



# Outlier Engineering

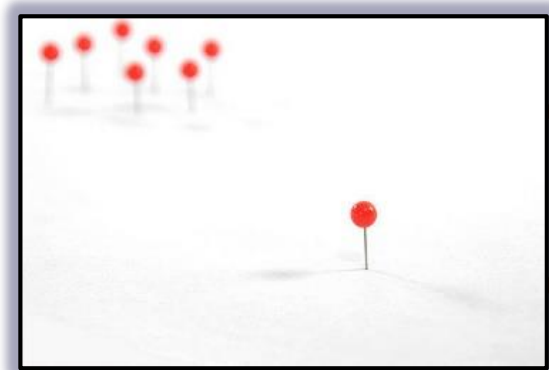


# Outliers

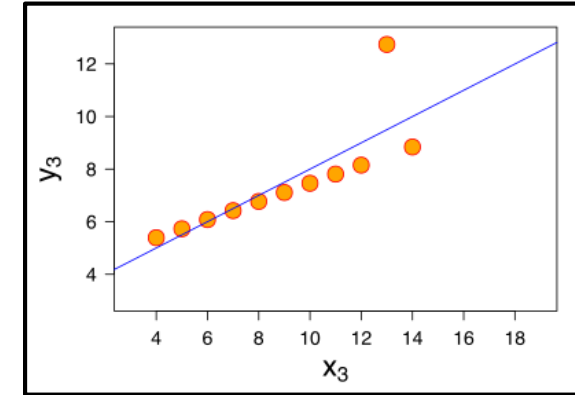
- An outlier is a data point which is significantly different from the remaining data.
- “An outlier is an observation which deviates so much from the other observations as to arouse suspicions that it was generated by a different mechanism.” [D. Hawkins. Identification of Outliers, Chapman and Hall , 1980.]



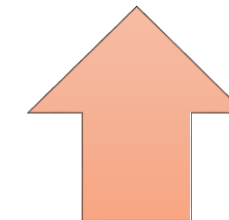
# Algorithms susceptible to outliers



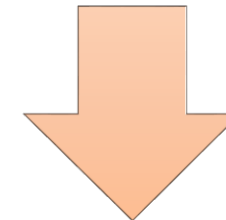
Linear  
models



Adaboost



Tremendous  
weights



Bad  
generalisation



# Handling Outliers

What can we do if we find outliers in our variables?



# Ways to engineer outliers

Trimming

- Removing outliers from the data set

Missing data

- Treat outliers as missing data and perform missing data imputation
- Any technique in section 4

Discretisation

- Put outliers into lower / upper bins
- Any technique in section 8

Censoring

- Capping
- Top / Bottom coding
- Winsorization

# Ways to engineer outliers

Trimming

- Removing outliers from the data set

Missing data

- Treat outliers as missing data and perform missing data imputation
- **Any technique in section 4**

Discretisation

- Put outliers into lower / upper bins
- **Any technique in section 8**

Censoring

- Capping
- Top / Bottom coding
- Winsorization

# Ways to engineer outliers

## Trimming

- Removing outliers from the data set

## Missing data

- Treat outliers as missing data and perform missing data imputation
- Any technique in section 4

