



Name	Harditya Shah
UID No.	2021300114
Batch	L

### Objectives:

- To explore and visualize a dataset related to Finance/Banking/Insurance/Credit using D3.js.
- To create basic visualizations (Bar chart, Pie chart, Histogram, Timeline chart, Scatter plot, Bubble plot) to understand data distribution and trends.
- To create advanced visualizations (Word chart, Box and Whisker plot, Violin plot, Regression plot, 3D chart, Jitter) for deeper insights and complex relationships.
- To perform hypothesis testing using the Pearson correlation coefficient to evaluate relationships between numerical variables in the dataset.

Dataset: Health Insurance Dataset

Link: [US Health Insurance Dataset](#)

### About Dataset:

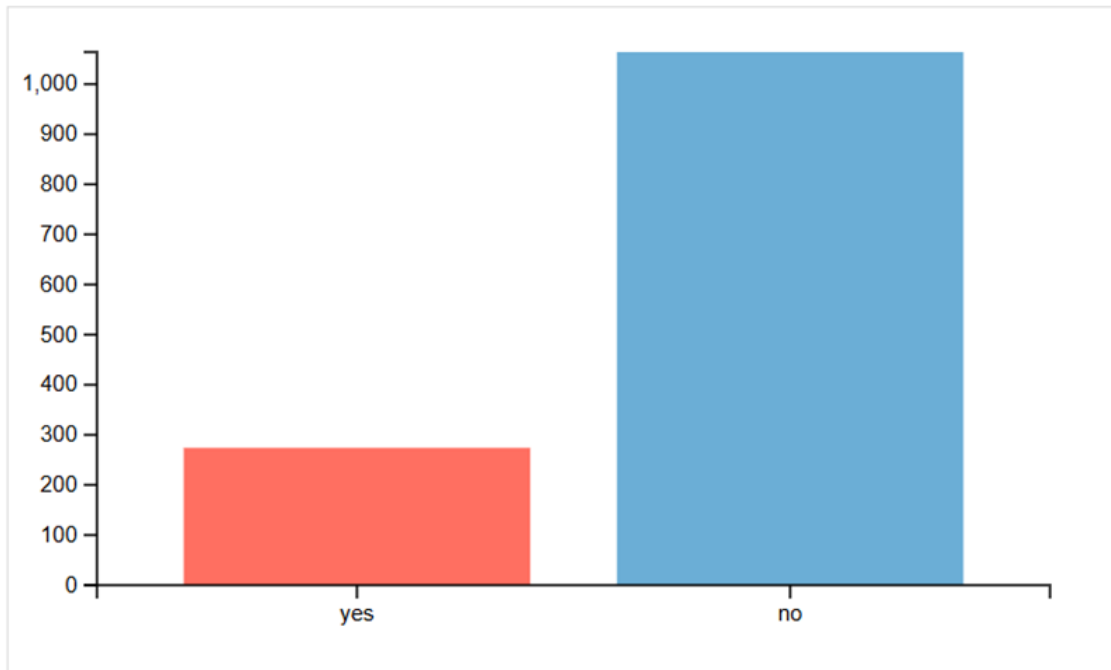
This dataset is of insured data, where the Insurance charges are given against the following attributes of the insured: Age, Sex, BMI, Number of Children, Smoker and Region. The attributes are a mix of numeric and categorical variables.

### Column Attributes:

- age: The age of the individual in years.
- sex: The gender of the individual (male or female).
- bmi: Body Mass Index, a measure of body fat based on height and weight.
- children: The number of children/dependents covered by the insurance.
- smoker: Indicates whether the individual is a smoker (yes or no).
- region: The individuals residential region in the U.S. (e.g., southwest, southeast, northwest, northeast).
- charges: The medical insurance charges billed to the individual.

