси абсцисс отложены даты: 8 янв., 18 янв., 28 янв., 7 февр., 17 февр., 27 февр., 9 мар., 19 мар., 29 мар. По оси ординат — значения индекса от 0 до 4. Индекс начинается на уровне 1,0, поднимается до 2,0 в середине января, падает до 1,0, затем резко возрастает до 3,0 в конце января. В начале февраля он снижается до 1,0, но в середине февраля снова достигает 3,0. В конце февраля и начале марта индекс продолжает колебаться на уровне 2,0-3,0, а к концу марта возвращается к значению 1,0.

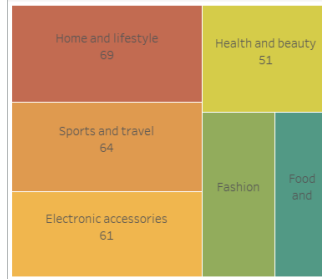
Дата	Индекс заболеваемости
8 янв.	1,0
18 янв.	2,0
28 янв.	3,0
7 февр.	1,0
17 февр.	3,0
27 февр.	2,0
9 мар.	2,0
19 мар.	1,0
29 мар.	1,0



6,820



Home and lifestyle 69	Health and beauty 51	
Sports and travel 64	Fashion 45	Food and drink 44
Electronic accessories 61		





06

Выводы

В процессе работы над дипломной работой было выполнено:

1. Обработка и анализ данных:
 - Проанализировали предоставленные данные и проверили их на ошибки.
2. Нормализованная схема данных (NDS):
 - Разработали структурированную нормализованную схему данных, обеспечивающую эффективное хранение и управление информацией.
3. Таблицы фактов и измерений (DDS):
 - Создали таблицы фактов и измерений, сформировав структуру схемы звезда для легкого доступа к ключевым показателям бизнеса.
4. ETL-процессы:
 - Разработали ETL-процессы для загрузки данных в NDS и DDS, обеспечивая эффективный и автоматизированный поток данных.
5. Дашборды в Tableau:
 - Построили дашборды в Tableau, визуализируя ключевые метрики, такие как общая выручка, прибыль, количество продаж, рейтинги и другие.

Решенные бизнес-задачи:

1. Управление финансами:
 - Мониторинг общей выручки компании и её филиалов.
 - Оценка общей прибыли и налога с выручки.
2. Анализ продаж:
 - Определение эффективности продаж в разрезе филиалов.
 - Изучение популярности продуктовых линеек.
3. Понимание клиентского поведения:
 - Анализ пола покупателей и их предпочтений в способах оплаты.
 - Выявление влияния наличия карты клиента.
4. Оптимизация ассортимента и рейтингов:
 - Ранжирование продуктовых линеек по прибыли и продажам.
 - Оценка среднего рейтинга от клиентов по категориям товаров и филиалам.

Итог работы:

Данная работа позволила закрепить полученные навыки в процессе прохождения курса.

Результат позволяет бизнесу эффективно управлять финансами, анализировать ключевые показатели продаж и клиентское поведение, принимать обоснованные решения для оптимизации бизнес-процессов и улучшения стратегии развития.



Спасибо за внимание!

Текст работы доступен на [GitHub](#)