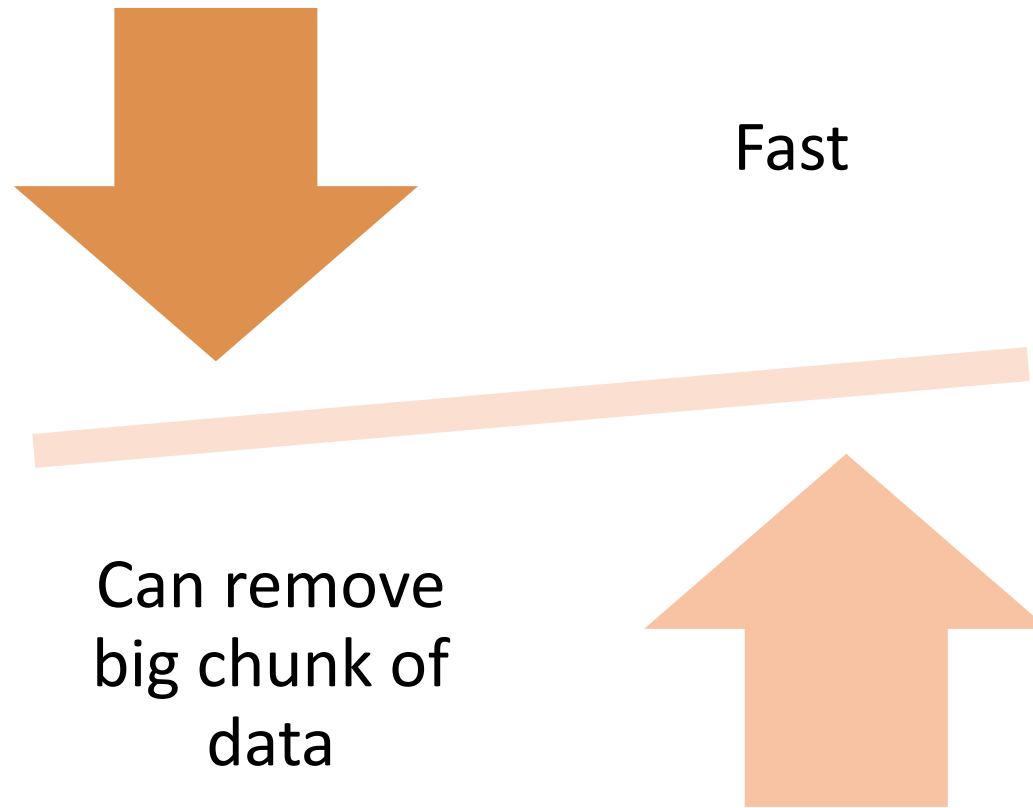
## Discretisation

- Put outliers into lower / upper bins
- Any technique in section 8

## Censoring

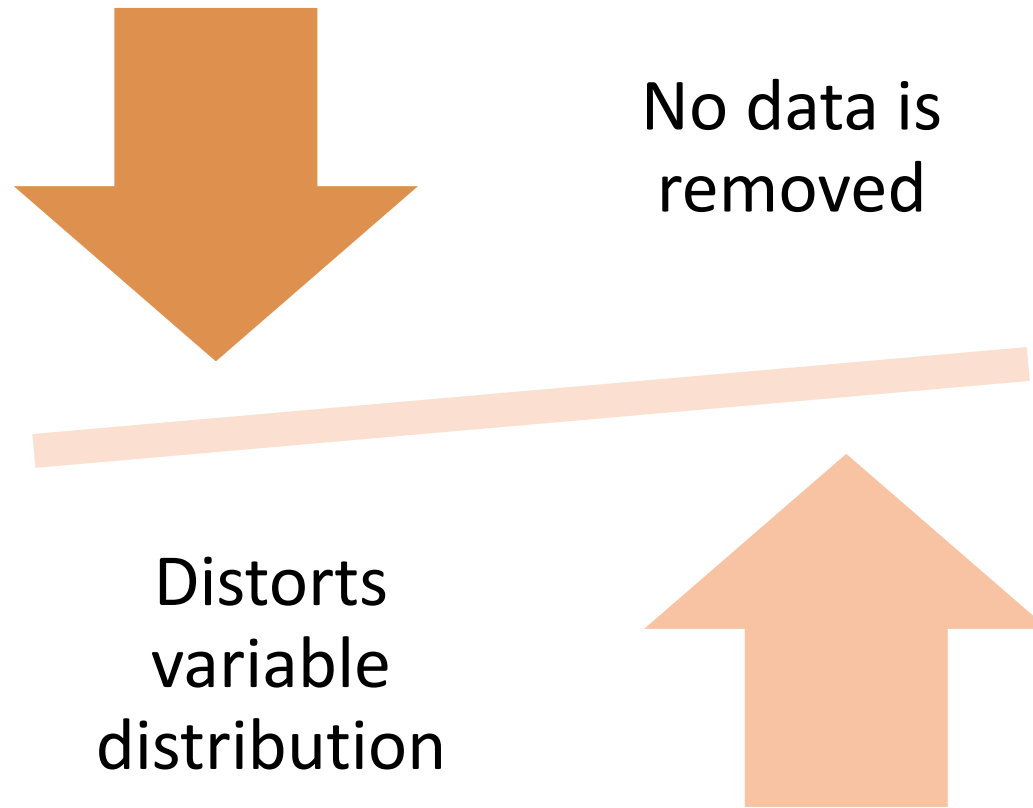
- Capping
- Top / Bottom coding
- Winsorization