### Basic Visualizations:

## Bar Chart (Smokers vs Non-Smokers)



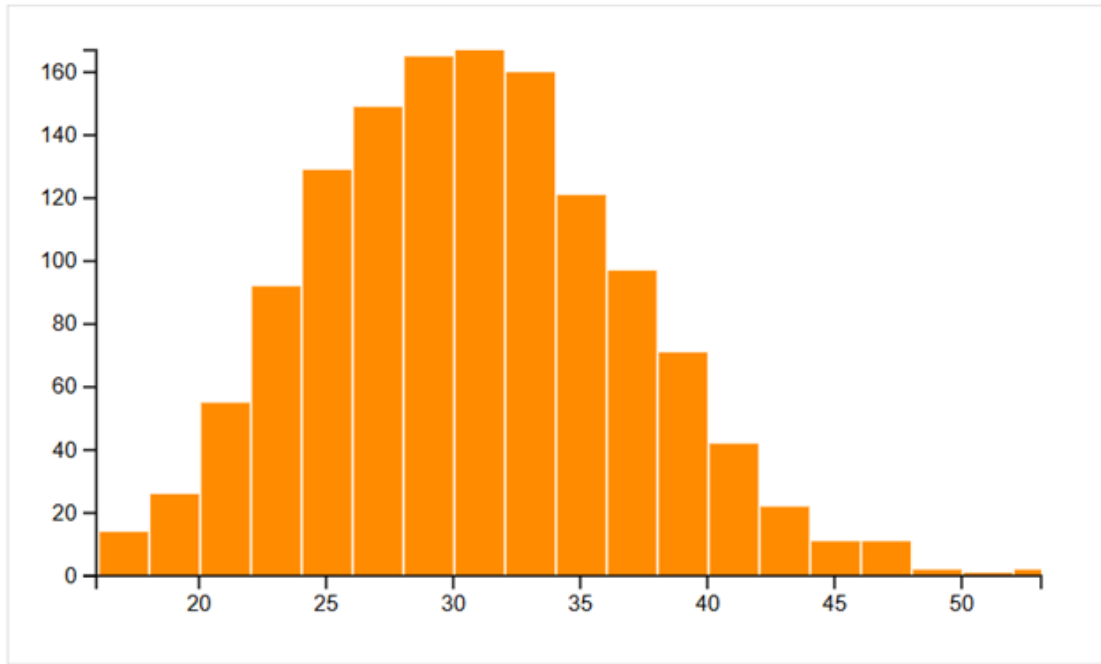
Non-Smokers Predominate: The “no” category has a significantly higher count compared to the “yes” category, indicating that there are many more non-smokers

## Pie Chart (Smokers vs Non-Smokers)



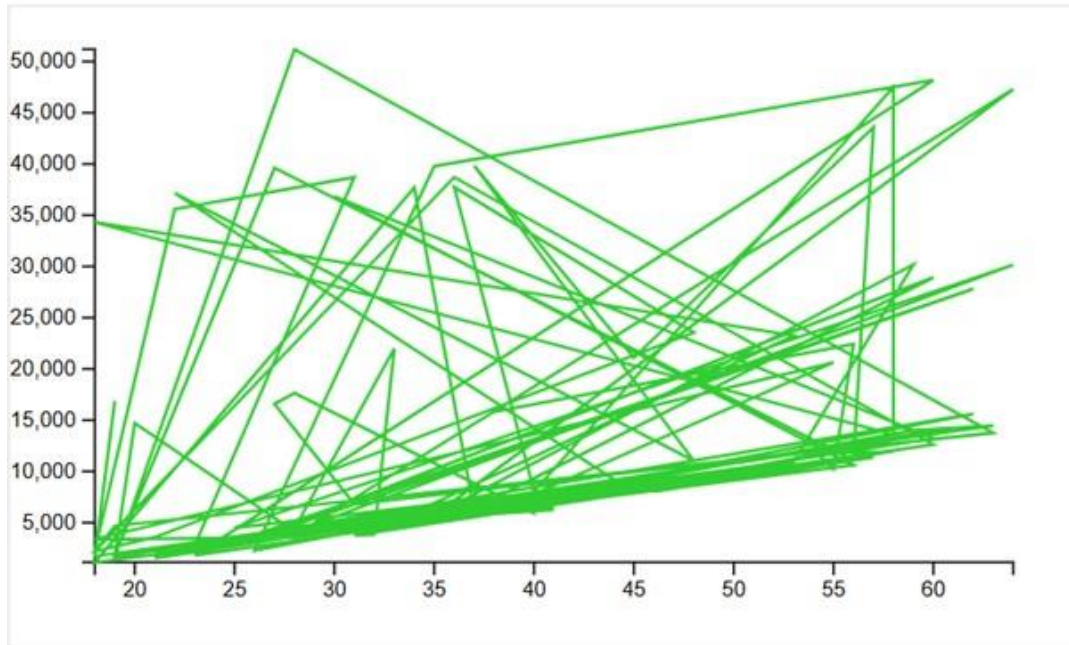
The visual disparity between the “yes” and “no” segments highlights the significant difference in smoking habits within the population.

