

Инжиниринг управления данными

Михайлов Роман

1 Определение области проекта

создание аналитической платформы, которая объединяет данные по странам (ВВП, сельскохозяйственный индекс, выбросы CO₂) создаёт локальную SQLite базу данных для хранения данных фильтрует данные и строит визуализацию по запросам к бд

2 Сбор данных

данные для базы данных берём из csv файлов с портала

<https://worlddataview.com/>

взяты 3 параметра

по каждому взят набор 21

csv файл по годам с 2000

по 2020

Copy	Print	CSV	Excel	PDF
Country	Crop production index (2004-2006 = 100)	Rank	Year	
Singapore	390.68	1	1961	
Seychelles	291.80	2	1961	
Saint Kitts and Nevis	290.61	3	1961	
Barbados	271.57	4	1961	
Saint Lucia	247.81	5	1961	
Trinidad and Tobago	221.84	6	1961	
Antigua and Barbuda	195.53	7	1961	
Tonga	152.15	8	1961	
Grenada	130.93	9	1961	
Norway	122.57	10	1961	
Japan	120.81	11	1961	
...

Исходные данные

для анализа взята статистика за год по параметрам:

1)gdp_growth 2)crop_index 3)co2_emissions

цель работы проанализировать как связана динамика параметров в разных странах

База данных

база данных создается локально с помощью библиотеки sqlite3 на питоне

2 этапа

1 первоначальный - база создается в первый раз с инпутом из папки инпут1 в которой хранится набор из 3х csv по 3м параметрам за один год - в данном случае 2000

2 этап в бд добавляются новые записи за следующие годы

1 этап

создаются 2 таблицы

Первая - `countries` хранит список всех стран, для которых есть данные, и присваивает каждой стране уникальный номер (ID) чтобы не хранить длинные названия стран в основной таблице данных. Вторая таблица `country_indicators` содержит значения показателей (ВВП, сельскохозяйственный индекс, выбросы CO₂) для каждой страны по годам. Здесь используется ID страны из таблицы `countries`, чтобы связать данные с конкретной страной.

1 этап

Также в этой таблице создан дополнительный столбец хранящий буквенный код страны (например, `RUS` для России) просто как пример доп данных предполагалось например хранить там флаги стран чтобы их можно было вызывать для визуализаций и выводов по запросам бд но некоторые страны не имеют своих флагов в библиотеках питон поэтому для примера заменён на буквенный код который не дублируется у каждой строки с данными по годам а хранится отдельно в таблице стран в одном экземпляре для экономии места

2 этап

следующие данные прошлых лет хранятся в отдельной папке инпут2 и уже позже добавляются к бд база данных строится из таблицы в которой хранится название страны из исходного csv файла, страна получает свой уникальный id при появлении новых стран в следующие года новые страны получают свои id id короче чем текстовое название страны что позволяет экономить память в бд если там будут записи за миллионы лет

db_create_and_update

1 база данных успешно создана и сохранена в файл

2 из отдельной папки с новыми данными успешно обновлена

db_analysis

анализ и визуализация данных из бд

Обеспечение качества данных

--- Запуск автоматизированных проверок качества данных ---

- ✅ Проверка 1 пройдена: Таблица 'countries' имеет ожидаемую структуру.
- ✅ Проверка 2 пройдена: Таблица 'country_indicators' имеет ожидаемую структуру.
- ✅ Проверка 3 пройдена: Таблица 'countries' не пуста. Записей: 218
- ✅ Проверка 4 пройдена: Таблица 'country_indicators' не пуста. Записей: 4099
- ✅ Проверка 5 пройдена: В таблице 'country_indicators' нет дубликатов по (country_id, year).
- ✅ Проверка 6 пройдена: Все country_id в 'country_indicators' ссылаются на существующие в 'countries'.
- ✅ Проверка 7 пройдена: Диапазон лет в 'country_indicators' (2000-2020) ожидаемый (2000-2020).
- ✅ Проверка 8 пройдена: В столбце 'co2_emissions' нет отрицательных значений.

--- Результат проверок ---

Пройдено: 8 / 8

- ✅ Все проверки качества данных пройдены успешно!

Анализ пропущенных данных

функция вызывает бд

```
conn = sqlite3.connect(db_path)

# Загрузим данные из country_indicators
query_all_data = """
SELECT c.country_name, ci.year, ci.gdp_growth, ci.crop_index, ci.co2_emissions
FROM country_indicators ci
JOIN countries c ON ci.country_id = c.country_id;
"""
df_full = pd.read_sql_query(query_all_data, conn)
conn.close()

print("Структура датафрейма из БД:")
print(df_full.info())

print("\nКоличество пропусков по столбцам:")
print(df_full.isnull().sum())

print("\nПроцент пропусков по столбцам:")
print((df_full.isnull().sum() / len(df_full)) * 100)
```

Структура датафрейма из БД:
<class 'pandas.core.frame.DataFrame'>

RangeIndex: 4099 entries, 0 to 4098

Data columns (total 5 columns):

