# Data Visualisation

DICE
ANALYTICS

# Visualising Numerical Data

# Scatterplot



**Income**

Response
Variable

**Age**

Explanatory
Variable

DICE ANALYTICS

# Characteristics of Relationship



Direction: +ve, -ve

Shape: curved, linear

Strength: strong, weak

Outliers

DICE ANALYTICS

# Correlation (example)



$r = 0.7$      $r = 0.3$      $r = 0$
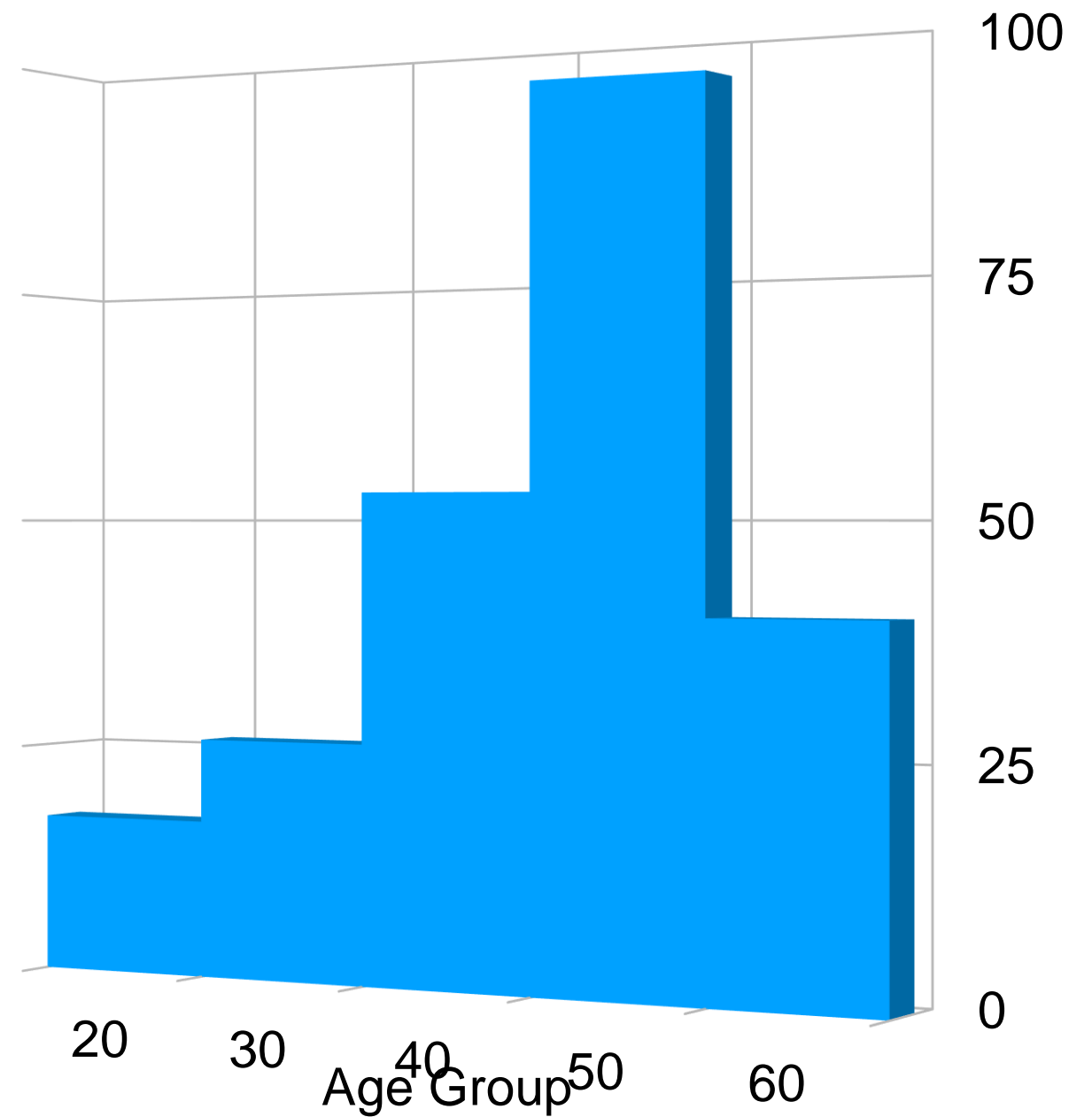
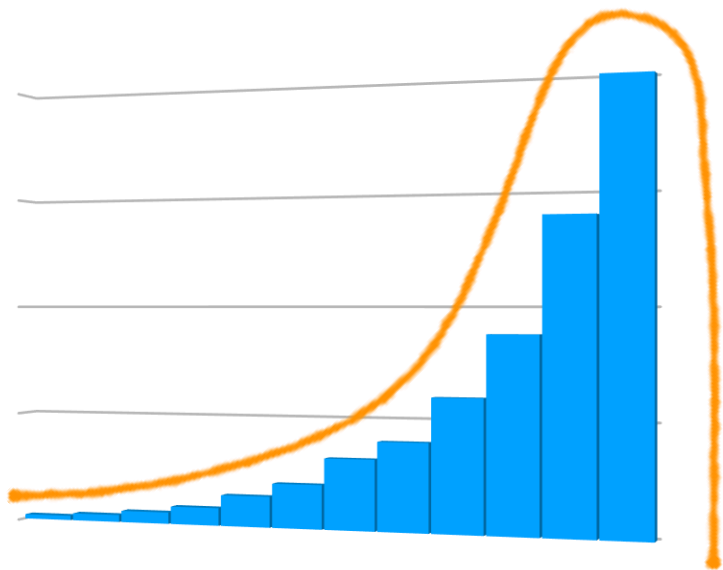$r = -0.7$      $r = -0.3$      $r = 0$

DICE ANALYTICS

# Histograms

- Help to view *data density*

- Help to see *shape of distribution*

  **1) Skewness**
  **2) Modality**

# Skewness



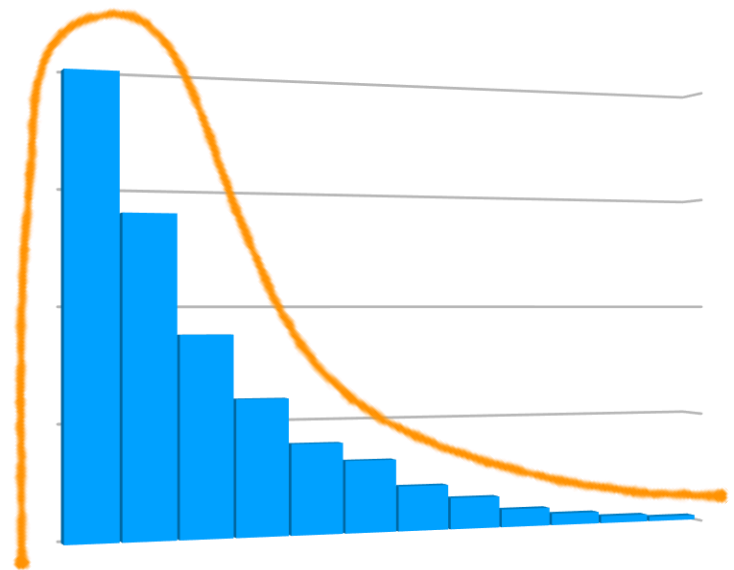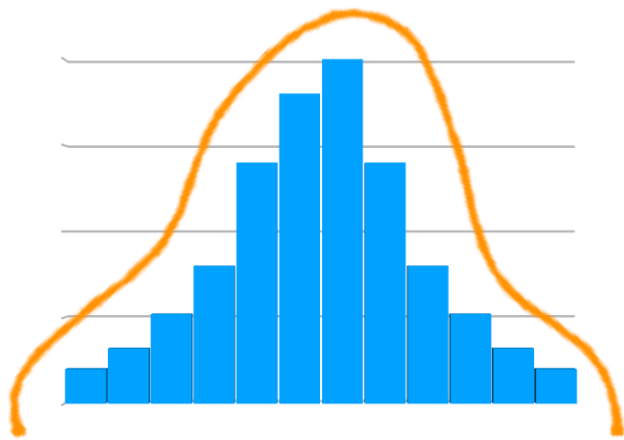Left Skewed — Symmetric — Right Skewed

