



STRATEGIC DECISION-MAKING USING POWER-BI

Submitted To : Dr. Arpit Yadav

Submitted By: Amartya Majumder

Table of Contents

1. PROBLEM STATEMENT	2
2. DATA REQUIREMENT.....	4
Data Requirements	4
3. DATA COLLECTION.....	6
4. DATA VALIDATION (Bias, Transparency, and Reliability)	9
5. DATA CLEANING.....	12
6. Tools Selection.....	16
1. Data Collection Tools	16
2. Data Validation Tools.....	16
3. Data Cleaning Tools	16
4. Data Analysis Tools.....	17
5. Data Visualization Tools.....	17
7. GRAPHS / CHART	18
1. Overview Charts	18
2. Distribution Charts	18
3. Relationship Charts	18
4. Group Comparisons.....	18
5. Categorical Analysis	19
6. Time-Based Analysis	19
7. Outcome Analysis	19
8. Correlation and Interaction	19
9. Advanced Interactive Dashboards	19
8. DASHBOARD	21
9. STORYTELLING/BUSINESS IMPACT.....	24

1. PROBLEM STATEMENT

The rise in diabetes prevalence globally poses a significant challenge for healthcare systems. By leveraging the dataset provided, which contains critical health parameters like glucose levels, BMI, insulin, age, and diabetes outcomes, we aim to uncover actionable insights to predict and prevent diabetes more effectively. The central question is:

"How can healthcare providers harness data analytics to identify diabetes risk factors, predict outcomes, and design targeted interventions for specific demographics?"

Objectives

1. Risk Factor Analysis:

- Explore how key health metrics such as glucose, BMI, insulin, and pregnancy frequency correlate with diabetes prevalence.
- Segment the data into meaningful categories (e.g., by BMI category, age group) to identify high-risk populations.

2. Demographic Insights:

- Analyze the impact of age, gender, and BMI category on diabetes outcomes to uncover vulnerable groups.
- Determine if specific demographic groups (e.g., older individuals with obesity) exhibit unique risk patterns.

3. Prediction and Early Intervention:

- Utilize the dataset to develop predictive insights, such as identifying early warning signs in prediabetic individuals.
- Provide healthcare providers with actionable recommendations for early interventions.

4. Personalized Healthcare Strategies:

- Design prevention and management plans tailored to each demographic group (e.g., dietary plans, exercise routines).
- Improve the efficacy of healthcare delivery by focusing on specific needs of individuals based on their risk profile.

5. Outcome Monitoring:

- Use diabetes outcome data (diabetic vs. non-diabetic) to evaluate the success of interventions.
- Create a framework for continuous monitoring and adjustment of healthcare strategies.

Expected Outcomes

1. Data-Driven Insights:

- A comprehensive understanding of the factors contributing to diabetes prevalence, enabling healthcare providers to focus on the most impactful areas.
 - Identification of high-risk segments based on BMI, glucose levels, and age.
2. **Predictive Tools:**
 - Development of data models or dashboards that can predict diabetes outcomes based on input health metrics.
 - Increased efficiency in identifying prediabetic individuals for early intervention.
 3. **Targeted Interventions:**
 - Creation of tailored programs addressing the specific needs of high-risk populations (e.g., educational campaigns, lifestyle modifications).
 - Better resource allocation by healthcare providers to areas with the highest potential impact.
 4. **Improved Health Outcomes:**
 - Reduction in the incidence of diabetes through early diagnosis and preventive strategies.
 - Enhanced quality of life for individuals at risk or diagnosed with diabetes by providing personalized care.
 5. **Scalable Solutions:**
 - A replicable model that can be applied to other chronic diseases, ensuring long-term benefits for healthcare systems.
 - A dashboard or reporting tool for continuous monitoring of diabetes metrics and outcomes.

2. DATA REQUIREMENT

Pregnancies	Glucose	BloodPressure	SkinThickness	Insulin	BMI	DiabetesPedigreeFunction	Age	Outcome	BMI Category	Age Group
8	183	64	0	0	23.3	0.672	32	Diabetic	Normal	21-40
5	116	74	0	0	25.6	0.201	30	Non-Diabetic	Overweight	21-40
10	115	0	0	0	35.3	0.134	29	Non-Diabetic	Obese	21-40
4	110	92	0	0	37.6	0.191	30	Non-Diabetic	Obese	21-40
10	168	74	0	0	38	0.537	34	Diabetic	Obese	21-40

To execute the analysis, the following columns are required:

Data Requirements

To address the problem statement and achieve the stated objectives, the following columns from the dataset are essential. Each column's relevance is aligned with the dataset and its role in the analysis:

Key Columns and Their Roles

1. Pregnancies:

- **Purpose:** Analyze how the number of pregnancies influences diabetes risk, particularly among women.
- **Objective Alignment:** Helps in identifying demographic-specific risk factors (e.g., maternal health impact).

2. Glucose:

- **Purpose:** Examine glucose levels as a primary indicator of diabetes risk.
- **Objective Alignment:** Central to predicting diabetes outcomes and identifying prediabetic individuals.