#	Column	Non-Null Count	Dtype
0	country_name	4099 non-null	object
1	year	4099 non-null	int64
2	gdp_growth	4005 non-null	float64
3	crop_index	3591 non-null	float64
4	co2_emissions	3634 non-null	float64

dtypes: float64(3), int64(1), object(1)

memory usage: 160.2+ KB

None

Количество пропусков по столбцам:

country_name	0
year	0
gdp_growth	94
crop_index	508
co2_emissions	465

dtype: int64

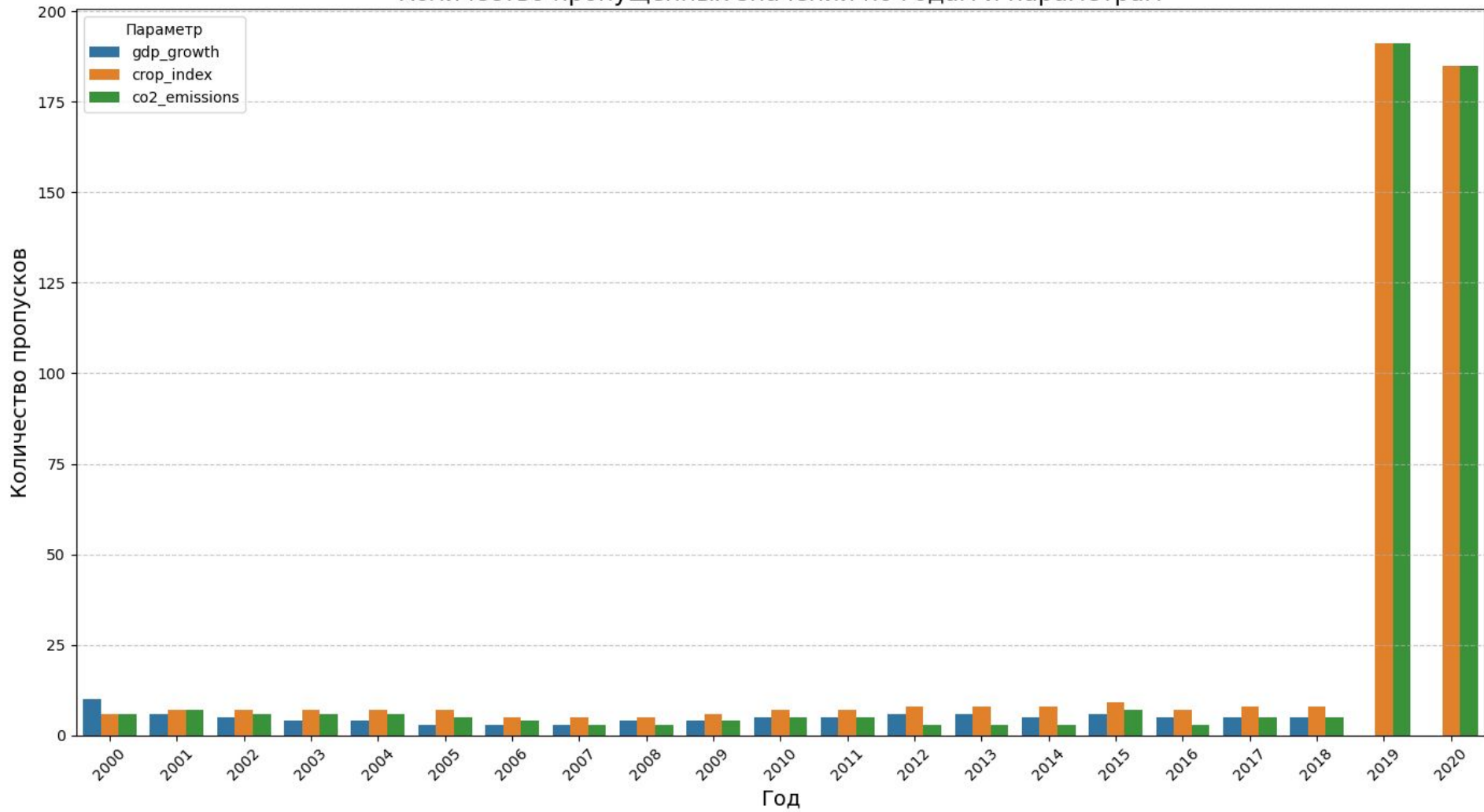
Процент пропусков по столбцам:

country_name	0.000000
year	0.000000
gdp_growth	2.293242
crop_index	12.393267
co2_emissions	11.344230

