

# Probability

## Lec 1

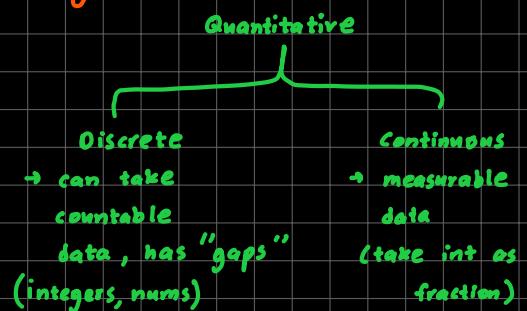
Data → facts, and figures

Dataset → Data collected in a particular study

Categorical Data → Data that can be grouped Qualitative into categories

- Primary Data: Data collected by oneself // Raw Data
- Secondary Data: Processed Data e.g.: news, TV

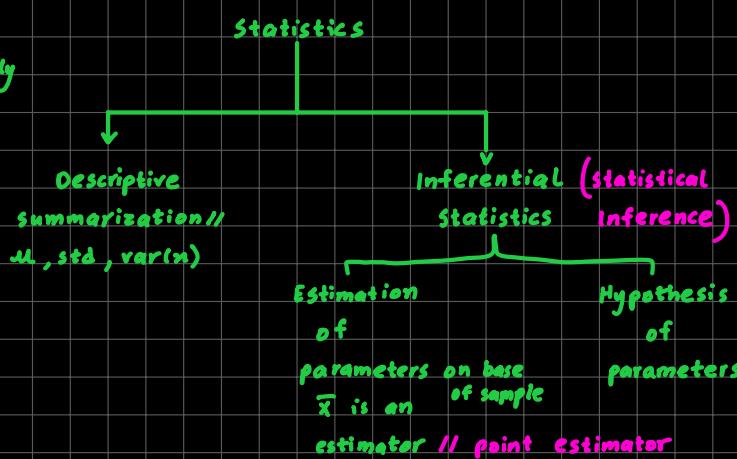
- Statistics: Numerical facts e.g. "statistics of sports"
- population: Set of all elements of interest in study → Survey to collect data on entire population → census
- sample: Subset of population → Survey to collect data on sample → Sample Survey
- Difficult to work on population, so sample is used
- observation: set of measurements obtained for a particular element parameter (measure) calc from pop // property
- $N = \text{pop size} \rightarrow \underline{\mu} = \text{pop mean}$
- $n = \text{sample size} \rightarrow \underline{x} = \text{sample mean}$
- statistic (measure) calc from sample
- Elements: Entities on which data are collected
- Sample Size: Num of Elements



- Variable: Characteristic of interest for the elements
- ↳ Qualitative: Variable that can not be expressed numerically
- Quantitative: Can be expressed numerically

Elements	Nominal	Ratio	2-Year Average Value (\$)	5-Year Average Return (%)	Expense Ratio (%)	Rating	Rank
Fund Name	Type						
American Century Infl. Disc	IE	14.37	30.53	1.41	3-Star		
American Century Tax-Free Bond	FI	10.73	3.84	0.49	4-Star		
American Century Ultra	DE	32.54	10.00	0.79	3-Star		
Artisan Small Cap	DE	16.92	15.67	1.18	3-Star		
Brown Cap Small Cap	DE	35.73	15.85	1.20	4-Star		
Brown Cap Small Cap	DE	11.47	12.26	0.53	3-Star		
Fidelity Contrafund	DE	73.11	17.99	0.89	5-Star		
Fidelity Overseas Fund	IE	48.39	23.46	0.90	4-Star		
Fidelity Select Electronics	DE	48.60	13.50	0.89	3-Star		
Fidelity St-Term Bond	FI	8.60	2.88	0.45	3-Star		
Gabelli Asset A	DE	49.81	16.70	1.36	4-Star		
Kidder Peabody Sm Cpt	DE	15.30	15.31	1.32	3-Star		
Kidder Peabody Sm Cpt	DE	17.14	15.31	1.31	3-Star		
Mathews Pacific Tiger	IE	27.86	32.70	1.16	3-Star		
Oakmark I	DE	40.37	9.51	1.01	2-Star		
Power Energ Mktgs Bd D	FI	10.00	13.50	1.25	3-Star		
RS Value Fund	DE	26.27	23.68	1.36	4-Star		
T. Rowe Price Latin Am.	IE	53.89	51.10	1.24	4-Star		
T. Rowe Price Mid Val	DE	22.46	16.91	0.80	4-Star		
Thomson Income A	DE	37.53	15.40	1.27	3-Star		
USAA Income	FI	12.10	4.31	0.62	3-Star		
Vanguard Equity Inc	DE	24.42	13.41	0.29	4-Star		
Vanguard Fund Inv	FI	15.06	2.37	0.16	3-Star		
Vanguard Sm Cpt Indx	DE	32.58	17.01	0.23	3-Star		
Wasatch Sm Cpt Growth	DE	35.41	13.98	1.19	4-Star		