3. BloodPressure:

- **Purpose:** Assess the correlation between blood pressure levels and diabetes prevalence.
- **Objective Alignment:** Helps in identifying additional health risks associated with diabetes.

4. SkinThickness:

- **Purpose:** Investigate its role as a proxy for body fat and its relationship with diabetes.
- **Objective Alignment:** Supports the analysis of obesity-related factors influencing diabetes.

5. Insulin:

- **Purpose:** Study insulin levels to detect signs of insulin resistance or deficiency, which are key in diabetes.
- **Objective Alignment:** Vital for identifying metabolic issues related to diabetes onset.

6. BMI (Body Mass Index):

- **Purpose:** Categorize individuals into BMI groups (e.g., obese, overweight) to study their impact on diabetes.

- **Objective Alignment:** Crucial for segmenting populations and tailoring preventive strategies.
- 7. **Diabetes Pedigree Function:**
 - **Purpose:** Analyze the genetic predisposition of individuals to diabetes.
 - **Objective Alignment:** Adds depth to the risk analysis by considering hereditary factors.
- 8. **Age:**
 - **Purpose:** Segment individuals into age groups to explore how age influences diabetes risk.
 - **Objective Alignment:** Key to designing age-specific preventive and management strategies.
- 9. **Outcome:**
 - **Purpose:** Define the binary classification (1 = Diabetic, 0 = Non-Diabetic) for the target variable.
 - **Objective Alignment:** Essential for analyzing and predicting diabetes outcomes.

3. DATA COLLECTION

1. Data Source

Data about diabetes can be sourced from multiple channels, depending on the purpose of the analysis. In this case, the dataset appears structured and is focused on individual health parameters associated with diabetes, indicating a reliable source. Potential sources include:

1. **Medical Institutions:**
 - Data from hospitals, clinics, and healthcare providers.
 - Includes records from patient diagnosis, laboratory tests, and ongoing treatment.
2. **Research Studies:**
 - Clinical studies specifically targeting diabetes, where participants are tested for key health metrics.
 - Typically includes high-quality, verified data.
3. **Open-Source Repositories:**
 - Platforms like **Kaggle** or the **UCI Machine Learning Repository** often host diabetes datasets for academic and research purposes.
 - The dataset provided may come from such a repository.
4. **National Health Databases:**
 - Health and demographic data collected through national programs, such as the **National Diabetes Information Clearinghouse (NDIC)** or similar organizations.
5. **Wearables and IoT Devices:**
 - Fitness trackers, glucometers, and continuous glucose monitors (CGMs) are capable of real-time data collection for parameters like glucose, BMI, and physical activity.

2. Data Collection Methods

The methods of data collection vary depending on the source and tools used. Key methods applicable to diabetes-related datasets include:

1. **Clinical Trials:**
 - In a controlled environment, participants are monitored for health parameters.
 - Involves direct measurements such as:
 - Blood glucose levels after fasting (Fasting Plasma Glucose Test).
 - Glucose tolerance test.
 - HbA1c tests to determine average blood sugar levels over 3 months.
2. **Electronic Health Records (EHR):**
 - Data stored digitally by hospitals or clinics during routine visits or treatments.
 - Includes longitudinal data for patients, such as their health history and progression of diabetes.
3. **Surveys and Questionnaires:**
 - Surveys conducted among targeted populations (e.g., diabetic patients or individuals at risk).

- Collect both self-reported information (age, pregnancy history) and clinically measured values.
- 4. **Wearable and Monitoring Devices:**
 - Devices like CGMs and smartwatches measure health metrics in real-time.
 - Track glucose, physical activity, and even insulin levels in some advanced devices.
- 5. **Community Screening Camps:**
 - Healthcare organizations or NGOs conduct free diabetes check-ups in rural and urban areas.
 - Useful for identifying undiagnosed diabetes cases.

3. Tools for Data Collection

Efficient data collection requires reliable tools and systems that ensure accuracy and minimize errors:

1. **Medical Instruments:**
 - **Glucometers:** Measure blood glucose levels instantly.
 - **Bioimpedance Scales:** Provide BMI, body fat percentage, and other metrics.
 - **Sphygmomanometers:** Record blood pressure.
2. **Digital Platforms:**
 - **REDCap, EpiData:** Software for securely managing survey data and clinical trial information.
 - **Mobile Apps and Portals:** Patient-reported data collection systems, ensuring accessibility.
3. **Data Warehousing:**
 - Cloud platforms or local servers to store and manage collected data securely.
 - Ensures integration of datasets from multiple sources.

4. Challenges in Data Collection

Despite the availability of robust systems, data collection faces several challenges:

1. **Accuracy and Reliability:**
 - Measurements such as glucose levels and BMI can vary due to calibration errors or human input errors.
2. **Incomplete Data:**
 - Missing entries for critical variables like insulin or glucose due to negligence or lack of follow-up.
3. **Data Privacy and Compliance:**
 - Collecting and managing sensitive health data requires adherence to legal frameworks such as:
 - **HIPAA (USA):** For protecting patient information.
 - **GDPR (EU):** For protecting personal data of EU citizens.
 - **Indian IT Act:** Applicable for healthcare data in India.
4. **Resource Constraints:**
 - Lack of access to advanced tools in rural areas or underfunded programs.