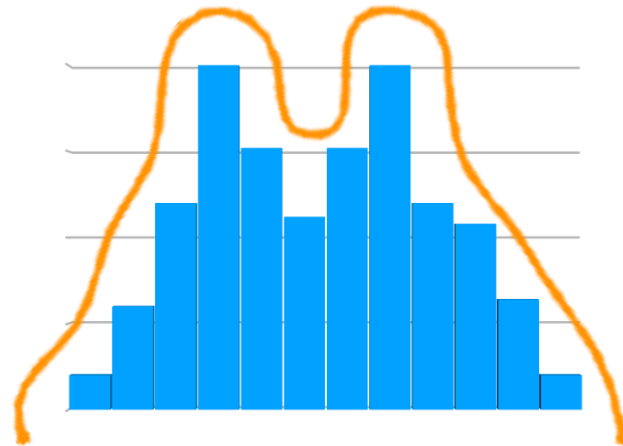-ve Skewness — Zero Skewness — +ve Skewness

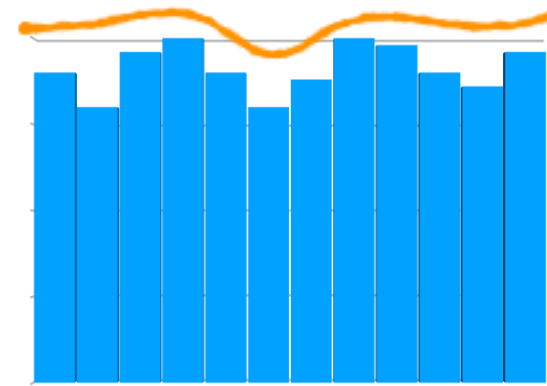- Draw a smooth curve to see skewness
- Don't rely on jagged edges

DICE ANALYTICS

# Modality



unimodal     bimodal     uniform     multimodal

DICE ANALYTICS
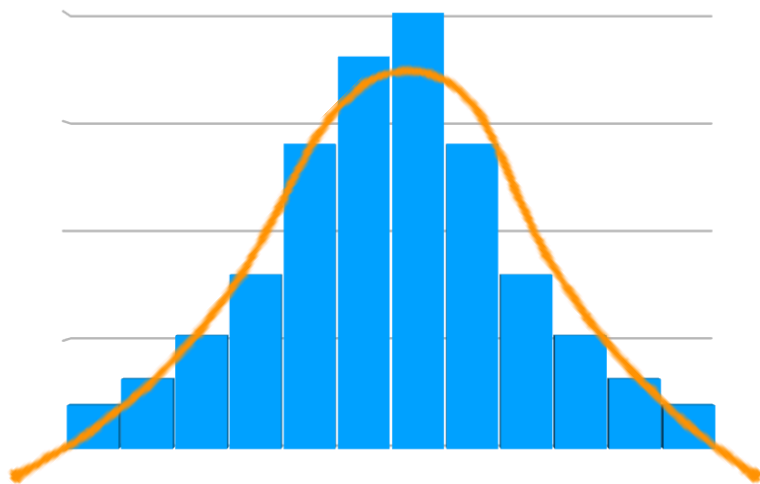
# Modality (Example)



Normal Distribution           Two separate groups           No trend

DICE
ANALYTICS

# Binwidth

# Measures of Center

**Data** : **56, 87, 34, 65, 77, 62, 90, 45, 77, 79**

| | |
|---|---|
| **Mean** | Arithmetic Average<br><br>$\text{Mean} = \dfrac{56 + 87 + 34 + 65 + 77 + 62 + 90 + 45 + 77 + 79}{10}$<br><br>$\text{Mean} = 67.2$ |
| **Mode** | Most frequent value/observation<br>Mode = 77 |
| **Median** | Midpoint of distribution ($50^{th}$ percentile)<br>$\text{Median} = \dfrac{77 + 62}{2} = 69.5$ |

**DICE** ANALYTICS

# Box Plots

**outliers**

IQR

Box

Whisker

Min. Value          Q1          Q2          Q3          Max. Value

**Min. Value** :Lower Extreme (that's not an outlier)
**Q1**          :Lower Quartile (25% of observations)
**Q2**          :Median (50% of observations)
**Q3**          :Upper Quartile (75% of observations)
**Max. Value** :Upper Extreme (that's not an outlier)
**IQR**          :Inter-Quartile Range = Q3 - Q1 (middle 50% of observations)

DICE ANALYTICS

# Box Plots & Skewness



Left Skewed

Symmetric

Right Skewed

DICE ANALYTICS

# Skewness vs Measures of Center



Mean < Median < Mode          Mean = Median = Mode          Mean > Median > Mode

Left Skewed          Symmetric          Right Skewed

DICE
ANALYTICS

# Intensity/Heat Maps

# Time Plots

# Measures of Spread

| Range | Variance |
|---|---|
| **Standard Deviation** | **Inter-quartile Range** |

DICE ANALYTICS

# Range

- Range = Max. Value - Min. Value

- **Data :    56, 87, 34, 65, 77, 62, 90, 45, 77, 79**

- Range = 90 - 34 = 56

**DICE** ANALYTICS

# Variance

- A measure of how much data (a variable) varies; how spread out a data set is about the mean.
- Average squared deviation from mean; has squared units of the variable

