# Brief Report of the Work Done

**Dataset Used**: Adoption of genetically modified crops, percentage of planting of different altered genes over the total plantings. Link: USDA ERS - Adoption of Genetically Engineered Crops in the U.S.

## 1. Solution Development:

- **Dataset Extraction**: The dataset was extracted directly from the source in .xlsx format, which contained several tables spread across three sheets.
- **Data Modeling**: I used Python and Jupyter Notebooks to process and unify the tables into a single coherent structure. The original file had distributed data, so I designed a custom script to consolidate this information into a single table. Due to the limited volume and structure of the dataset, I opted not to use a Kimball schema for dimension and fact tables.

**Python Script Used**:

```python
import pandas as pd
from datetime import datetime
import re
```

```python
# input the file of the dataset
file_path = 'BiotechCropsAllTables2024.xlsx'
xls = pd.ExcelFile(file_path)

# List to store the unified data
data = []
```

```python
# Iterate over each sheet
for sheet_name in xls.sheet_names:
    df = pd.read_excel(xls, sheet_name=sheet_name, header=None)

    # Find the rows where each table starts
    table_starts = df[df.apply(lambda row: row.str.contains('State/Year', na=False).any(), axis=1)].index

    for start in table_starts:
        # Get and clean the gene name
        gene = df.iloc[start-1, 0].split('(')[0].strip()

        # Get and clean the years
        years = df.iloc[start, 1:].values
        years = [str(year).strip() for year in years if pd.notna(year)]
        try:
            years = [int(float(year)) for year in years]
        except ValueError:
            continue  # Skip if conversion fails

        # Get the states and adoption percentages
        table_data = df.iloc[start+1:]
        for i, row in table_data.iterrows():
            state = str(row[0])

            # Skip empty or invalid rows
            if pd.isna(state) or state.strip() == "":
                break