# Trimming: pros and cons





# Capping: pros and cons





# Detecting Outliers

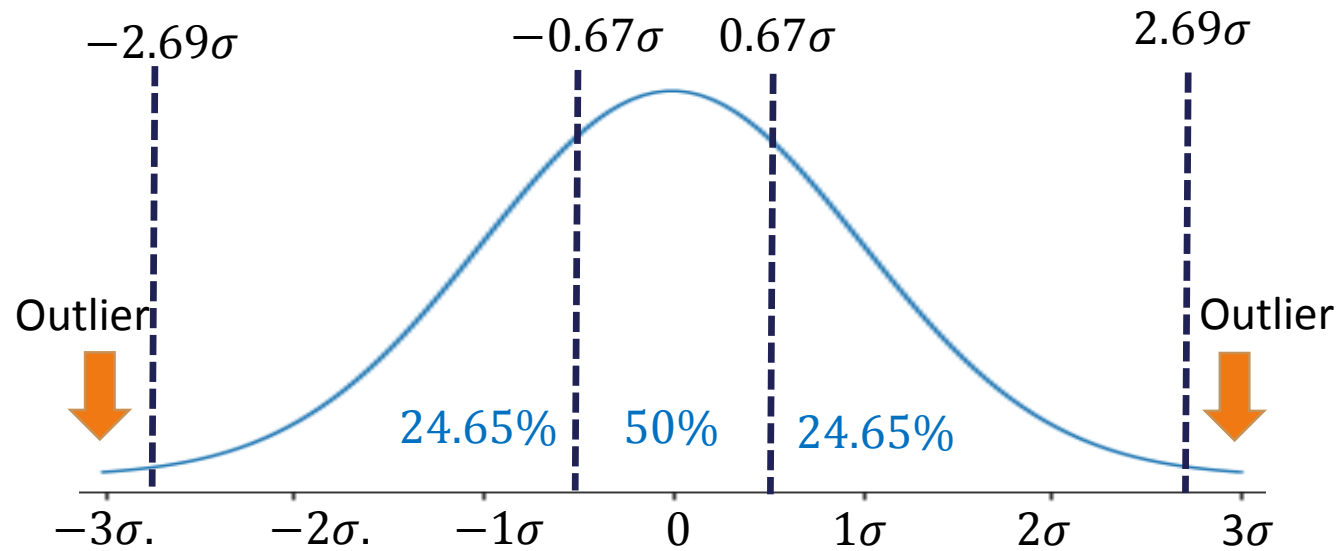
Extreme Value Analysis



# Detecting outliers

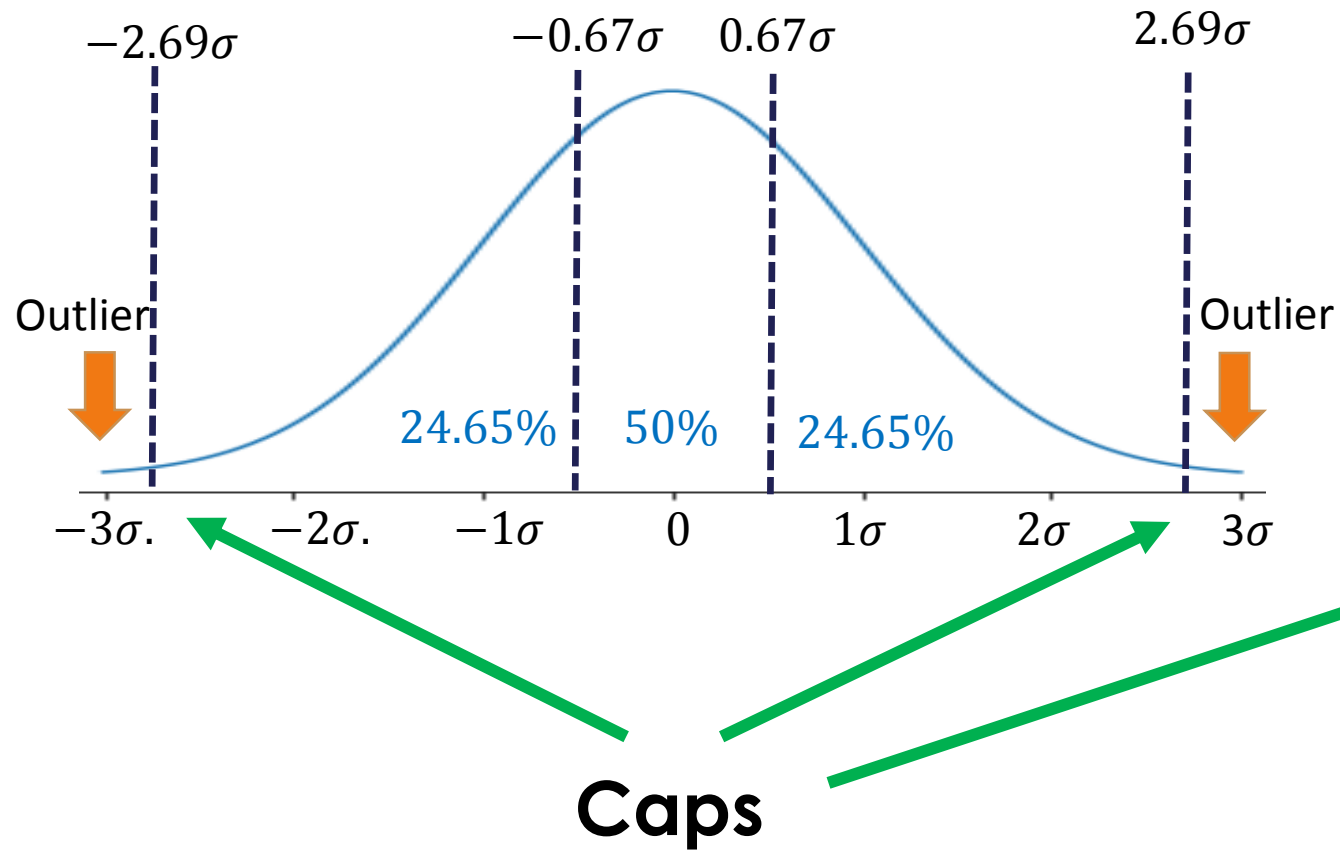
- Gaussian distribution (mean and std)
- Inter-quantal range proximity rule
