# Sayak Dey 19BOE10072

# VIT Bhopal

# SmartInternz

## Assignment-4

## Dataset : Insurance.csv

## Overview:



## Uploaded Dataset:

## Data refinery Flow:



## Data Profile:

## Data Visualization:

## Bmi vs premium:

# Age vs premium:



# Smoker vs premium:

## Data Audit:





View Output: Data Audit of [7 fields]

| | | | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|
| 7 | premium | | Continuous | 1121.874 | 63770.428 | 13270.422 | 12110.011 | 1.516 | -- | 1338 |

| | Field | Measurement | Outliers | Extremes | Action | Impute Missing | Method | % Complete | Valid Records | Null Value |
|---|---|---|---|---|---|---|---|---|---|---|
| 1 | age | Continuous | 0 | 0 | None | Never | Fixed | 100.000 | 1338 | 0 |
| 2 | sex | Categorical | -- | -- | -- | Never | Fixed | 100.000 | 1338 | 0 |
| 3 | bmi | Continuous | 4 | 0 | None | Never | Fixed | 100.000 | 1338 | 0 |
| 4 | children | Continuous | 18 | 0 | None | Never | Fixed | 100.000 | 1338 | 0 |
| 5 | smoker | Categorical | -- | -- | -- | Never | Fixed | 100.000 | 1338 | 0 |
| 6 | region | Categorical | -- | -- | -- | Never | Fixed | 100.000 | 1338 | 0 |
| 7 | premium | Continuous | 7 | 0 | None | Never | Fixed | 100.000 | 1338 | 0 |

## Building K-Means Clustering Model:



## Model Output:

## View Model: K-Means

K-Means Clustering Model ⓘ

**EVALUATION**

**Cluster Quality**

MODEL VIEWER

Model Information

Feature Importance

Cluster Sizes

Cluster Comparison

Clusters

Cell Distributions (Absolute)

Cell Distributions (Relative)

Build Settings

### Cluster Quality ⓘ

```
 -1.0        -0.5        0.0        0.5        1.0
        Silhouette Measure of Cohesion and Separation
```

**Cluster Quality Parameters**

| | |
|---|---|
| Overall Clustering Quality (Avg. Silhouette) | 0.260 |
| Total Within Clusters Sum of Squares | 0.145 |
| Average Within Cluster Sum of Squares | 0.029 |
| Average SSB (Between ss) | 0.073 |

---

## View Model: K-Means

K-Means Clustering Model ⓘ

**EVALUATION**

Cluster Quality

MODEL VIEWER

**Model Information**

Feature Importance

Cluster Sizes

Cluster Comparison

Clusters

Cell Distributions (Absolute)

Cell Distributions (Relative)

Build Settings

### Model Information ⓘ

| | | |
|---|---|---|
| Algorithm | | K-Means |
| Model Class | | Center Based |
| Number of Features | | 6 |
| Distance Measure | | Euclidean |
| Number of Clusters | | 5 |
| | Cluster 1 | 94 (8.79%) |
| | Cluster 2 | 411 (38.45%) |

# View Model: K-Means

## K-Means Clustering Model ⓘ

### Model Information ⓘ

| | | |
|---|---|---|
| Number of Clusters | | 5 |
| Number of instances in each cluster | Cluster 1 | 94 (8.79%) |
| | Cluster 2 | 411 (38.45%) |
| | Cluster 3 | 220 (20.58%) |
| | Cluster 4 | 130 (12.16%) |
| | Cluster 5 | 214 (20.02%) |
| Ratio of sizes (Largest to smallest) | | 4.372 |

**EVALUATION**

- Cluster Quality

**MODEL VIEWER**

- **Model Information**
- Feature Importance
- Cluster Sizes
- Cluster Comparison
- Clusters
- Cell Distributions (Absolute)
- Cell Distributions (Relative)
- Build Settings

---

# View Model: K-Means

## K-Means Clustering Model ⓘ

### Feature Importance ⓘ



| Feature | Importance |
|---|---|
| sex | 1.00 |
| smoker | 1.00 |
| region | 0.38 |
| bmi | 0.01 |
| children | 0.00 |
| age | 0.00 |

