# Stone – Data-Driven Insights & Business Solutions

**Machine Learning & Analytics | SQL, Python, R | Healthcare, Banking, Forecasting**
Budgeting | Healthcare & Insurance Analysis | Bank & Customer Insights
Data Visualization | Power BI | Excel Modeling | Jupyter Notebook

# SQL

**Invalid Pickup Locations**

- Some trips have `PULocationID` values that do not exist in the official `taxi_zone_lookup`.
- **Outcome:** Highlights data inconsistencies and missing location mappings.

**Unusual Trip Distances**

- Some trips recorded zero, negative, or excessively high distances (e.g., over 50 miles).
- **Outcome:** Identifies potential errors in trip data affecting fare calculations.

**Mismatched Fare Amounts**

- Some fares are significantly too low (<$2) or too high relative to trip distance.
- **Outcome:** Detects incorrect fare entries or potential fraudulent activity.

**Service Zone Anomalies**

- Certain service zones are linked to multiple boroughs, which should not occur.
- **Outcome:** Flags inconsistencies in zone mapping that could impact reporting.

```
        tpep_dropoff_datetime AS DropoffTime,
        ████████
        payment_type,
        fare_amount,
        tip_amount,
        total_amount
FROM
        ████████
ORDER BY
        tpep_dropoff_datetime,
        trip_distance DESC,
        passenger_count;

-- Question 2:
SELECT
        VendorID,
        COUNT(*) AS RecordCount
```

Results | Messages

| | | | | | |
|---|---|---|---|---|---|
| 2023-07-27 | 2023-07-27 | 9.98 | 1 | 42.9 | 12.54 | 75.24 |

| | RecordCount |
|---|---|
| 2 | 29003 |
| 1 | 8982 |
| 6 | 3 |

| service_zone | BoroughCount |
|---|---|
| Airports | 1 |
| Boro Zone | 5 |
| EWR | 1 |
| N/A | 1 |
| Yellow Zone | 1 |

| PickupDate | DropoffDate | passenger_count | total_amount | DropoffBorough | DropoffServiceZone |
|---|---|---|---|---|---|
| 2023-07-27 | 2023-07-27 | 2.0 | 30.4 | | Yellow Zone |
| 2023-07-27 | 2023-07-27 | 1.0 | 26.2 | Queens | |
| 2023-07-27 | 2023-07-27 | 1.0 | 10.8 | Manhattan | |

# SQL

**App Count per Category**

- Groups apps by category and counts the number of apps in each.
- **Outcome:** Identifies which categories have the most or least apps.

**Most Reviewed App**

- Retrieves the app with the highest number of reviews.
- **Outcome:** Highlights the most engaged and popular app.

**Game Genre Statistics**

- Finds the number of apps, max installs, and min reviews per game genre.
- **Outcome:** Provides insights into game market trends and user engagement.

**Unique Genres per Category**

- Counts distinct genres within each app category.
- **Outcome:** Shows category diversity and market segmentation.

```sql
-- count the numberof apps per category
SELECT
    Category,
    COUNT(AppName) AS AppCount
FROM
    googleplaystore1
GROUP BY
    Category
ORDER BY Category;

-- Retrieve
SELECT
    AppName,
    Reviews
FROM
    googleplaystore1
WHERE
    CAST(Reviews AS INT) = (SELECT MAX(CAST(Reviews AS INT)) FROM googleplaystore1);

-- #3
SELECT
    Genres,
    COUNT(AppName) AS AppCount,
    MAX(CAST(Replace(REPLACE(Installs, '+',''),',','') AS INT)) AS MaxInstalls,
    MIN(CAST(Reviews AS INT)) AS MinReviews
FROM
    googleplaystore1
WHERE
    Category = 'GAME'
GROUP BY
    Genres
ORDER BY
    Genres;
```