            # Clean the state name by removing disclaimers "X/X"
            if state.strip() != "United States":
                state = re.split(r'[\s]', state)[0].strip()
            else:
                state = state.strip()

            # Extract adoption percentages
            adoption_percentages = row[1:].values
            #adoption_percentages = int(adoption_percentages)
            for year, adoption_percentage in zip(years, adoption_percentages):
                if pd.notna(adoption_percentage):
                    data.append({
                        'gene': gene,
                        'state': state,
                        'crop': sheet_name,
                        'year': year,
                        'adoption_percentage': adoption_percentage
                    })
```
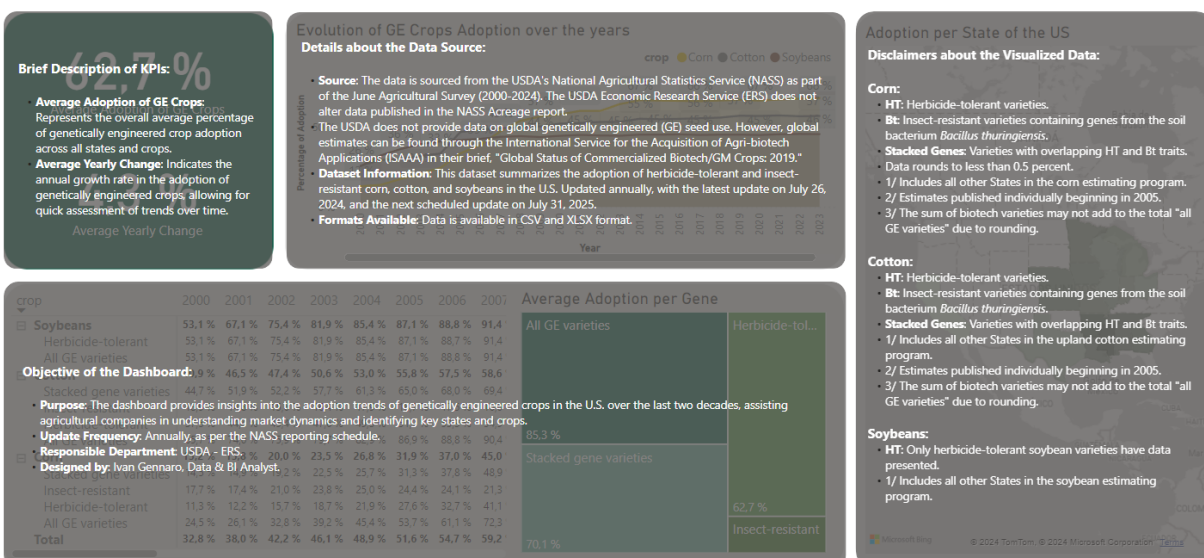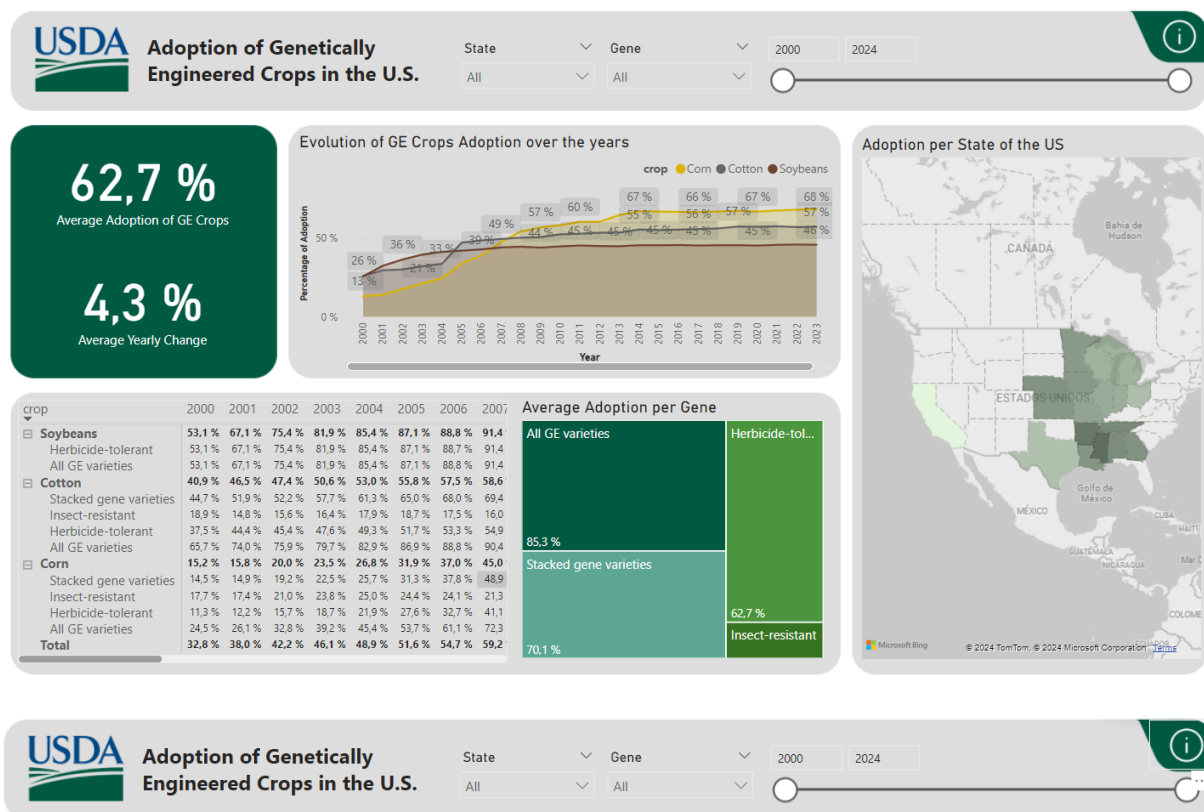
- **Dashboard Design in Power BI**: Importing the modeled data (.xlsx) into Power BI and designing an interactive dashboard with DAX measures such as:
  - Average adoption percentage
  - Average annual growth rate

## 2. Dashboard Objective

The dashboard aims to provide valuable insights into trends in genetically modified crop adoption over the past two decades by state in the U.S. This is particularly useful for a company in the agricultural sector focused on genetic modification, enabling them to identify:

- **Adoption Trends**: Understanding how the adoption of different genes has evolved over the years and in which states they have been most prevalent.
- **Geographical Information**: Visualizing which regions have the most significant adoption of these technologies, aiding strategic business decisions.

**3. Conclusion**

Given the problem presented, I opted for an efficient solution planned over three days:

- **Day 1**: Dataset search, BI planning, and data modeling development.
- **Day 2**: Preliminary dashboard design and data importation.
- **Day 3**: Visualization refinement, creation of DAX measures, and documentation.

**Total time invested**: Approximately 5 hours.