

MEDICAL INSURANCE COST ANALYSIS REPORT



Table of Contents

Sl. No	Content	Page No.
1.	Problem Statement	3
2.	Data Requirement	4
3.	Data Collection	5
4.	Data Validation	8
5.	Data Preprocessing	8
6.	Tools used	11
7.	Dashboard	14
8.	Presentation	16

Problem Statement

The medical insurance dataset includes demographic, lifestyle, and health-related data, along with insurance charges for individuals. This data presents an opportunity to analyze the factors influencing insurance costs.

Objective:

To understand the primary factors affecting medical insurance charges and develop predictive models for estimating costs based on customer profiles. Specifically, the goal is to analyze:

- **Demographic Impact:** Examining the role of age, sex, and region in determining insurance costs.
- **Lifestyle and Health Factors:** Understanding the influence of smoking status, BMI, and number of children on charges.
- **Predictive Modeling:** Creating a model to predict insurance costs based on the demographic and lifestyle characteristics.

Key Questions:

1. How do factors like age, BMI, and smoking status influence insurance charges?
2. Are there significant regional differences in insurance costs?
3. What demographic groups are associated with higher insurance charges?
4. How accurately can we predict insurance costs based on the available features?

By addressing these questions, healthcare and insurance companies can better understand risk factors, set fairer premiums, and create personalized insurance plans.

Unveiling Factors Influencing Medical Insurance Costs

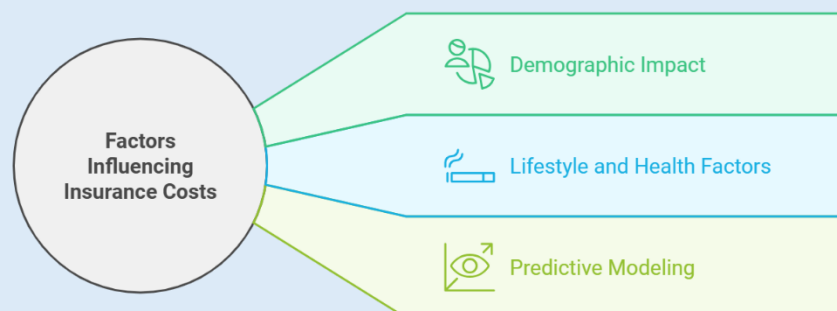


Figure 1.1

Data Requirements

To utilize this dataset effectively for medical insurance cost analysis, the following requirements should be met:

1. **Age:** Integer values, ideally within a realistic range (e.g., 18–100 years), as insurance costs are typically age-dependent.
2. **Sex:** Categorical variable with two values (male, female). It should be encoded or standardized for modeling purposes if necessary.
3. **BMI:** Continuous variable representing the individual's BMI. Values should ideally fall within a reasonable range for adults (e.g., 15–50).
4. **Children:** Integer representing the number of dependents. It should be non-negative, with typical values ranging from 0 to a few children.
5. **Smoker:** Binary categorical variable (yes or no) indicating smoking status, as this greatly influences health risk.
6. **Region:** Categorical variable with predefined geographic locations (southwest, southeast, northwest, northeast). Values should be consistent for analysis.
7. **Charges:** Continuous variable representing insurance charges in currency units. This should be a non-negative float, as it represents costs billed to the individual.

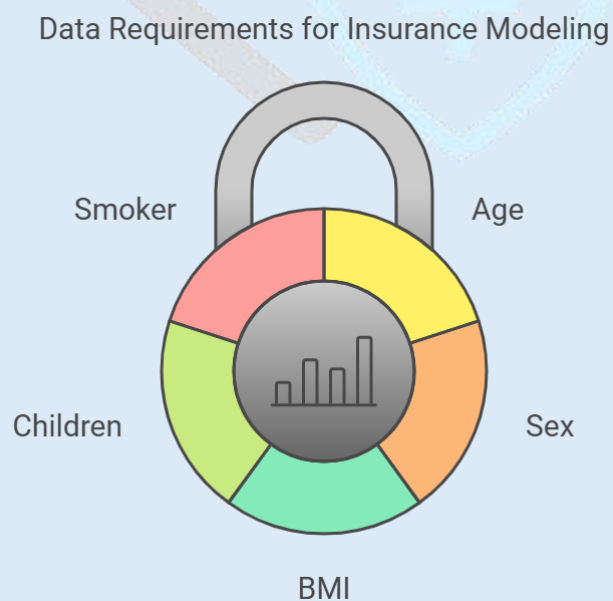


Figure 1.2

Data Collection

Data Collection Categories

1. Personal Demographics

- **Age:** Directly ask the individual's age or calculate it based on their date of birth.
- **Sex:** Collect information on gender, typically through self-reporting (male, female, or other, if applicable).

2. Health Information

- **BMI (Body Mass Index):**
 - Collect the height and weight of each individual to calculate BMI (BMI = weight in kg / (height in meters)²).
 - This could be collected via medical records, self-reporting, or direct measurement if feasible.
- **Smoking Status:**
 - Directly ask individuals if they currently smoke or have smoked in the past.
 - If collecting from health records, ensure the data indicates whether the individual is a current smoker, a former smoker, or has never smoked.

