



# Diabetes Detection: From Data to Insights

This presentation explores a data analytics project focused on the early detection of diabetes. We will cover the entire process, from data acquisition to generating actionable insights. Our goal is to leverage data for improved healthcare outcomes.

# Identifying & Sourcing Relevant Datasets



## Public Datasets

Pima Indians Diabetes Database is a common public resource.



## Private Datasets

Electronic Health Records (EHR) and patient records offer detailed insights.



## Data Privacy

Adherence to HIPAA regulations is crucial for patient data.

We sourced both public and private datasets for this project. Key attributes included glucose, hbA1c level, BMI, and age. Strict adherence to HIPAA regulations ensured data privacy.





# Cleaning & Handling Missing Values



## Identify Issues

Missing values, outliers, and inconsistencies are common.



## Apply Techniques

Imputation, removal, and modeling address data gaps.

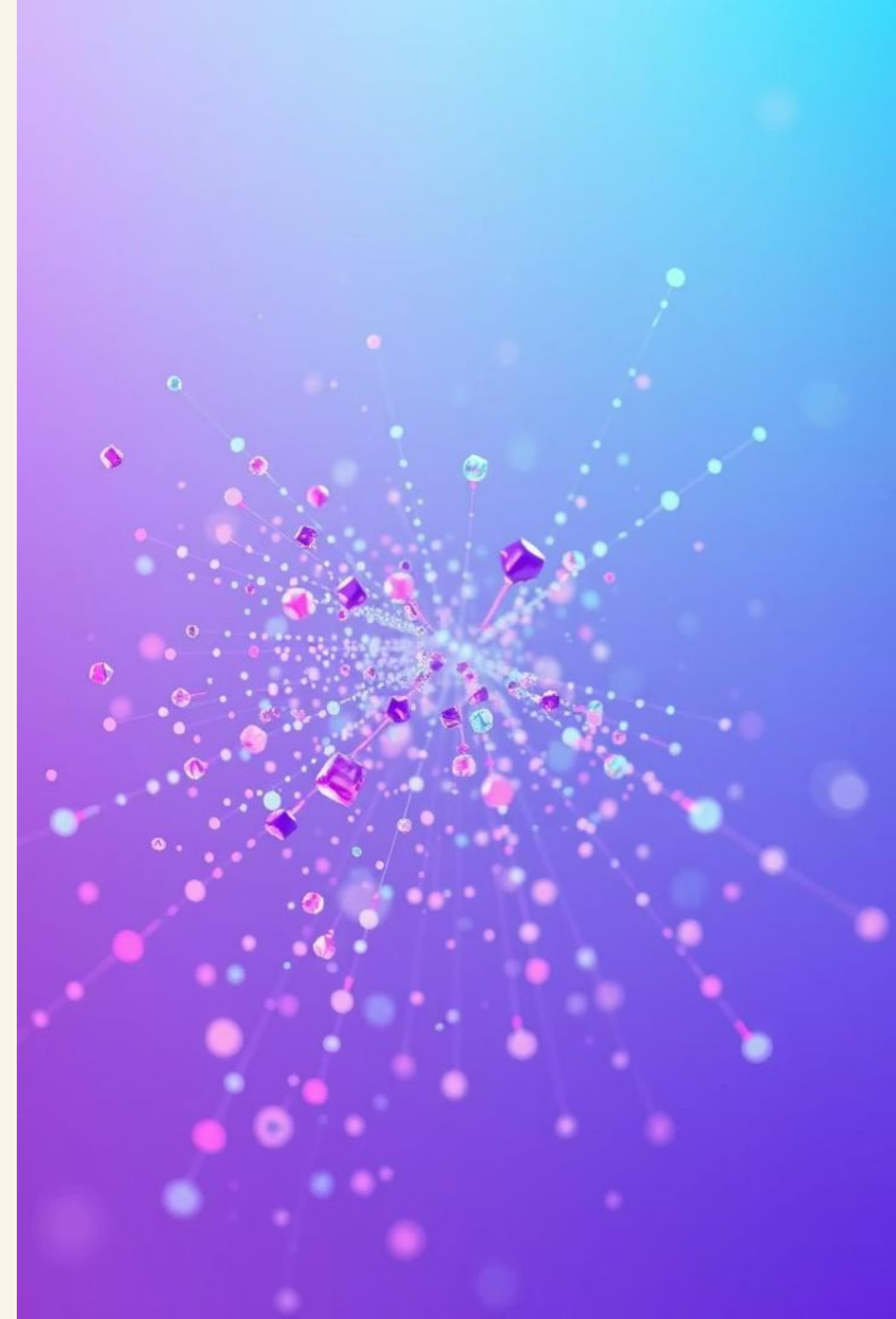


## Assess Impact

Missing values can significantly reduce model accuracy.

