# Handling Missing Data

① Remove all missing rows

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

③ <u>Imputation</u> :- Estimate a probability model \ or most frequent for the missing variable & replace the missing value with one or more samples from probability model. \ (categorical if value is

→ <u>Types of Missing Data</u>

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 = $\boxed{NOT\ MCAR}$

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

Ext

| Y | Gender | Race | Income |
|---|--------|------|--------|
|   | M | Asian | 888 |
|   | F | India~ |   |
|   | M | American |   |

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

If not its MCAR.

<u>Missing Data & R</u>

Dealing with Outliers → Winsorizing common outliers → Winsorization — Shrink 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 1st technique (when we assume data is gaussian)

→ Issue: outliers can affect mean and other calculations

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

(from -8) DVA

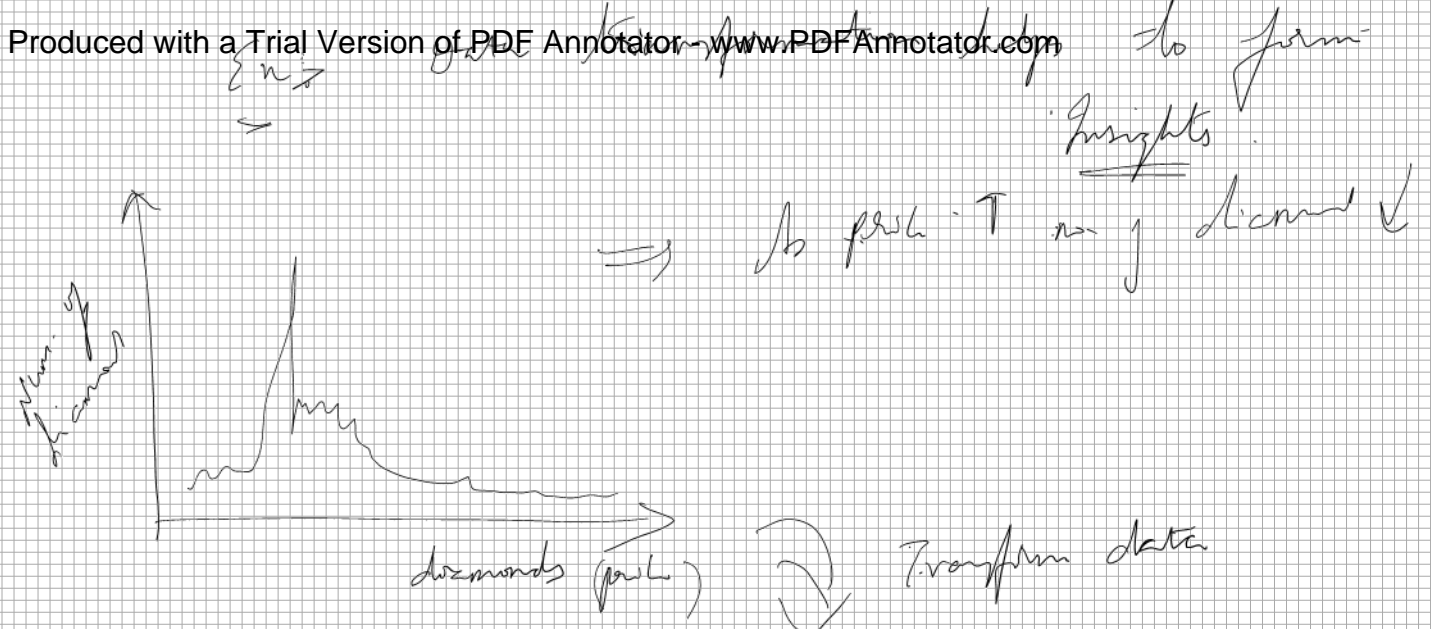## Data Transformations : Skewness & Power Transformations

Data is generally drawn from a highly skewed distribution and that is not well described by a common distribution.

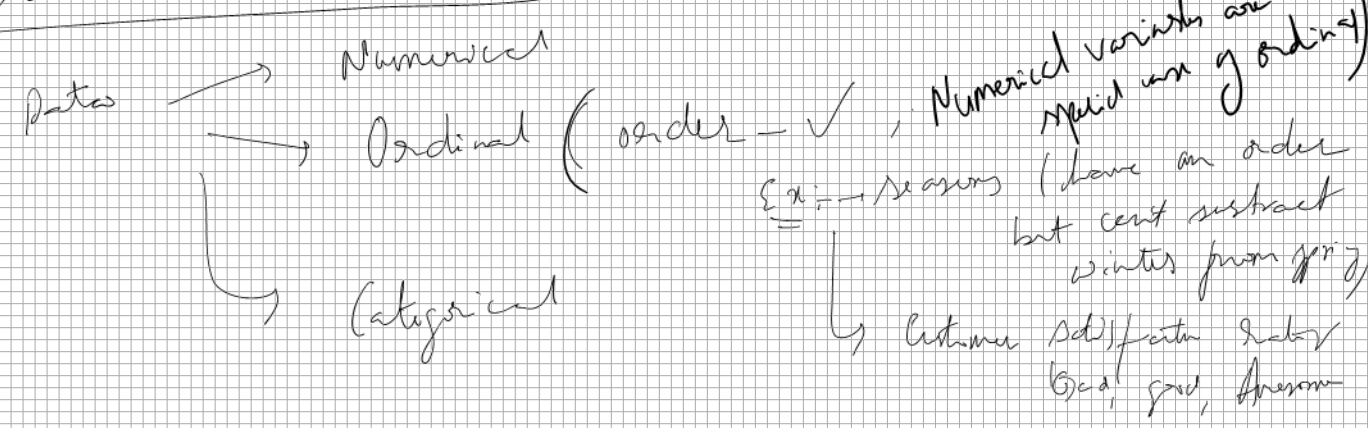→ A single transformation may map the data to a form that is well described by common distributions.