5. Sample Bias:

- Overrepresentation of specific demographic groups, such as urban populations or specific age groups, leading to skewed analysis.

5. Future Data Requirements

To improve the dataset's comprehensiveness and address gaps, future data collection efforts should focus on:

1. Lifestyle Metrics:

- Data on physical activity, diet patterns, and sleep habits.
- Smoking and alcohol consumption as potential diabetes risk factors.

2. Socioeconomic Variables:

- Information on income, education level, and occupation.
- Provides insights into the social determinants of diabetes risk.

3. Psychological Factors:

- Stress levels and mental health indicators.
- Explore the link between chronic stress and diabetes.

4. Longitudinal Data Collection:

- Monitor the same individuals over a long period to identify trends, patterns, and progression of diabetes.

5. Geographical and Cultural Data:

- Region-specific factors, such as urban vs. rural disparities in healthcare access.
- Dietary preferences based on cultural habits.

6. Integration of Wearables:

- Leveraging IoT devices for continuous, real-time data monitoring.
- Enhances data granularity and accuracy.

4. DATA VALIDATION (Bias, Transparency, and Reliability)

Data validation is a critical step in ensuring the dataset's accuracy, fairness, and trustworthiness. For diabetes-related data, validation focuses on eliminating errors, identifying biases, ensuring transparency in data collection and processing, and establishing reliability to generate meaningful insights.

1. Addressing Bias

Bias in a dataset occurs when certain groups, categories, or variables are overrepresented, underrepresented, or inaccurately represented, leading to skewed insights.

Types of Bias and Mitigation:

1. Sampling Bias:

- **Problem:** The dataset might overrepresent certain age groups, BMI categories, or geographical locations (e.g., urban areas over rural).
- **Mitigation:**
 - Ensure proportional representation across age groups, genders, BMI categories, and geographical regions.
 - Use stratified sampling techniques to ensure the inclusion of diverse subpopulations.

2. Measurement Bias:

- **Problem:** Inaccurate or inconsistent measurements of health metrics like glucose, insulin levels, or BMI.
- **Mitigation:**
 - Use standardized measurement procedures (e.g., calibrated glucometers and weight scales).
 - Train data collectors to reduce human error.

3. Response Bias:

- **Problem:** Self-reported data (e.g., pregnancies or age) might be inaccurate due to recall issues or intentional misreporting.
- **Mitigation:**
 - Cross-validate self-reported data with medical records or clinical measurements where possible.
 - Use direct measurement tools instead of relying solely on self-reported data.

4. Algorithmic Bias:

- **Problem:** If the dataset is used for predictive models, training on biased data may result in unfair predictions (e.g., favoring certain BMI categories or age groups).
- **Mitigation:**
 - Perform fairness checks on models to ensure unbiased outputs.
 - Use diverse datasets to train machine learning models.

2. Ensuring Transparency

Transparency is crucial to building trust in the dataset and the insights derived from it. It involves clearly documenting how data is collected, processed, and validated.

Steps for Transparency:

1. Documentation of Data Collection:

- Clearly document the methodology for data collection, including:
 - Tools and techniques used (e.g., glucometers, BMI calculations).
 - Locations, demographics, and sample sizes.

2. Data Processing and Cleaning:

- Provide detailed logs of data preprocessing steps:
 - Handling of missing values (e.g., imputation methods used).
 - Removal of duplicates or anomalies.
 - Categorization criteria for BMI and age groups.

3. Access to Metadata:

- Ensure the availability of metadata that explains:
 - Column definitions and units (e.g., glucose measured in mg/dL).
 - Data sources (e.g., clinical trials, surveys).

4. Stakeholder Communication:

- Share the dataset and insights with stakeholders, explaining any limitations, assumptions, or potential biases.

3. Establishing Reliability

Reliability ensures that the dataset is consistent, reproducible, and valid for deriving meaningful conclusions.

Key Steps for Reliability:

1. Consistency Checks:

- Verify the uniformity of data collection procedures across different samples or regions.
- Ensure all glucose, insulin, and BMI values align with known clinical thresholds.

2. Outlier Detection:

- Identify and address outliers in the dataset:
 - For example, unusually high insulin or glucose values that might indicate data entry errors.

3. Handling Missing Data:

- Determine the extent of missing values in critical columns (e.g., glucose, BMI).
- Use appropriate imputation techniques:
 - Mean or median substitution for numerical variables.
 - Predictive models for imputing missing data.

4. Validation Against Benchmarks:

- Compare data trends with established benchmarks from healthcare organizations:
 - Average glucose levels across age groups.
 - BMI distributions in diabetic and non-diabetic populations.