| | Category | AppCount |
|---|---|---|
| 1 | 1.9 | 1 |
| 2 | ART_AND_DESIGN | 6 |
| 3 | AUTO_AND_VEHICLES | 8 |
| 4 | BEAUTY | 5 |
| 5 | BOOKS_AND_REFERENCE | 2 |
| 6 | BUSINESS | 4 |
| 7 | COMICS | |
| 8 | COMMUNICATION | 7 |

| | AppName | Reviews |
|---|---|---|
| 1 | Facebook | 78 |

| | Genres | AppCount | MaxInstalls | MinReviews |
|---|---|---|---|---|
| 1 | Action | 5 | | |
| 2 | Action;Action & Adventure | | | |
| 3 | Adventure | | | |
| 4 | Adventure;Action & Adv... | | | |
| 5 | Arcade | | | |

# SQL

```
1   SELECT e.emp_no,
2       e.first_name,
3       e.last_name,
4           e.birth_date,
5           d.from_date,
6           d.to_date
7           SELECT e.emp_no,
8       e.first_name,
9       e.last_name,
10          t.title,
11          t.from_date,
12          t.to_date
13  INTO retirement_titles
14  FROM employees as e
15  LEFT JOIN titles as t
16  ON (e.emp_no = t.emp_no)
17  WHERE (birth_date BETWEEN '1952-01-01' AND '1955-12-31')
18  ORDER BY e.emp_no;
19
20
21  -- Use Dictinct with Orderby to remove duplicate rows
22  SELECT DISTINCT ON (rt.emp_no) rt.emp_no,
23  rt.first_name,
24  rt.last_name,
25  rt.title
26
27  INTO unique_titles
```

| | count bigint | title character varying |
|---|---|---|
| 1 | 29414 | Senior Engineer |
| 2 | 28254 | Senior Staff |
| 3 | 14222 | Engineer |
| 4 | 12243 | Staff |
| 5 | 4502 | Technique Leader |
| 6 | 1761 | Assistant Engineer |
| 7 | 2 | Manager |

# SQL

```
-- Creating tables for PH-EmployeeDB
CREATE TABLE departments (
        dept_no VARCHAR(4) NOT NULL,
        dept_name VARCHAR(40) NOT NULL,
        PRIMARY KEY (dept_no),
        UNIQUE (dept_name)
);

CREATE TABLE employees (
        emp_no INT NOT NULL,
        birth_date DATE NOT NULL,
        first_name VARCHAR NOT NULL,
        last_name VARCHAR NOT NULL,
        gender VARCHAR NOT NULL,
        hire_date DATE NOT NULL,
        PRIMARY KEY (emp_no)
);

CREATE TABLE dept_manager (
dept_no VARCHAR(4) NOT NULL,
        emp_no INT NOT NULL,
        from_date DATE NOT NULL,
        to_date DATE NOT NULL,
FOREIGN KEY (emp_no) REFERENCES employees (emp_no),
FOREIGN KEY (dept_no) REFERENCES departments (dept_no),
```

| emp_no integer | first_name character varying | last_name character varying | title character varying |
|---|---|---|---|
| 10001 | | | Senior Engineer |
| 10004 | | | Senior Engineer |
| 10005 | | | Senior Staff |
| 10006 | | | Senior Engineer |
| 10009 | | | Senior Engineer |
| 10011 | | | Staff |
| 10018 | | | Senior Engineer |

# SAS- 1st part (Insurance Data)

**Health Classification & Insurance Premiums (Insurance Data)**

- Logistic regression categorizes individuals into "Healthy," "Sick," or "Severely Sick" to analyze premium adjustments.
- **Outcome:** Enhances accuracy in pricing insurance policies.