- Sample Variance

$$s^2 = \frac{\sum (X - \bar{X})^2}{N - 1}$$

- Population Variance

$$\sigma^2 = \frac{\sum (x - \mu)^2}{N}$$

DICE ANALYTICS

# Variance (Example)

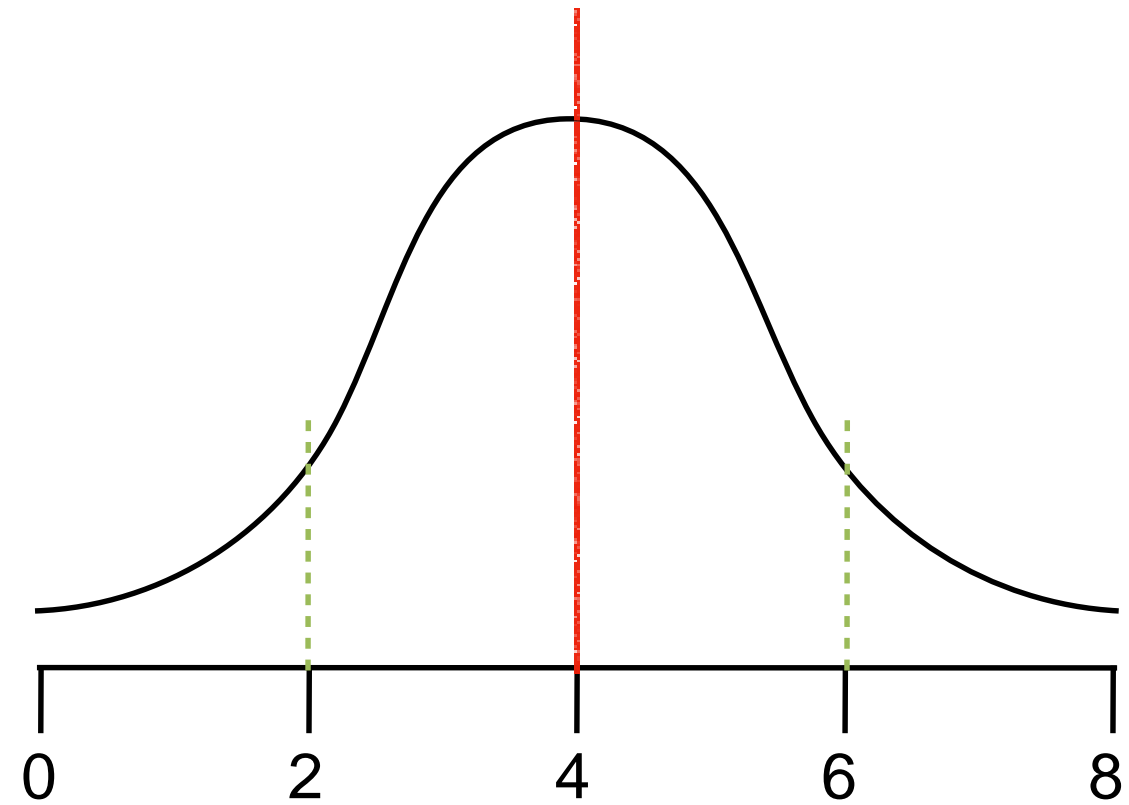- **Data : 56, 87, 34, 65, 77, 62, 90, 45, 77, 79**

$$s^2 = \frac{\sum (X - \bar{X})^2}{N-1} = \frac{(56 - 67.2)^2 + (87 - 67.2)^2 + \ldots\ldots + (79 - 67.2)^2}{10-1}$$

Sum of Squares

$$= \frac{2995.6}{9}$$

$$= 332.8$$

DICE ANALYTICS

# Why Square The Differences?

- Get rid of negatives, so that the negatives and positives do not cancel each other during addition.

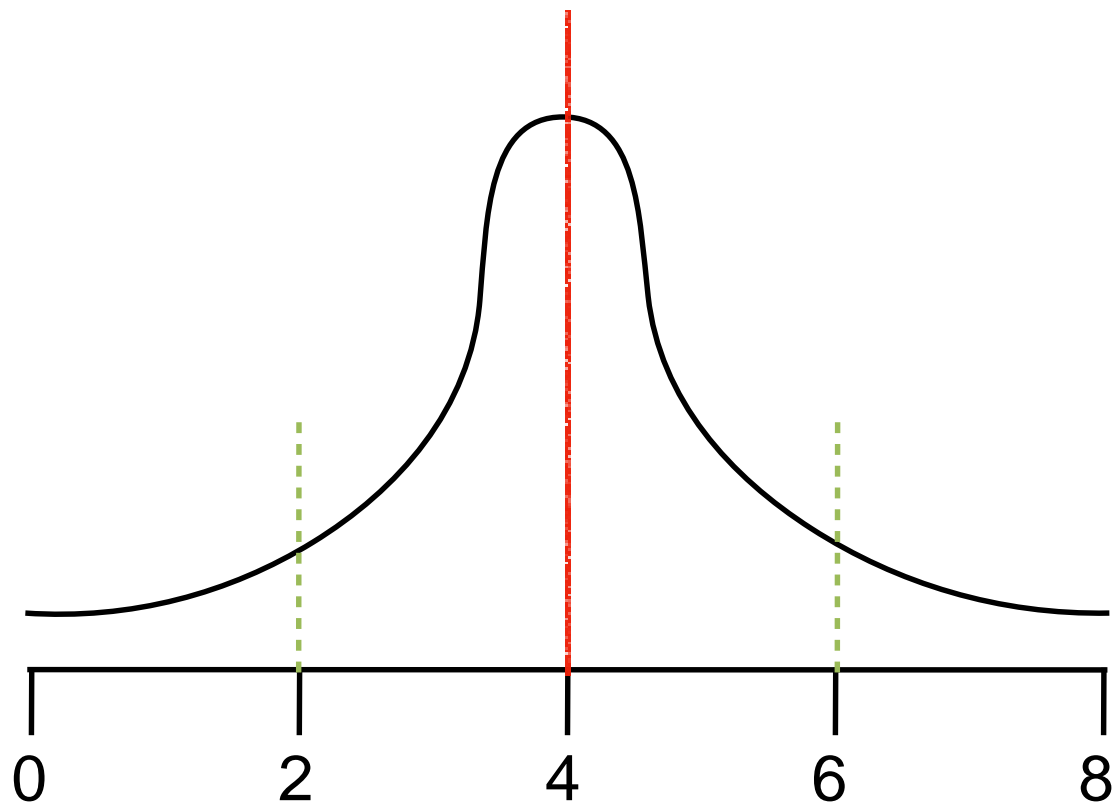- Increase larger deviations more than smaller ones so that they are weighed more heavily.



$(2-4) + (6-4) = -2 + 2 = 0$

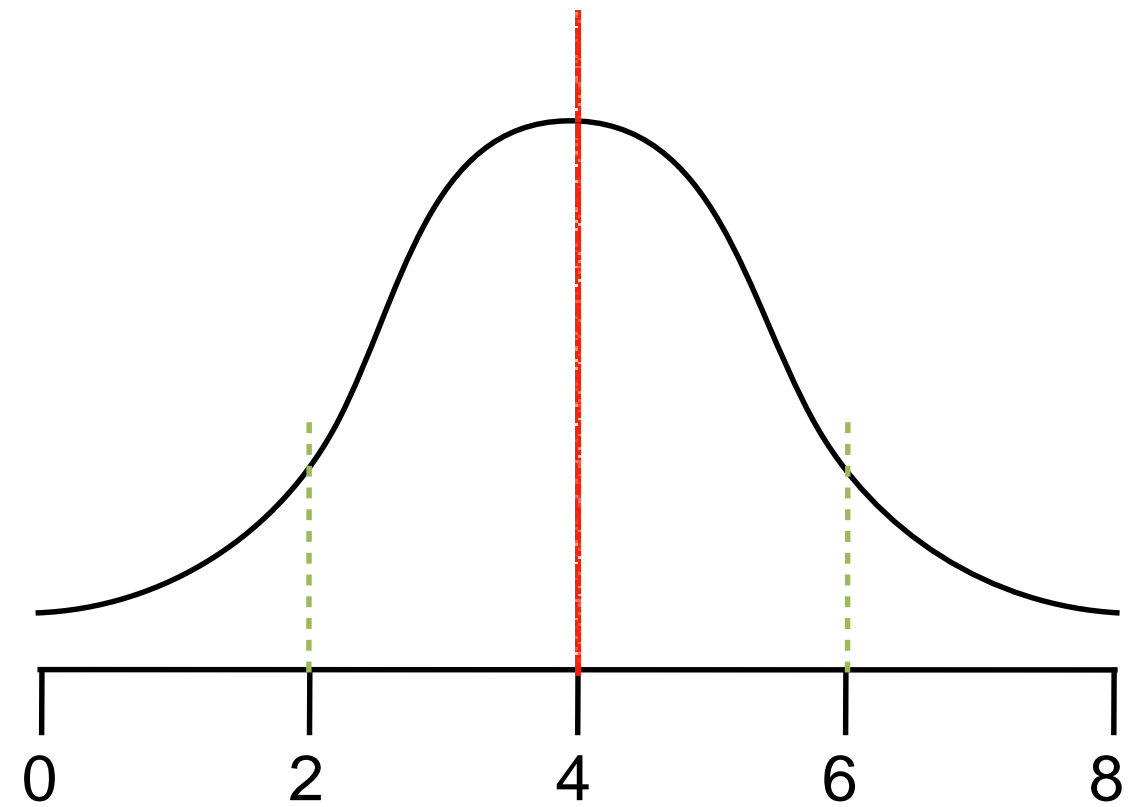DICE ANALYTICS

# Standard Deviation (SD)

- Square root of Variance
- It has the same units as the variable, which makes it useful in comparisons and calculations