5. Reproducibility:

- Ensure that the data analysis pipeline produces consistent results when re-executed.
- Use version-controlled systems to manage data cleaning and processing steps.

4. Tools and Techniques for Validation

1. Statistical Analysis:

- Perform distribution checks to ensure normalcy in key variables (e.g., glucose, BMI).
- Use correlation matrices to validate relationships between variables (e.g., age vs. BMI).

2. Automated Data Validation Tools:

- Leverage tools like **OpenRefine**, **Talend**, or **DataRobot** for automated error detection and correction.

3. Human Review:

- Conduct manual audits of a subset of the dataset to verify automated results.

4. Bias Detection Algorithms:

- Use fairness and bias detection algorithms (e.g., SHAP, Fairlearn) to identify and mitigate potential biases.

Expected Outcome of Data Validation

By addressing bias, ensuring transparency, and establishing reliability, the validated dataset will enable:

1. Accurate and Unbiased Insights:

- Reliable conclusions on diabetes trends and risk factors.
- Meaningful correlations between BMI, age, glucose, and other variables.

2. Fair and Ethical Usage:

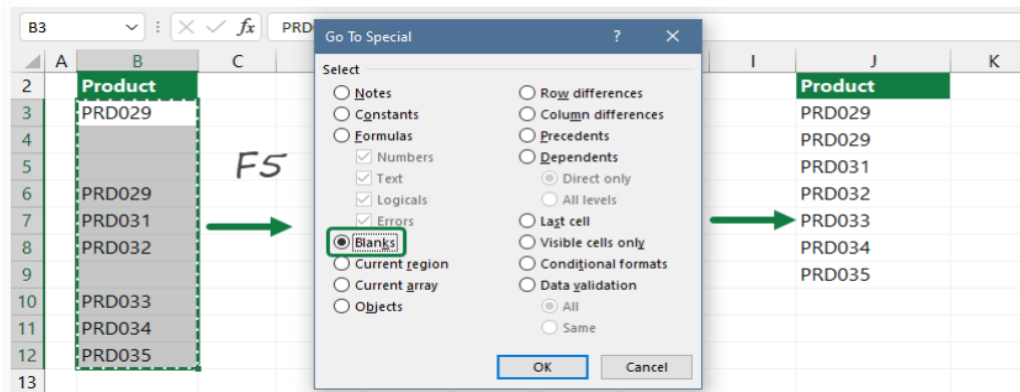
- Equitable representation of all demographic groups.
- Reduced risk of skewed predictions or conclusions.

3. Improved Decision-Making:

- Empower healthcare providers, researchers, and policymakers to make data-driven decisions for diabetes prevention.

5. DATA CLEANING

Purpose: To prepare the dataset for meaningful analysis by addressing errors, inconsistencies, and irrelevant data.



Data

cleaning involves detecting and rectifying errors, handling missing values, and standardizing the dataset.

Data cleaning is the process of preparing raw data for analysis by identifying and rectifying errors, inconsistencies, and inaccuracies. It ensures the dataset is accurate, complete, and ready for insightful analysis. Cleaning the diabetes dataset is crucial for producing reliable insights that drive meaningful healthcare decisions.

In-Depth Steps for Data Cleaning

1. Handling Missing Values

Missing values in critical columns (e.g., glucose, BMI, insulin) can reduce the effectiveness of the analysis if not addressed.

- **Identification:**
 - Use statistical tools to identify missing or null values in the dataset. For instance:
 - Check for blank cells or "NaN" values.
 - Compute the percentage of missing data per column.
- **Approach to Handle Missing Data:**
 - **Imputation:**
 - For numerical columns (e.g., glucose, insulin): Use the mean, median, or regression-based predictions.
 - For categorical columns (e.g., BMI category, outcome): Use the mode or predictive modeling.
 - **Dropping Rows:**
 - If a significant portion of a row's data is missing, consider removing it (ensuring this doesn't lead to sample bias).
 - **Flagging Missing Data:**
 - Add a new column to indicate rows with imputed or missing data for transparency.

2. Removing Duplicates

Duplicate records can distort analysis and produce skewed results.

- **Detection:**
 - Check for duplicate rows based on unique identifiers or combinations of key attributes (e.g., age, pregnancies, BMI, and glucose).
- **Resolution:**
 - Remove duplicates while ensuring important records are retained.
- **Tool:**
 - Use software like Python's pandas library (`.drop_duplicates()`) or Excel's "Remove Duplicates" function.

3. Handling Outliers

Outliers are extreme values that can distort statistical analysis and skew results.

- **Detection:**
 - Use visualizations (e.g., boxplots, histograms) to identify outliers in variables like glucose, BMI, or insulin levels.
 - Use statistical methods (e.g., Z-scores or the Interquartile Range [IQR]):
 - Values beyond 1.5 times the IQR from the 1st or 3rd quartile are considered outliers.
- **Approach:**
 - Verify the accuracy of the outlier:
 - For valid clinical outliers (e.g., extremely high glucose in diabetic patients), retain them.
 - For erroneous outliers (e.g., negative insulin levels), correct or remove them.
 - Transform data if required:
 - Use logarithmic or square root transformations for highly skewed distributions.