dtype: float64

видим что разные параметры имеют разное число пропусков

Количество пропущенных значений по годам и параметрам



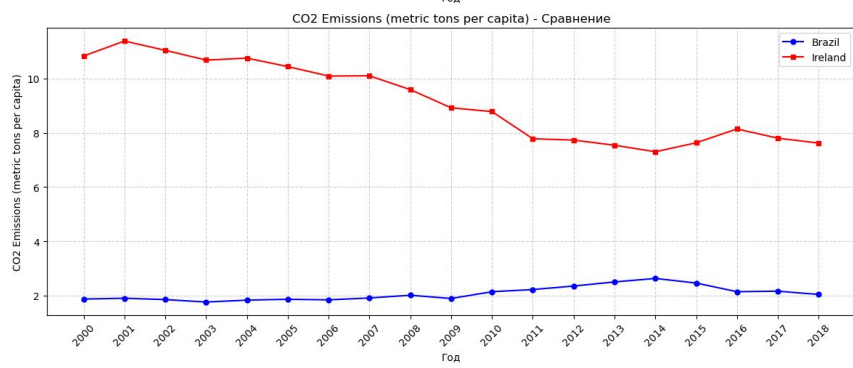
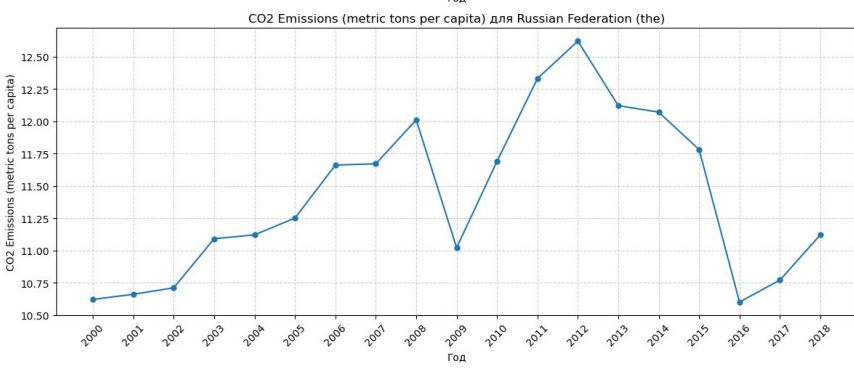
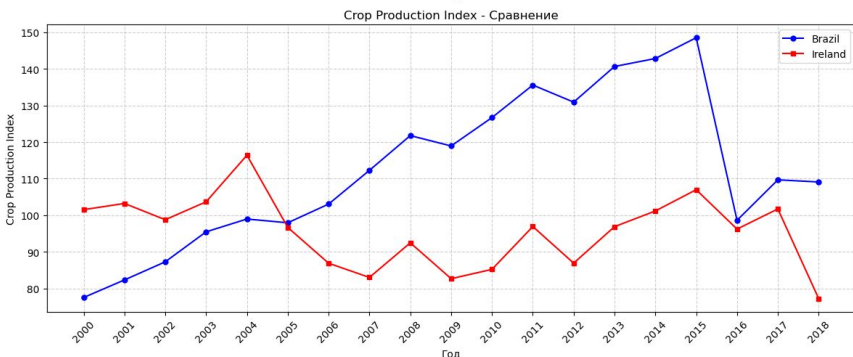
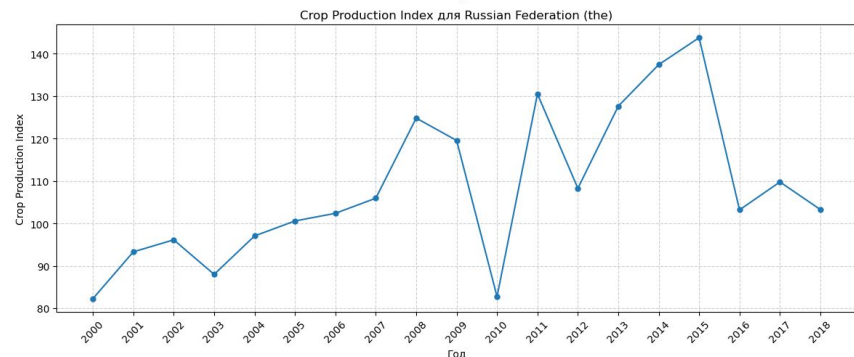
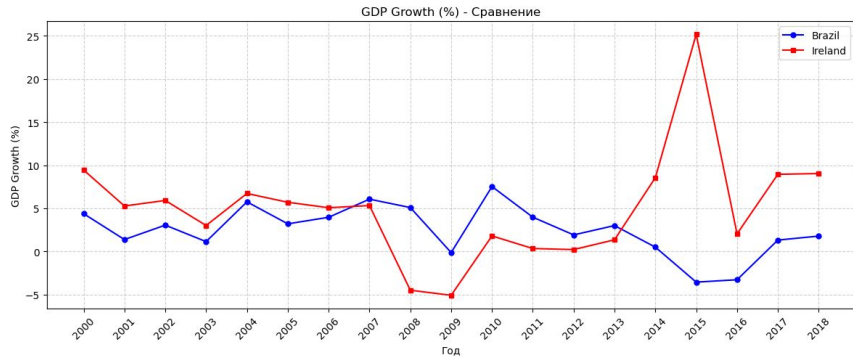
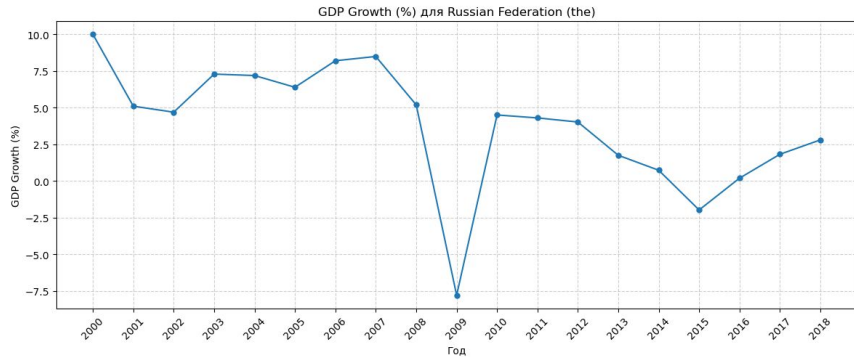
Фильтрация

