# Linear Regression-2

Prof. Asim Tewari
IIT Bombay

# Linear Regression

$$Y = \beta_0 + \beta_1 X + \epsilon$$

0 mean
$\sigma^2$ var

Data
X , Y

| $x_1$ | $y_1$ |
| $x_2$ | $y_2$ |
| ⋮ | ⋮ |
| $x_n$ | $y_n$ |

Sample

Limited
Data →
$(y_i, x_i)$

$\hat{\beta}_0$ , $\hat{\beta}_1$

$\hat{\beta}_0$ , $\hat{\beta}_1$

# Characteristic Function

The characteristic function of a random variable $X$ is

$$\phi_X(t) \equiv E(e^{itx}) = \int_{-\infty}^{\infty} e^{itx} \underbrace{f_X(x)}_{\text{Pdf of } X} \, dx$$

$$e^{itx} = \frac{1}{\angle 0} + \frac{itx}{\angle 1} + \frac{(itx)^2}{\angle 2} + \cdots$$

$$\therefore \phi_X(t) = E\left[ \frac{1}{\angle 0} + \frac{itx}{\angle 1} + \frac{(itx)^2}{\angle 2} + \cdots \right]$$

$$= \frac{1}{\angle 0} + \frac{it \, E[x]}{\angle 1} + \frac{(it)^2 E[x^2]}{\angle 2} + \frac{(it)^3 E[x^3]}{\angle 3} + \cdots$$

# Characteristic Function

$$\therefore \ \phi_x(t) = \frac{1}{\angle 0} + \frac{it\,m_1}{\angle 1} + \frac{(it)^2 m_2}{\angle 2} + \cdots \cdots$$

where $m_n$ is the $n^{th}$ moment of the r.v. $X$

$$\text{ie } m_n \equiv E[X^n]$$

$$\phi_x(t)\Big|_{t=0} = 1 \quad ; \quad \frac{d}{dt}\phi_x(t)\Big|_{t=0} = im \; ; \quad \frac{d^n}{dt^n}\phi_x(t)\Big|_{t=0} = (i)^n m_n$$

# Moment generating Function

A moment generating function of a r.v. $X$ is

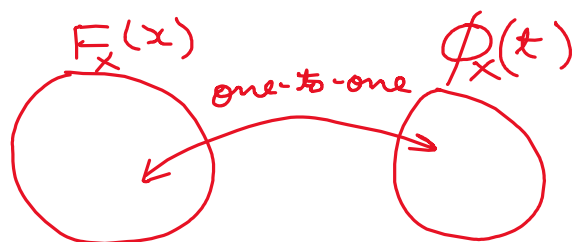$$M_x(t) = \phi_x(-it) = \int_{-\infty}^{\infty} e^{tx} f_x(x) \, dx$$

Now

$$\therefore \quad \frac{d^n}{dt^n} \phi_x(t) \bigg|_{t=0} = (i)^n m_n \implies \frac{d^n}{dt^n} M_x(t) \bigg|_{t=0} = m_n$$

$= n^{th}$ moment of $X$

# Characteristic Function

$$\therefore \phi_X(t) = \frac{1}{\angle 0} + \frac{it}{\angle 1} m_1 + \frac{(it)^2}{\angle 2} m_2 + \frac{(it)^3}{\angle 3} m_3 + \cdots$$

There is a one-to-one correspondance between the cumulative distribution function and the characteristic function.