**EVALUATION**

- Cluster Quality

**MODEL VIEWER**

- Model Information
- **Feature Importance**
- Cluster Sizes
- Cluster Comparison
- Clusters
- Cell Distributions (Absolute)
- Cell Distributions (Relative)
- Build Settings

## View Model: K-Means



**K-Means Clustering Model** ⓘ

EVALUATION

Cluster Quality

MODEL VIEWER

Model Information

Feature Importance

Cluster Sizes

Cluster Comparison

**Clusters**

Cell Distributions (Absolute)

Cell Distributions (Relative)

Build Settings

### Clusters ⓘ

Input Importance

■ 1.0  ■ 0.8  ■ 0.6  ■ 0.4  □ 0.2  □ 0.0

| Cluster | cluster_1 | cluster_2 | cluster_3 | cluster_4 | cluster_5 |
|---------|-----------|-----------|-----------|-----------|-----------|
| Size | | | | | |
| Inputs | sex **female (100.00%)** | sex **male (100.00%)** | sex **female (100.00%)** | sex **male (100.00%)** | sex **female (100.00%)** |
| | smoker **yes (100.00%)** | smoker **no (100.00%)** | smoker **no (100.00%)** | smoker **yes (100.00%)** | smoker **no (100.00%)** |

---

## View Model: K-Means

**K-Means Clustering Model** ⓘ

EVALUATION

Cluster Quality

MODEL VIEWER

Model Information

Feature Importance
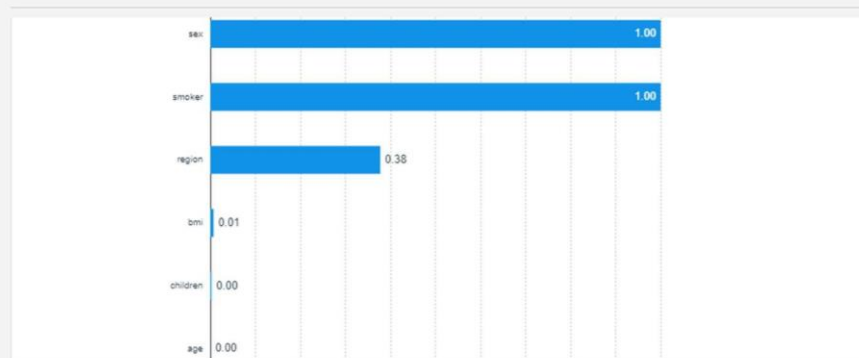
Cluster Sizes

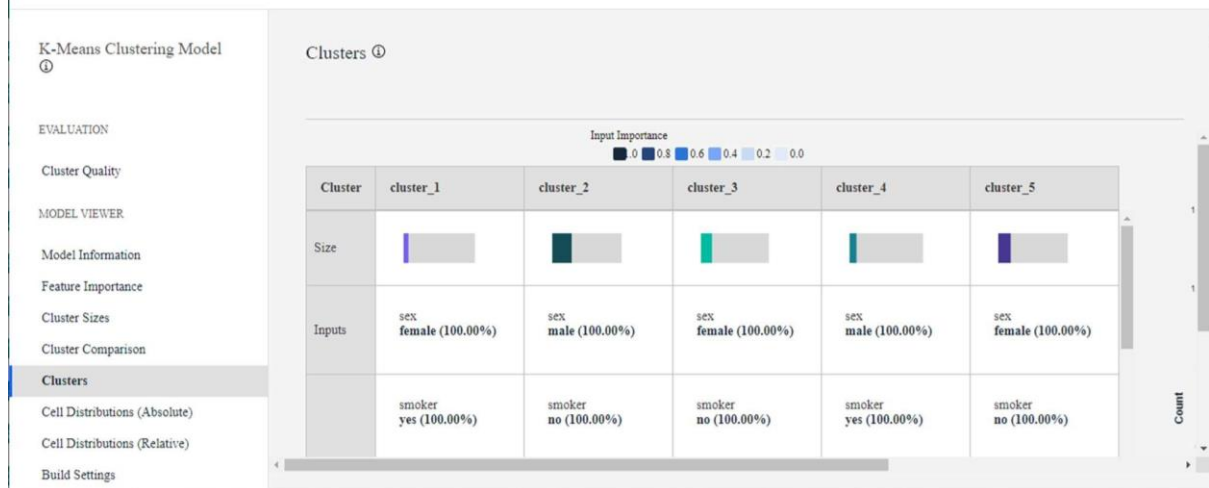Cluster Comparison

Clusters

**Cell Distributions (Absolute)**

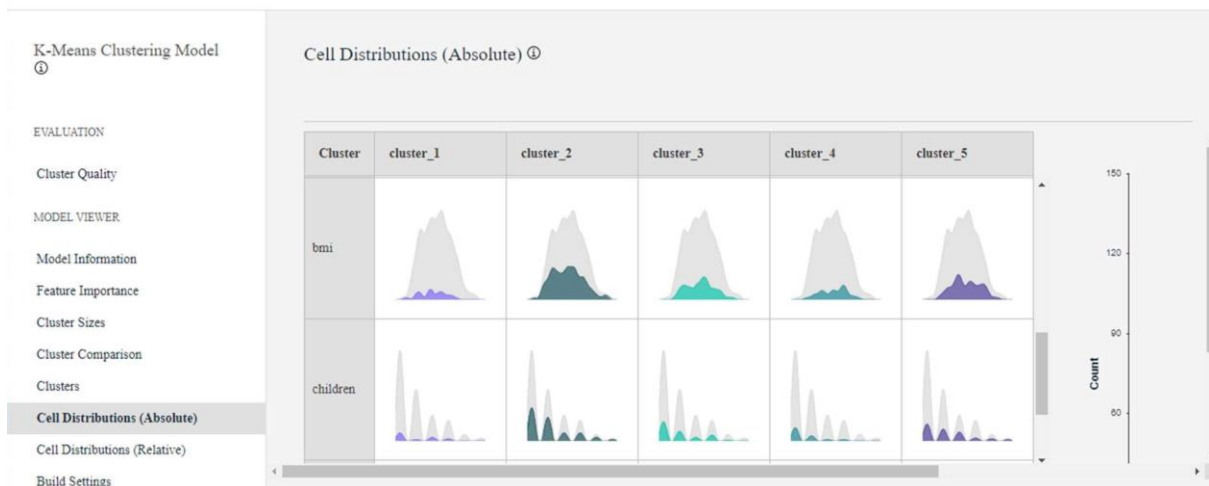Cell Distributions (Relative)