3. Family and Dependents

- **Number of Children:** Ask individuals how many dependents or children are covered under their insurance plan. This can be collected directly from the insured individual or from the insurance application form.

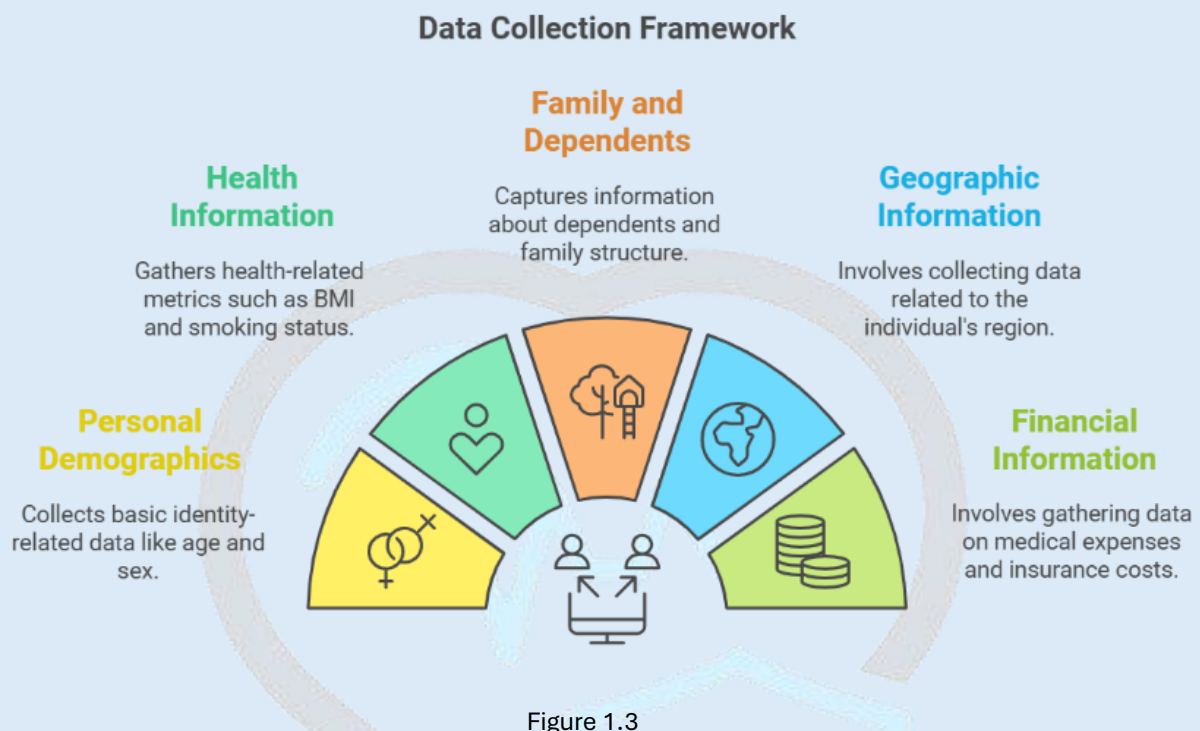
4. Geographic Information

- **Region:** Collect the individual's location data, ideally limited to the region (such as Southwest, Southeast, Northwest, or Northeast).
- You may collect this via the individual's home address, but store only the broader region to protect privacy.

5. Financial Information

- **Charges (Insurance Cost):**

- Collect the total medical expenses charged to the individual by insurance over a specified period.
- This information may be available from insurance companies, hospital billing departments, or directly from the individual if they have access to their billing records.



Data Collection Methods

1. Surveys and Questionnaires

- Conduct surveys or questionnaires that individuals can fill out to provide data on age, sex, BMI (or height and weight for calculation), smoking status, children, and region.
- Surveys can be administered online, over the phone, or in person, depending on accessibility.

2. Medical and Insurance Records

- Collaborate with healthcare providers or insurance companies to access anonymized or consented individual records.

- This approach can ensure data accuracy, especially for BMI, smoking status, and charges.

3. Self-Reporting and Interviews

- Interview individuals directly for their smoking status, number of children, and other demographic information.
- This may work well for data that people are comfortable sharing directly.

4. Government or Public Health Data Sources

- Use publicly available health or demographic datasets if they align with your data requirements (e.g., government health surveys that capture smoking prevalence, BMI, or regional health costs).
- Note that some variables may require additional validation or filtering.