- Quantiles

# Normal distribution



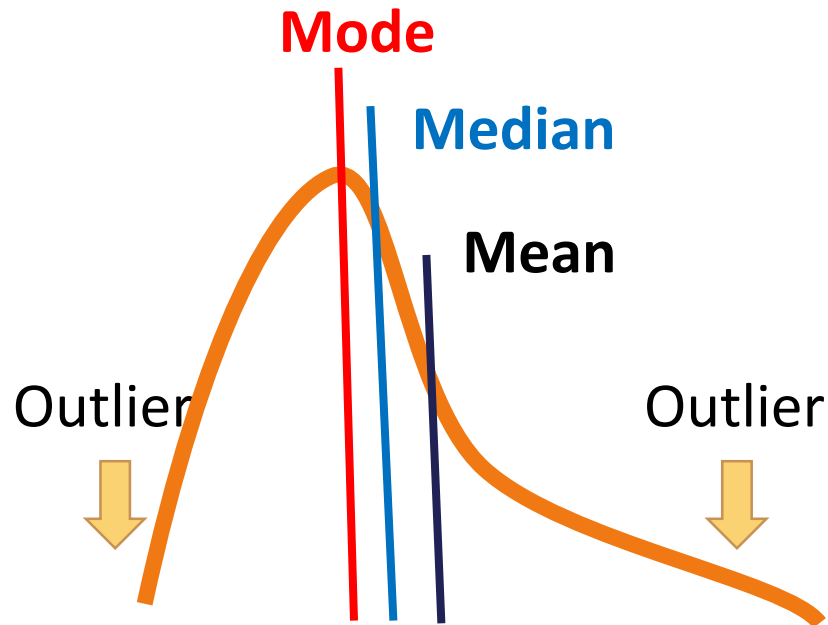
- ~99% of the observations of a normally distributed variable lie within the mean  $\pm 3 \times$  standard deviations.
- Values outside **mean  $\pm 3 \times$  standard deviations** are considered outliers