$F_X(x)$ ——one-to-one——→ $\phi_X(t)$

If the r.v. has a probability density function $f_X(x)$ then

$$f_X(x) = F_X'(x) = \frac{1}{2\pi} \int e^{-itx} \phi_X(t) \, dt$$

# Characteristic Function

$$\therefore \phi_X(t) = \frac{1}{\angle 0} + \frac{it\, m_1}{\angle 1} + \frac{(it)^2}{\angle 2} m_2 + \frac{(it)^3}{\angle 3} m_3 + \cdots$$

If a r.v. $X$ has $\mu = 0$ and $\sigma^2 = 1$, i.e $X \sim (0, 1)$

$$\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\uparrow\quad\uparrow$$
$$\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\mu\quad\sigma^2$$

then $\phi_X(t) = 1 + 0 - \dfrac{t^2}{2} + O(t^2)$

For a normal distribution $N(\mu, \sigma^2)$

$$\phi_X(t) = e^{it\mu - \frac{1}{2}\sigma^2 t^2}$$

and for $N(0,1)$ ; $\boxed{\phi_X(t) = e^{-\frac{t^2}{2}}}$

For $N(\mu, \sigma^2)$

$$f_X = \frac{1}{\sigma\sqrt{2\pi}} e^{-\frac{1}{2}\left(\frac{x-\mu}{\sigma}\right)^2}$$

For $N(0,1)$

$$f_X = \frac{1}{\sqrt{2\pi}} e^{-x^2/2}$$

# Central limit theorem

$$X_i \overset{iid}{\sim} (\mu, \sigma^2)$$

Sample mean $\quad \overline{X}_n = \frac{1}{n} \sum_{i=1}^{n} X_i$

$\therefore$ Expected value of the sample mean is

$$E[\overline{X}_n] = E\left[\frac{1}{n} \sum_{i=1}^{n} X_i\right] = \frac{1}{n}\left(\sum_{i=1}^{n} E[X_i]\right)$$

$$= \frac{1}{n}\left(\sum_{i=1}^{n} \mu\right) = \frac{1}{n} n\mu = \mu$$

# Central limit theorem

Variance of the sample mean $Var(\bar{X}_n)$

$$Var(\bar{X}_n) = Var\left(\frac{1}{n}\sum_{i=1}^{n} X_i\right) = \frac{1}{n^2}\sum_{i=1}^{n} Var(X_i) = \frac{1}{n^2}\sum_{i=1}^{n}\sigma^2$$

$$\Rightarrow Var(\bar{X}_n) = \frac{n}{n^2}\sigma^2 = \frac{\sigma^2}{n}$$

Mean of sample mean $E(\bar{X}_n) = \mu$

and Variance of sample mean $Var(\bar{X}_n) = \frac{\sigma^2}{n}$

# Central limit theorem

Now we define $Z_n \equiv \dfrac{n\,\bar{X}_n - n\mu}{\sigma\sqrt{n}}$

$\Rightarrow Z_n = \dfrac{n\,\frac{1}{n}\sum\limits_{i=1}^{n} X_i - n\mu}{\sigma\sqrt{n}} = \dfrac{\sum\limits_{i=1}^{n} X_i - n\mu}{\sigma\sqrt{n}}$

$= \dfrac{\sum\limits_{i=1}^{n} (X_i - \mu)}{\sigma\sqrt{n}}$

If now we define

$Y_i = \dfrac{X_i - \mu}{\sigma}$ then

$Z_n = \sum\limits_{i=1}^{n} \dfrac{Y_i}{\sqrt{n}}$

# Central limit theorem

$$\therefore \quad Y_i = \frac{X_i - \mu}{\sigma} \quad ; \quad E(Y_i) = 0 \quad \text{and}$$

$$Var(Y_i) = Var\left(\frac{X_i}{\sigma}\right) = \frac{1}{\sigma^2} \sigma^2 = 1$$

$$\therefore \quad \phi_y(t) = 1 - \frac{t^2}{2} + O(t^2)$$

$$\boxed{Z_n = \sum_{i=1}^{n} \frac{Y_i}{\sqrt{n}}}$$

$$\therefore \quad \phi_{Z_n}(t) = E\left[e^{it\left(\frac{Y_1 + Y_2 + \cdots Y_n}{\sqrt{n}}\right)}\right]$$