→ Once transformed a suitable model can then be fitted to data.

→ Power Transformation family ( Neat from video lecture )

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

$\Sigma n_{?}$     to form

Insights

$\Longrightarrow$   to predict T on g diamond

diamonds (price)    Transform data

Transformed data

looks like it is

Bi-modal in nature

log (price)

$\longrightarrow$ Data Transformation : Binning

Data $\longrightarrow$ Numerical

$\longrightarrow$ Ordinal ( order - ✓ , Numerical variables are a special use of ordinal)

$\Sigma x$: — seasons ( have an order but can't subtract winter from spring )

$\longrightarrow$ Categorical

$\longmapsto$ Customer satisfaction rating Bad, good, Awesome

$\longrightarrow$ Data Transformation $\longrightarrow$ Indicator Variable

$\longrightarrow$ Indicator variables are used in conjunction with Binning.

Bin — 1-5 $\Rightarrow$ 0 ; 5-16 $\Rightarrow$ 1 , 16→ $\Rightarrow$ 2

$\Sigma x$:

| data | Bins | Indicator |
|------|------|-----------|
| 1 | 0 | 000 |
| 4 | 0 | 000 |
| 15 | 1 | 010 |
| 20 | 2 | 100 |

# MCAR Vs MAR

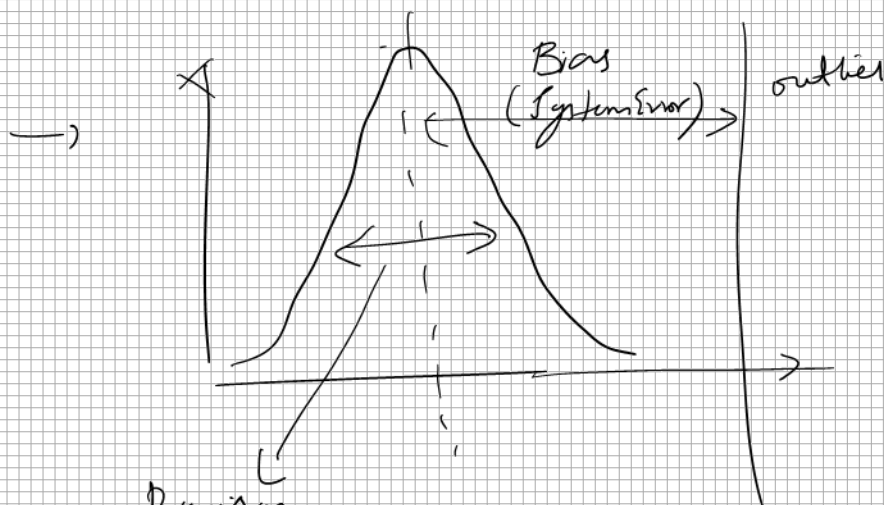↳ Any of the 3 techniques used for handling missing data are reasonable, though some are better.

For MAR : The 1st two techniques may introduce some bias, the 3rd technique is maybe reasonable depending on how well the probability model is built.

Outliers Types
- → Corrupted values ( human Err )
- → Unlikely values
  ↳ are substantially unlikely values given our modelling assumptions.

Robustness :- lack of sensitivity of data analysis procedures to outliers.

"mean" gets affected by outlier, median doesn't.



Bias (System Error) → outlier

Precision (uncertainty) also referred to as random Error.

→ Handling outliers
- → Truncating
- → Winsorization - replace it with acceptable Extreme values
- → Robustness :- keep outlier in data, but using robust procedure

... frequency of value,

we may use median rather than mean.

i.e median is more robust to outliers.

> * winsorize ( df $ x, [1::100] )
> ↳ library (robustHD)

---

Tall Data :   one or more columns are Keys
              Atleast one column is value

↳ is convenient for adding new records incrementally
  & for removing old records.
                  i.e simply adding additional rows.

→ melt() :  Converts wide data to tall data.

↳ Not Easy for summarizing.

---

Wide Data :  Represents in multiple columns the information
             that tall data holds in multiple rows.

↳ simpler to analyze.
    → no need to pass over all data, jump to
specific location.

↳          to add/remove Entries.

Ex :

| Date       | apples | oranges |
|------------|--------|---------|
| 2015/01/01 | 200    | 150     |
| 2015/01/02 | 220    | 130     |

Tall data:

| 2015/01/01 | apples  | 200 |
| 2015/01/01 | oranges | 150 |
| 2015/01/02 | apples  | 220 |
| 2015/01/02 | oranges | 130 |

Wide Data:

Tall data

→ melt - wide ⟶ tall (reshape2 package)

→ [ acast / dcast inverse of melt ]

↳ returns an array   ↳ returns a data frame.

↳ :-   tall data to wide data.