**Predicting Patient Recovery (Insurance Data)**

- Regression models assess how age, condition, blood sugar, and disease classes influence recovery rates.
- **Outcome:** Improves risk assessment for insurance and healthcare planning.



```
/* Step 2: Check and Recode Variables if Necessary */
data insurance_data;
    set insurance_data;

    /* Recode class to descriptive categories */
    if class = 1 then class_cat = "Healthy";
    else if class = 2 then class_cat = "Sick";
    else if class = 3 then class_cat = "Severely Sick";

    /* Standardize Sex values */
    if upcase(Sex) = "MALE" then sex_cat = "Male";
    else if upcase(Sex) = "FEMALE" then sex_cat = "Female";
run;

/* Verify the Recoded Variables */
proc freq data=insurance_data;
    tables class_cat sex_cat;
run;

/* Step 3: Logistic Regression Analysis */
proc logistic data=insurance_data;
    class disease_classes(ref="S") class_cat(ref="Healthy") sex_cat(ref="Male") / param=ref;
    model cured(event='1') = class_cat sex_cat age current_condition blood_sugar revival_days disease_classes;
run;

/* Step 4: Assess Model Fit with ROC Curve */
proc logistic data=insurance_data plots(only)=roc;
    class disease_classes(ref="S") class_cat(ref="Healthy") sex_cat(ref="Male") / param=ref;
    model cured(event='1') = class_cat sex_cat age current_condition blood_sugar revival_days disease_classes;
run;
```
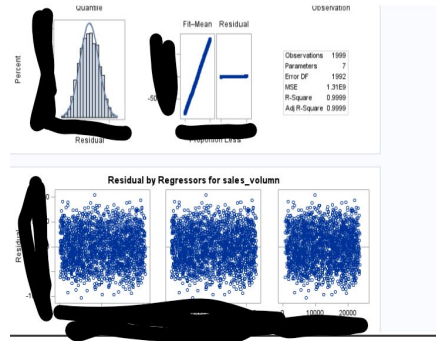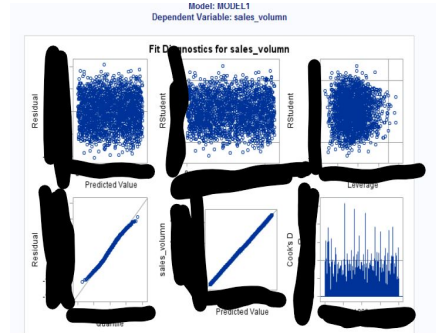
# SAS - 2nd (Forecast)

**Sales Volume Impact Analysis (Forecast Consulting Data)**

- Regression analysis identifies key factors (S1, S2, S4, S5, S6, occasional channel) affecting sales volume.
- **Outcome:** Supports strategic decision-making in marketing and resource allocation.

**Correlation Between Factors (Forecast Consulting Data)**

- Examines the relationships between sales volume and sales channels.
- **Outcome:** Provides insights into key drivers of sales performance.





```
/* Step 1: Importing the Dataset */
proc import datafile="C:\Users\leiker-s\Desktop\Data Case Analysis\Forecast Consulting Data-1.csv"
    out=forecast_data
    dbms=csv
    replace;
    getnames=yes;
run;


/* Step 1.1: Exploring the Data Structure */
proc contents data=forecast_data;
run;


/* Step 2: Descriptive Statistics */
proc means data=forecast_data;
    var s1 s2 s4 S5 S6          ____ ___ es_volumn;
run;


/* Step 3: Regression Analysis to Determine Channel Effects */
proc reg data=forecast_data;
    model sales_volumn = s1 s2 s4 S5 S6 _____;
run;


/* Step 4: Correlation Analysis to Identify Relationships */
proc corr data=forecast_data;
    var sales_volumn s1 s2 s4 S5 S6 _____;
run;


/* Step 5: Stepwise Regression for Reallocation Strategy */
proc glmselect data=forecast_data;
    model sales_volumn = s1 s2 s4 S5 S6 _____ ___ ____ / selection=stepwise;
run;
```

# Python

**Simulating Dice Rolls**

- Rolls two six-sided dice repeatedly and tracks the outcomes.
- **Outcome:** Generates realistic dice roll distributions for probability estimation.

**Checking for Double Sixes**

- Simulates 24 dice rolls in a game and determines if at least one double six appears.
- **Outcome:** Evaluates the likelihood of rolling double sixes in a session.