Build Settings

### Cell Distributions (Absolute) ⓘ

| Cluster | cluster_1 | cluster_2 | cluster_3 | cluster_4 | cluster_5 |
|---------|-----------|-----------|-----------|-----------|-----------|
| bmi | | | | | |
| children | | | | | |

## View Model: K-Means

**K-Means Clustering Model** ⓘ

Cluster Quality

MODEL VIEWER

Model Information
Feature Importance
Cluster Sizes
Cluster Comparison
Clusters
Cell Distributions (Absolute)
**Cell Distributions (Relative)**
Build Settings
Training Summary

### Cell Distributions (Relative) ⓘ



## View Model: K-Means

**K-Means Clustering Model** ⓘ

Cluster Quality

MODEL VIEWER

Model Information
Feature Importance
Cluster Sizes
Cluster Comparison
Clusters
Cell Distributions (Absolute)
Cell Distributions (Relative)
**Build Settings**
Training Summary

### Build Settings ⓘ

| | |
|---|---|
| Use partitioned data | true |
| Calculate raw propensity scores | false |
| Calculate adjusted propensity scores | false |
| Number of clusters | 5 |
| Generate distance field | false |
| Cluster label | String |
| Label prefix | cluster |

**Building a Plot node with BMI vs Premium vs Age:**

**Output:**


View Output: bmi v. premium v. age

**Changing number of Clusters from 5 to 3:**

**Output:**



View Output: bmi v. premium v. age