4. Standardizing Data Formats

Inconsistent formats can hinder data analysis and lead to errors.

- **Age:**
 - Ensure all age values are numeric and within a valid range (e.g., 0–120 years).
- **BMI, Glucose, and Insulin:**
 - Ensure all values are in consistent units (e.g., mg/dL for glucose).
- **Date Formats:**
 - Standardize any date columns to a uniform format (e.g., YYYY-MM-DD).

5. Correcting Errors

Errors in data entry or recording can lead to inaccuracies.

- **Common Errors:**
 - Negative values for metrics like age, BMI, or glucose.
 - Impossible pregnancy values (e.g., negative or extremely high numbers).
- **Resolution:**
 - Replace erroneous values with accurate entries if available or impute them based on statistical methods.

6. Consistency in Categorical Data

Categorical variables like BMI category or outcome must be consistent.

- **Detection:**
 - Identify spelling errors or inconsistent capitalization (e.g., "Normal" vs. "normal").
- **Resolution:**
 - Standardize all categories using a predefined list (e.g., Normal, Obese, Overweight, Underweight).
- **Tool:**
 - Use Python's pandas library (`.str.lower()` or `.str.capitalize()`) or Excel's "Find and Replace."

7. Encoding Categorical Variables

To make the dataset ready for machine learning models, categorical variables may need to be encoded.

- **Approaches:**
 - **One-Hot Encoding:** Create binary columns for each BMI category (e.g., Normal, Obese).
 - **Label Encoding:** Assign numerical values to categories (e.g., Normal = 0, Obese = 1).

8. Resolving Inconsistencies Across Related Variables

Some variables may have logical dependencies (e.g., BMI and glucose levels should correlate to some extent).

- **Cross-Validation:**
 - Check for inconsistencies between related columns.
 - Example: Pregnancies should be zero for males or individuals under 12 years old.

9. Creating Derived Variables

Derived variables can provide additional insights or simplify analysis.

- **Examples:**

- BMI Category: Derived from BMI values based on predefined ranges.
- Risk Score: Create a composite score combining glucose, BMI, and age.

10. Final Dataset Validation

Once data cleaning is complete:

- Recheck for missing values, duplicates, and inconsistencies.
- Validate data distributions using summary statistics (mean, median, standard deviation).

6. Tools Selection

Selecting the right tools is critical to ensure efficient data processing, analysis, and visualization. The choice of tools should align with the dataset's structure, volume, and the objectives of the analysis. For the diabetes dataset, the following categories of tools are recommended:

1. Data Collection Tools

Microsoft Excel or Google Sheets:

- Best for smaller datasets or manual data entry.
- Provides quick validation, formatting, and data previews.

Database Management Systems (DBMS):

- **MySQL or PostgreSQL:** Ideal for managing structured data at scale.
- **MongoDB:** Useful if unstructured or semi-structured data is also involved.

APIs or Web Scraping Tools:

- **Python Libraries:**
 - requests and BeautifulSoup for web scraping.
 - Pandas to import data from APIs or JSON files.

2. Data Validation Tools

- **Python with Libraries:**
 - Pandas: For detecting missing values, duplicates, and outliers.
 - SciPy or Statsmodels: For statistical validation and checking bias in distributions.
 - Fairlearn: Evaluates bias and fairness in datasets.
- **R Programming:**
 - Useful for statistical analysis and bias detection.
- **Excel/Google Sheets:**
 - Suitable for quick cross-validation, filtering, and creating pivot tables.

3. Data Cleaning Tools

- **Python Libraries:**
 - **Pandas:** For handling missing data, duplicates, and outliers.
 - **Numpy:** For data transformations and numerical calculations.
 - **Scikit-learn:** For preprocessing steps like scaling, encoding, and imputation.
 - **Matplotlib and Seaborn:** For identifying outliers visually.
- **R Programming:**
 - Packages like tidyverse or dplyr for data cleaning and preprocessing.

- **Data Preparation Tools:**
 - **Trifacta or Alteryx:** Advanced tools for data wrangling and transformation.

4. Data Analysis Tools

- **Python:**
 - Libraries like Scikit-learn for modeling and predictive analysis.
 - Matplotlib, Seaborn, and Plotly for visualization.
- **R Programming:**
 - ggplot2: For data visualization.
 - caret: For machine learning models.
- **Tableau or Power BI:**
 - For dynamic dashboards and exploratory data analysis.
- **Excel/Google Sheets:**
 - For descriptive statistics, summary analysis, and pivot tables.

5. Data Visualization Tools

- **Tableau:**
 - For creating interactive dashboards and visual storytelling.
- **Microsoft Power BI:**
 - For building business-friendly visualizations.
- **Python Visualization Libraries:**
 - **Matplotlib and Seaborn:** Ideal for basic plots like histograms, scatter plots, and boxplots.
 - **Plotly or Dash:** For interactive and web-based visualizations.
- **Excel:**
 - Quick charts for smaller datasets.