$$= E\left[\prod_{k=1}^{n} e^{i\frac{t\, Y_k}{\sqrt{n}}}\right] = \prod_{k=1}^{n} \underbrace{E\left[e^{i\frac{t\, Y_k}{\sqrt{n}}}\right]}_{\phi_y(t/\sqrt{n})}$$

$$= \left[\phi_y(t/\sqrt{n})\right]^n$$

# Central limit theorem

$$X_i \overset{iid}{\sim} (\mu, \sigma^2); \quad Y_i \equiv \frac{X_i - \mu}{\sigma}$$

$$Z_n = \sum_{i=1}^{n} \frac{Y_i}{\sqrt{n}}$$

$$\therefore \quad \phi_{Z_n}(t) = \left[ \phi_Y\left(t/\sqrt{n}\right) \right]^n = \left[ 1 - \frac{t^2}{2n} + O\left(\frac{t^2}{n}\right) \right]^n$$

As we increase the sample size $n$, we get the limit

$$\lim_{n \to \infty} \phi_{Z_n}(t) = \lim_{n \to \infty} \left[ 1 - \frac{t^2}{2n} + O\left(\frac{t^2}{n}\right) \right]^n$$

$$= e^{-t^2/2}$$

$\}$ This is same as the characteristic function for $N(0,1)$

Hence, $\lim_{n \to \infty} Z_n = N(0,1)$

# Central limit theorem

$$X_i \overset{iid}{\sim} (\mu, \sigma^2) \quad ; \quad \bar{X}_n = \frac{1}{n} \sum_{i=1}^{n} X_i$$

$$Y_i = \frac{X_i - \mu}{\sigma}$$

$$Z_n = \sum_{i=1}^{n} \frac{Y_i}{\sqrt{n}}$$

$$\lim_{n \to \infty} Z_n \sim N(0,1) \implies \lim_{n \to \infty} \sqrt{n} \frac{(\bar{X}_n - \mu)}{\sigma} \sim N(0,1)$$

$$\implies \lim_{n \to \infty} \sqrt{n} (\bar{X}_n - \mu) \sim N(0, \sigma^2)$$

$$\implies \lim_{n \to \infty} (\bar{X}_n - \mu) \sim N\left(0, \frac{\sigma^2}{n}\right)$$

$$\implies \lim_{n \to \infty} \bar{X}_n \sim \mu + N\left(0, \frac{\sigma^2}{n}\right) \sim N\left(\mu, \frac{\sigma^2}{n}\right)$$

# Central limit theorem

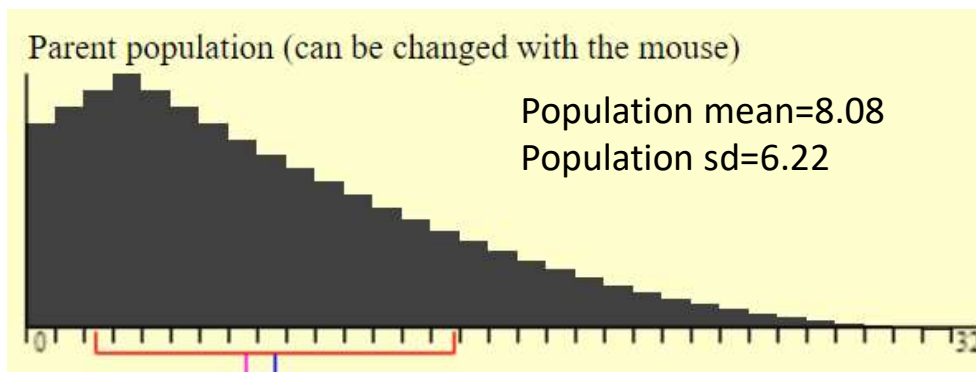For $X_i \overset{iid}{\sim} (\mu, \sigma^2)$ if we define sample mean $\overline{X}_n = \frac{1}{n} \sum_{i=1}^{n} X_i$

then mean of sample mean $(\overline{X}_n)$ is $E[\overline{X}_n] = \mu$

and variance of sample mean $(\overline{X}_n)$ is $Var(\overline{X}_n) = \frac{\sigma^2}{n}$