Data cleaning addressed missing values, outliers, and inconsistencies.

Techniques like imputation were used, for example, replacing missing BMI values with the mean. This process ensures data quality and improves model reliability.



# Feature Selection & Engineering

## Improve Performance

Relevant features enhance model accuracy and efficiency.

## Combine Data

We combined hbA1c level and glucose measurements for richer insights.



Feature selection and engineering are vital for strong model performance. We employed statistical tests and domain knowledge. Creating new features, like combining insulin and glucose levels, provides deeper insights into diabetes indicators.





# Ensuring Data Integrity & Consistency



## Data Validation

Range checks and cross-field validation ensure accuracy.



## Data Quality

High data quality is critical for model reliability.

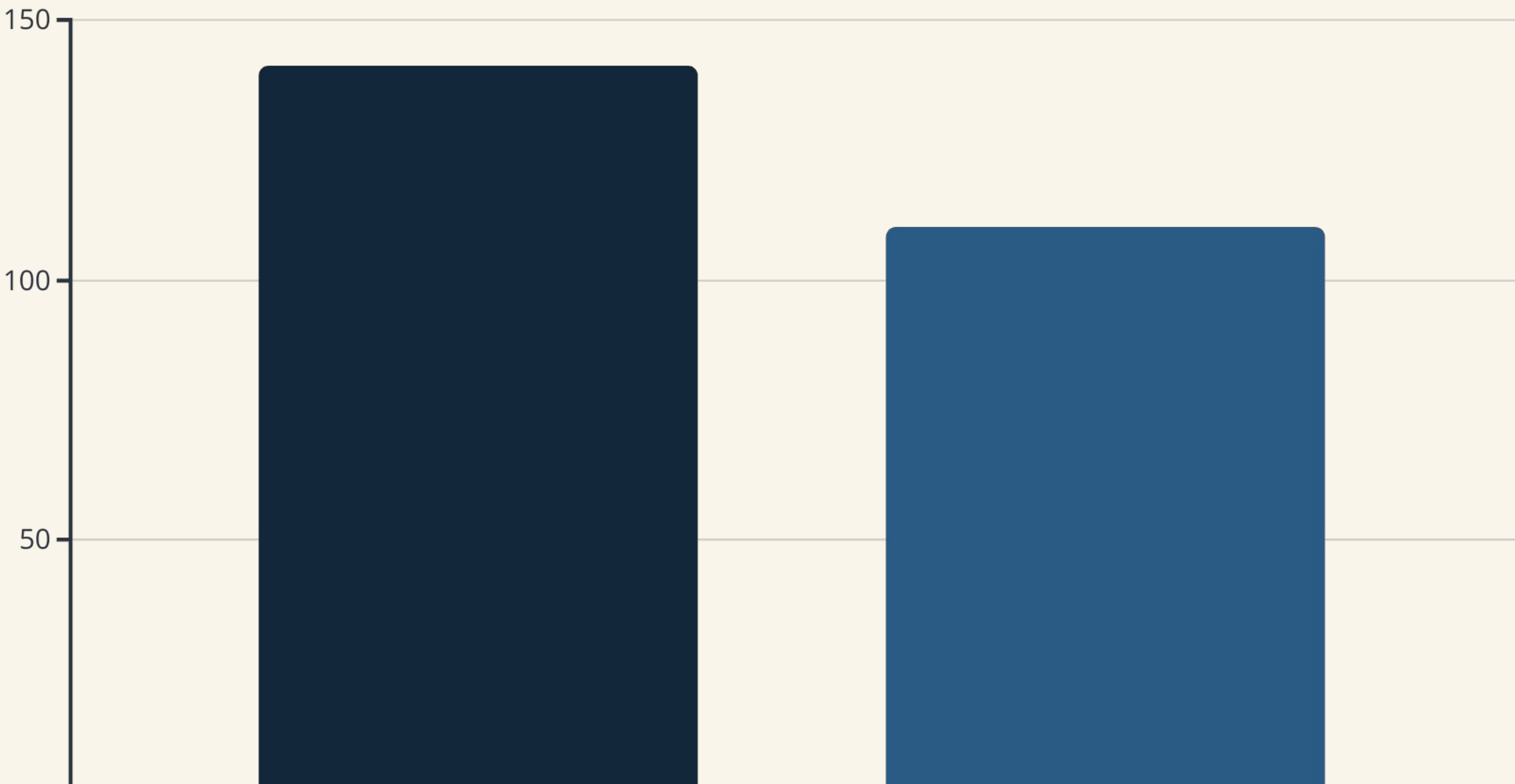


## Integrity Tools

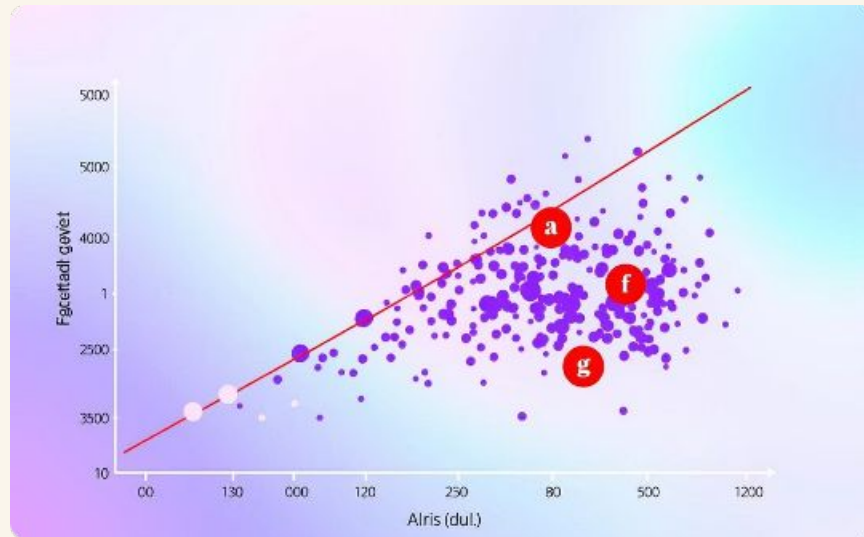
SQL constraints and data profiling maintain consistency.

Data integrity is paramount for reliable models. We used validation techniques like range checks to address inconsistencies. Tools like SQL constraints and data profiling were essential in maintaining high data quality throughout the project.