- Sample SD

$$s = \sqrt{s^2} = \sqrt{\frac{\sum (X - \overline{X})^2}{N-1}}$$

- Population SD

$$\sigma = \sqrt{\sigma^2} = \sqrt{\frac{\sum (x - \mu)^2}{N}}$$

**DICE** ANALYTICS

# Spread



Less Spread

Low Variance

Low Deviation

More Spread

High Variance

High Deviation

# Robust Statistics

- Measures on which extreme observations or outliers have little effect

|  | Robust | Non-Robust |
|---|---|---|
| Spread | IQR | SD, Range |
| Center | Median | Mean |

**Skewed**　　　　　**Symmetric**

DICE ANALYTICS
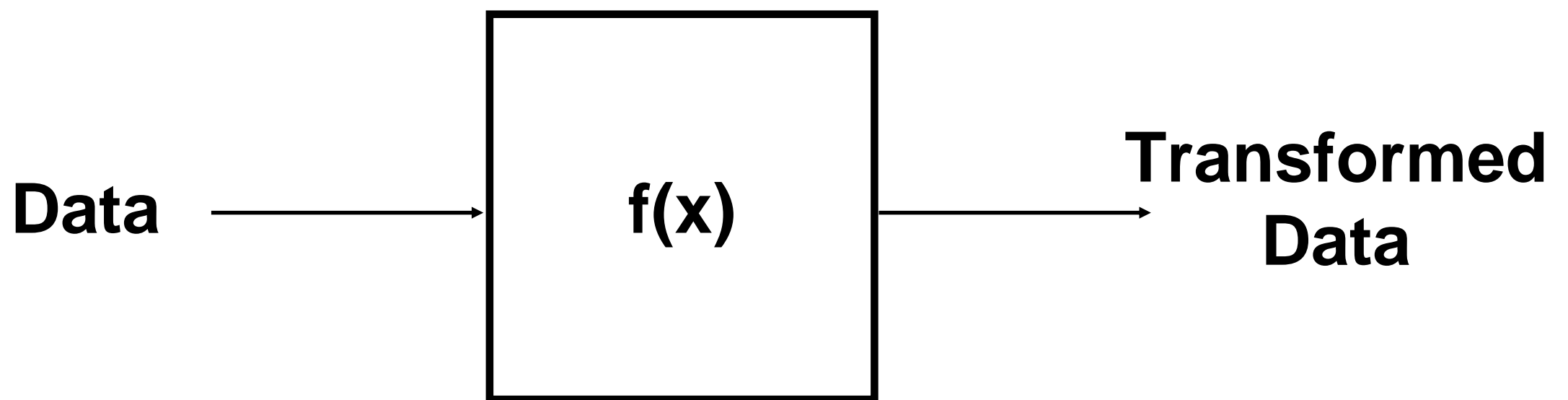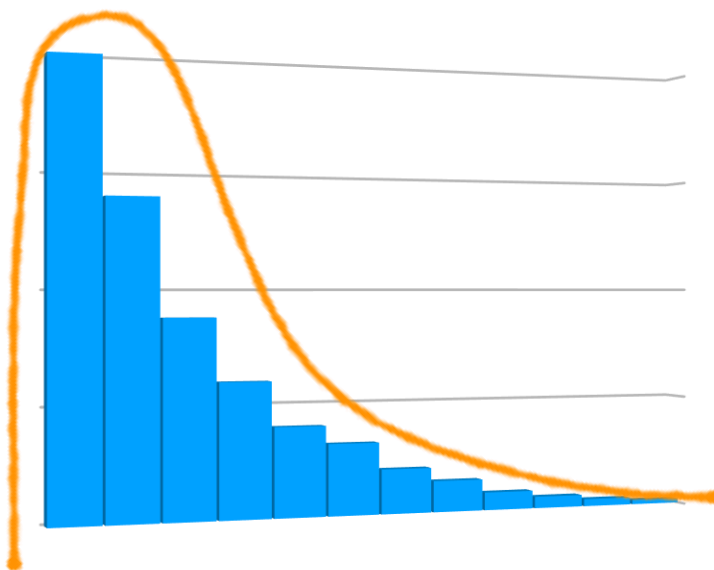
# Data Transformations

- Applying a Function f(x) to adjust scales of data.
- Done usually when data is skewed, so that it becomes easier to perform *modelling.*
- Done to convert non-linear relationship into a linear relationship.

