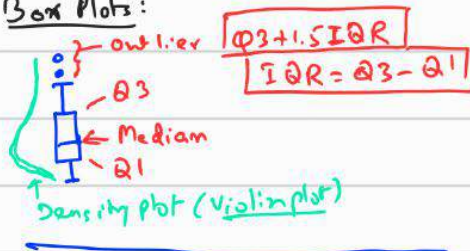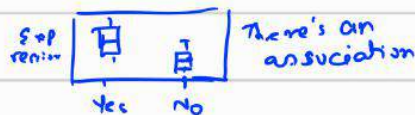- Bar Chart – Categorical
- Histogram – Numerical / Continuous
- Density Plots – smoothed histogram (kernel)

Box Plots:

- outlier $Q3 + 1.5 IQR$
- $IQR = Q3 - Q1$
- $Q3$
- Median
- $Q1$

Density Plot (Violin plot)

2 Variables:
① 1 cat, 1 numeric side by side box plt


S+P region    There's an association
Yes    No

② 2 cat
a) side by side bar chart

%    y    no
b) Stacked %

c) Mosaic          d) scatter:

% Yes    % No

f) istribution:
① Normal Symmetric

Uniform → Symmetric

③ Right skewed (Positively skewed)


② 


④ Left skewed (-ve)


⑤ Exponential
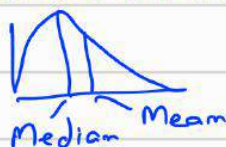

① Trimmed mean: Mean after removing top/bottom $\alpha$% of data.

② Median – Not sensitive to outliers.
   – Non parametric meanure


Mean = Median          Median   Mean

③ Standard Dev. : $\sqrt{\dfrac{\sum_{i=1}^{n}(x_i - \bar{x})^2}{n-1}}$   ④ Average Abs. Dev. $\sqrt{\dfrac{\sum_{i=1}^{n}|x_i - \bar{x}|}{n-1}}$

why $(n-1)$
If we have n obs. we loose one d.f for estimation.
For ep. $n = 1$ & $x_1 = 86$, s.D = ?     no s.D as we have only one obs.
let's say $n = 4$ and given a statistic $\bar{x} = 80$, we can choose $X_1, X_2, X_3$ freely but $X_4$
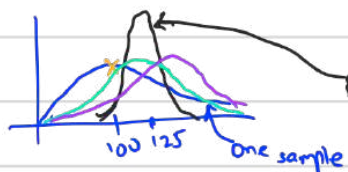has to be some number to make $\bar{x} = 80$. so d.f = 3

Normal Distⁿ : 68/95/99.7

$$Z = \frac{x - \mu}{\sigma}$$

Inferential Statistics –

Inference of Population from sample.

Sampling distribution – Theoretical set of ALL POSSIBLE statistics (eg. $\bar{X}$) we could get. Eg. Taking sample size of 25 over and over again. Assume pop. mean $\mu = 125$, one sample mean = 100
$\sigma = 20$ , $n = 25$

Expected value of $\bar{X} = \mu = 125$

S.D of $\bar{X}$ = Standard Error $(SE \bar{X}) = \frac{\sigma}{\sqrt{n}} = \frac{20}{\sqrt{25}} = 4$

↑ on average how far estimated $\bar{X}$ deviate from $\mu$
and as it is $\frac{\sigma}{\sqrt{n}}$ ie, as sample size increases The
S.E is smaller.