# Normal distribution



- ~99% of the observations of a normally distributed variable lie within the mean  $\pm 3 \times$  standard deviations.
- Values outside **mean  $\pm 3 \times$  standard deviations** are considered outliers

# Skewed distributions

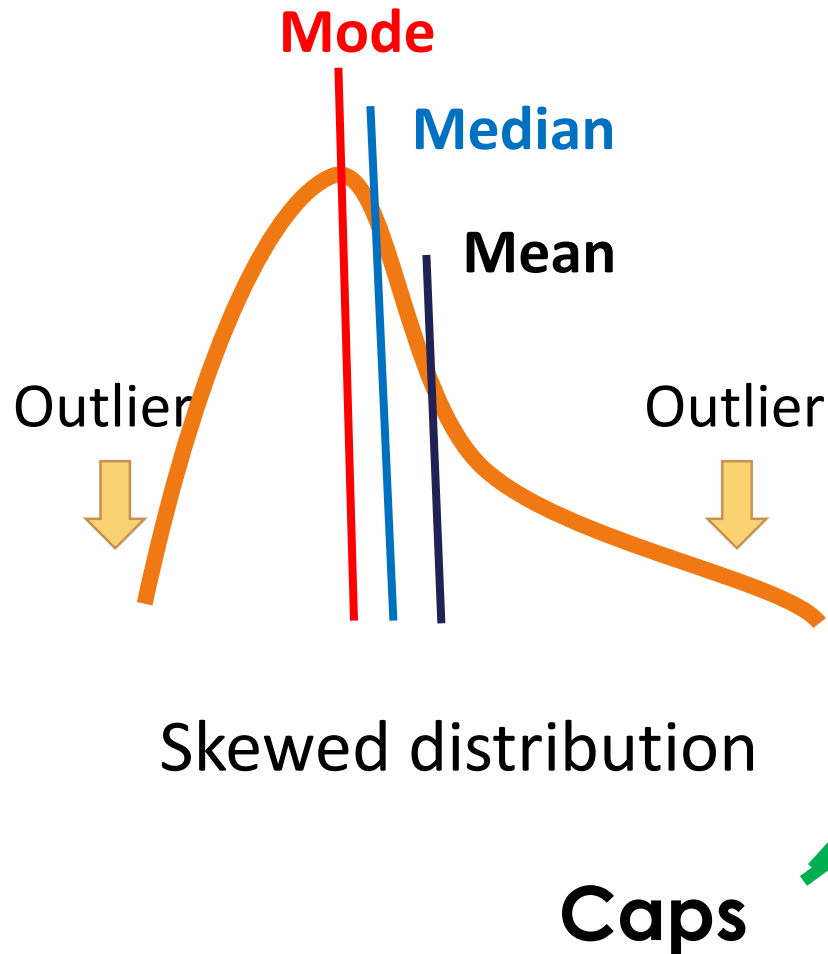


Skewed distribution

- The general approach is to calculate the quantiles, and then the inter-quantile range (IQR), as follows:
- $IQR = 75^{th} \text{ Quantile} - 25^{th} \text{ Quantile}$
- **Upper limit =  $75^{th} \text{ Quantile} + IQR \times 1.5$**
- **Lower limit =  $25^{th} \text{ Quantile} - IQR \times 1.5$**