**Monte Carlo Probability Estimation**

- Runs multiple trials (e.g., 100,000) to estimate the probability of rolling at least one double six.
- **Outcome:** Provides a reliable statistical probability based on large-scale simulations.

**Final Probability Calculation**

- Computes the success rate from all trials and outputs the estimated probability.
- **Outcome:** Delivers an accurate probability approximation based on empirical data.

```python
import random

def roll_dice():
    """Simulate rolling two six-sided dice."""
    die1 = random.randint(1, 6)
    die2 = random.randint(1, 6)
    return die1, die2

def simulate_game(num_rolls=24):
    """Simulate rolling the dice 'num_rolls' times."""
    for _ in range(num_rolls):
        die1, die2 = roll_dice()
        if die1 == 6 and die2 == 6:
            return True  # At least one double six rolled
    return False  # No double six rolled

def monte_carlo_simulation(num_trials=100000):
    """Run the simulation for 'num_trials' and calculate the probability of rolling at least one double six."""
    successful_trials = 0

    for _ in range(num_trials):
        if simulate_game():
            successful_trials += 1

    probability = successful_trials / num_trials
    return probability

if __name__ == "__main__":
    num_trials = 100000  # Number of trials for the Monte Carlo simulation
    probability = monte_carlo_simulation(num_trials)
    print(f"The estimated probability of rolling at least one double six in 24 rolls is approximately {probability:.4f}")
```

# Python

**Reading Rose Bowl Data**

- Loads team names from a file listing Rose Bowl winners from 1902 to 2020.
- **Outcome:** Extracts historical game results for analysis.

**Counting Wins per Team**

- Uses a counter to count how many times each team has won the Rose Bowl.
- **Outcome:** Identifies teams with the highest number of victories.

**Saving Win Counts to CSV**

- Writes the team names and their total wins to a new CSV file.
- **Outcome:** Creates a structured dataset for further analysis.

**Displaying Teams with More Than Four Wins**

- Filters and prints only teams that have won more than four times.
- **Outcome:** Highlights the most successful teams in Rose Bowl history.

```python
#Initialization
import csv
from collections import Counter

#INPUT - Read the file and get teams' data
def read_rosebowl(filename):
    with open(filename, 'r') as file:
        teams = file.read().splitlines()
    return teams


#PROCESS: Count the number of Rose Bowl for each team
def win_team_count(teams):
    return Counter(teams)


# Write the results to a new CSV file
def write_wins_csv(wins, output_filename):
    with open(output_filename, 'w', newline='') as csvfile:
        writer = csv.writer(csvfile)
        writer.writerow(['Team', 'Wins'])
        for team, win_count in wins.items():
            writer.writerow([team, win_count])

# OUTPUT: Display teams with more than 4 wins
def display_teams_more_than_4wins(wins):
    print("Teams with more than 4 wins:")
    for team, win_count in wins.items():
        if win_count > 4:
            print(f"{team}: {win_count} wins")


def main():
    filename = r'C:\Users\leiker-s\Desktop\ANLY 615\Python\Module 3\Rosebowl.txt'
    output_filename = 'Rosebowl_Wins.csv'

    teams = read_rosebowl(filename)

    teams_wins = win_team_count(teams)
    write_wins_csv(teams_wins, output_filename)

    display_teams_more_than_4wins(teams_wins)

if __name__ == "__main__":
    main()
```
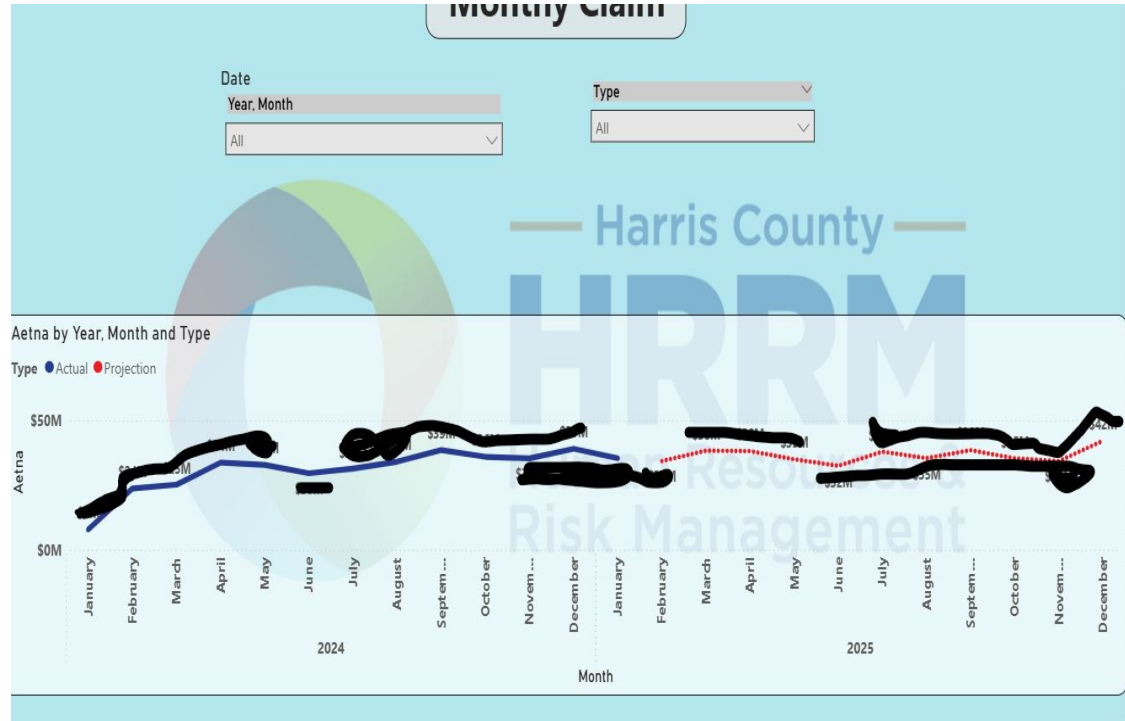
# Power BI