Now by Central limit theorem, we get that

$$\underset{n \to \infty}{Lim} \quad \overline{X}_n \sim N\left(\mu, \frac{\sigma^2}{n}\right)$$

# Central limit theorem



Parent population (can be changed with the mouse)

Population mean=8.08
Population sd=6.22

Sample size=2

Distribution of Means, N=2
Sample mean=8.14
Sample sd=4.48

Sample size = 10

Distribution of Means, N=10
Sample mean=8.09
Sample sd=1.24

Sample size =25

Distribution of Means, N=25
Sample mean=8.08
Sample sd=1.24

As n↑ the distribution of sample mean approaches a normal distribution.

Ref: http://onlinestatbook.com/stat_sim/sampling_dist/index.html

# Central limit theorem



Parent population (can be changed with the mouse)
Population mean=8.08
Population sd=6.22

Distribution of Means, N=2
Sample mean=8.14
Sample sd=4.48

Distribution of Means, N=10
Sample mean=8.09
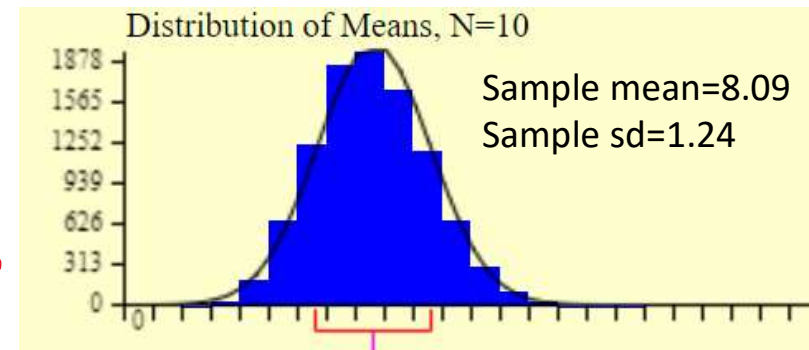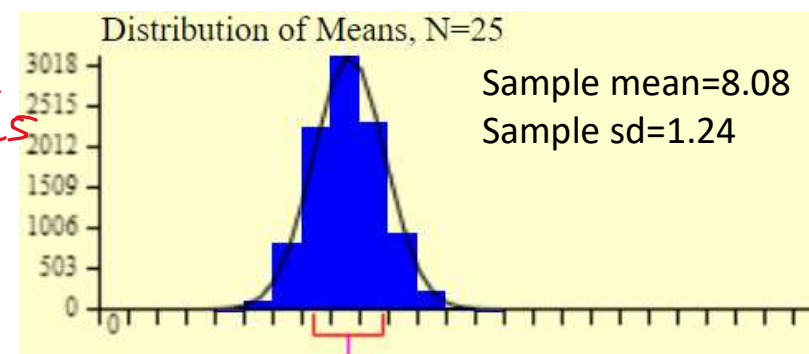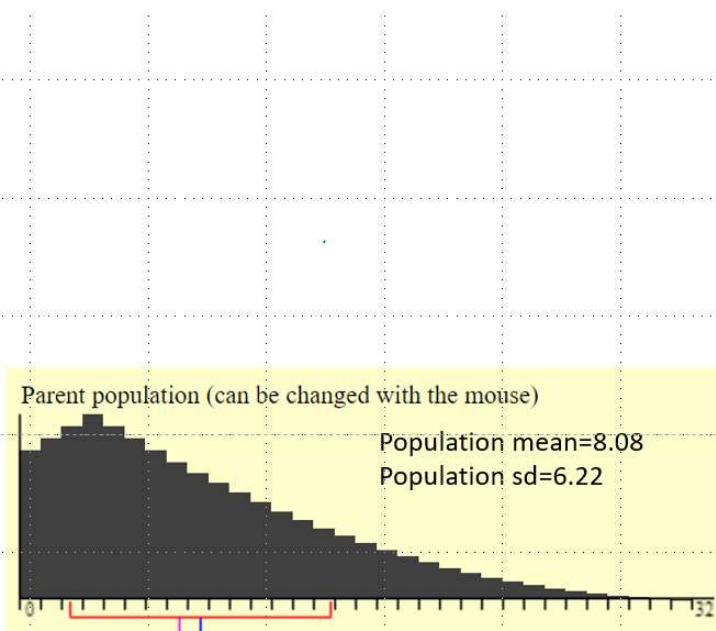Sample sd=1.97