**Data** → **f(x)** → **Transformed Data**

# (Natural) Log Transformation

- To transform data that is positively skewed
- Usually done when data is concentrated near Zero (relative to the few large values in data)
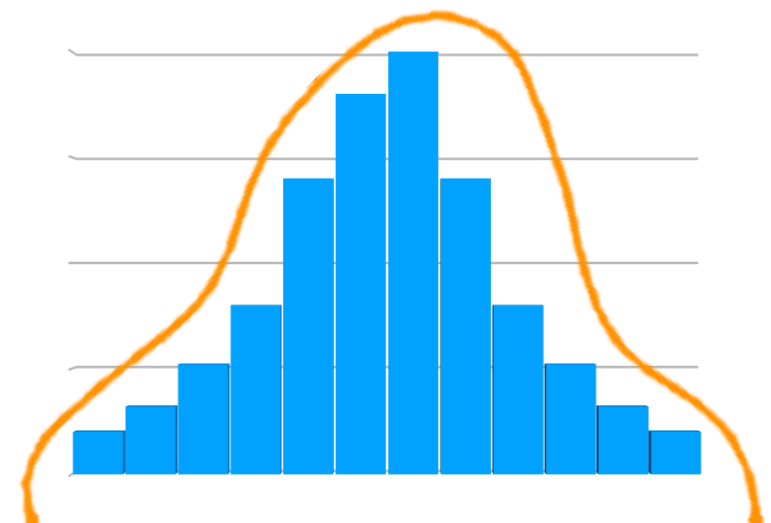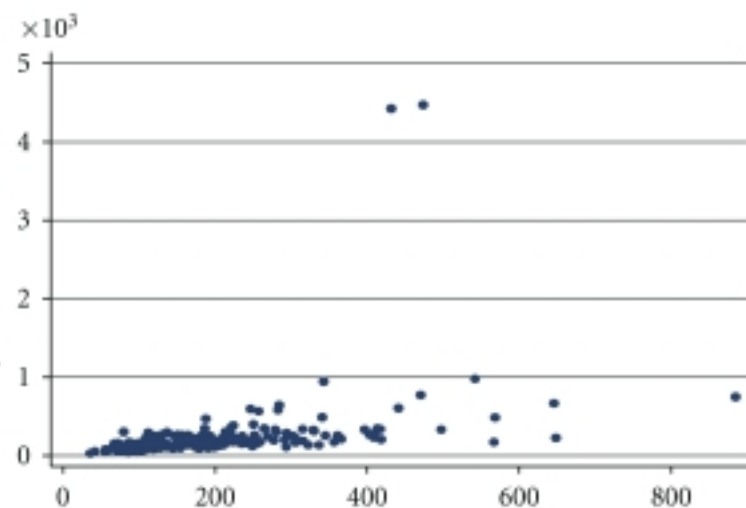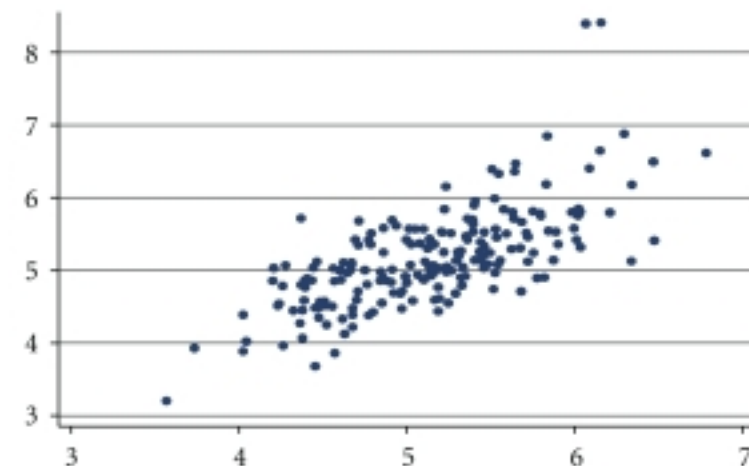
Right Skewed

Natural
Log

Symmetric

DICE ANALYTICS

# Log Transformation

- To make the relationship between two variable more linear
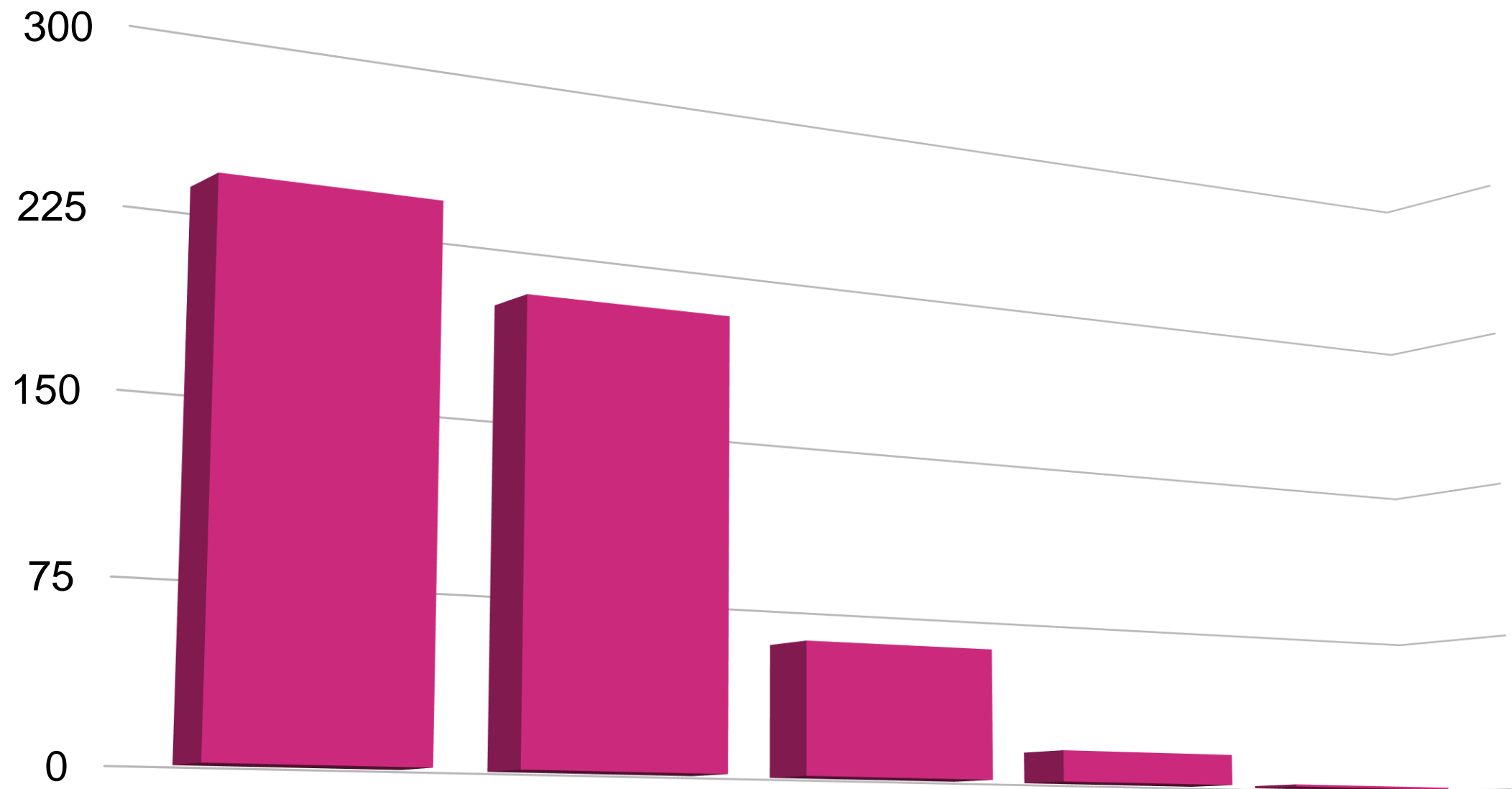- Most of the simple methods for modelling work only when relationship is linear