**Actual Financial Reports – Weekly, Monthly, Quarterly, Yearly**

- Tracks real-time spending and revenue in healthcare and pharmacy operations.
- **Outcome:** Provides a structured financial overview for performance evaluation.

**Future Projections – Medical & Pharmacy (Grants, Rebates, etc.)**

- Uses historical data to forecast upcoming expenses, revenue, and funding sources.
- **Outcome:** Supports budgeting decisions and strategic planning for cost management

# Power BI

**Drug Name Identification & Inflation Impact**

- Analyzes a list of medications, highlighting those affected by price inflation.
- **Outcome:** Identifies cost trends to support budgeting and policy adjustments.

**Seasonal Claims Increase & Dashboard Insights**

- Tracks rising claims, especially during flu season in winter.
- **Outcome:** Uses charts and multiple dashboards to drive data-backed decisions.

Utilization Detail by Medical Cost C...  ● Ambulatory ...  ● Emergen...  ● Home He...  ● Inpatient ...  ● Lab  ● Medical ...  ● Mental ...  ▶

Sum of Value

6M

4M

2M

0M

2          4          6          8          10          12

Month

0.00M                                                                492.74M

| Utilization Detail by Medical Cost Category | January | February | March | April | May | June | July | August | September | October | November | Decembe |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Ambulatory Facility | | | | | | | | | | | | |
| Emergency Room | | | | | | | | | | | | |
| Home Health | | | | | | | | | | | | |
| Inpatient Facility | | | | | | | | | | | | |
| Lab | | | | | | | | | | | | |
| Medical Pharmacy | | | | | | | | | | | | |
| Mental Health | | | | | | | | | | | | |
| Primary Physician | | | | | | | | | | | | |
| Radiology | | | | | | | | | | | | |
| Specialist Physician | | | | | | | | | | | | |
| **Total** | | | | | | | | | | | | |

# VLOOKUP

I use the VLOOKUP formula to match IDs and check if retirees and employees are still here in our company. This is just one of many ways I use VLOOKUP regularly.