100 125   one sample

―――――― X ――――――

S.D of $\bar{X}$ : $\sigma_{\bar{x}} = \frac{\sigma}{\sqrt{n}}$

S.E of $\bar{X}$ : $SE_{\bar{x}} = \frac{s}{\sqrt{n}}$

s is sample S.D

Some Properties:
① $Var(aX) = a^2 Var(X)$   a is constant

② $Var(x_1 + x_2) = Var(x_1) + Var(x_2)$
if $x_1, x_2$ independent

$x_i$ has mean $\mu$, var $= \sigma^2$
$\bar{X} = \frac{x_1 + x_2 \cdots + x_n}{n}$

Now,
$$Var(\bar{X}) = Var\left(\frac{x_1 + x_2 + \cdots x_n}{n}\right)$$

$$= \frac{1}{n^2} Var(x_1 + x_2 + \cdots x_n)$$

$$= \frac{1}{n^2}\left[Var(x_1) + Var(x_1) + \cdots Var(x_n)\right]$$

$$= \frac{1}{n^2}\left[\sigma^2 + \sigma^2 + \cdots + \sigma^2\right]$$

$$= \frac{n\sigma^2}{n^2} = \frac{\sigma^2}{n} \Rightarrow S.D_{\bar{x}} = \frac{\sigma}{\sqrt{n}}$$

Confidence Interval
――――――― X ―――――――

Let's look at a skewed distⁿ $\mu = 125, \sigma = 20$ ; Sample $\bar{X} \approx$ Normal $\mu_{\bar{x}} = 125, \sigma_{\bar{x}} = \frac{20}{\sqrt{25}} = 4$
$n = 25$

→ 95% of Time The $\bar{X}$ will be within
$\mu \pm 2 SD_{\bar{x}}$
or $\mu \pm 2\sigma/\sqrt{n}$
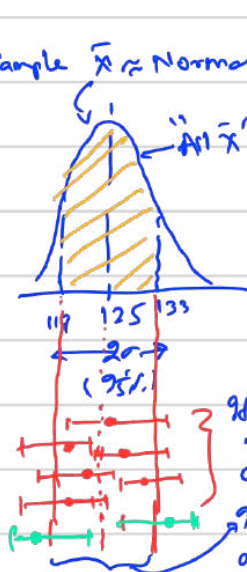
→ conversely 95% time The $\mu$ will be within
$\bar{X} \pm 2 \cdot SD_{\bar{x}}$ or $\boxed{\bar{X} \pm 2 \frac{\sigma}{\sqrt{n}}}$

In reality we don't know $\sigma$, replace with s
$\boxed{\bar{X} \pm t \cdot \frac{s}{\sqrt{n}}}$ where $\frac{s}{\sqrt{n}}$ is $S.E_{\bar{x}}$
"t" instead of 2


"All $\bar{x}$"
117  125  133
←20→
(95%)

If we drew a sample
n = 25, it could be
any of these points
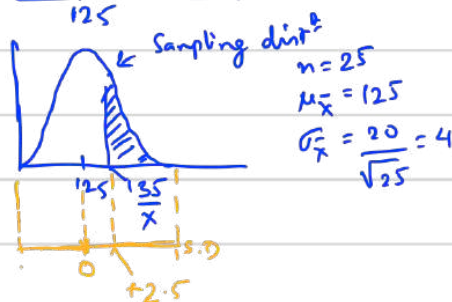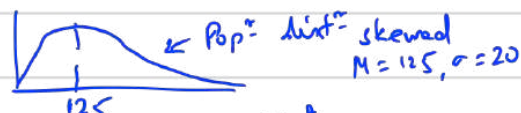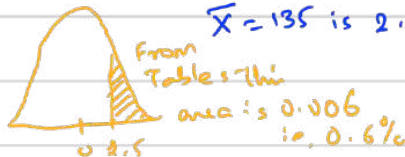95% times in this region
and true mean lies with

# Hypothesis Testing:

What's the prob. of $\bar{X} \geq 135$?

← Pop$^\underline{n}$ dist$^\underline{n}$ skewed $M = 125, \sigma = 20$

125

let's standardize

$$Z = \frac{\bar{X} - M}{\sigma/\sqrt{n}} = \frac{135 - 125}{4} = 2.5$$

← Sampling dist$^\underline{n}$
$n = 25$
$\mu_{\bar{x}} = 125$
$\sigma_{\bar{x}} = \frac{20}{\sqrt{25}} = 4$

This mean

$\bar{X} = 135$ is 2.5 S.D above the mean

From Tables This area is 0.006 i.e. 0.6%

0 2.5

125 135
$\bar{X}$

S.D

0    +2.5

→ Pop$^\underline{n}$ Mean 125, Prob. of estimating $\bar{X} \geq 135$ is 0.6%.

