

Central limit thm:

Population view	Population
1	5%
2	10%
3	15%
4	20%
5	
6	20%
7	10%
8	10%
9	5%
10	3%
	2%
	100%

any distribution

$$\begin{aligned} \text{true mean} &= 1(0.05) + 2(0.01) + \dots + 10(0.02) \\ &= 6.67 \end{aligned}$$

Sample 1000 times, 100 people each

$[\mu_1, \mu_2, \dots, \mu_{1000}] \rightarrow$ normal dist
around mean of 6.67

law of large numbers

Population size	Population
1	5%
2	10%
3	15%
4	20%
5	20%
6	20%
7	10%
8	10%
9	5%
10	2%
$\frac{1}{100} \times 100\% = 100\%$	

true mean = $1(0.05) + 2(0.01) + \dots + 10(0.02)$
 $= 0.67$

Sample $\frac{n}{10}$ people, mean $\neq 0.67$

Sample 100 people, mean $\rightarrow 0.67$

Sample 1000 people, mean ≈ 0.67

Sample entire people, mean = 0.67

$$\text{Sample size } n = \left(\frac{Z \sigma}{E} \right)^2$$

$Z \Rightarrow$ z-score of confidence level
 1.96 for 95%

$\sigma \Rightarrow$ sd of population
 $E \Rightarrow$ margin of error desired

Measures of Central Tendency

$$\text{Mean}(x) \mu = E[x] = \int_{-\infty}^{\infty} x f(x) dx$$

$$\text{Variance}(x) \text{Var}(x) = E[(x-\mu)^2] = E(x^2) - (E[x])^2$$

$$\sigma = \sqrt{\text{Var}(x)}$$

$$E(x/x=y) = \int_{-\infty}^{\infty} x f_{x|y}(x|y) dx$$

$$\text{Cov}(x,y) = E((x-E[x])(y-E[y])) = E(xy) - E[x]E[y]$$

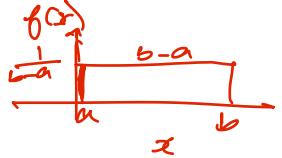
$$\rho(x,y) = \frac{\text{Cov}(x,y)}{\sigma_x \sigma_y}$$

Uniform Dist

$$f(x) = \frac{1}{b-a}$$

$$E[X] = \int_a^b x f(x) dx$$

$$\begin{aligned}
 &= \int_a^b x \cdot \frac{1}{b-a} dx \\
 &= \frac{1}{b-a} \int_a^b x dx \\
 &= \frac{1}{b-a} \left[\frac{x^2}{2} \right]_a^b = \frac{b^2 - a^2}{2(b-a)}
 \end{aligned}$$



Hypothesis Testing

H_0 :

H_a : reject or ^{not}reject

test statistic ✓ p-value ✓ $P_{\text{val}} < \alpha$ ✓

(z, t, Chi-sq)

one tailed:

$H_0: \mu = \mu_0$

$H_a: \mu < \mu_0$ or $\mu > \mu_0$

more powerful \rightarrow directional
 in campaign more effective?
 has higher voter turnout.

two tailed:

$H_0: \mu = \mu_0$

$H_a: \mu \neq \mu_0$

non-directional
 new policy subject effect on unemployment

Z-test

assumed test statistic follows
normal distribution under
null hypothesis.

$$Z\text{-statistic} \quad z = \frac{\bar{x} - \mu_0}{\sigma / \sqrt{n}} \sim N(0, 1)$$

e.g.: significant diff b/w means of
2 groups or mean diff sig. if
known value

Sample size is large ($n > 30^+$)
 σ is known (pop variance)

$\bar{x} \rightarrow$ sample mean
 $\mu_0 \rightarrow$ pop mean
 $\sigma \rightarrow$ pop sd
 $n \rightarrow$ sample size

campaign effectiveness:

$$\begin{aligned} \text{no. of contacts voted for DFK} &= \bar{x} = 0.10 \\ \text{sd of part dicty result} &= 0.05 = \sigma \end{aligned}$$

New strategy (sample size 1000 voters)

$$H_0: \text{New strategy does not support } (\mu_0) \quad \bar{x} = \mu_0$$

$$\frac{n}{1000} \text{ voters} \equiv 45\% \text{ support} = \bar{x}$$

$$H_a: \quad \bar{x} > \mu_0$$

$$z = \frac{0.10 - 0.10}{0.05 / \sqrt{1000}} = \frac{0.05}{0.05} \times \sqrt{1000} = 31.62$$

$$31.62 > 1.645 \quad \left[\begin{array}{l} \text{tailed test} \\ \text{5% significance level} \end{array} \right]$$

reject null hypo

\rightarrow Prob of observing 5% larger increase
by chance alone is low

1.645 \rightarrow cumulative area under normal curve from left up to 1.645

t + t-df

$$t\text{-statistic} \quad t = \frac{\bar{x} - \mu_0}{\sigma / \sqrt{n}} \sim t_{n-1}$$