7. GRAPHS / CHART

Visualizations play a vital role in understanding the dataset, uncovering trends, and communicating insights effectively. Here's a detailed list of recommended charts/graphs aligned with the diabetes dataset and objectives:

1. Overview Charts

- **Total Individuals (KPI):**
 - **Type:** Numeric Display or Card.
 - **Purpose:** Shows the total number of individuals in the dataset, providing a snapshot of the dataset size.
- **Average Metrics (KPI):**
 - **Type:** Numeric Displays for Average BMI, Age, Glucose, etc.
 - **Purpose:** Highlights key dataset averages for high-level understanding.

2. Distribution Charts

- **Age Distribution:**
 - **Type:** Histogram or Density Plot.
 - **Purpose:** Understand the age spread and identify the most common age groups.
- **BMI Distribution:**
 - **Type:** Histogram or Boxplot.
 - **Purpose:** Identify outliers and trends in BMI.
- **Glucose Levels Distribution:**
 - **Type:** Violin Plot or Histogram.
 - **Purpose:** Examine glucose levels and check for skewness or abnormalities.

3. Relationship Charts

- **Pregnancies vs. Outcome:**
 - **Type:** Stacked Bar Chart.
 - **Purpose:** Visualize the relationship between the number of pregnancies and diabetes outcomes.
- **BMI vs. Glucose Levels:**
 - **Type:** Scatter Plot.
 - **Purpose:** Analyze if higher BMI correlates with higher glucose levels or risk of diabetes.
- **Age vs. Insulin Levels:**
 - **Type:** Bubble Chart.
 - **Purpose:** Explore patterns across age groups and insulin levels.

4. Group Comparisons

- **Average Metrics by Age Group:**
 - **Type:** Bar Chart with Categories (e.g., 21–40, 41–60, 61+).

- **Purpose:** Compare averages like BMI, Glucose, and Insulin across age groups.
- **Outcome by Age and BMI Category:**
 - **Type:** Heatmap.
 - **Purpose:** Identify combinations of age and BMI categories with the highest prevalence of diabetes.

5. Categorical Analysis

- **BMI Categories:**
 - **Type:** Pie Chart or Donut Chart.
 - **Purpose:** Show the proportion of individuals in each BMI category (Normal, Obese, Overweight, Underweight).
- **Outcome by BMI Category:**
 - **Type:** Stacked Bar Chart.
 - **Purpose:** Understand diabetes prevalence in each BMI category.

6. Time-Based Analysis

(If the dataset includes a timeline or sequential data)

- **Trend of Glucose Levels Over Time:**
 - **Type:** Line Chart.
 - **Purpose:** Analyze how glucose levels change over time for individuals or groups.

7. Outcome Analysis

- **Diabetes Outcome by Age:**
 - **Type:** Line Chart.
 - **Purpose:** Show trends in diabetes occurrence across different age groups.
- **Outcome Proportion:**
 - **Type:** Pie Chart or Bar Chart.
 - **Purpose:** Display the proportion of diabetic vs. non-diabetic individuals.

8. Correlation and Interaction

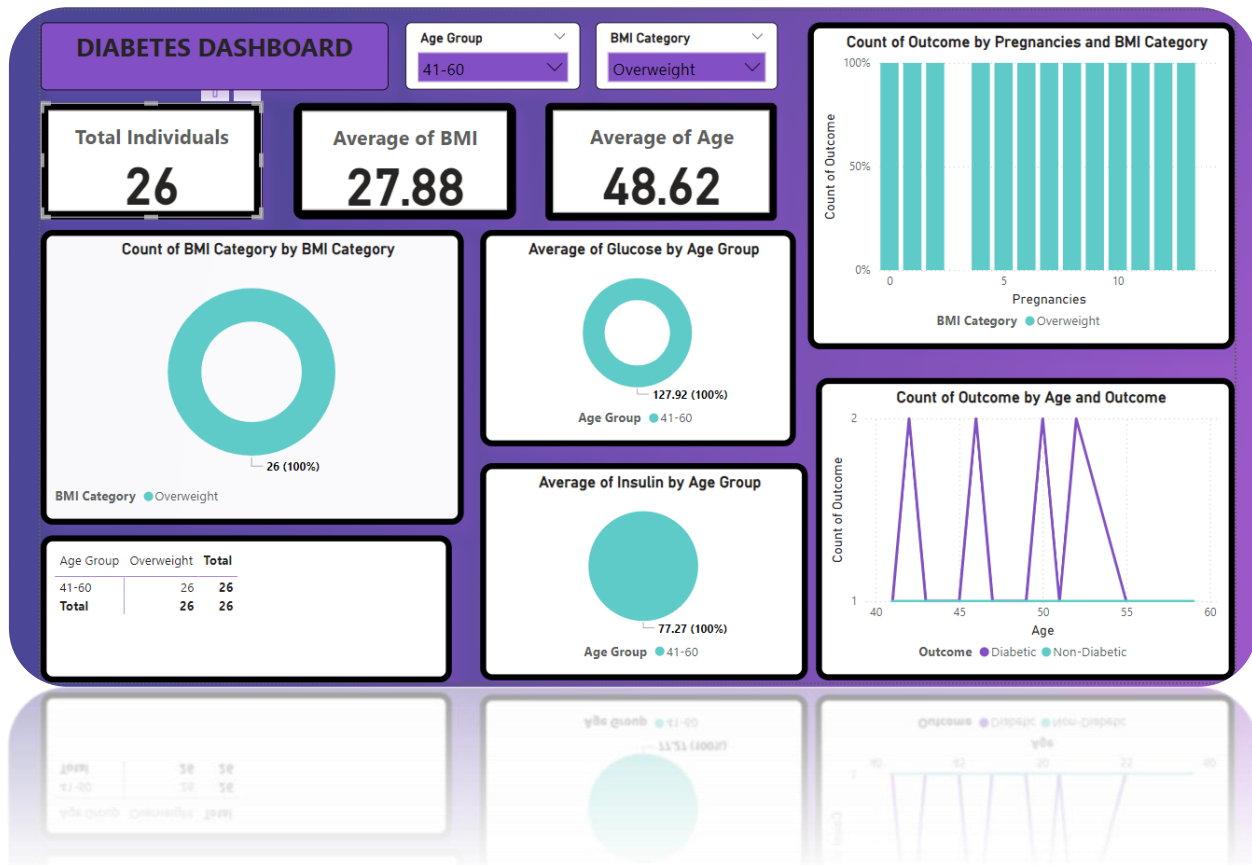
- **Correlation Matrix:**
 - **Type:** Heatmap.
 - **Purpose:** Understand relationships between variables like glucose, BMI, age, and insulin.
- **3D Scatter Plot (Optional):**
 - **Type:** 3D Plot (e.g., Age vs. Glucose vs. Insulin).
 - **Purpose:** Explore multidimensional relationships.