видим что есть сильные проблемы с данными с начала пандемии - 2 года
можно отсеивать как недостоверные

Количество стран без пропусков по всем параметрам (2000-2018): 170

Анализ данных и визуализация

написана функция для визуализации уровней параметров по годам для одной страны по выбору и для двух стран VS по выбору



Инсайты

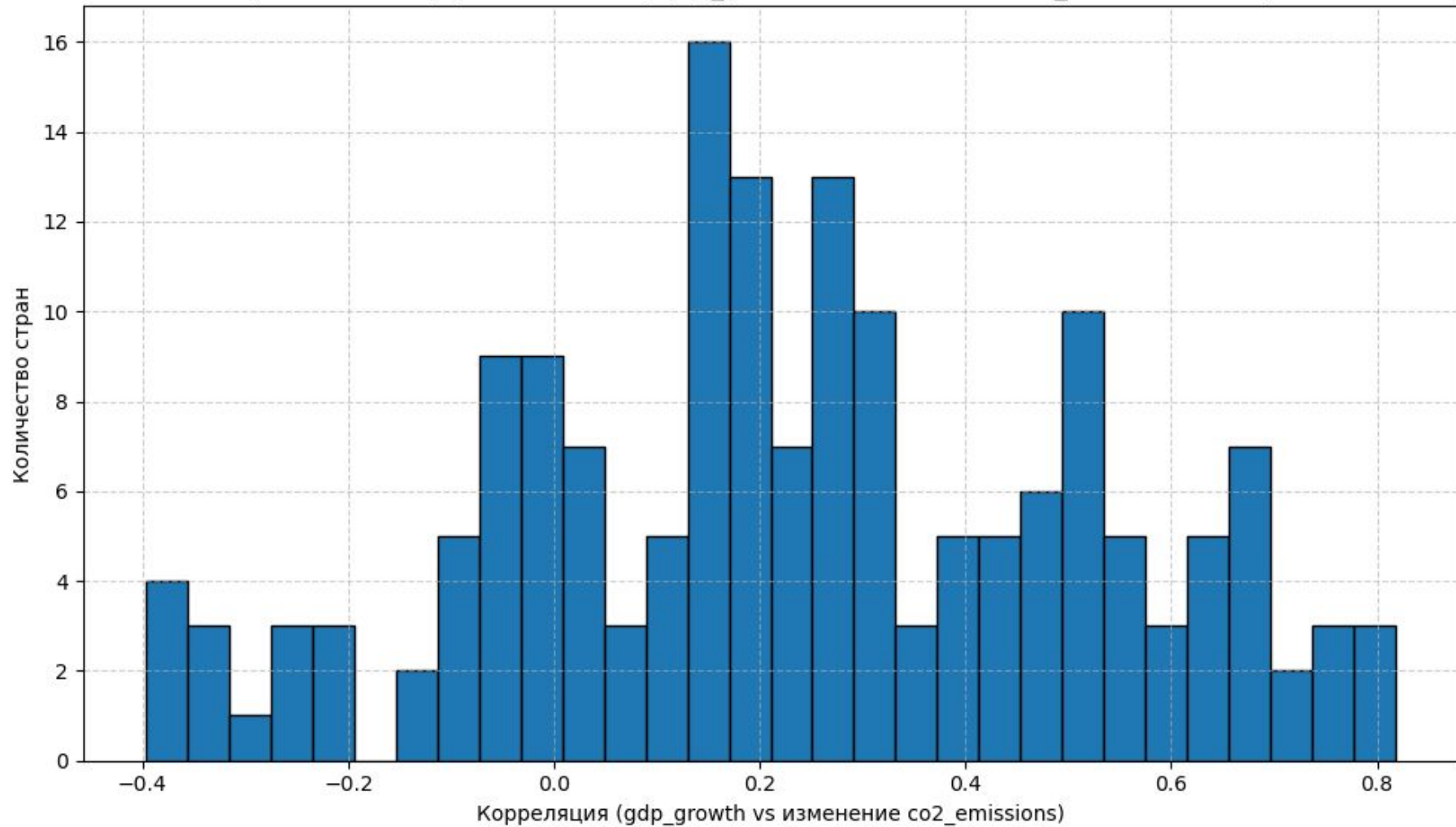
видим что в разных странах синхронизация отличается вероятно из за того что кредиты в банковской сфере берутся не на ближайший сезон посева и тд

наглядно видно что например в ирландии резкий экономический спад кризиса 2008 начался на год раньше чем в бразилии при этом бразилия не сокращала производство сх а наращивала