Note, for extreme outliers, multiply the IQR by 3 instead of 1.5

# Skewed distributions



- The general approach is to calculate the quantiles, and then the inter-quantile range (IQR), as follows:

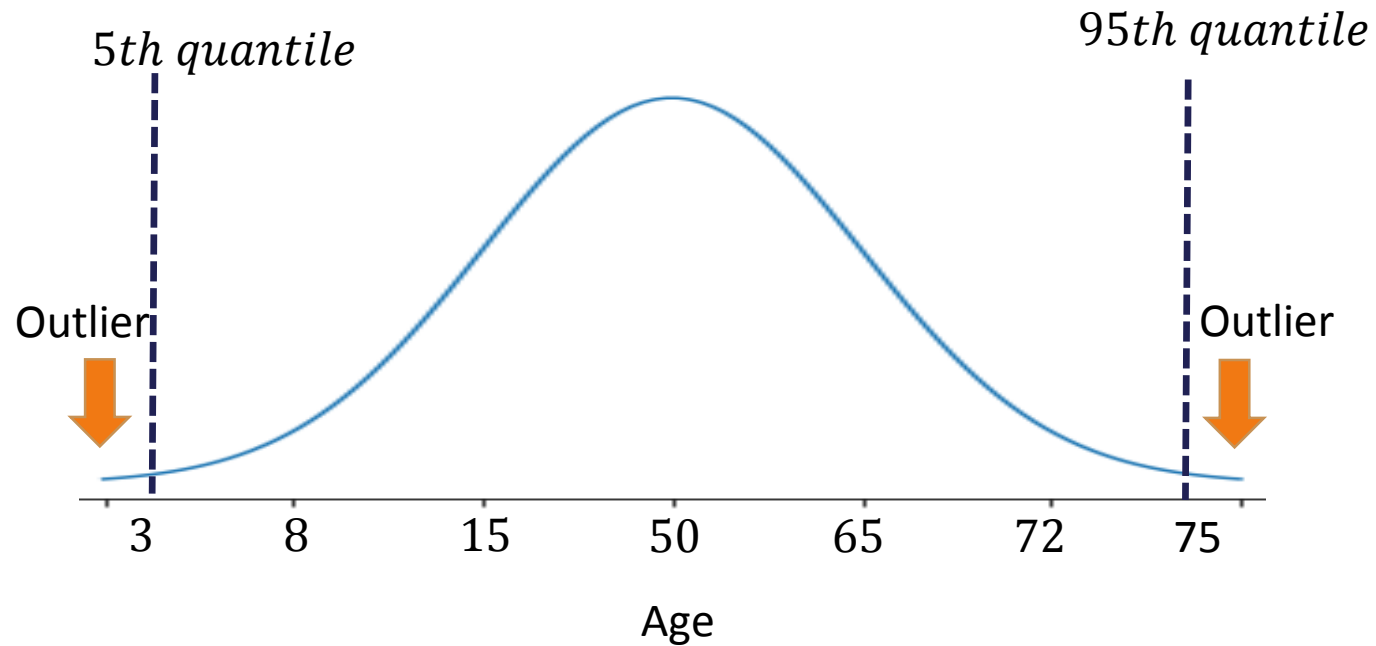
- $IQR = 75^{\text{th}} \text{ Quantile} - 25^{\text{th}} \text{ Quantile}$