Distribution of Means, N=25
Sample mean=8.08
Sample sd=1.24

Sample Size =2

Sample size = 10

Sample size =25

As $n\uparrow$ the distribution of sample mean approaches a normal distribution.

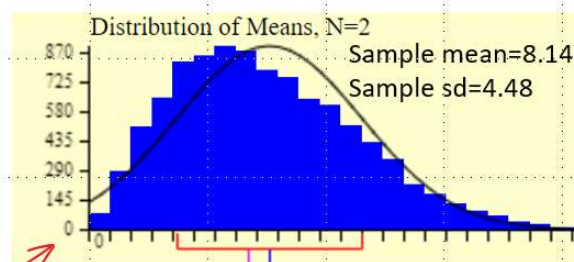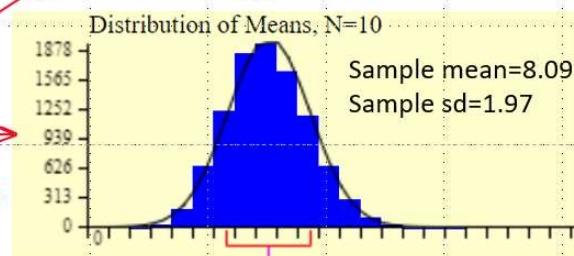$E(\bar{X}_2) = 8.14 \simeq \mu = 8.08$