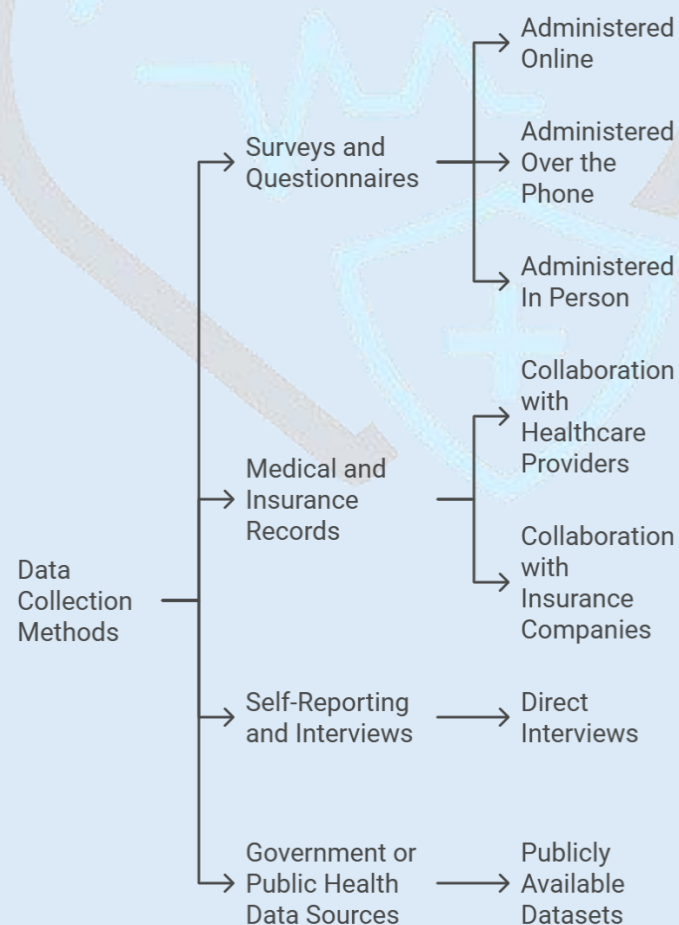


Figure 1.4

Ethical Considerations and Privacy

- **Informed Consent:** Always obtain informed consent from individuals, explaining the purpose of data collection, how their data will be used, and their right to withdraw.
- **Data Anonymization:** Anonymize personal identifiers to protect individual privacy, especially for sensitive information like health records and insurance charges.
- **Data Security:** Store collected data securely, using encryption and access controls to prevent unauthorized access.

Compliance with Regulations: Adhere to relevant privacy regulations (such as GDPR or HIPAA in the United States) when collecting health and demographic information.

Data Validation

After data collection, perform checks for:

- **Completeness:** Ensure no missing data, especially in key variables (age, BMI, smoking status, region, and charges).
- **Accuracy:** Cross-check data where possible (e.g., confirm BMI by recalculating from reported height and weight).
- **Consistency:** Verify that categorical values (e.g., region, sex) align with the expected categories and formats.

This structured approach will help ensure the data collected is reliable, relevant, and ethically gathered for medical insurance cost analysis.

Data Preprocessing

1. Missing Values:

- **Check for missing values:** Use functions like `isnull()` or `isna()` to identify missing data.
- **Handle missing values:**
 - a) **Numerical:** Impute missing values with mean, median, mode, or predict using machine learning models.

- b) **Categorical:** Impute with the most frequent category or create a separate category for missing values.

2.Outliers:

- **Identify outliers:** Use box plots, z-scores, or interquartile range (IQR) to detect outliers.
- **Handle outliers:**
 - a) **Cap outliers:** Set a maximum or minimum value.
 - b) **Winsorization:** Replace outliers with the nearest non-outlier value.
 - c) **Trimming:** Remove outliers.

3.Data Normalization:

- **Scale numerical features:** Use techniques like Min-Max scaling or Standardization to bring features to a common scale. This is especially important for features like Age, BMI, and Charges.

4.Categorical Data Encoding:

- **Label encoding:** Assign numerical labels to categorical variables with ordinal relationships (e.g., low, medium, high).
- **One-hot encoding:** Create binary features for each category, suitable for nominal categorical variables (e.g., Sex, Smoker, Region).

5.Data Consistency:

- **Check for inconsistencies:** Ensure data types are correct, units are consistent, and there are no logical errors.
- **Clean inconsistencies:** Correct errors, standardize units, and ensure data integrity.