Source: Morningstar Funds500 (2008)



## Lec 2

### Measurement Scales:

① Nominal → Data consists of labels/names used to identify an attribute of the element e.g.: FundType

→ numeric/non-numeric codes may be used e.g.: 1 → coke, 2 → sprite, 3 → Fanta

② Ordinal → Data has properties of nominal scale and order/rank of data is meaningful  
e.g.: Morning-star rank → 1-star // describes the rank, + has an order 1st, 2nd, 3rd

③ Interval → Data has properties of ordinal scale

→ interval b/w values expressed in terms of fixed unit of measure  
→ Always numeric

e.g.: SAT Scores → can be ranked in terms of performance

→ Differences b/w intervals are meaningful → st 1 (620-550) 70 more points than st2.

④ Ratio → Data has properties of interval scale + ratio of 2 values is meaningful

→ Zero value indicates that nothing exists for the variable at zero point

e.g.: cost, distance, height

• Interval / Ratio → Quantitative

• Nominal / Ordinal → Qualitative

## • Study Design:

- ① Time Series: Data collected over several periods of time  
 ↳ trends + forecast + multivariate time series (more than 2 interested var)  
 + univariate " " (only 1 interested var)

time	Profit
2000	\$ 175k
2001	:
:	:
2010	:

- ② Cross-sectional Study / Design: Data collected at the same/ approximately the same point in time

	Temp	Humidity
city 1	:	:
:	:	:
city 2	:	:
:	:	:
city 3	:	:

- Sampling error: processing sample instead of population
- Non-sampling error: Mistakes that happen during data collection, recording, analysis
- Underestimation: prediction/calculation too low 45 min required for a task, but prediction of 30 min for completion
- Overestimation: " too high 45 min, but 70 min prediction

## LEC 3 Summarization of Qualitative variable

### ① Relative and % relative Frequency distribution

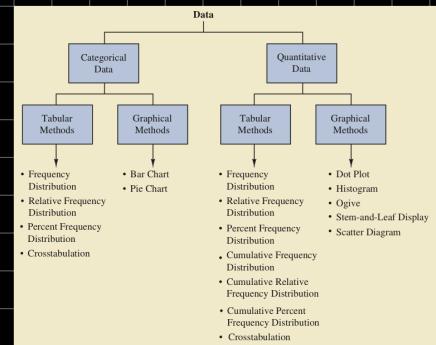
None	f	R.F	PRF(%)	C.F	(%)
	19	:	:	19	
	8	:	:	19+8	27
	5	:	:	19+8+5	32
total ( )					

- Minimum
- Maximum
- Most Preferred
- Outliers
- Minimum Preferred

• R.F = Relative frequency =  $\frac{\text{frequency of the class}}{n}$  ns total frequency

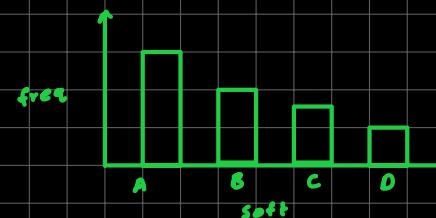
• PRF = % relative frequency = " \* 100

• C.F = Cumulative frequency / percentile = prev + next; order does not matter



### ② Visualization

- Bar chart / pie chart → Qualitative
- n-axis: categories → trend with time
- y-axis: frequency
- note of scaling



time: n-axis  
 var of interest: y-axis

- Pie chart / Sector diagram / Pizza Sliced diagram

- mean = measure for central location / avg
- median = " " " / middle value when sorted in ascending order  $\frac{n+1}{2}$
- mode = most occurring value
- Range = Largest - Smallest value / measure of variability dependance on extreme values
- Interquartile Range (IQR) = Q3 - Q1 / " " " not dependant on extreme values

$$\text{Variance} = s^2 = \frac{\sum (x_i - \bar{x})^2}{n-1}, \quad \frac{2f(x_i - \bar{x})^2}{n-1}, \quad \frac{1}{n-1} \left( \sum x^2 - \frac{(\sum x)^2}{n} \right) \quad \cdot 6sf \quad s^2 = \frac{\sum x_i^2 - n\bar{x}^2}{n-1}$$

$$\text{Standard Deviation} = \sqrt{\frac{\sum (x_i - \bar{x})^2}{n-1}}, \quad \sqrt{\frac{2f(x_i - \bar{x})^2}{n-1}}, \quad \sqrt{\frac{1}{n-1} \left( \sum x^2 - \frac{(\sum x)^2}{n} \right)}$$