# Other Transformation

- You may use other transformations or create of your own

- For instance: Square Root, Square, Inverse

DICE
ANALYTICS

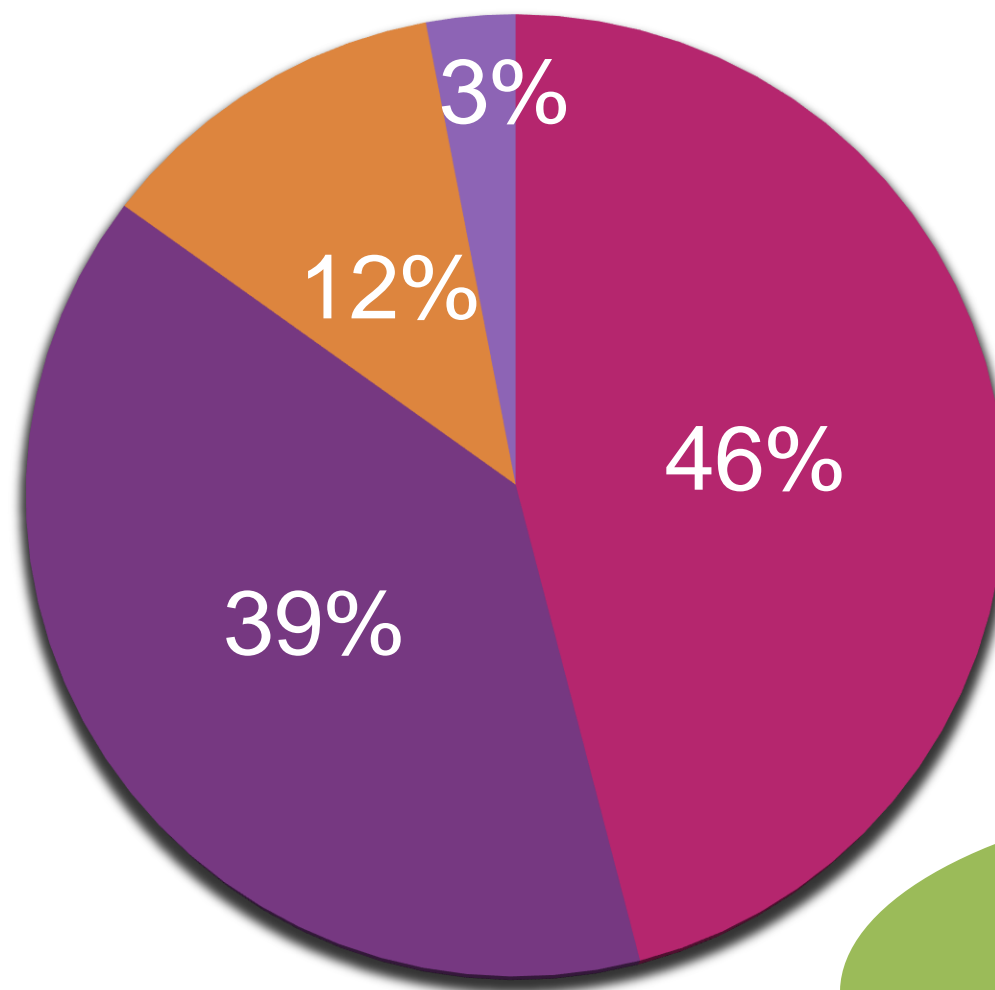# Visualising Categorical Data

# Bar Plot



Frequency

DICE ANALYTICS

# Bar Plot vs Histogram

- Bar Plot for Categorical Variables, Histogram for Numerical Variables

- X-axis in Histogram must be a Number Line

- Ordering of bars is not interchangeable in Histogram as compared to Bar Plot

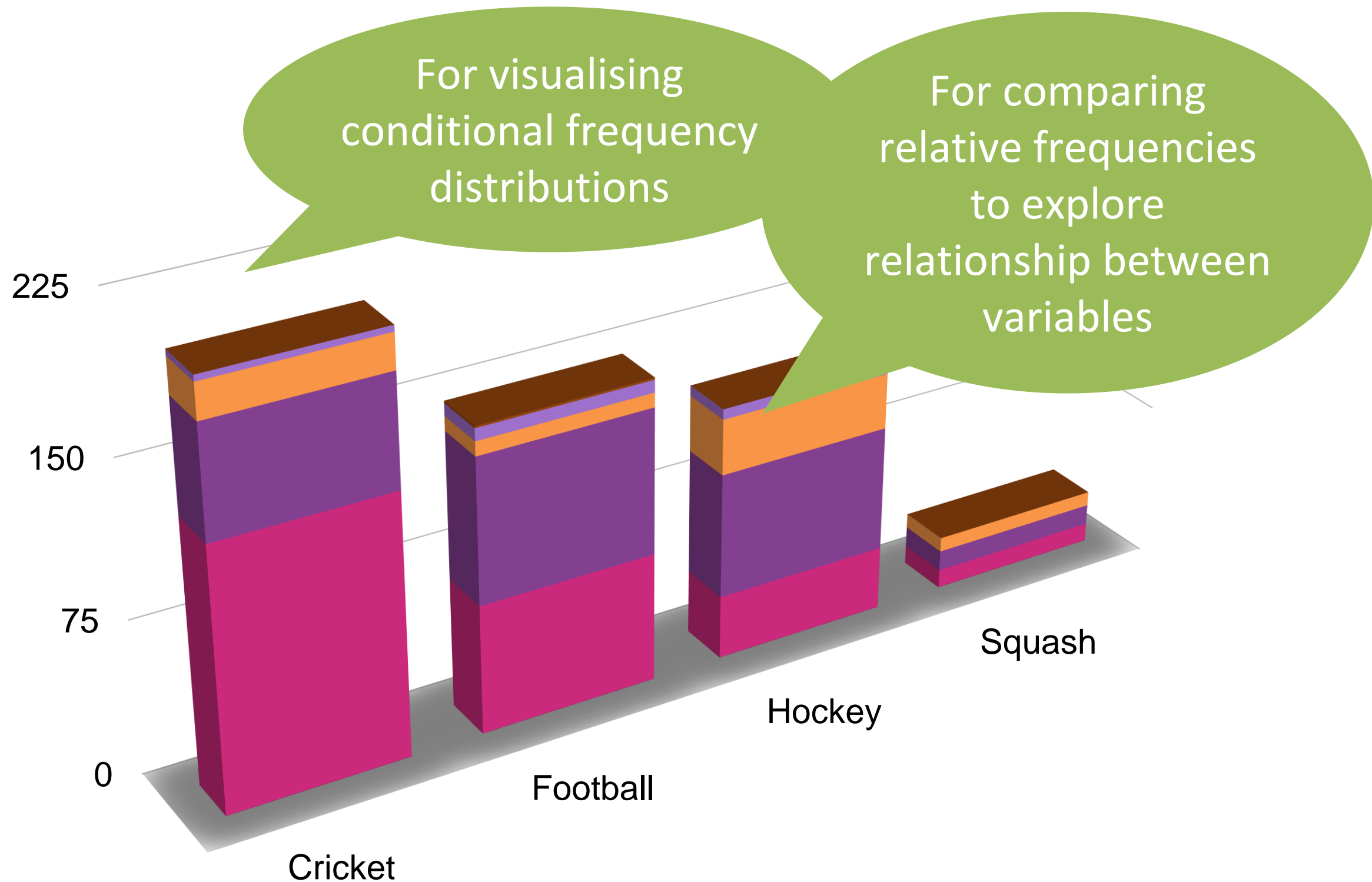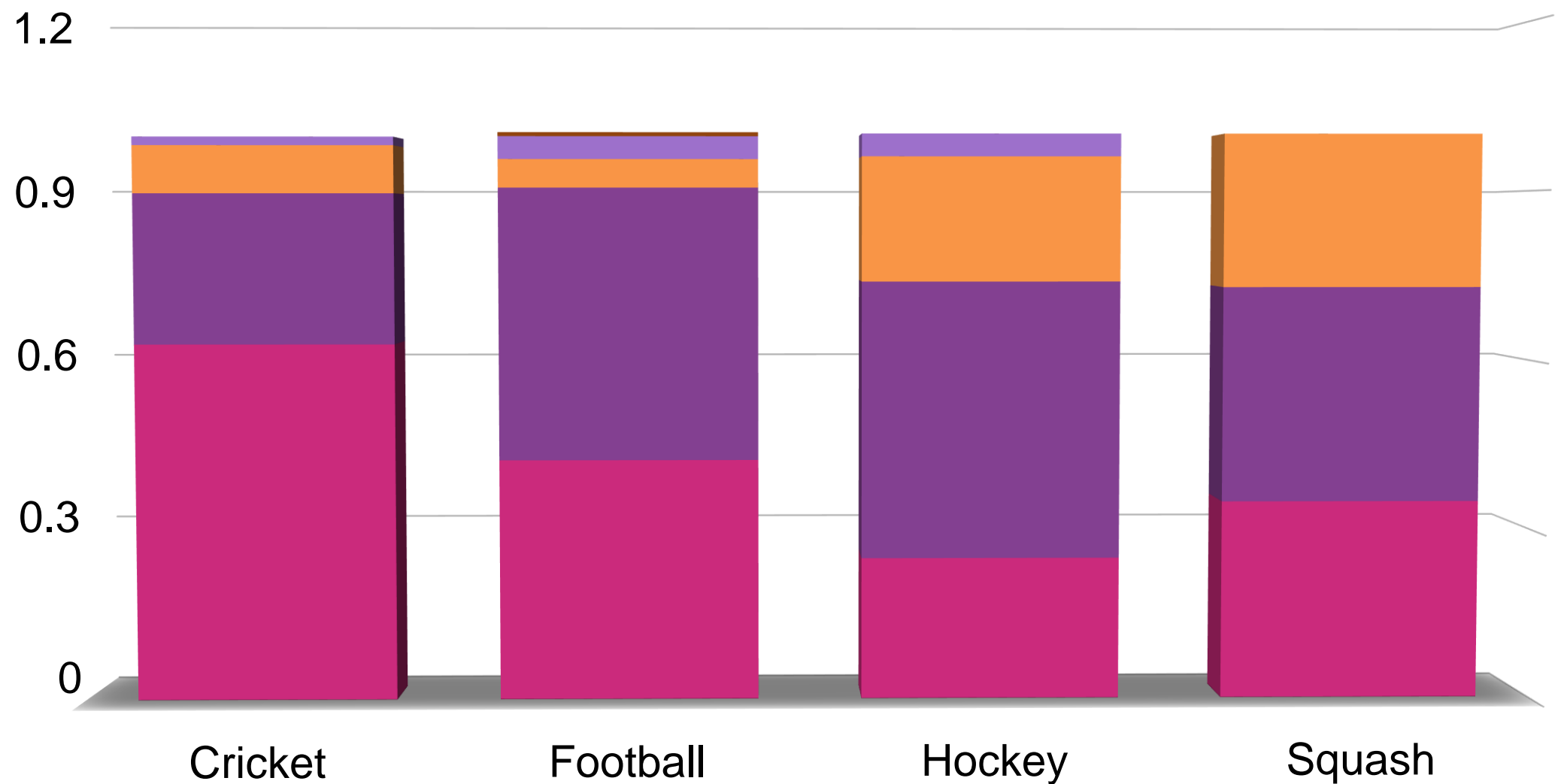**DICE** ANALYTICS