- **Upper limit =  $75^{\text{th}} \text{ Quantile} + IQR \times 1.5$**

- **Lower limit =  $25^{\text{th}} \text{ Quantile} - IQR \times 1.5$**

Note, for extreme outliers, multiply the IQR by 3 instead of 1.5

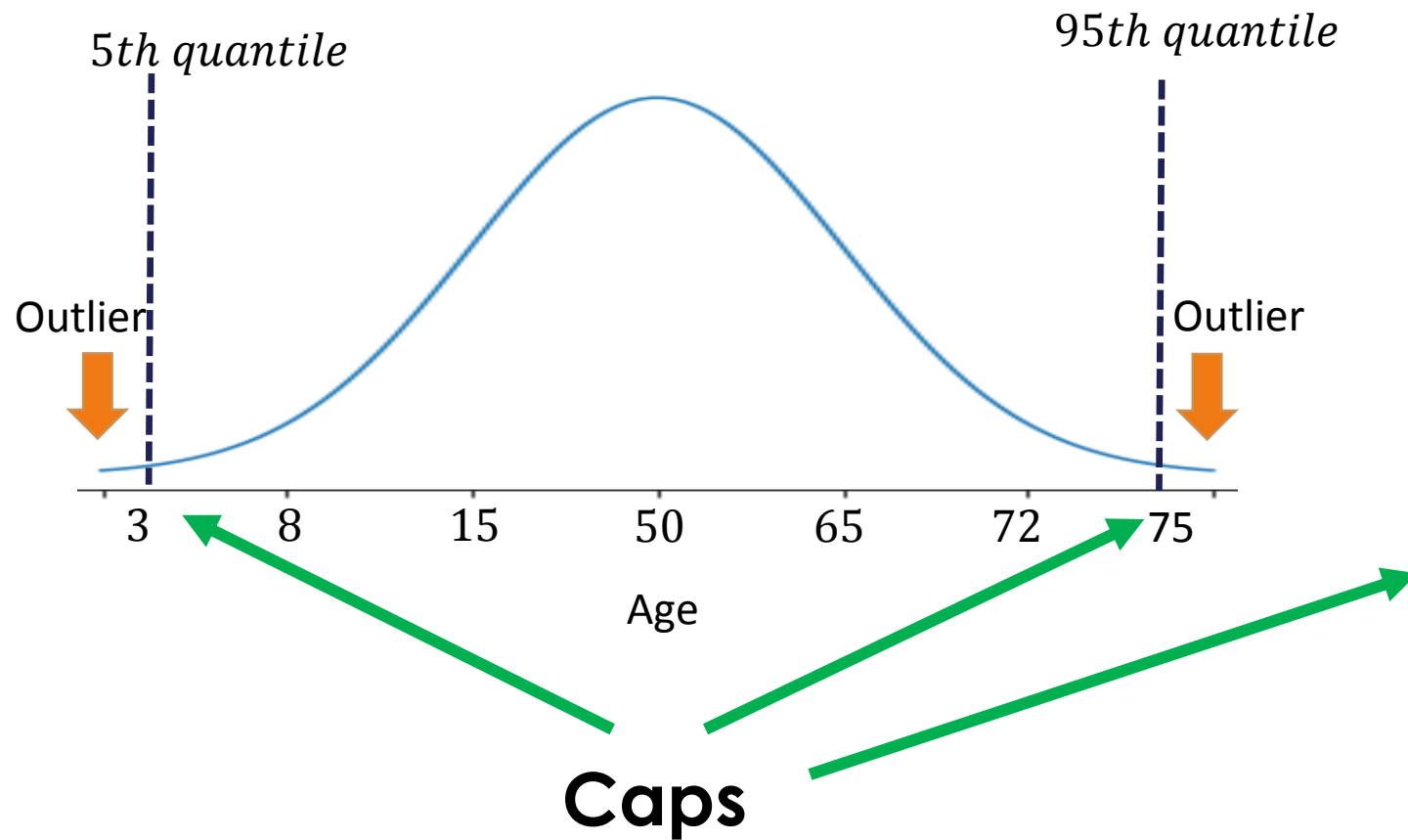
# Normal distribution



- ~95% of the observations above the 5<sup>th</sup> quantile
- ~95% of the observations below the 95<sup>th</sup> quantile
- Values **above or below the 95<sup>th</sup> or 5<sup>th</sup> quantile** are considered outliers



# Normal distribution



- ~95% of the observations above the 5<sup>th</sup> quantile
- ~95% of the observations below the 95<sup>th</sup> quantile
- Values **above or below the 95<sup>th</sup> or 5<sup>th</sup> quantile** are considered outliers

# Accompanying Jupyter Notebook



- Accompanying Jupyter Notebook
- How to perform outlier engineering:
  - Gaussian approximation
  - IQR
  - Quantiles
  - Pandas and Feature-engine

# THANK YOU

[www.trainindata.com](http://www.trainindata.com)