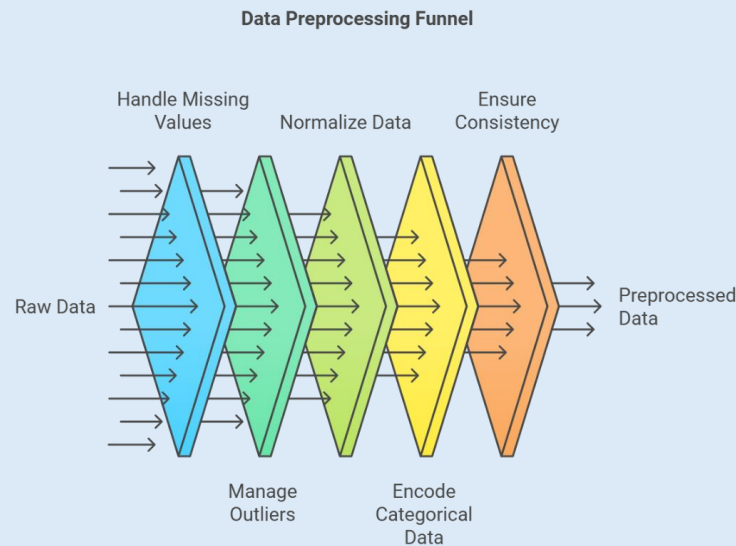


Figure 1.5

Tools

1. Tableau

- **Description:** Tableau is a powerful data visualization tool focused on creating interactive, shareable dashboards that allow for a deep dive into data.
- **Key Uses:**
 - Easily connect to various data sources, including CSV files, databases, and cloud-based data.
 - Create detailed visualizations, like bar charts, scatter plots, and heat maps, to analyze the influence of different factors (e.g., age, BMI, smoker status) on insurance charges.
 - Use Tableau's drag-and-drop interface to create and customize complex charts with minimal coding knowledge.
 - Interactive dashboards allow users to explore the data by filtering or drilling down by factors like age group, region, or smoker status.
- **Best Suited For:** Interactive, real-time dashboarding and exploratory data analysis.

2. Power BI

- **Description:** Power BI, by Microsoft, is a business analytics tool that provides interactive visualizations and robust business intelligence capabilities.
- **Key Uses:**
 - Connect easily to both Excel and SQL databases to import and transform data.
 - Similar to Tableau, Power BI allows users to create interactive dashboards and reports, ideal for showing patterns and trends in insurance costs.
 - Power BI's AI-driven insights can highlight unusual trends or outliers (like the effect of smoking on charges).
 - Offers integration with other Microsoft services, which is beneficial for organizations already using Microsoft's ecosystem.
- **Best Suited For:** Business intelligence reporting, particularly in environments that use Microsoft Office products.

3. MySQL

- **Description:** MySQL is a relational database management system used primarily for data storage and querying.
- **Key Uses:**
 - Though not primarily a visualization tool, MySQL is essential for organizing, querying, and preparing large datasets for analysis.
 - Using SQL queries, users can filter, group, and aggregate data efficiently, like calculating average charges by age group or region.
 - It's often the backend for applications that pull data for visualizations in tools like Tableau or Power BI.
 - Data can be pre-aggregated in MySQL for faster loading and visualization in other tools.
- **Best Suited For:** Data storage, processing, and initial data preparation for large datasets.

4. Advanced Excel

- **Description:** Excel is a widely-used spreadsheet tool, and its advanced features allow for both analysis and visualization.
- **Key Uses:**
 - With pivot tables and pivot charts, Excel can summarize data, such as calculating total charges for different regions or smoker status.
 - Conditional formatting can highlight patterns or outliers, making it easy to spot trends visually (e.g., higher charges for smokers).
 - Advanced chart types and data analysis tools (such as data tables, sparklines, and scenario analysis) allow for creating a range of visualizations, from basic histograms to complex combo charts.
 - Excel's Data Analysis ToolPak can be used for more advanced statistical analysis if needed.
- **Best Suited For:** Quick, ad-hoc data analysis and visualizations for smaller datasets.

5. Python

- **Description:** Python, a versatile programming language, offers a broad range of libraries for data analysis and visualization, such as Matplotlib, Seaborn, and Plotly.
- **Key Uses:**
 - Matplotlib and Seaborn can create detailed and highly customized plots, such as scatter plots, box plots, and regression plots, to show relationships between variables.
 - Plotly and Bokeh enable the creation of interactive, web-ready visualizations.
 - Python also supports complex data analysis and machine learning, allowing users to visualize insights from models (e.g., showing the impact of BMI on charges).
 - For large datasets, Python is effective in handling and processing data before visualization.

- **Best Suited For:** Customizable, in-depth data visualization, particularly for data scientists and analysts familiar with coding. Ideal for predictive modeling and advanced statistical analysis.

In summary, while Tableau and Power BI excel at interactive dashboarding, MySQL is crucial for data storage and pre-processing, Advanced Excel is best for quick analyses, and Python offers extensive customization for both visualization and in-depth statistical analysis.