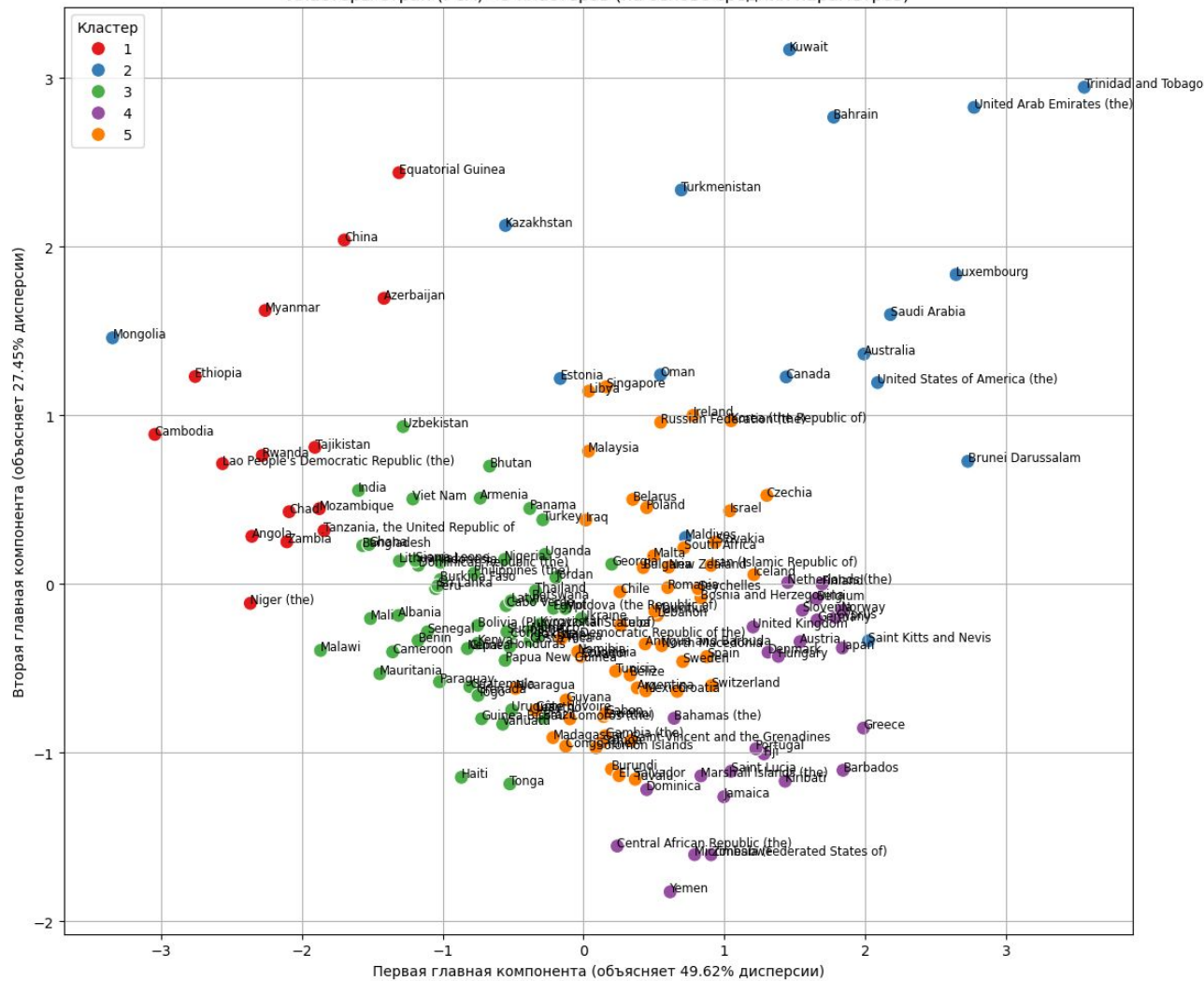
Поиск и визуализация паттернов и корреляций