## Summarization of Quantitative Variables

Q. The following scores represent the final exam scores for an elementary statistics course:

23, 60, 79, 32, 57, 74, 52, 78, 82, 36, 80, 73, 81, 95, 41, 65, 92, 85, 55, 76, 52, 10, 64, 75, 78, 64, 84, 98, 81, 67, 41, 71, 83, 54, 74, 82, 88, 62, 79, 43, 60, 78, 89, 17, 48, 84, 90, 15, 79, 34, 67, 80, 69, 74, 63, 82, 85, 61, 25, 72

① Create classes + intervals 5 sy km aur 20 say zindah nhi if generalizing

$$\cdot k = \lceil 1 + 3.3 \log n \rceil \quad \text{e.g. } k = 6.867$$

$$\approx 7 \rightarrow \text{N.o. of classes}$$

$$\cdot \left\lceil \frac{R}{k} \right\rceil = \text{interval} \quad ; \quad R = \text{range}$$

· Range: span of data

→ class boundary creation (S1)

→ construct Group frequency

→ calculate r.f., prf, CF

→ Draw r.f. histogram, and comment on distribution of data.

→ calc. avg score and std.

## LEC 4

- Histogram → Overall commenting on data
- Class limits have units of original data
- Class boundaries → Horizontal / vertical diff, find common point

lower → sub

upper → add

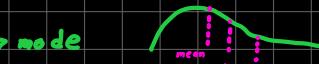
· Create CB

· Behaviours:

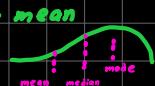
① Symmetric / Normal: Pattern like a bell shaped curve.  
Same freq for equidistant points

→ Skewness:

+vely skewed: mean > median > mode



-vely skewed: mode > median > mean



$$② S_k = \frac{\text{mean} - \text{mode}}{\text{SD}}$$

$$= \frac{3(\text{mean} - \text{median})}{\text{SD}}$$

- $S_k < 0$  : -vely skewed      moderate/ high according to value
- $S_k > 0$  : +vely skewed
- $S_k = 0$  : Symmetric

Range: -3 — 3

③ Histogram

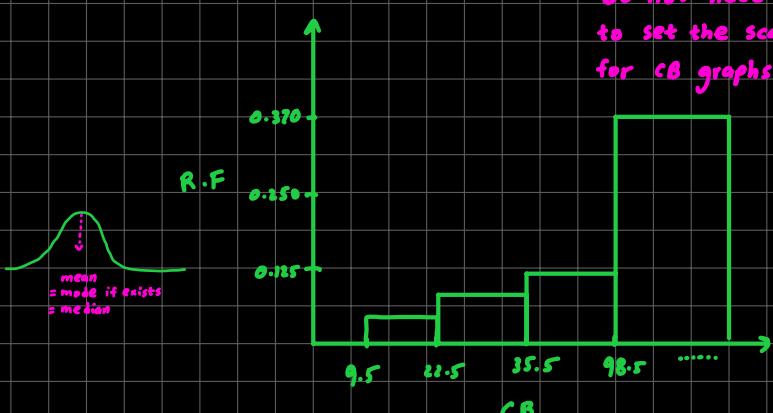
$$\cdot \bar{x} = \frac{\sum x}{n}$$

$$= \frac{\sum f x}{\sum f} \quad (\text{group mean})$$

·  $\bar{x} = 65.83 \rightarrow \text{class limits matching}$



Do not need  
to set the scale  
for CB graphs



## LEC 5

### Countable

- Discrete variables  $\rightarrow$  classes not recommended  $\rightarrow$  act as a qualitative variable.
  - $\rightarrow$  range and intervals too small
  - $\rightarrow$  make ungrouped freq.
  - $\rightarrow$  class intervals (grouped freq)

• Walpole Ex 1.17 , pg # 31

a. calc. 10% trimmed mean for smokers

b. Make dot-plot of datasets A and B on the same line

c. calculate coefficient of variation (cv) and comment on consistency of the datasets

• Trimmed mean:

- ① Sort
- ② "Round" the num of values
- ③ Take mean after discarding values

• std = dispersion from mean

• trimmed mean for removal of outliers

• Dot-plot:

- plot each point on a number line
- make key
- Dispersion, interval density
- Plot mean



• if overall, then



else,



$$\begin{aligned} \cdot CV &= \frac{\sigma}{\mu} \times 100 \quad ] \text{ pop} \\ &= \frac{s}{\bar{x}} \times 100 \quad ] \text{ sample} \end{aligned}$$