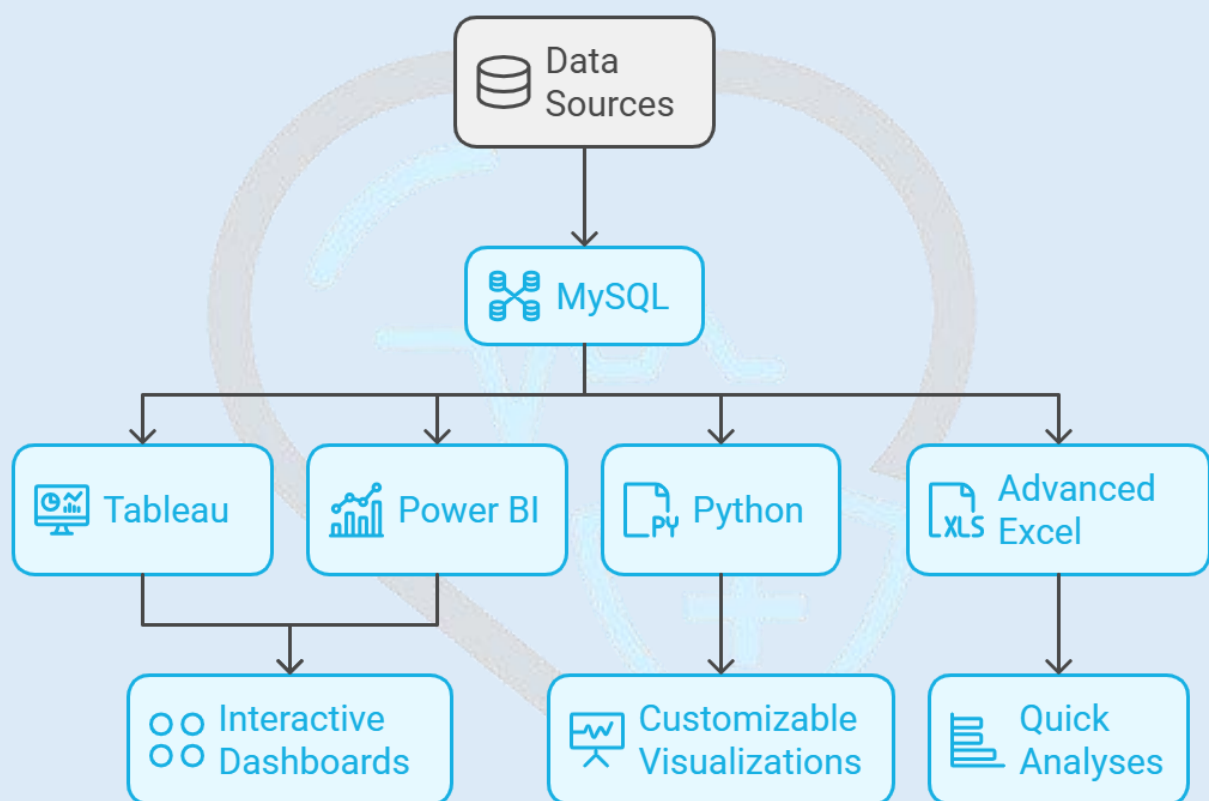


Figure 1.6

Dashboard

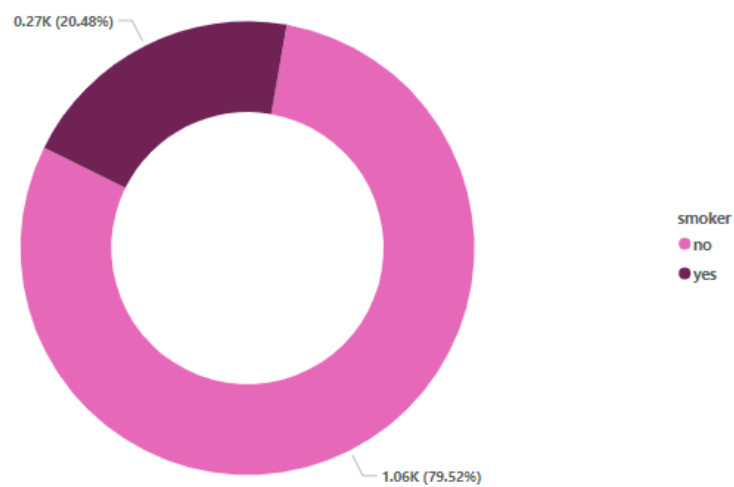
13.27K

Average of charges