в ноутбуке `db_analysis` находятся функции для анализа корреляций и визуализации

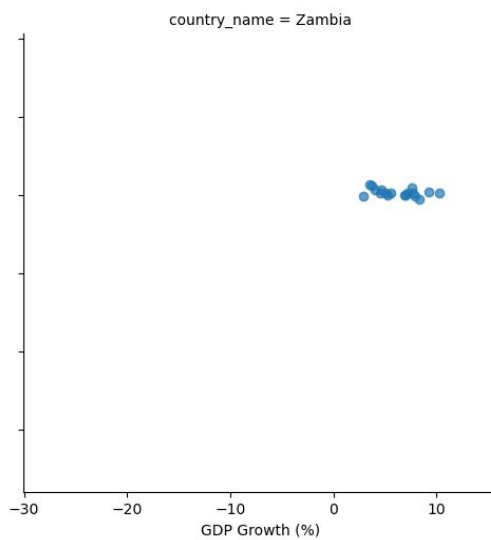
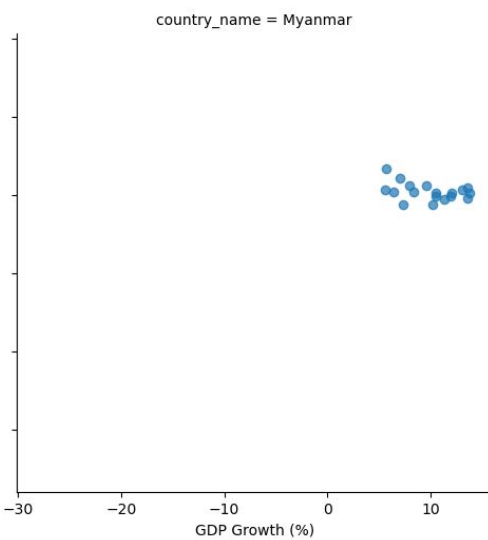
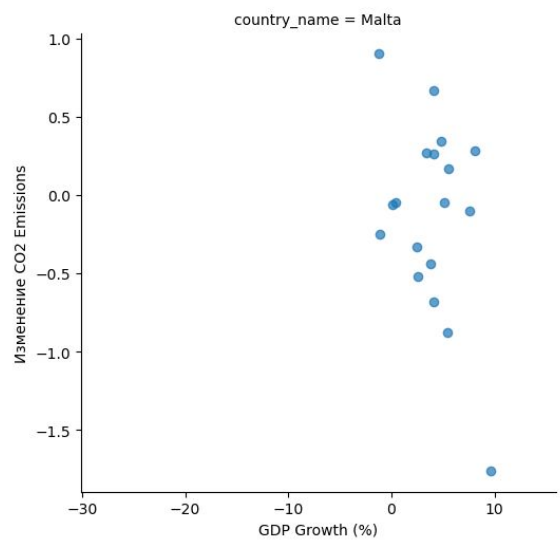
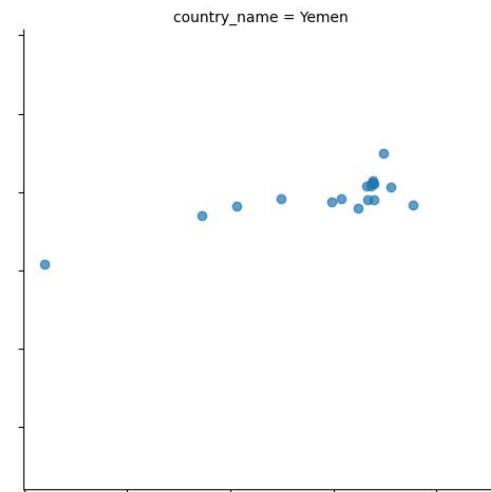
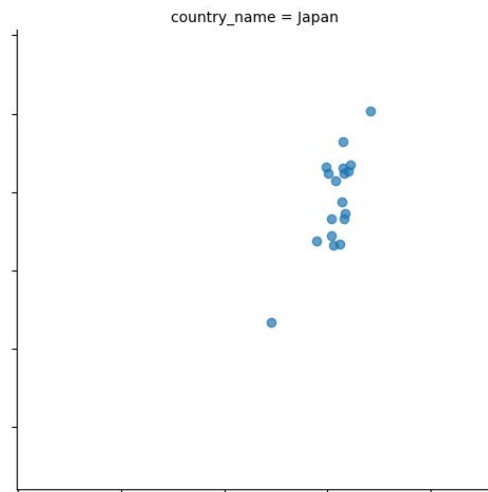
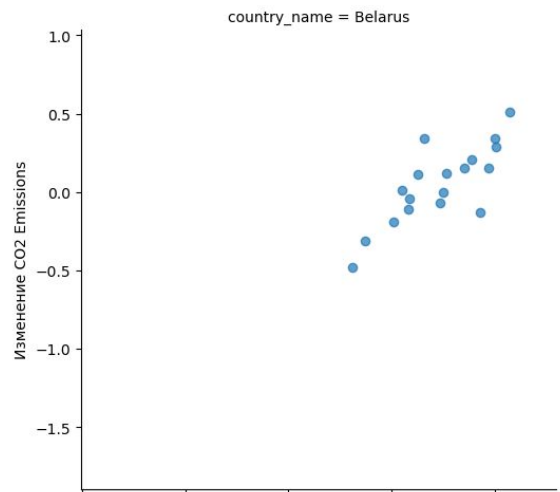
Распределение корреляций между gdp_growth и изменением co2_emissions по странам



Первая главная компонента (объясняет 49.62% дисперсии)



Сравнение: Высокая vs Низкая корреляция (gdp_growth vs изменение co2_emiissions)



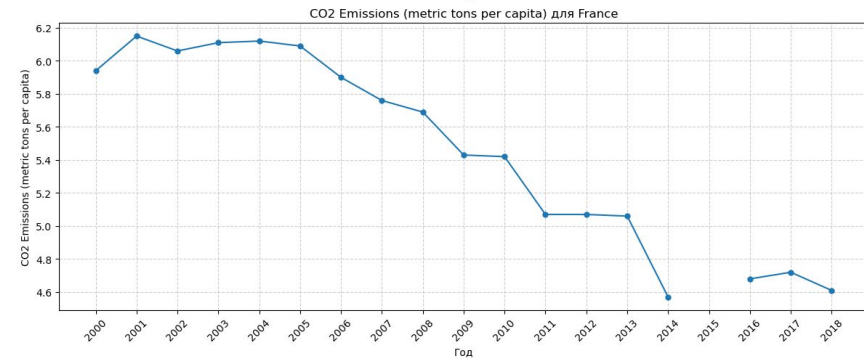
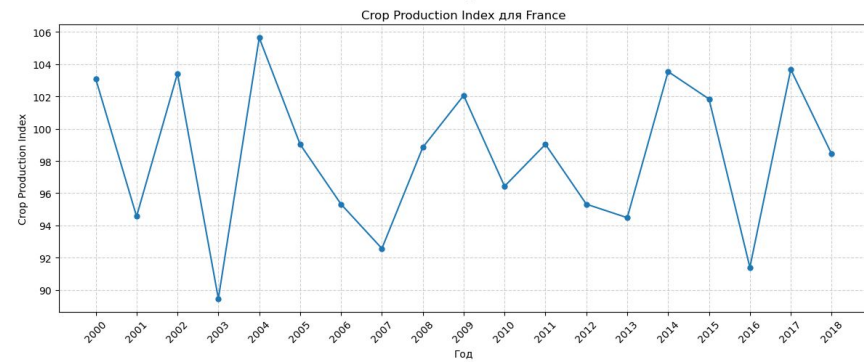
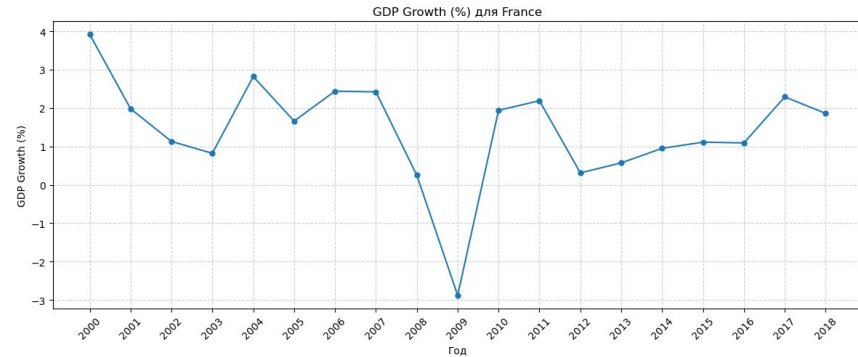
Анализ

