# Handling Missing Data

① Remove all missing rows

② (a) Replace with generic substitute values (mean/median/s.t.e)

③ Imputation: Estimate a probability model for the missing variable & replace the missing value with one or more samples from probability model.

(or most frequent if value is categorical)

→ ## Types of Missing Data

Data may be missing for a variety of reasons.
→ corrupt during its transfer or storage

① MCAR — Missing Completely At Random

— pros of an observation being missing does not depend on observed or unobserved measurements.

Ext movie rating from users, since some movies are more popular than others, some movies may not have ratings = | NOT MCAR |

② MAR → given the observed data, the probability that data is missing does not depend on unobserved data

Ext

| Y | Gender | Race | Income |
|---|--------|-------|--------|
|   | M | Asian | 800 |
|   | F | Indian |  |
|   | M | American |  |

If missing data depends on Gender/Race then it is MAR

If not its MCAR.

## Missing Data & R

Dealing with Outliers     Trandical common outliers

→ Winsorization — Shorten outliers

→ Robustness — Keep the outliers & analyze data using a robust procedure.

Detecting Outliers
① values below the alpha percentile (or 100-α percentile)
② values more than c times std. dev^n from the mean.

→ follows 1^st technique (when we assume data is gaussian)

→ Issue :- outliers can affect mean and other calculations

→ So to avoid this remove most extreme
(or) just use percentiles (more robust)

(From - 8) DVA

Data Transformations : Skewness & Power Transformations

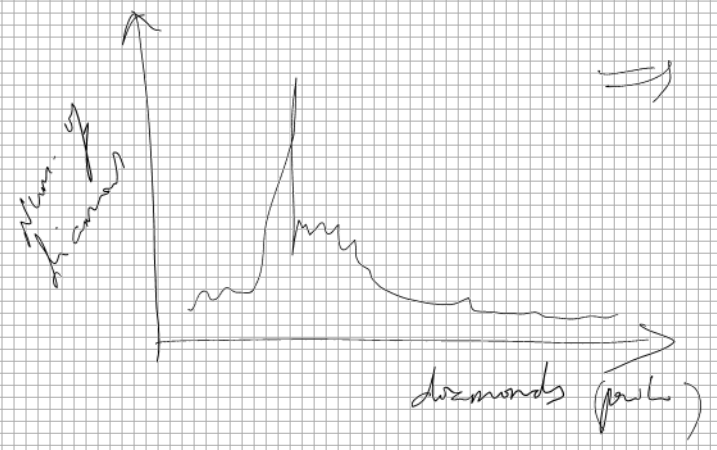Data is generally drawn from a highly skewed distribut^n and that is not well described by a common distribution.
— A single transformation may map the data to a form that is well described by common distribution.

→ One transformed a suitable model can then be fitted to data.
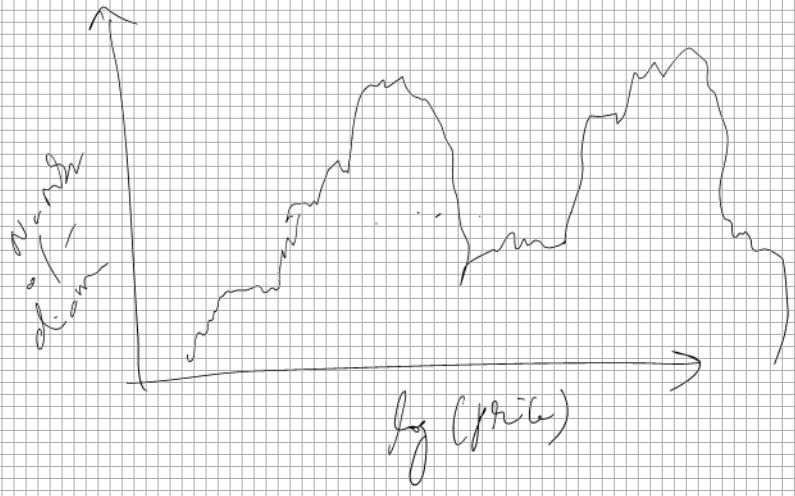
→ Power Transformation family ( Next from video lecture )

$$f_\lambda(x) = \begin{cases} (x^\lambda - 1)/\lambda & \lambda > 0 \\ \log x & \lambda = 0 \quad x > 0 \\ -(x^\lambda - 1)/\lambda & \lambda < 0 \end{cases}$$

$\Longrightarrow$ to predict T price of diamond

to form *Insights*

Number of diamonds

diamonds (price) $\curvearrowright$ Transform data

Number of diamonds

log (price)

Transformed data $\Longrightarrow$ looks like it is

*Bi-modal* in nature.