$\sqrt{Var(\bar{X}_2)} = 4.48 \simeq \dfrac{\sigma}{\sqrt{n}} = 4.40$

$E(\bar{X}_{10}) = 8.09 \simeq \mu = 8.08$

$\sqrt{Var(\bar{X}_{10})} = 1.97 \simeq \dfrac{\sigma}{\sqrt{n}} = 1.97$

$E(\bar{X}_{25}) = 8.08 = \mu = 8.08$

$\sqrt{Var(\bar{X}_{25})} = 1.24 = \dfrac{\sigma}{\sqrt{n}} = 1.24$

# Application of Central limit theorem

Population $\xrightarrow{\quad n \quad}$ Sample

$$\overline{X}_n = \frac{1}{n} \sum_{i=1}^{n} X_i$$

$$X_i \overset{iid}{\sim} (\mu, \sigma^2)$$

Find population mean from a sample of size $n$.
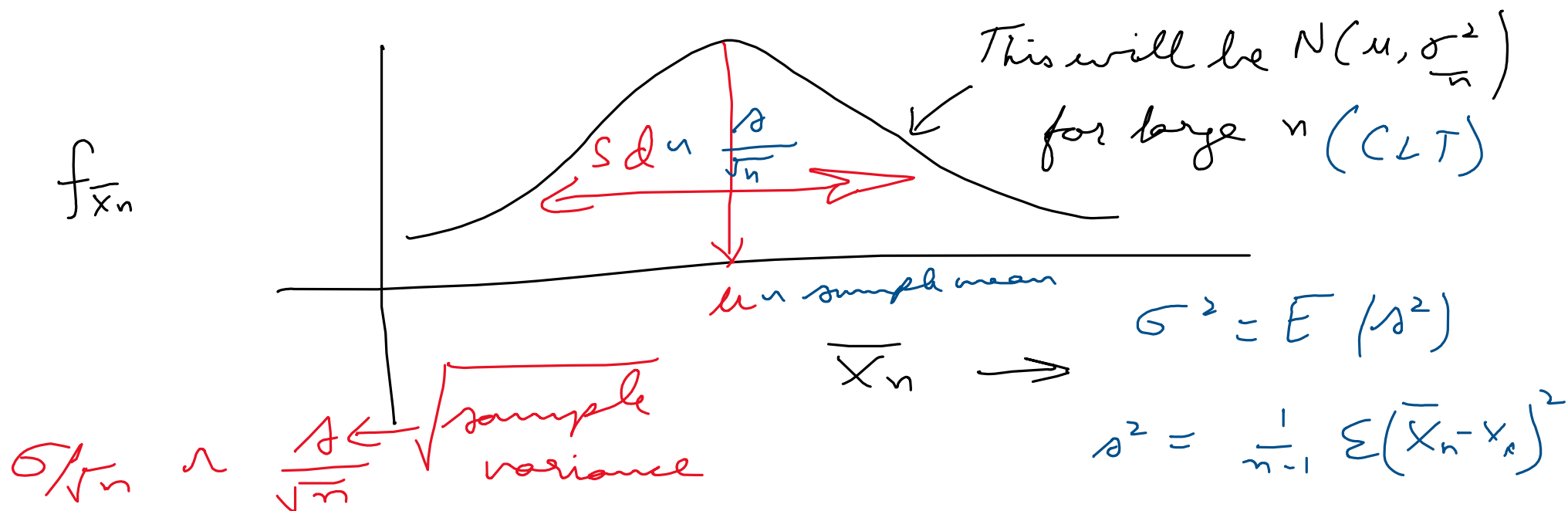
Population mean $\mu = E(\text{sample mean}) = E\left(\frac{1}{n} \sum_{i=1}^{n} X_i\right)$

Population variance $\sigma^2 = E(\text{sample variance}) = E\left(\frac{1}{n-1} \sum_{i=1}^{n} (\overline{X}_n - X_i)^2\right)$

CLT $\lim_{n \to \infty} \overline{X}_n = N\left(\mu, \frac{\sigma^2}{n}\right)$

# Application of Central limit theorem

$$X_1 \overset{iid}{\sim} (\mu, \sigma^2) \qquad \Big| \qquad \overline{X}_n$$



$f_{\overline{X}_n}$

$Sd \sim \dfrac{s}{\sqrt{n}}$

This will be $N\left(\mu, \dfrac{\sigma^2}{n}\right)$
for large $n$ (CLT)

$\mu \sim$ sample mean

$\overline{X}_n \longrightarrow$

$\sigma^2 = E(s^2)$

$s^2 = \dfrac{1}{n-1} \Sigma \left(\overline{X}_n - X_i\right)^2$

$\dfrac{\sigma}{\sqrt{n}} \sim \dfrac{s}{\sqrt{n}} \leftarrow \sqrt{\text{sample variance}}$

# Application of Central limit theorem

$$X_i \overset{iid}{\sim} (\mu, \sigma^2) \qquad \bar{X}_n = \frac{1}{n} \sum_{i=1}^{n} X_i$$

$$CLT \qquad \lim_{n \to \infty} \bar{X}_n = N\left(\mu, \frac{\sigma^2}{n}\right)$$

$$\mu = \bar{X}_n \pm Z^x \frac{\sigma}{\sqrt{n}} \qquad \sim N\left(\bar{X}_n, \frac{\sigma^2}{n}\right)$$

| C  | $Z^x$ |
|----|-------|
| 95 | 1.96  |
| 98 | 2.326 |
| 99 | 2.576 |

# Application of Central limit theorem

There is a 95% prob that $\mu$ is within $\pm 1.96$ Sd of $\overline{x}_n$.