наглядно видно что в разных странах параметры изменяются не одинаково
кластеризация по корреляциям не даёт однозначного объединения
богатейших стран или преимущественно сельскохозяйственных стран что
указывает на отсутствие единой модели волшебной палочки

Работа с пропусками

написаны функции для поиска стран имеющих пропуска в данных

для примера взята франция с 1 пропуском по отчётам выбросов углекислого газа за 2015й год



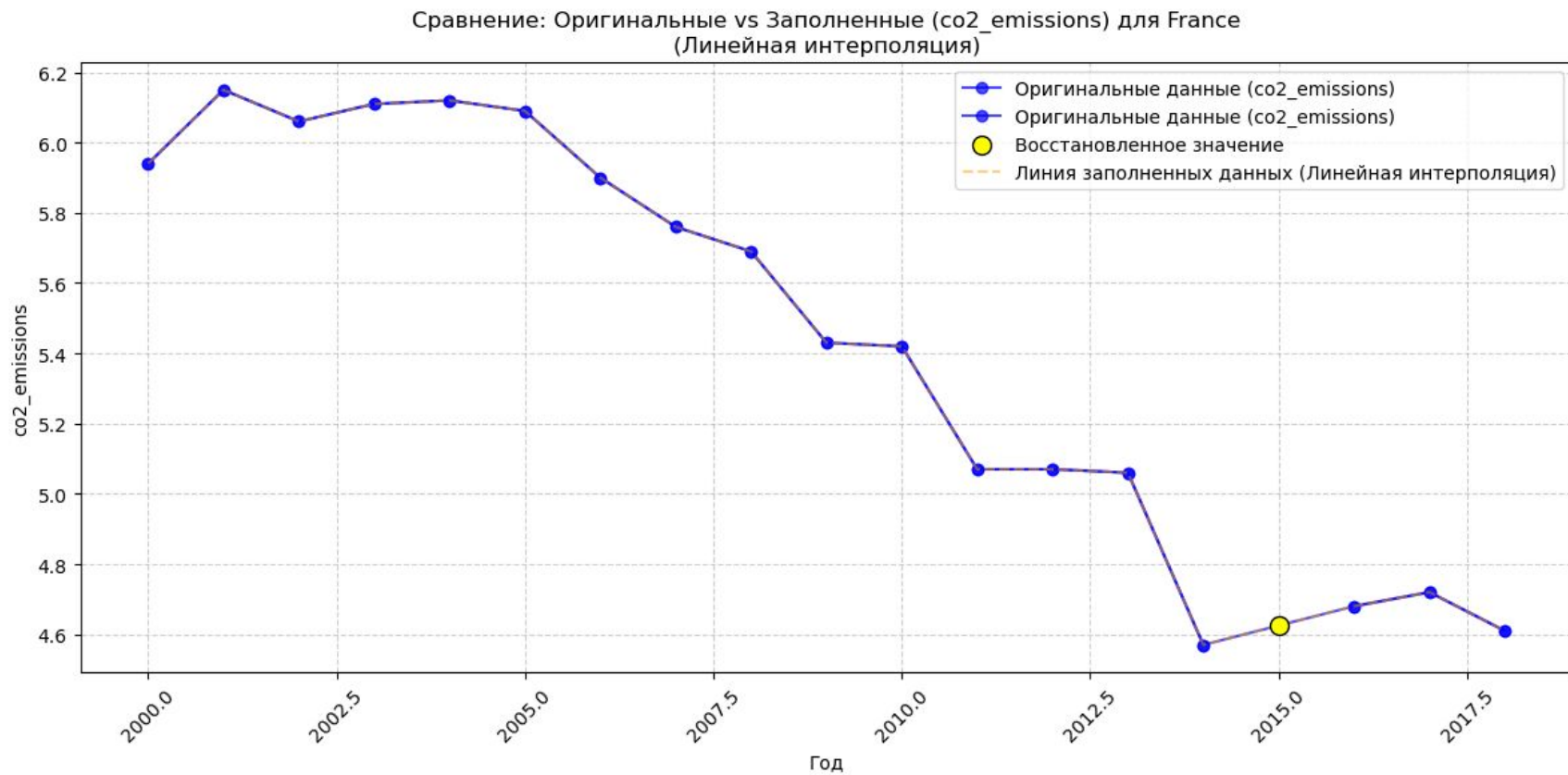
Заполнение пропусков

написаны 2 разные функции которые восстанавливают пропущенные значения разными способами