## Histogram (BMI Distribution)



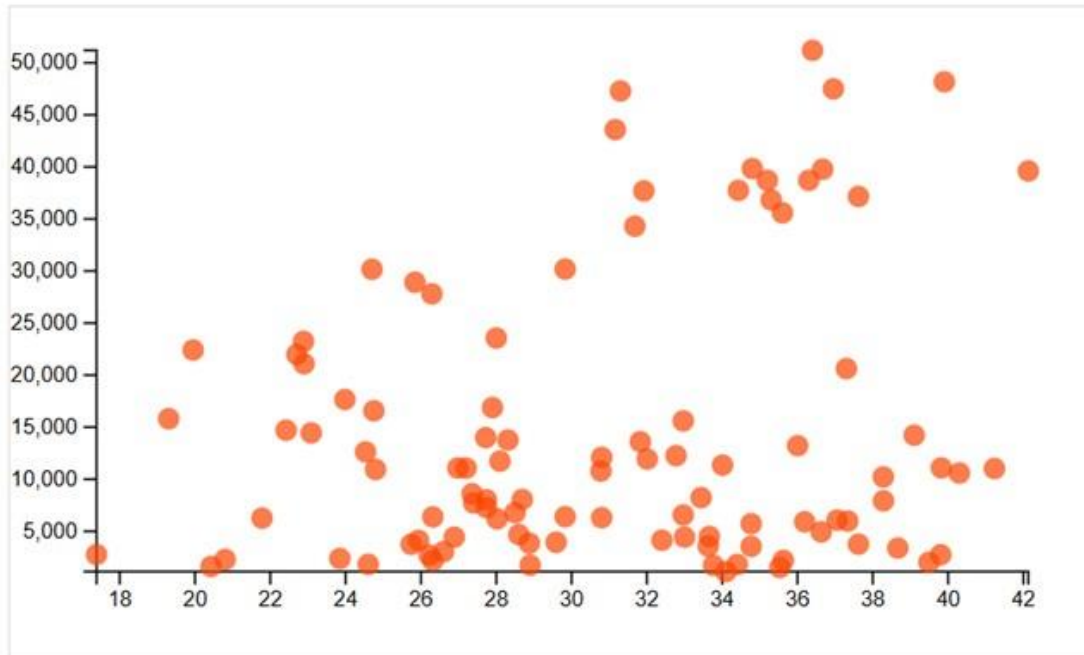
The highest frequency of BMI values is centered around 30, which indicates that most individuals in this dataset have a BMI close to this value.