# Line Chart



The trend of sum of Spend (actual & forecast) for Transaction Date Month. Color shows details about Category.

# Healthcare Costs (Managerial Accounting)

Managerial accounting- Analyzing FY25 Healthcare Costs to improve budgeting and forecasting. Tracking enrollment trends helps ensure sustainable fundings.

# PivotChart + Pivot Table

| Sum of Claim Amt | Column Labels ▾ | | | |
|---|---|---|---|---|
| | ⊞ Oct | ⊞ Nov | ⊞ Dec | Grand Total |
| Row Labels ▾ | | | | |
| Medical - In-Network | ████████████████████████████ | | | |
| Medical - Out-of-Network | | | | |
| Rx | | | | |
| Grand Total | | | | |

| Count of Claim Amt | Column Labels ▾ | | | |
|---|---|---|---|---|
| | ⊞ Oct | ⊞ Nov | ⊞ Dec | Grand Total |
| Count of Claim Amoun ▾ | | | | |
| Medical - In-Network | ████████████████████████████ | | | |
| Medical - Out-of-Network | | | | |
| Rx | | | | |
| Grand Total | | | | |

| | | Volume Claim | | |
|---|---|---|---|---|
| | | Nov | Dec | |
| | Medical | 104,470 | 111,775 | |
| | Rx | | | |
| | Grand Total | | | |

# Pivot Table & Chart

## Profit and Sales

| Region | Category | Profit | Sales |
|--------|----------|--------|-------|
| AsiaPac | Telephones and Com.. | $916,003 | $1,639,4! |
| | Tables | $865,879 | $1,367,8: |
| | Chairs & Chairmats | $691,200 | $1,186,9! |
| | Office Machines | $566,976 | $1,677,1: |
| | Storage & Organization | $429,069 | $772,1: |
| | Bookcases | $361,681 | $576,1: |
| | Computer Peripherals | $340,243 | $596,6: |
| | Appliances | $296,971 | $573,0! |
| | Copiers and Fax | $268,432 | $722,7: |
| | Office Furnishings | $215,290 | $420,7: |
| | Binders and Binder A.. | $168,273 | $654,3: |
| | Paper | $97,297 | $294,8: |
| | Pens & Art Supplies | $58,044 | $113,1: |
| | Scissors, Rulers and .. | $45,479 | $63,4: |
| | Envelopes | $34,556 | $100,5: |
| | Labels | $9,866 | $31,1: |
| | Rubber Bands | $5,671 | $12,1: |
| EMEA | Tables | $643,292 | $997,3! |
| | Telephones and Com.. | $405,059 | $725,8: |
| | Chairs & Chairmats | $348,065 | $593,6: |
| | Office Machines | $305,828 | $783,8: |
| | Storage & Organization | $251,128 | $444,3: |
| | Office Furnishings | $159,444 | $298,2: |
| | Appliances | $131,333 | $254,6: |
| | Copiers and Fax | $130,829 | $360,8: |
| | Computer Peripherals | $129,759 | $229,3: |
| | Bookcases | $129,584 | $206,2: |
| | Binders and Binder A.. | $96,492 | $342,7: |