1 метод - усредняется 1 предыдущий и 1 следующий

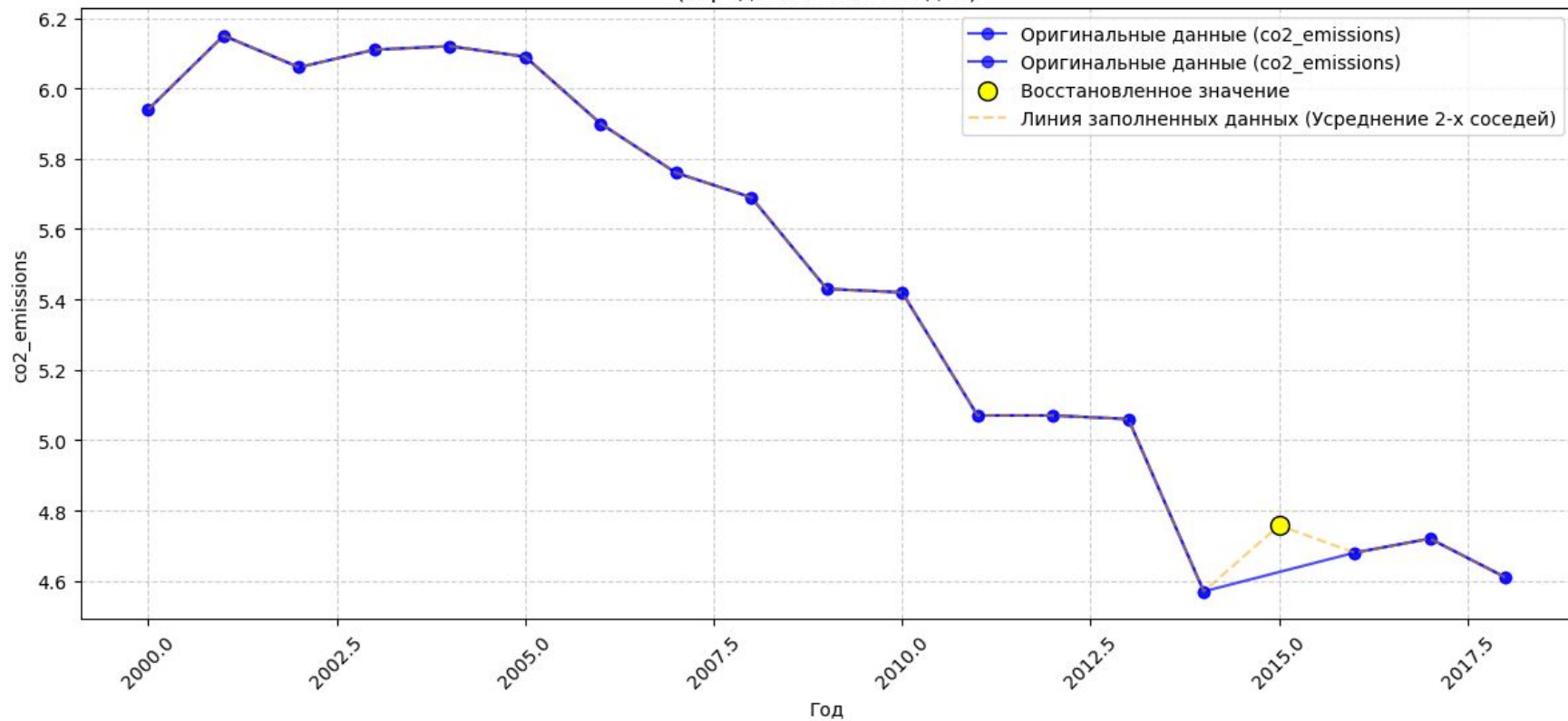
2 метод - усредняется 2 предыдущих и 2 следующих

метод 1



метод 2

Сравнение: Оригинальные vs Заполненные (co2_emissions) для France
(Усреднение 2-х соседей)



сравнение

наглядно видно что разные подходы возвращают разные значения

это важно учитывать для анализа и прогнозов

Итоги

Создана БД

В бд добавляются новые записи

Созданы функции анализа качества данных

Протестированы методы заполнения пропущенных данных

Созданы функции визуализации и анализа корреляций в данных

Проанализированы корреляции и тенденции на основе данных

Выводы

нет единой модели зависимости между ввп , уровнем производств сх культур и выбросами CO₂

различные страны демонстрируют разные корреляции и различные паттерны по годам что указывает на невозможность единого курса на эко френдли производства в странах с разной структурой экономики и культурами потребления

нельзя просто воспроизводить производственные практики и ожидать пропорционального снижения выбросов и так далее