## Timeline Chart (Charges vs Age)



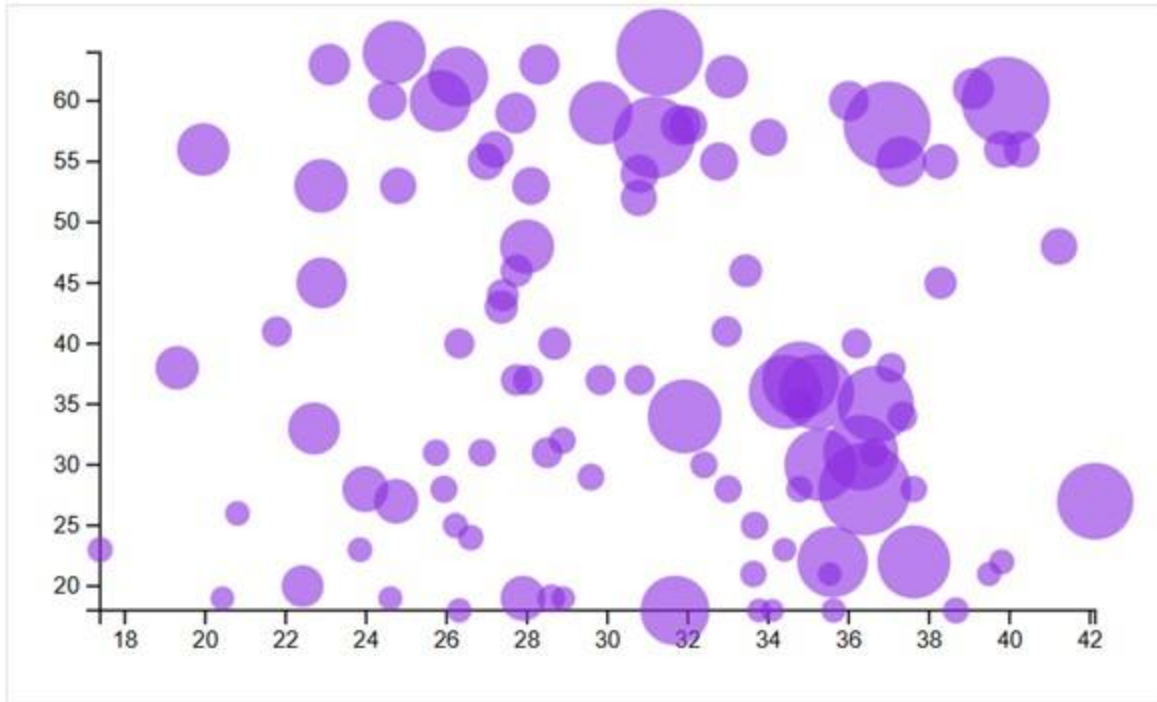
Trend: There is a general trend where charges increase with age, though there are fluctuations.

## Scatter Plot (BMI vs Charges)



There is a positive correlation between BMI and charges, meaning that as BMI increases, the charges also tend to increase.

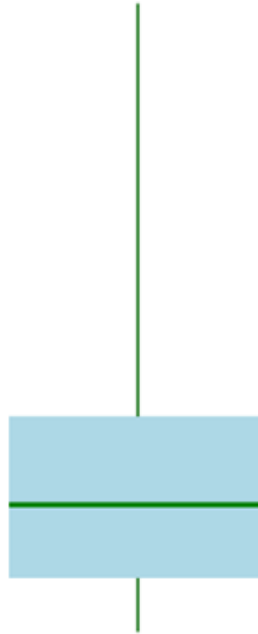
## Bubble Plot (BMI, Age, and Charges)



There are clusters of larger bubbles (higher BMI) at various age points, suggesting that higher BMI individuals incur higher charges.

Advanced Visualizations:

## Box Plot (Charges)



The box plot summarizes the distribution of insurance charges. It shows the median, quartiles, and potential outliers in the charges data