Normal dist of $\overline{x}_n$

$Sd = \sigma/\sqrt{n}$

Normal dist

95%

$f_{\overline{x}_n}$

1.96 Sd     1.96 Sd

$\overline{x}_n$     $\mu$     $\overline{x}_n$     $\overline{X}_n$

$s^2$

2.5%     2.5%

| C | $Z^x$ |
|-----|-------|
| 95% | 1.96 |
| 98% | 2.326 |
| 99% | 2.576 |

$$Sd = \frac{\sigma}{\sqrt{n}} \sim \frac{\sqrt{s^2}}{\sqrt{n}}$$

$$\mu = \overline{x}_n \pm \boxed{Z^x} \; Sd \Big\}\; \begin{array}{l} \text{with} \\ C\% \text{ prob} \end{array}$$

# Application of Central limit theorem

$$X_i \overset{iid}{\sim} (\mu, \sigma^2) \qquad \text{If we take a sample of singe } n \text{ then } \bar{X}_n \equiv \sum_{i=1}^{n} \frac{x_i}{n}$$

Then for large $n$, CLT $\Rightarrow \qquad \bar{X}_n \rightarrow N\left(\mu, \frac{\sigma^2}{n}\right)$

| C | Z* |
|------|-------|
| 99% | 2.576 |
| 98% | 2.326 |
| 95% | 1.96 |

Therefore we can say with C confidence

that $\quad \mu = \bar{X}_n \pm Z^*(c)\,(Sd)$

Sample variance
$= \dfrac{\sigma}{\sqrt{n}} = \dfrac{E(S)}{\sqrt{n}}$

$$\therefore \quad \mu = \frac{1}{n}\sum_{i=1}^{n} x_i \pm Z^*(c)\,\frac{\frac{1}{n-1}\sum(\bar{X}_n - x_i)^2}{\sqrt{n}}$$

# t-distribution

If $X_i \overset{iid}{\sim} N(\mu, \sigma^2)$ and $\overline{X}_n \equiv \frac{1}{n} \sum_{i=1}^{n} X_i$

$$\overline{Z}_n \equiv \frac{\overline{X}_n - \mu}{\sigma/\sqrt{n}}$$
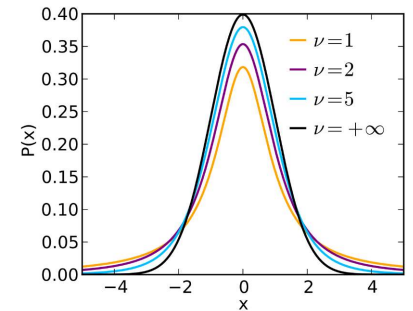
CLT as $n \to \infty$  $Z_n \to N(0,1)$

$t \equiv \dfrac{\overline{X}_n - \mu}{S/\sqrt{n}}$, then the r.v. 't' follows a distribution

known as $t$- distribution with

$\nu = n - 1$ degrees of freedom

pdf of t-distribution $= \dfrac{\Gamma\left(\frac{\nu+1}{2}\right)}{\sqrt{\nu\pi}\ \Gamma(\nu/2)} \left(1 + \dfrac{x^2}{\nu}\right)^{-\frac{\nu+1}{2}}$

# Confidence Interval

CLT based: If $X_i \overset{iid}{\sim} (\mu, \sigma^2)$

then $\qquad \mu = \overline{X} \pm Z^*(C) \dfrac{\sigma}{\sqrt{n}}$

Sample variance $S$ is an unbiased estimate of $\sigma$

| C | $Z^*$ |
|------|-------|
| 99% | 2.676 |
| 98% | 2.326 |
| 95% | 1.96 |

t-distribution: If $X_i \overset{iid}{\sim} N(\mu, \sigma^2)$

then $\qquad \mu = \overline{X} \pm t^*(C) \dfrac{S}{\sqrt{n}}$

If $n > 30$ use CLT, else use t-distribution

For 95% C

| n | n-1 | $t^*$ |
|-----|-----|-------|
| 6 | 5 | 2.571 |
| 11 | 10 | 2.228 |
| 31 | 30 | 2.042 |
| $\infty$ | $\infty$ | 1.960 |