• Shows how large the standard deviation is from the mean

• std has same unit of measurement

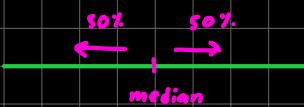
- Relative comparison + unitless
- std = abs measure of dispersion
- CV = relative measure //
- ideal when comparing variability of variables having different standard deviation and different mean

- CV ↑ less consistent
- CV ↓ more consistent

Heaps  $\rightarrow$  median

## LEC 6

• median - cuts data into 2 halves



### Percentile Points:

- ① Sort the data
- ② Calculate Index
- ③  $i = \frac{P}{100} \times n ; 1-99$

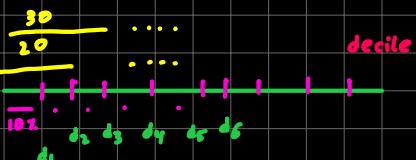
③ If  $i$  is an integer, then  $p$ th percentile will be the average of values placed at  $i^{\text{th}}$  and  $(i+1)^{\text{th}}$  position



$\rightarrow$  If  $i$  is non integer, round up, and pick the value at that position



$\rightarrow$  groups data in segments e.g:



Sorted data

$Q_2 = \text{median}$

$d_5 = Q_3 = \text{median}$

Q- Considering data of monthly starting salaries of 12 business graduates

3450  
3550  
3650  
3480  
3355  
3310  
3490  
3730  
3540  
3920  
3520  
3480

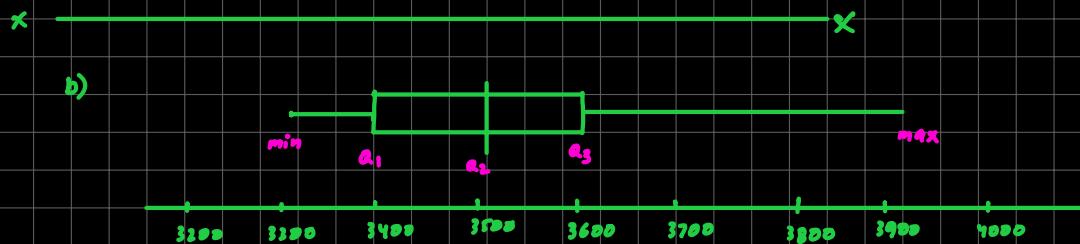
- a) Calc five-point summary of the data and IQR  
 b) Make a box-plot and detect if there exists any outlier, also comment on the shape of the distribution  
 c) Calc z-scores and detect any outlier if exists  
 measure of relative location of a value from mean

$$\text{IQR} = \frac{Q_3 - Q_1}{2}$$

- a) Minimum point  
 Maximum point  
 Quartiles points

IQR  $\rightarrow$  middle span

① Apply procedure for finding percentiles



$\rightarrow$  take diff to confirm the larger box.  $Q_3 - Q_2$ ,  $Q_2 - Q_1$ ,  
 right box large  $\rightarrow$  + very skewed  
 left box large  $\rightarrow$  - very skewed

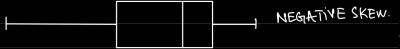


$$\text{Lower fence} = Q_1 - 1.5(\text{IQR})$$

SYMMETRICAL

$$\text{Upper fence} = Q_3 + 1.5(\text{IQR})$$

POSITIVE SKEW



NEGATIVE SKEW

c)  $z = \frac{x - \bar{x}}{s}$

;  $-3 \dots 3$

if any point exists outside the range, then the  $x$  value corresponding to it is the outlier.

$\cdot z = -1.5$ ,  $x_i$  is 1.5 standard deviations below the mean

$\cdot z = 1.5$ , // above the mean

## LEC 7

- anomalies = outliers
- Weights = frequency
- cross tabulation / Joint frequency table

\* combining 2 individual tabs

2 Qualitative

	Fuel $\rightarrow$	Premium	Regular	Diesel	total
compact					$\rightarrow$
midsize		frequency of each			$\rightarrow$
large			$\downarrow$	$\downarrow$	$\rightarrow$
total					

column wise total

marginal frequency/partial

Grand total

## 2 Quantitative

→ median if deviated data

cylinder	city MPG			
	$\leq 15$	16 - 20	21 - 25	total
1 - 5	0	0	4	4
6 - 10	2	3	0	5
11 - 15	1	0	0	1
total	3	3	4	10

## 1 Qualitative, 1 Quantitative

size	cylinder			
	1 - 5	6 - 10	11 - 15	total
compact	2	2	0	4
midsize	2	2	0	4
large	0	1	1	2
total	4	5	1	10