

## Chapter - 5

### Univariate Statistical Analysis

- 5.1 → Data Mining Tasks in Discovering knowledge in Data
- It involves six tasks (already covered in Chapter - 1)
- Description (understanding existing data)
  - Estimation (from sample, compute values like their true value)
  - Prediction (predict the future value to estimate its true value)
  - Classification (clustering the data objects into predefined class)
  - Clustering (grouping)
  - Association ("go together")

- 5.2 → Statistical Approaches to Estimation and Prediction

#### Prediction

- Definition → predicts future or unknown values based on pattern in existing data.
- Time oriented → future-oriented
- Type of target → can be numerical or categorical

#### Estimation

- Assigns a current unknown numerical value to a target variable.

→ present-oriented.

- Always numerical (continuous)

Technique used → Regression, DT, NN classification model

- Regression, KNN, Bayesian model.

Opp Nature → Future number or class label

- Produces numeric estimate.

### S.3 Statistical Inference

- It is a process of drawing conclusions about a population based on sample data;
- In data mining, it's used to understand patterns, relationship and predictions using statistical techniques.
- Statistical inference consists of methods like estimation and tests of hypothesis about population characteristics based on the information contained in the sample.
- Population is the collection of all elements (persons, item or data) of interest on a particular study.

→ Sample Statistics to estimate unknown population  
sample statistics Estimation population parameters (known)

Mean → $\bar{x}$	→	$\mu$
SD → s	→	$\sigma$
Proportion → p	→	$\pi$

- 5.4 How confident are we in our estimates
- estimation
- we can't check the churn behaviour of all future cell phone customers.
  - So we take a sample, compute a value & use it to estimate the true value.

Ex: If your churn sample, of 600 out of 3333 customers churned  
 $\text{estimated churn rate} = \frac{600}{3333} = 0.18 = 18\%$ .

Point Estimation (provides a single value not a range)

It is a single population that number calculated from the sample that is used to guess the Population Parameter.

- But)
- To find confidence calculate accuracy of the estimate.
  - Accuracy of the estimate can be calculate through sampling error?
  - why we need sampling error?  
 because a sample never perfectly represents the population.  
 (subset of the population)

$$\text{so sampling error} = |\text{sample estimate} - \frac{\text{Population}}{\text{Parameter}}|$$

i.e. the population parameter  
 is usually unknown. so we can't compute sampling error exactly. that's why

We estimate SE using standard error (SE)

→ Ex If the population parameter is known (rare case) then sampling error = [sample estimate (SE) - population parameter]

$$\text{Ex} \quad \text{Population mean} = 50$$

$$\text{Sample mean} = 47 \quad SE = |47 - 50| = |-3|$$

→ Ex If the population parameter is unknown (common case) then we estimate the sampling error using standard error (SE)

i.e standard error = estimated sampling error.

$$\Rightarrow SE = \frac{s}{\sqrt{n}} \quad \text{where } s = \text{Sample Standard deviation}$$

$$\text{Ex} \quad SD = 6, \text{ Sample size } n = 100$$

$$SE = \frac{6}{\sqrt{100}} = \frac{6}{10} = 0.6.$$

∴ Sampling error is  $\pm 0.6$ .

## 5.5 Confidence interval estimation of the mean

### confidence interval (CI)

It gives a range of values within which the true population mean is likely to lie.

It is based on

- Sample mean
- sample standard deviation
- sample size
- level of confidence (90%, 95%, 99%)

### confidence interval estimate

It consists:

(I) Point estimate → A single number from the sample

ex. Sample mean ( $\bar{x}$ )

(II) Margin of error = how far the estimate might be from the true value. (It measures the precision of the interval estimate)  
ex. Smaller margin of error = more precision

(III) confidence level (90%, 95%, 99%)

↳ this tells how confident we are that the interval contains the true parameter.

- or setting Standard error
- $\hat{x}$  if the population
- general form of a confidence interval
- $CI = \text{point estimate} \pm \text{margin error}$
- It creates two bounds.
- lower limit = point estimate - ME
- upper limit = " + ME

- the  $t$ -interval for estimating a population mean when:
- (i) the SD is unknown (usually true)
  - (ii) & either the population is normal OR sample size is large ( $n \geq 30$ )
- then use  $t$ -distribution.

so, using  $t$ -interval

$$CI = \bar{x} \pm t_{\frac{\alpha}{2}, n-1} \times \frac{s}{\sqrt{n}}$$

where

$\alpha$  = significance level  
(It represents the probability of making a mistake by rejecting a true statement about the population)

confidence level	$\alpha$
90%	0.10
95%	0.05
99%	0.01

$\bar{x}$  = sample mean (point estimate)

$t_{\frac{\alpha}{2}, n-1}$  = ( $t$  critical value at chosen confidence level)

$s$  = sample standard deviation

$n$  = sample size

$\frac{s}{\sqrt{n}}$  = standard error

$t \times SE$  = margin of error

\* why we use  $t$ - instead of  $Z$ .  
 →  $Z$  is used when population SD is known (rare)  
 →  $t$  " " " SD must be estimated from  
 the sample (common case)

-: sample size increases,  $\rightarrow t$ -distribution  
 or normal distribution, no difference

95% confidence interval means if we take  
 many samples & build intervals, 95% of those  
 intervals will contain the true population mean.