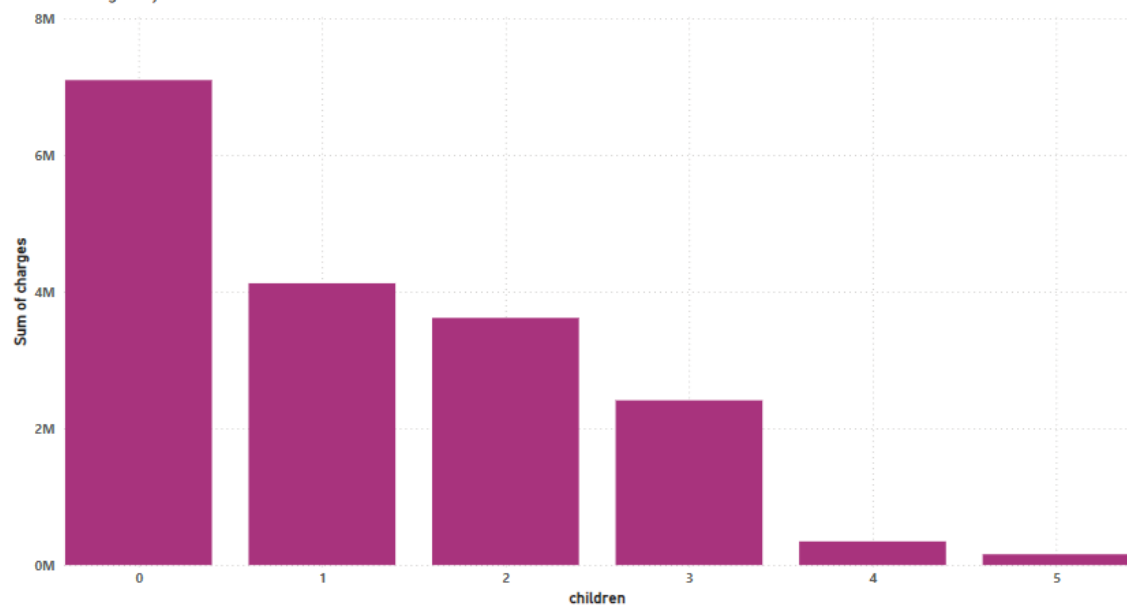
63.77K

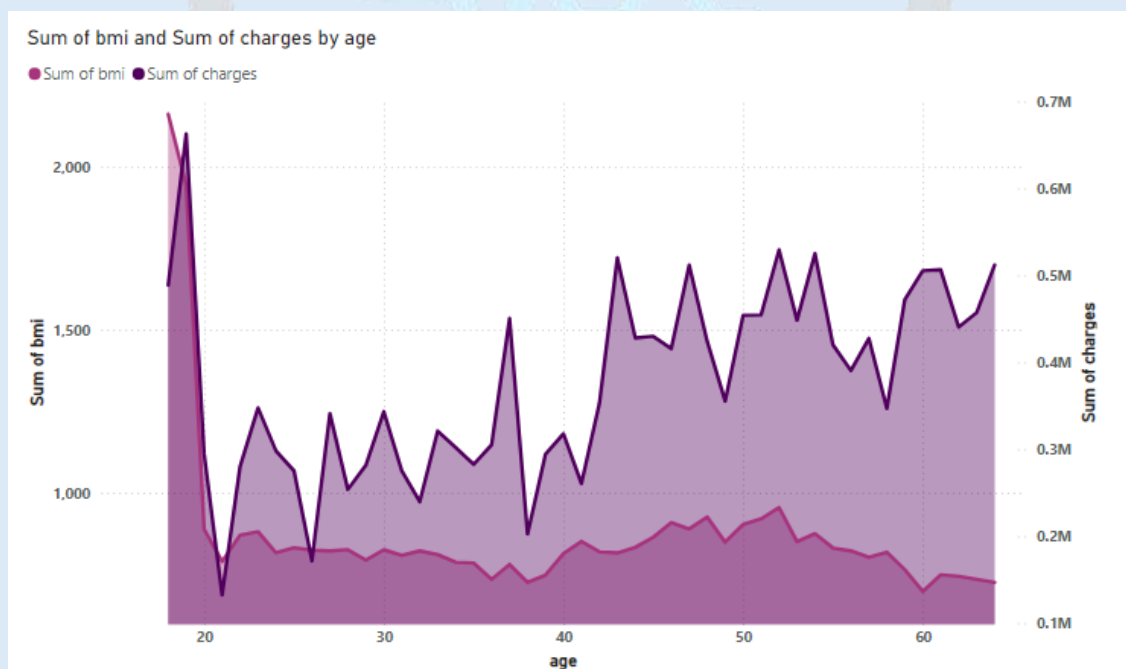
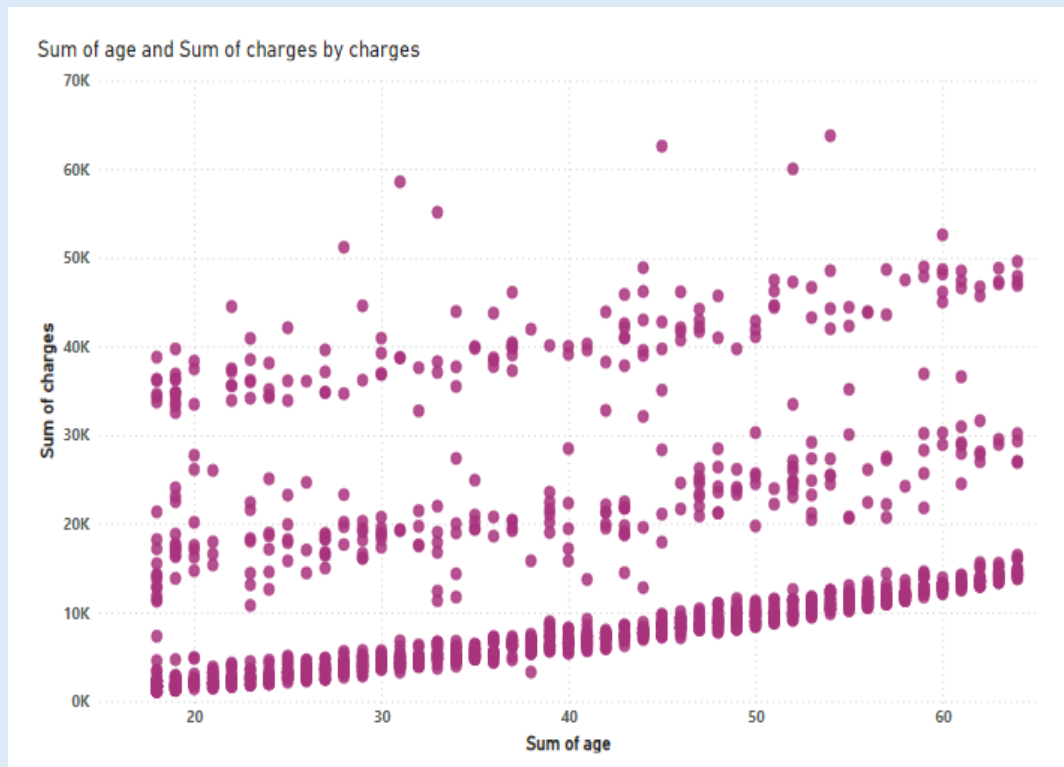
Max of charges