Eg. For smokers we believe that $\mu > 125$; if 125 is for healthy pop$^\underline{n}$.
We take a sample (n=25) of smokers – their mean $\bar{X} = 135$ & $S = 20$.

$H_0$: $\mu_{smokers} = 125$

$H_1$: $\mu_{smokers} > 125$

For this we use t-dist$^\underline{n}$ (same as Z dist$^\underline{n}$)

$$t = \frac{\bar{X} - M_0}{S/\sqrt{n}}$$

(can be thought of as Normal dist$^\underline{n}$ for samples)

now in this scenario 0.6% is p-value

⟹ Assuming "Pop$^\underline{n}$ mean is 125"

## t-dist$^\underline{n}$:

used instead of z-dist$^\underline{n}$ - b/c don't know pop$^\underline{n}$ std. dev. and must use sample std. dev. to estimate Std. Error.
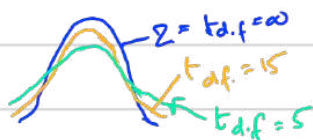
Sample std. dev. is not a true estimate of $\sqrt{\sigma}$

so as $n \to \infty$
$\sqrt{S} \to \sqrt{\sigma}$

Conf. interval
$$\bar{X} \pm t * \frac{S}{\sqrt{n}}$$

Hyp. testing
$$t = \frac{\bar{X} - M_0}{S/\sqrt{n}}$$

$Z = t_{d.f = \infty}$
$t_{d.f. = 15}$
$t_{d.f = 5}$

## confidence Interval for Mean

Eg. Mean BMI of pop$^\underline{n}$ $M = 27.2$

take sample $n = 16$, $\bar{X} = 25.2$, $s = 5$

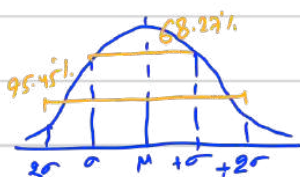$$\bar{X} \pm t_{n-1} * \frac{S}{\sqrt{n}}$$

95% C.I
2-sided

$= 25.2 \pm 2 * \frac{5}{\sqrt{16}} = 25.2 \pm 2.5 \quad (22.7, 27.7)$

$t_{15} = 2$ (just for example)

if we take repeated samples of n=16 for every 100 c.I, about 95% will have the true mean within the intervals.

one-sided
  95%,
we interested in upper bound

$$\bar{X} + (1.64)\frac{s}{\sqrt{16}} = 25.2 + 1.64\left(\frac{s}{\sqrt{16}}\right) = 27.25$$



68.27%
95.45%

$2\sigma$  $\sigma$  $M$  $+\sigma$  $+2\sigma$

this qtty refers to value
That says 95% of observation lie below 1.64 s.d above the
mean or less

## Controlling the margin of error :

$$\bar{X} \pm t\left(\frac{s}{\sqrt{n}}\right)$$
$$\underbrace{\phantom{t\left(\frac{s}{\sqrt{n}}\right)}}_{\text{margin of error}}$$

Reducing ME : $\downarrow t$ – reducing $t$ – lower CI – $\chi$
             $\downarrow s$ – can't be reduced
             $\uparrow n$
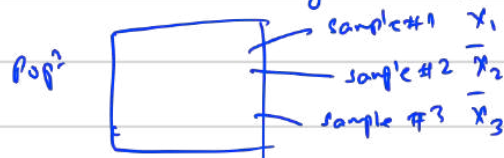
Suppose we want M.E 0.5 ie. $\bar{X} \pm 0.5$
$$n = \left(\frac{t * s}{ME}\right)^2 = \frac{2 \times 5}{0.5} = 400 \qquad | \text{value from previous page}$$