## Discount and Profit Margin

Order Date



## Profit

$2,393 — $916,003

## Profit

$2,393 ▭ $916,003

## Measure Names

- ▪ Discount
- ▪ Product Base Margin

## Region Profit and Sales

Region

# Pivot Table

# Power BI

**Actual Sales Volume**

- Analyzes total sales transactions over a specific period.
- **Outcome:** Measures real-time business performance based on actual sales data.
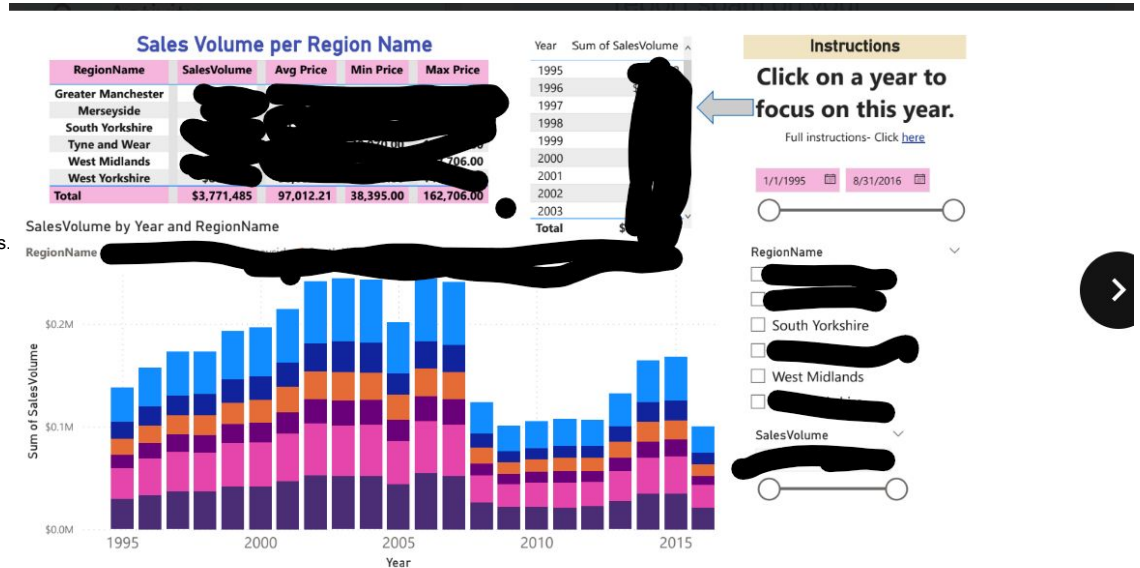
**Monthly Earnings Calculation**

- Aggregates revenue earned per month from recorded sales.
- **Outcome:** Tracks financial trends and identifies seasonal variations.

**Yearly Earnings Analysis**

- Sums up total revenue for each year to compare long-term growth.
- **Outcome:** Evaluates annual performance and supports strategic planning.

**Forecasting Future Revenue**

- Uses past sales data to predict upcoming earnings.
- **Outcome:** Helps structure financial goals and resource allocation.

# Power BI

**Actual Sales Volume**

- Analyzes total sales transactions over a specific period.
- **Outcome:** Measures real-time business performance based on actual sales data.

**Monthly Earnings Calculation**

- Aggregates revenue earned per month from recorded sales.
- **Outcome:** Tracks financial trends and identifies seasonal variations.

**Yearly Earnings Analysis**

- Sums up total revenue for each year to compare long-term growth.
- **Outcome:** Evaluates annual performance and supports strategic planning.

**Forecasting Future Revenue**

- Uses past sales data to predict upcoming earnings.
- **Outcome:** Helps structure financial goals and resource allocation.

# Power BI

**Top Sales Locations**

- Identifies the locations with the highest total sales.
- **Outcome:** Pinpoints the most profitable areas for business growth.
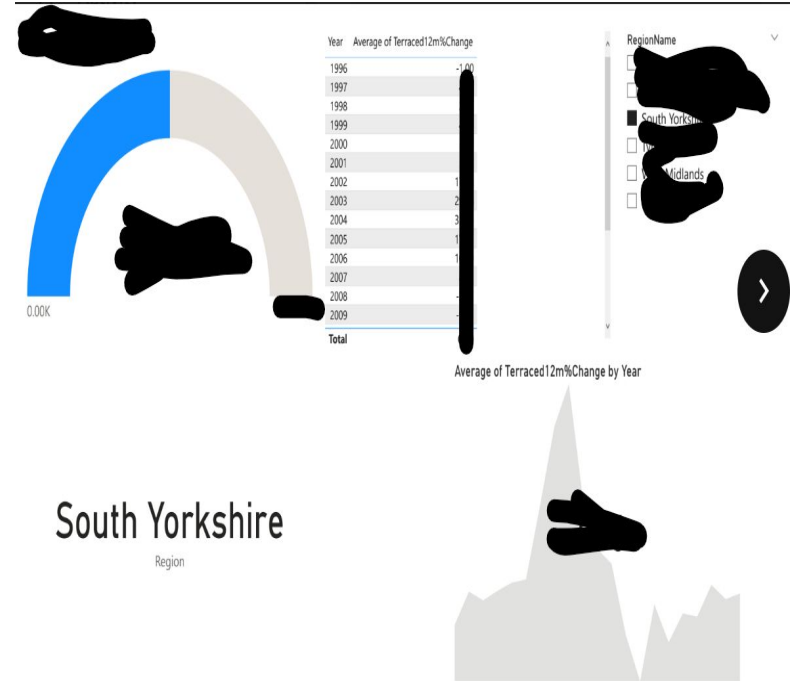
**Yearly Sales Performance**