## Violin Plot (Charges Distribution)

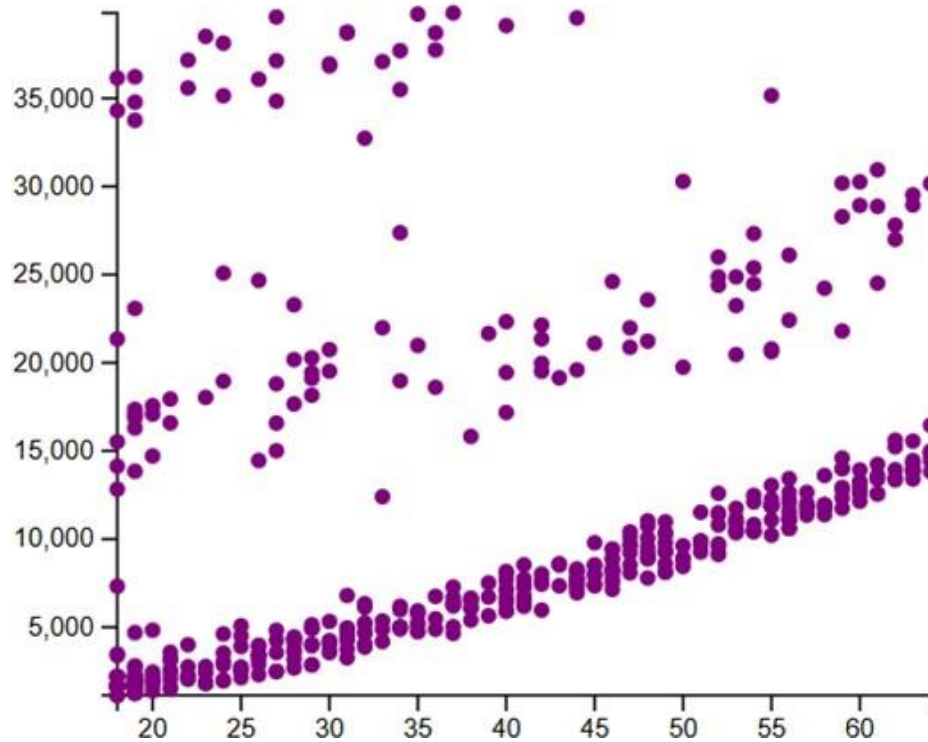


Violin Plot of Charges



The violin plot provides a detailed view of the distribution of charges, combining the box plot's summary statistics with a density estimation. This helps visualize where the bulk of the charges lie and whether there are multiple modes in the data.

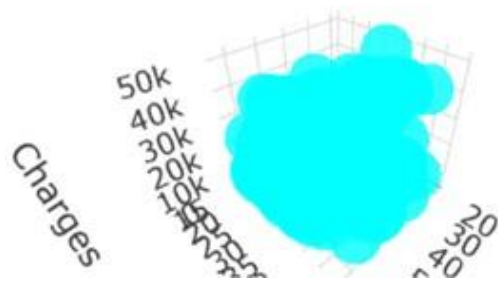
**Regression Plot (Age vs Charges)**



This scatter plot with a regression line illustrates the relationship between age and insurance charges. The slope of the regression line indicates how charges increase with age.

## 3D Scatter Plot (Age, BMI, Charges)

3D Scatter Plot



The 3D scatter plot allows for the visualization of how age, Body Mass Index (BMI), and insurance charges interact. Each point represents an individual in the dataset, with three dimensions of data being plotted.

### Jitter Plot (Charges)



The jitter plot adds a slight random displacement to points in a scatter plot to prevent overlap, thus visualizing the distribution of charges more clearly. It can reveal clustering of points and density in specific areas of charge values.

#### Hypothesis Testing

Perform hypothesis testing to evaluate the correlation between age and charges in the dataset.

#### Hypothesis

1. Null Hypothesis ( $H_0$ ): There is no correlation between age and charges.
2. Alternative Hypothesis ( $H_1$ ): There is a correlation between age and charges.

CODE AND OUTPUT:

```
sma > (No subject) > ⚡ adv7.py > ...
1  from scipy.stats import pearsonr
2  import pandas as pd
3
4  # Load dataset
5  file_path = 'sma\No subject\insurance.csv'
6  data = pd.read_csv(file_path)
7
8  # Calculate Pearson correlation coefficient between 'age' and 'charges'
9  corr_age_charges, p_value_age_charges = pearsonr(data['age'], data['charges'])
10
11 # Print the results
12 print(f"--- Hypothesis Testing for Age vs Charges ---")
13 print(f"Null Hypothesis (H0): There is no correlation between age and charges.")
14 print(f"Alternative Hypothesis (H1): There is a correlation between age and charges.")
15 print(f"\nPearson Correlation Coefficient: {corr_age_charges:.4f}")
16 print(f"P-Value: {p_value_age_charges:.4f}")
17
18 alpha = 0.05 # Significance level
19 if p_value_age_charges < alpha:
20     print("Reject the null hypothesis (H0). There is evidence to suggest a correlation between age and charges.")
21 else:
22     print("Fail to reject the null hypothesis (H0). There is no evidence to suggest a correlation between age and charges.")
23
```

PROBLEMS OUTPUT DEBUG CONSOLE PORTS TERMINAL JUPYTER

```
PS D:\vscode> python -u "d:\vscode\sma\No subject\adv7.py"
--- Hypothesis Testing for Age vs Charges ---
Null Hypothesis (H0): There is no correlation between age and charges.
Alternative Hypothesis (H1): There is a correlation between age and charges.

Pearson Correlation Coefficient: 0.2990
P-Value: 0.0000
Reject the null hypothesis (H0). There is evidence to suggest a correlation between age and charges.
PS D:\vscode>
```

Conclusion: Using the health insurance dataset, visualizations created with D3.js revealed significant insights. These interactive and dynamic visualizations, such as regression plots and custom box plots, facilitated deeper exploration of data patterns and trends. These tools are invaluable for understanding complex relationships within the dataset and enhancing data analysis capabilities.