Count of smoker by smoker



Sum of charges by children





Presentation (Insight)

The dashboard presents key metrics and visualizations related to a medical insurance cost analysis. Below are detailed insights from the data, with implications for healthcare and insurance businesses.

1. Summary Statistics of Medical Charges

- **Average Medical Charge:** \$13.27K
- **Maximum Medical Charge:** \$63.77K

Insight:

- These statistics indicate a substantial range in medical insurance costs, with some individuals incurring exceptionally high charges. This variability is likely due to factors such as health status, age, and lifestyle choices.
- **Business Impact:** The wide disparity in charges highlights the need for data-driven pricing models that can adjust premiums based on individual risk factors. For insurers, understanding this cost distribution is essential for setting premiums that cover high-cost outliers without overpricing low-risk individuals.

2. Age and Insurance Charges

- The scatter plot of **Sum of Age** vs. **Sum of Charges** shows that while charges tend to increase with age, there is significant dispersion, indicating that not all older individuals have high charges and vice versa.
- The **Sum of BMI and Sum of Charges by Age** line chart further illustrates this trend, showing that while BMI and charges fluctuate, older age groups generally incur higher medical costs.

Insight:

- While age is a contributing factor to higher charges, it is not the sole predictor, as there is considerable variation across age groups.
- **Business Impact:** Insurance providers can use age as one factor in risk assessment but should incorporate other data points (like health status and lifestyle) for more accurate pricing. This approach enables more personalized premiums, reducing the financial burden on younger, lower-risk individuals while covering costs for high-risk age groups.

3. Smoking Status and Medical Costs

- The **Donut Chart** shows that 20.48% of individuals in the dataset are smokers, while 79.52% are non-smokers.
- Smokers generally incur higher charges, suggesting a strong correlation between smoking status and medical costs.

Insight:

- The data reinforces that smoking is a major health risk and a significant driver of medical costs. Smokers represent a substantial cost burden for insurers, as they are at higher risk for diseases and health complications.
- **Business Impact:** Insurers can leverage this insight to design policies with targeted premiums for smokers and non-smokers. Additionally, introducing or expanding wellness programs focused on smoking cessation could help reduce costs over time, benefiting both insurers and policyholders. By encouraging healthier lifestyles, insurers can lower the risk pool and manage claims more effectively.

4. BMI and Medical Charges by Age

- The **Line Chart** comparing **Sum of BMI** and **Sum of Charges by Age** suggests that high BMI levels, especially in older age groups, correspond to higher insurance costs.
- Younger individuals with high BMI may still incur lower charges, but the trend shifts as age increases, with older adults with high BMI facing significantly higher charges.

Insight:

- High BMI levels, particularly in older adults, correlate with increased healthcare costs due to the higher likelihood of associated health issues (e.g., cardiovascular disease, diabetes).
- **Business Impact:** This data highlights the potential value of BMI-based wellness incentives in health insurance. Insurers might offer discounts or incentives for individuals who maintain a healthy BMI, promoting preventive healthcare practices. Encouraging policyholders to adopt healthier lifestyles can reduce long-term healthcare costs for insurers.

5. Medical Charges by Number of Children

- The **Bar Chart** displaying **Sum of Charges by Children** reveals that medical charges are highest for individuals with one or two children and gradually decrease with larger family sizes (three or more children).
- The peak at one or two children may indicate a high demand for healthcare services in smaller families but shows a stabilizing trend as the number of children increases.

Insight:

- Families with one or two children appear to incur the highest medical costs, which may be due to various factors, including the specific healthcare needs of young families.
- **Business Impact:** Insurers could design family-oriented plans that offer flexible pricing based on family size, making healthcare more affordable for families with different numbers of dependents. Tailoring plans based on average family medical costs allows insurers to cater more effectively to this demographic, potentially enhancing customer satisfaction and retention.