If you are planning expt. – to get s $\begin{cases} \text{literature} \\ \text{small pilot study} \end{cases}$

## Bootstraping & Resampling :

Parametric / large "n"



Pop²

sample #1  $\bar{X}_1$
sample #2  $\bar{X}_2$
sample #3  $\bar{X}_3$

sampling dist² $\approx$ Normal
s.D of $\bar{X} = \sigma_{\bar{X}} = \frac{s}{\sqrt{n}}$

• Bootstrap : Instead of large sample.
             & when SE estimate is difficult.

sample → $\bar{X}$

Resample with replacement $\bar{X}_2^*$
   "        "        :    $\bar{X}_3^*$
                         $\vdots$
                         $\bar{X}_B^*$
$\left.\right\}$ Bootstrap
Sampling dist²

Eg.  60, 75, 80, 85, 90, with $\bar{X} = 78$, $s = 11.51$, $S.E = \frac{11.51}{\sqrt{5}} = 5.15$

Random sample #1 : Pick a sample 75, then another w. replacement 90, ...... we get
                   75, 90, 80, 90, 85  $\bar{X}_1^* = 84$
         #2        85, 60, 75, 85, 60  $\bar{X}_2^* = 73$
         $\vdots$

$SE_{\bar{X}}^* = 5.57$ (quite close to 5.15)
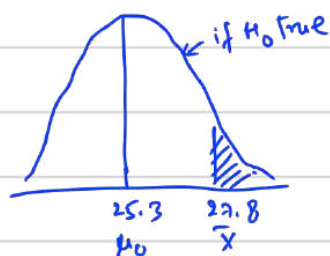
## One-Sample t-test:

known    Mean BMI USA 2008 = 25.3 ... has it increased

2018 $n=25$, $\bar{x} = 27.8$, $s = 6$

$H_O$   $\mu_{208} = 25.3$  — why $H_O$ is stated this way? because we know it is true. if we make $H_A$ as null we don't what is true, what do we expect
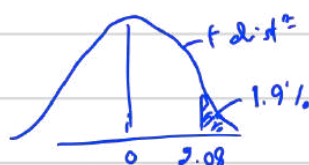
$H_A$   $\mu_{2018} > 25.3$



if $H_O$ true

$25.3$   $27.8$
$\mu_O$   $\bar{x}$

$t_{statistic} = \dfrac{\bar{x} - \mu_O}{s/\sqrt{n}} = \dfrac{27.8 - 25.3}{6/\sqrt{25}} = 2.08$

27.8 is 2.08 stand. errors above what we'd expect if $H_O$ true.



t dist$^n$

1.9%

0   2.08

pvalue 0.019   ie, if $H_O$ true, prob. of $\bar{x} \geq 27.8$ is $\simeq 1.9\%$

so $H_O$ not true

95% CI   $27.8 \pm 2\left(\dfrac{6}{\sqrt{25}}\right) \rightarrow (25.4, 30.2)$

## One Vs. Two sided t-test:   $t_{STA} = 2.08$

| One sided:   $H_A$ $\mu_{2018} > 25.3$ | Two sided:   $H_A$ $\mu_{2018} \neq 25.3$ |
|---|---|
| here we're looking at prob. of getting estimate 2.08 or more above $25.3$ ($\mu_O$) | Prob. of getting estimate that is 2.08 or more std. error away from $\mu_O$ |



1.9%

0

this will be 3.8%.



-2.08   0   +2.08

what if One-sided p-value is 3%, then two would be 6% - so reject $H_O$
well depends on the problem.

## Hypothesis test & CI:

$H_O: \mu = \mu_O$   $H_A: \mu \neq \mu_O$

$t_{STA} = \dfrac{\bar{x} - \mu_O}{s/\sqrt{n}}$   CI: $\bar{x} \pm t * \dfrac{s}{\sqrt{n}}$

$\hookrightarrow$ given pvalue $-\alpha$



$\mu_O$   $\bar{x}$

Reference pt for p-value

Reference pt for CI

CI