- Aggregates total sales for each year.
- **Outcome:** Evaluates long-term trends and revenue consistency.

**Sales Distribution by Region**

- Compares sales performance across different locations.
- **Outcome:** Helps optimize resource allocation and marketing strategies.

**Future Sales Forecasting**

- Uses past yearly sales data to predict future revenue trends.
- **Outcome:** Supports strategic decision-making for expansion and investment.

# Tableau

**Customer Demographics & Average Age**

- Analyzes customer age distribution and key demographic trends.
- **Outcome:** Provides insights into the bank's target audience and customer segments.
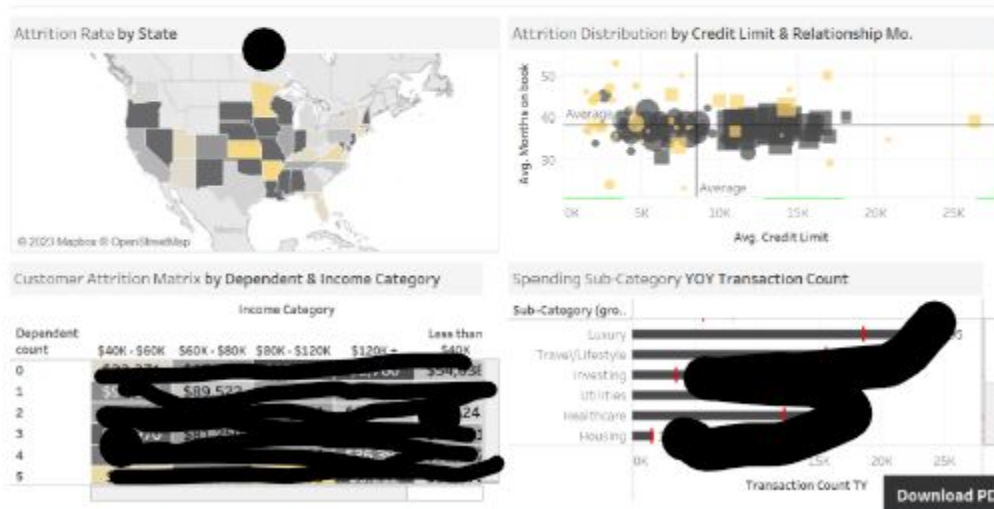
**Lifestyle & Financial Product Usage**

- Examines customer spending habits, travel expenses, and loan preferences.
- **Outcome:** Identifies which financial services (loans, credit cards, travel perks) are most popular.

**Top States with the Most Customers**

- Ranks states based on customer concentration and banking activity.
- **Outcome:** Helps in regional expansion and targeted financial services.

**Credit Usage & Debt Patterns**

- Evaluates how much credit customers use and their repayment behaviors.
- **Outcome:** Assists in risk assessment and customized financial product offerings.

# Tableau

**Historical Data Collection**

- Analyzes Rose Bowl winners from **1902 to 2020** to track long-term trends.
- **Outcome:** Establishes a structured historical dataset for performance analysis.

**Forecasting Future Trends**

- Uses past wins to identify dominant teams and predict future success.
- **Outcome:** Helps measure patterns and build a structured forecast for upcoming years.