9. Advanced Interactive Dashboards

- **Filters for Custom Views:**

- Include dropdowns to filter by BMI category, age group, or outcome.
 - Dynamic charts update based on user selections.
- **Dynamic KPI Cards:**
 - Change values based on applied filters (e.g., showing average glucose levels for a selected age group).

8. DASHBOARD



This dashboard provides insights into diabetes-related data segmented by **BMI Category**, **Age Group**, and various metrics like **Glucose Levels**, **Insulin Levels**, and **Outcomes** (diabetic or non-diabetic). Here's a breakdown of the analysis using univariate, bivariate, and multivariate approaches:

1. Univariate Analysis

Univariate analysis examines individual variables to identify their distribution and central tendencies.

Key Observations:

1. Total Individuals:

- The dataset includes **26 individuals**, all within the Overweight BMI category and the 41-60 age group.
- There is no diversity in BMI Category or Age Group in this subset, limiting further univariate segmentation.

2. Average of BMI:

- The **mean BMI** for this subset is **27.88**, indicating overweight individuals based on WHO BMI thresholds.
- Since all individuals fall into this category, BMI is uniform across the dataset.
- 3. **Average of Age:**
 - The **average age** is **48.62 years**, reflecting middle-aged individuals (41-60 years age group).
- 4. **Glucose Levels:**
 - The **average glucose level** is **127.92**.
 - This value might be near or above the threshold for concern, depending on diabetic diagnostic standards.
- 5. **Insulin Levels:**
 - The **average insulin level** is **77.27**, which is likely significant for identifying insulin resistance or diabetes.

2. Bivariate Analysis

Bivariate analysis explores relationships between two variables, offering insights into correlations or trends.

Key Observations:

1. **Count of BMI Category by BMI Category:**
 - A redundant chart that reiterates the entire population (100%) is Overweight. It could be omitted for better clarity.
2. **Count of Outcome by Pregnancies and BMI Category:**
 - This bar chart demonstrates the **distribution of diabetic outcomes** across the number of pregnancies within the Overweight BMI category.
 - All pregnancies fall uniformly into the Overweight category with no variation.
 - Further analysis of how pregnancy impacts diabetic outcomes would be insightful but isn't evident here due to the uniform BMI grouping.
3. **Count of Outcome by Age and Outcome:**
 - This chart shows the count of **diabetic vs. non-diabetic individuals** segmented by age.
 - The data reveals variability in diabetic status with age, suggesting age might be a contributing factor to outcomes.
 - Peaks indicate specific age groups (e.g., 45-55) are more likely to exhibit diabetic outcomes.

3. Multivariate Analysis

Multivariate analysis examines interactions between three or more variables, uncovering deeper relationships.