Q. From a customer service calls statistics as given  
 sample mean  $\bar{x} = 1.563$ , sample SD  $= s = 1.315$   
 sample size  $n = 3333$  & confidence level = 95%. use  
 $t$  interval ( $t_{\alpha/2}, n-1 \approx 1.96$ )

s-1 complete standard error,  $SE = \frac{s}{\sqrt{n}} = \frac{1.315}{\sqrt{3333}} \approx 0.02278$

s-2 compute margin of error  $ME = t \times SE$   
 $= 1.96 \times 0.02278 \approx 0.045$

s-3 form the confidence interval  
 Lower bound =  $\bar{x} - ME = 1.563 - 0.045 = 1.518$   
 upper bound =  $\bar{x} + ME = 1.563 + 0.045 = 1.608$

Interpretation

$$Q \quad \bar{x} = 1.607, n = 28 \text{ for } 95\% \text{ t control}$$

$df = n - 1 = 27$

$$\text{interval estimate} \quad \bar{x} \pm t_{12} \left( \frac{s}{\sqrt{n}} \right) = 1.607 \pm 2.052 \left( \frac{1.82}{\sqrt{27}} \right) = 1.607 \pm 0.095$$

we are 95% confident that the true mean number of customer service calls for the population lies between 1.518 & 1.608

Q The contents of seven similar containers of sulfuric acid are 9.8, 10.2, 10.4, 9.8, 10.0, 10.2 & 9.6 g.

Find 95% CI for the mean contents of all such containers, assuming an approximately normal distribution,  $\bar{x} = 10.0$  &  $s = 0.283$  &  $t_{0.925} = 2.447$

$$\text{Hence } CI = 10.0 - (2.447) \left( \frac{0.283}{\sqrt{7}} \right) < \mu < 10.0 + 2.447 \left( \frac{0.283}{\sqrt{7}} \right)$$

Q The average zinc concentration recovered from a sample of measurements taken in 36 different locations in a river is found to be 2.6 grams per millimeter. Find the 95% & 99% confidence intervals for the mean zinc concentration in the river. Assume that the population standard deviation is 0.3 gram per millimeter.

Sol

Sample mean,  $\bar{x} = 2.6 \text{ g/ml}$

Population standard deviation  $\sigma = 0.3 \text{ g/ml}$  (known)

Sample size  $n = 36$

so we use  $Z$ -based interval because  $\sigma$  is known.

S-1 compute Standard Error (SE) =  $\frac{\sigma}{\sqrt{n}} = \frac{0.3}{\sqrt{36}} = 0.05$

S-2 Z-critical values

→ For 95% confidence = 1.96

99% " = 2.576

### -3 Margin of error (ME)

$$95\% \rightarrow ME_{95} = Z_{0.025} \times SE = 1.96 \times 0.05 = 0.098$$

$$99\% \rightarrow ME_{99} = Z_{0.005} \times SE = 2.576 \times 0.05 = 0.1288$$

### -4 confidence intervals

#### 95% CI

$$\text{lower} = \bar{x} - ME = 2.6 - 0.098 = 2.502$$

$$\text{upper} = \bar{x} + ME = 2.6 + 0.098 = 2.698$$

for 95% CI the true mean zinc conc' lies bet' 2.502

#### 99% CI

$$\text{lower} = 2.6 - 0.1288 = 2.4712$$

$$\text{upper} = 2.6 + 0.1288 = 2.7288$$

for 99% CI the true mean zinc concentration lies  
between 2.4712 & 2.7288 g/ml



<u>Situation</u>
$\sigma$ known
$\sigma$ unknown, $n < 30$
$\sigma$ unknown, $n \geq 30$
Population normal + $\sigma$ unknown

#### use

Z-interval

t-interval

t-interval (preferred)

but z-interval can be used

t-interval

but

$CI = \bar{x} \pm Z \frac{\sigma}{\sqrt{n}}$

## S.6 How to reduce the margin of error

the margin of error (ME) for a 95% confidence interval for the population mean  $\mu$  is

$$ME = \sqrt{t_{\alpha/2}^2 \left(\frac{s}{\sqrt{n}}\right)}$$

$$\approx 0.045$$

so now smaller is the ME, the more precise our estimation.

so how to reduce

ME contains 3 quantities.

(1)  $t_{\alpha/2}$ , which depends upon confidence level & sample size.

(2) Sample standard deviation ( $s$ ), which is the characteristic of the data.

may not be changed as follows  
if the sample size, we can reduce then

(i) If confidence level is reduced then  
 $t_{\alpha/2}$  then reduce ME, (Not recommended)

(ii) By increasing the sample size, we can reduce ME while maintaining a constant level of confidence.

Ex Sample 5000 customers,

with same  $S = 1.315$ . the ME.

$$ME = t_{\frac{\alpha}{2}} \left( \frac{S}{\sqrt{n}} \right) = 1.96 \left( \frac{1.315}{\sqrt{5000}} \right)$$
$$= 0.036$$

increase in

with sample size reduce the ME.