(\approx) sample size is small
pop mean is unknown
(\Rightarrow) ~~normal~~

when $s^2 = \frac{\sum (x_i - \bar{x})^2}{n-1}$

1) indep/unpaired t-test

compare means of 2 indep groups
(like test scores)

t-dist similar to normal
but longer tail

2) paired / depn t-test

same group, diff times.
(test scores before & after extra doses)

assumptions

indep: var should be similar
indep: indep samples

example
in free age diff b/w Dem & Rep

H₀: No diff

H₁: there is diff

$$t = \frac{\bar{x}_1 - \bar{x}_2}{\sqrt{\frac{s_1^2}{n_1} + \frac{s_2^2}{n_2}}} = \frac{-5}{\sqrt{11/30}} = -3.5$$

indep + t	
Dem Sample	Rep Sample
30 voters	30 voters
25 age mean	50 age mean
5 sd age	6 age sd

$$df = n_1 + n_2 - 2$$

$$p\text{-value} = 0.00077 < 0.05$$

reject H₀

= area under f-distr that is more extreme than -3.5

2-tailed: both tails

→ prob of obs a value more extreme than -3.5 if H₀ is true

→ how is p calculated
t-distribution (~~assumed~~ population is unknown)

$n \rightarrow$ sample size
⇒ sample dist of mean
follows normal
but t-dist

Chi-squared test

(assess goodness of fit)

$$\chi^2 = \sum \frac{(O_i - E_i)^2}{E_i}$$

$O_i \rightarrow$ observed value

$E_i \rightarrow$ expected value

$df = \# \text{ categories in distribution}$

Determines if there is a significant association b/w categorical variables

Categorical
variables are indep

	E_1	E_2
J_1	10	10
J_2	15	30

no depn.

e.g.

age	Dow		Repub
	M	F	
M	100	150	250
F	180	120	300
	280	270	550

H_0 : No assoc.

H_a : Yes assoc.

$$E_{ij} = \frac{\text{Row total}_i \times \text{Col total}_j}{\text{Grand total}} =$$

$$\begin{bmatrix} 250 & 200 \\ 550 & 550 \end{bmatrix} + \begin{bmatrix} 180 & 300 \\ 550 & 550 \end{bmatrix} = \begin{bmatrix} 127.27 & 122.73 \\ 152.73 & 147.27 \end{bmatrix}$$

$$\chi^2 = \frac{(100 - 127.27)^2}{127.27} + \frac{(150 - 122.73)^2}{122.73} + \frac{(180 - 152.73)^2}{152.73} + \frac{(300 - 147.27)^2}{147.27}$$

$$= 21.03$$

$$DF = (n_{\text{row}} - 1)(n_{\text{col}} - 1) = (2 - 1)(2 - 1) = 1$$

P-value = 0.0000045 < 0.05
Reject the null hypothesis

Chi-squared dist

P-value: If the null hypo is true,
what is the prob of observing
value this extreme (or more extreme)

Z stat cont

55% with part A

new policy affects or
not

n = 1000 voters

520 support part A = 0.52

480 not support part A = 0.48

H₀: unchanged at 55%

H_a: changes from 55%

Z found test

$$Z = \frac{\hat{p} - p_0}{\sqrt{\frac{p_0(1-p_0)}{n}}} = \frac{0.52 - 0.55}{\sqrt{0.52(0.48)}} \times \sqrt{1000} \\ = -1.91$$

$$-1.96 < -1.91 < 1.96$$

within range

can't reject null hypo

assumption

pop is normal dist
esp for sample size ≥ 30

2 80

Distributions

Normal

histogram, Q-Q plot

Shapiro-Wilk test

Kolmogorov-Smirnov test

pdf

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

Symmetric, mean=median=mode

Binomial

Normal (Gaussian)

Uniform

Binomial

Poisson

Exponential

chi-sq

t-distr

F-distr

Uniform



$$f(x) = \frac{1}{b-a}$$

dorf

drf

Don't

t-test example 2 (dep/paired)

Same ppl. diff time (before + after policy)

SD \bar{d} , before = $\frac{\text{mean}}{6}$, sd 1.5
 after = $\frac{\text{mean}}{7}$ sd 1.2

H_0 : no effect.

H_a : significant effect

$$t\text{-test} = \frac{\bar{d}}{S_d / \sqrt{n}} = \frac{17-6}{2.06 / \sqrt{50}} = 2.63$$

p-value = 0.0111 < 0.05 (statistically significant)

Do not reject H_0