# Pie Chart

Cricket  Football  Hockey  Squash  Not Sure

3%

12%

46%

39%

Use Bar Plot instead

DICE ANALYTICS
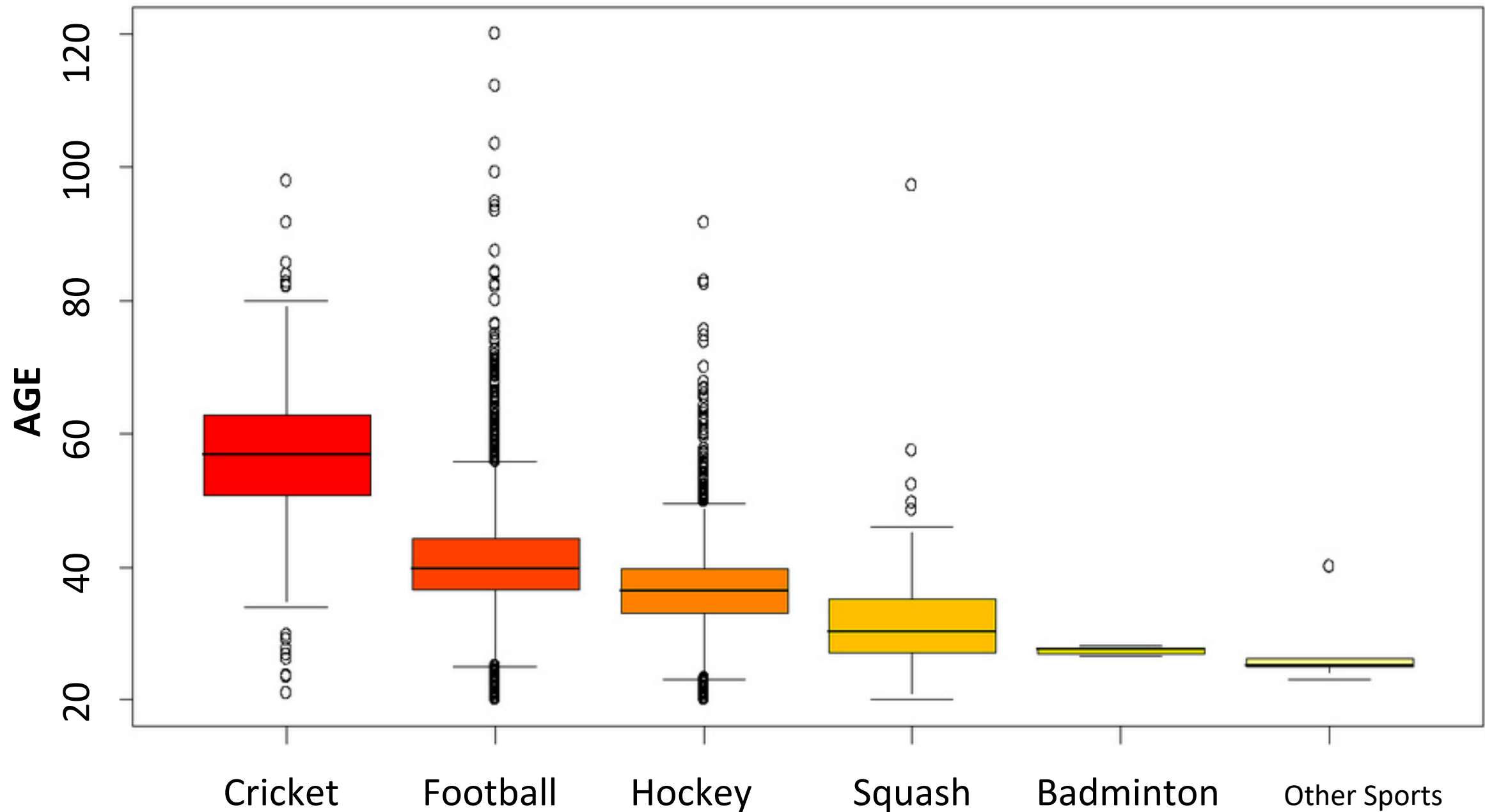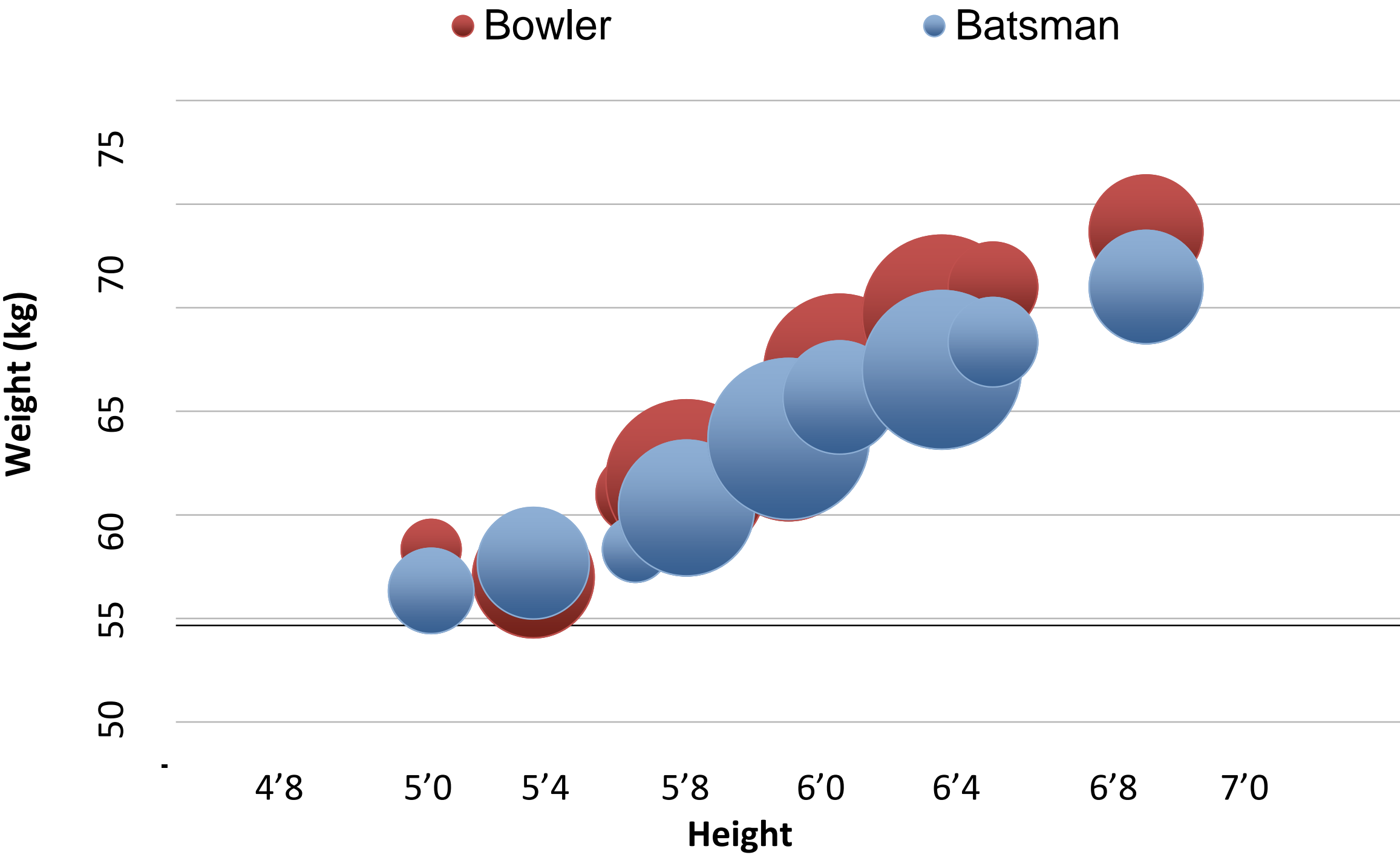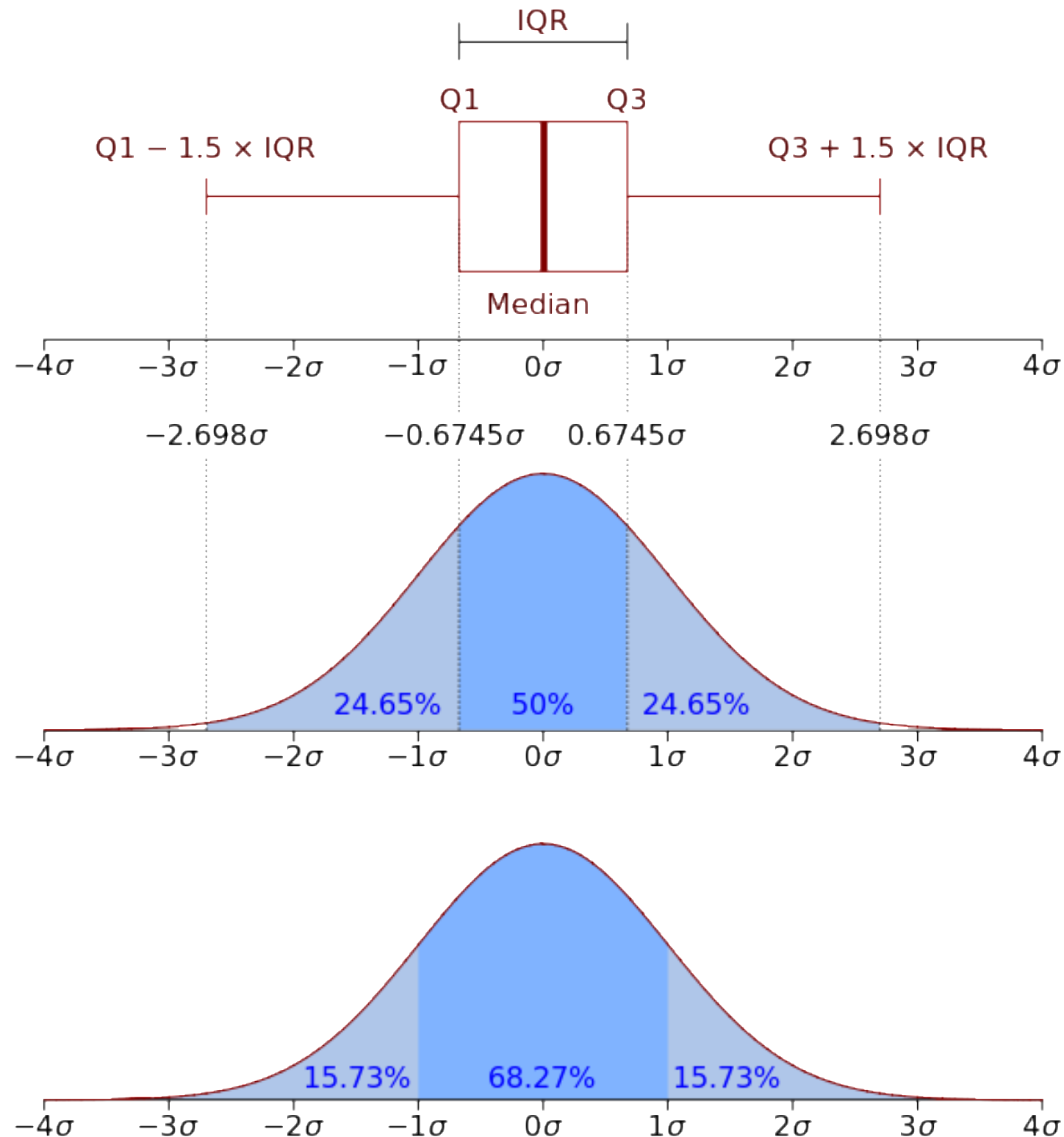
Relative Frequency Segmented Bar Plot

# Side-by-Side Box Plots

Building density against Urban Atlas code

# Outliers

# Why do EDA

- To understand data properties

- To find patterns in data

- To suggest modelling strategies

- To "debug" analyses

- To communicate results

(From JHU)

DICE
ANALYTICS

# Why do EDA

https://www.youtube.com/watch?v=jbkSRLYSojo

DICE
ANALYTICS