Overall Business Implications and Recommendations

1. **Data-Driven Pricing Models:** The variation in medical costs across different demographics (age, smoking status, BMI, and family size) suggests that insurers should adopt more granular, data-driven models to determine premiums. By segmenting risk profiles more precisely, insurers can create fair and sustainable pricing structures.
2. **Wellness and Preventive Health Programs:** The higher costs associated with smokers and individuals with higher BMIs indicate a strong potential for wellness programs. Insurers could introduce incentives for smoking cessation, healthy BMI maintenance, and regular health screenings. Such programs not only help reduce claims but also improve long-term health outcomes, benefiting both insurers and policyholders.
3. **3. Customization of Family Plans:** With insights on cost variations by the number of children, insurers can create customized family plans that cater to the unique needs of small and large families. This flexibility can improve customer satisfaction by aligning coverage options with actual needs and costs.
4. **4. Proactive Risk Management:** Age remains a significant factor in predicting insurance costs. By understanding the age distribution of their customer base and identifying high-cost age groups, insurers can proactively manage risk by adjusting reserves and planning for higher claims in certain demographics.

5. **5. Leveraging Insights for Customer Retention:** Personalized and flexible insurance plans based on insights from data analytics can boost customer retention by making healthcare more affordable and accessible. Customized policies, health incentives, and family-friendly plans can enhance customer loyalty and attract new policyholders in a competitive market.

These insights provide a roadmap for insurance and healthcare providers looking to optimize costs, improve customer satisfaction, and encourage healthier lifestyles among their customers. Data-driven decision-making in insurance pricing and policy design is not only a competitive advantage but also a pathway to better health outcomes for all.

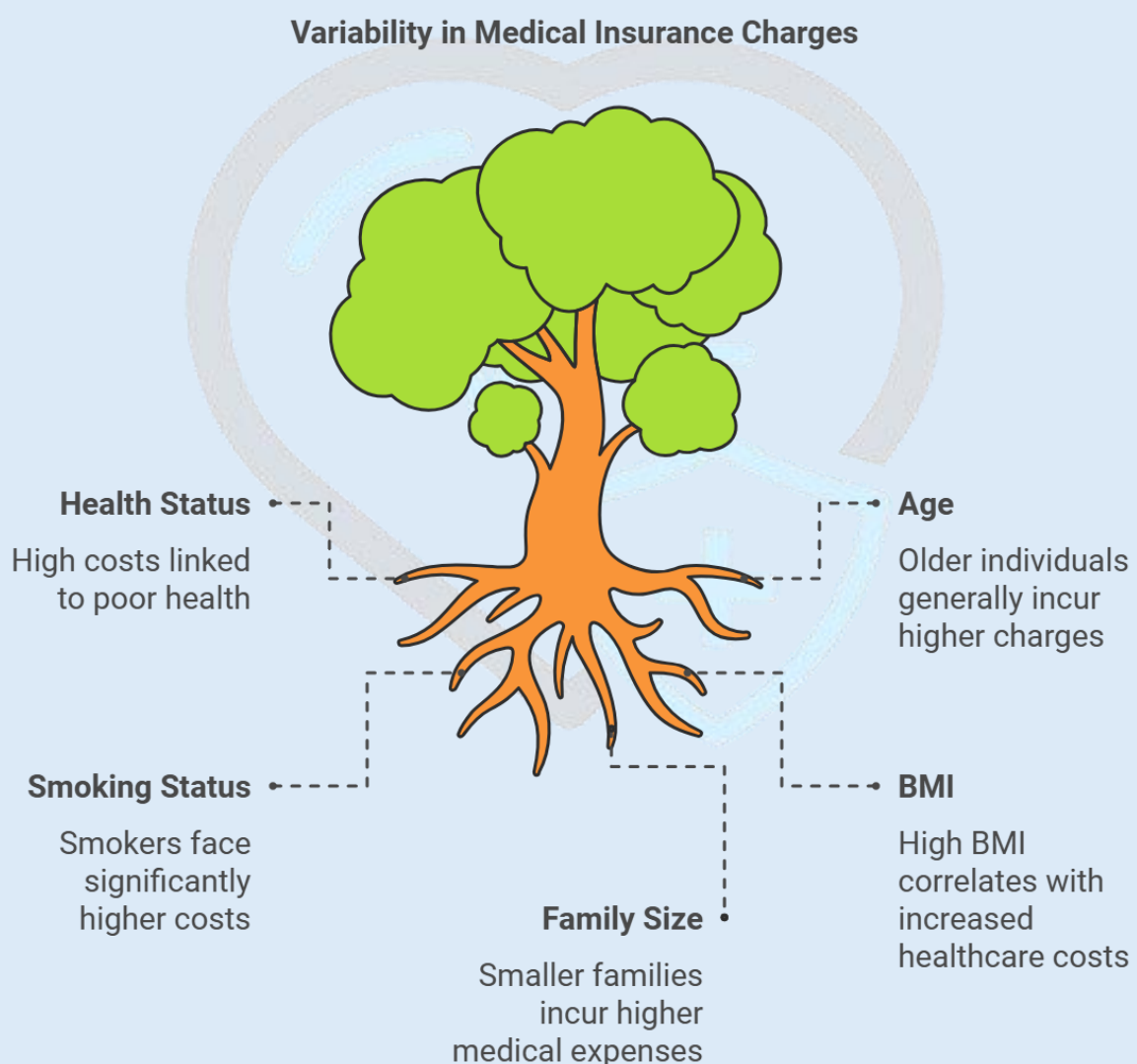


Figure 1.7