# Summary Statistics & Insights

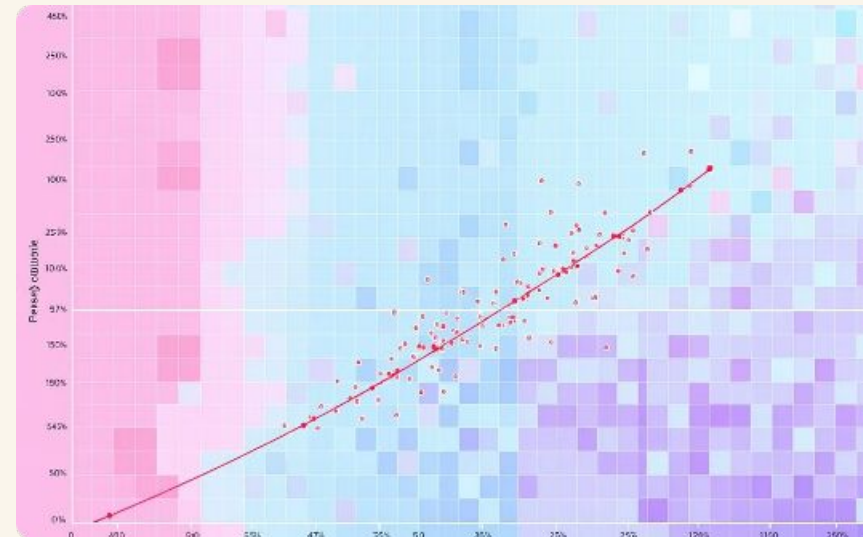


# Identifying Patterns, Trends, & Anomalies



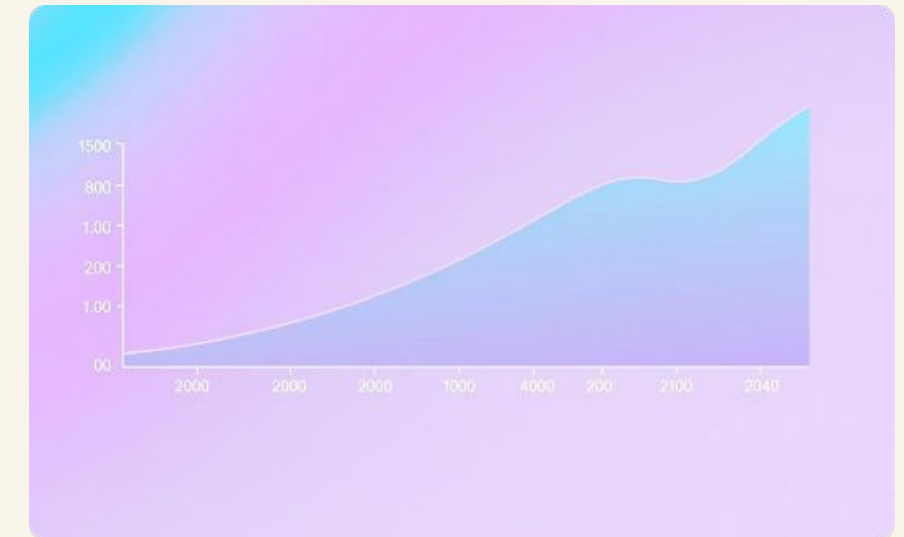
## Scatter Plots

Revealed relationships between glucose levels and BMI.



## Heatmaps

Displayed correlations across various medical features.

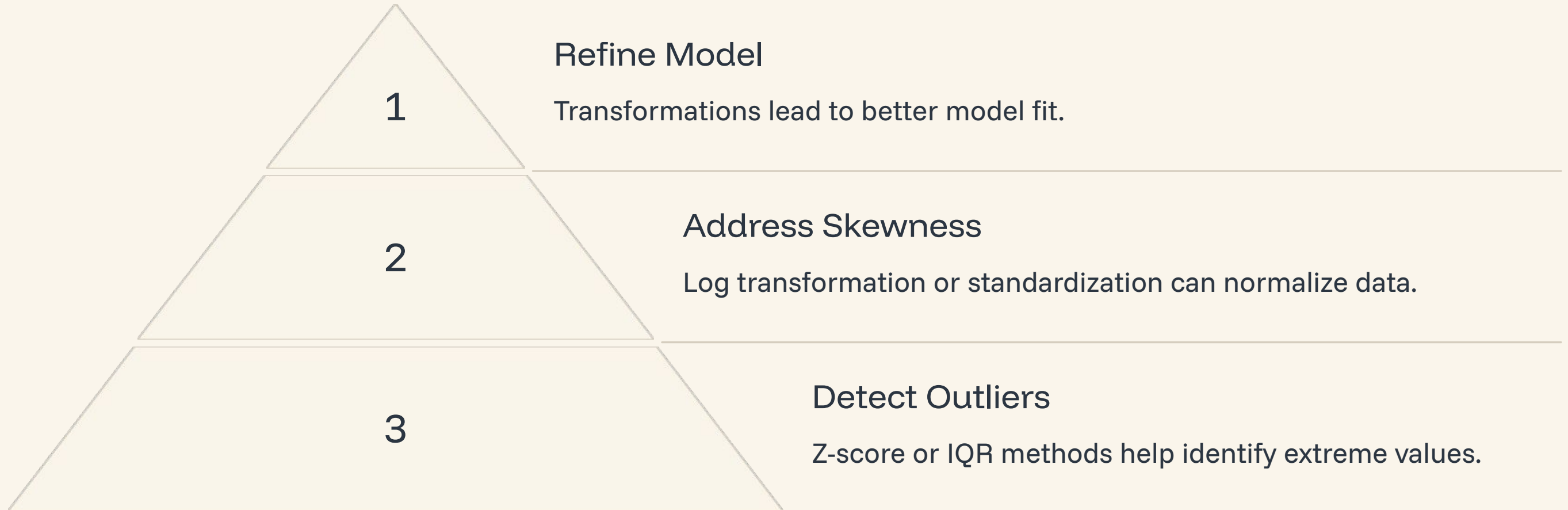


## Trend Analysis

Illustrated changes in diabetes incidence over time.

Data visualization, including scatter plots and heatmaps, helped us identify patterns. We analyzed trends in diabetes incidence and detected anomalies like high glucose levels with unusually low BMI, indicating unique patient profiles.

# Handling Outliers & Data Transformations



Outliers were identified using methods like Z-score and IQR. These extreme values can significantly impact model performance. Data transformations, such as log transformation, were applied to address skewness and improve model fit.



# Initial Visual Representation of Key Findings



Clear and effective visualizations are crucial for communicating insights. We created charts and graphs using tools like Python's Matplotlib and Seaborn. Interactive elements, such as tooltips and filters, enhanced data exploration and storytelling.

# Interpretation & Storytelling with Data



## Actionable Recommendations

Translate insights into concrete steps.



## Target Groups

Identify specific populations for preventative measures.



## Ethical Considerations

Address privacy, bias, and fairness in data use.

The final step involves translating data insights into actionable recommendations. We identified target groups for preventative measures. Ethical considerations, including patient privacy, bias, and fairness, were paramount. Our findings hold potential for significantly reducing diabetes incidence.