Key Observations:

1. **Age Group, Glucose, and Insulin Levels:**

- All data points are within the 41-60 age group. However:
 - The average glucose (127.92) and insulin levels (77.27) could be compared with diabetic and non-diabetic outcomes for insights into the physiological profiles.
 - This correlation could suggest a predictive model for identifying diabetes risk.
- 2. **Pregnancies, BMI, and Outcome:**
 - The Pregnancies dimension, combined with BMI Category and Outcome, suggests:
 - Higher pregnancy counts do not necessarily correspond to varying diabetic outcomes in this subset.
 - If combined with additional BMI categories or age brackets, this might yield stronger patterns.
- 3. **Age, Outcome, and Insulin/Glucose Levels:**
 - The **interaction between age, outcome, and insulin/glucose levels** could explain how these metrics interact to influence diabetic status.
 - For example, individuals aged 50-55 with glucose levels above 127.92 and insulin around 77.27 may exhibit trends correlating with diabetes.

Key Insights and Next Steps

- **Uniformity of BMI and Age Group:**
 - The data is restricted to Overweight individuals aged 41-60, which limits broader generalizations.
 - Expanding the dataset to include diverse BMI categories and age groups is crucial.
- **Outcome Trends:**
 - Age appears to have a potential correlation with diabetic outcomes, as seen in the age vs. outcome chart.
 - Further analysis into other variables (e.g., glucose, insulin) influencing outcomes is recommended.
- **Data Utilization:**
 - Pregnancies and BMI outcomes require richer datasets to yield meaningful patterns.
 - Incorporating additional features like Activity Levels, Diet, or Family History would enhance analysis.

Recommended Enhancements to the Dashboard

1. **Diversity in Variables:**
 - Include more BMI categories (Underweight, Normal, Obese).
 - Add more age groups or consider continuous age data.
2. **Advanced Visualizations:**
 - Introduce heatmaps for correlations (e.g., Glucose, Insulin, Outcome).
 - Use scatter plots for relationships between glucose/insulin levels and outcomes.
3. **Drill-Down Filters:**
 - Add interactive filters for deeper segmentation (e.g., by gender, lifestyle factors, or medical history).

9. STORYTELLING/BUSINESS IMPACT

1. Storytelling Framework

The goal of storytelling in this context is to transform raw data insights into a cohesive narrative that highlights the underlying trends, relationships, and potential actions. The storytelling framework for the diabetes dashboard focuses on three key aspects: **problem, analysis, and action.**

1. Problem Identification:

- **Context:** Diabetes is a global health challenge, and early identification of risk factors can significantly improve patient outcomes.
- **Target Population:** Middle-aged, overweight individuals are at a higher risk of developing diabetes. The dataset focuses on individuals aged 41-60 who fall into the Overweight BMI category.
- **Key Concern:** What factors—such as glucose levels, insulin levels, and pregnancy history—contribute to diabetic or non-diabetic outcomes in this specific population?

2. Analysis:

- **Insights from Data:**
 - The average glucose level (127.92) and insulin level (77.27) suggest potential markers for identifying diabetic risk.
 - Age appears to influence diabetic outcomes, with certain subgroups showing a higher likelihood of diabetes.
 - Uniformity in BMI category and age limits the variability for broader generalizations.
- **Patterns and Trends:**
 - Correlations between insulin and glucose levels with diabetic outcomes provide actionable insights.
 - High glucose levels in combination with age-specific factors increase the likelihood of diabetes.

3. Action:

- **Engagement:** Equip healthcare professionals with these insights to focus interventions on high-risk age groups within the overweight category.
- **Impact:** Improved diagnostics and personalized healthcare interventions targeting middle-aged individuals can mitigate the risk of diabetes onset and reduce long-term healthcare costs.

2. Key Messages

- **The Problem:** Diabetes is prevalent among middle-aged, overweight individuals, and early intervention is critical.
- **Insights:**
 - Glucose and insulin levels are key predictors of diabetic outcomes.
 - Age plays a significant role in determining outcomes within this population group.
- **Call to Action:**

- Focus on high-risk individuals with elevated glucose and insulin levels.
- Expand the dataset to include diverse BMI categories and age groups to improve generalizability.

3. Business Impact

Leveraging insights from this analysis can have far-reaching implications in healthcare delivery, preventive measures, and overall well-being.

1. Improved Diagnostics:

- By identifying patterns in glucose, insulin, and age, healthcare professionals can develop **predictive models** for early diabetes detection.

2. Personalized Interventions:

- Tailoring interventions for specific risk groups (e.g., middle-aged, overweight individuals with high glucose levels) enhances effectiveness and efficiency.
- Introduce targeted health campaigns for this group to promote healthy lifestyle changes.

3. Cost Reduction in Healthcare:

- Early detection and management of diabetes can prevent long-term complications, reducing healthcare costs.
- Resources can be allocated more effectively, focusing on high-risk populations.

4. Policy Development:

- The insights can guide public health policies aimed at mitigating diabetes risk in specific demographic groups.

4. Storytelling for Stakeholders

The insights from the analysis can be presented as follows:

• Healthcare Providers:

- Highlight actionable insights, such as glucose and insulin thresholds, to refine diagnosis and treatment protocols.

• Public Health Officials:

- Advocate for policies targeting middle-aged, overweight individuals with specific glucose and insulin levels.

• Patients:

- Share visual and intuitive dashboards to raise awareness and encourage proactive health measures.