

COVARIANCE CONTINUATION:

$\text{covar}(x, y)$ is not independent of the unit in which x and y are measured.

eg: Suppose x & y are rvs measured in cm have
 $\text{cov}(x, y) = 0.15$

If we change the unit to mm, then $X_1 = 10x$, $Y_1 = 10y$ are new rvs $\text{covar}(X_1, Y_1) =$
 $\text{covar}(10x, 10y) = 150$ ($K^2 \text{var}(x, y)$)

Solution So we normalize covariance to measure the correlation in absolute scale

$$\rho_{xy} \text{ always } * \rho_{xy} = \frac{\text{cov}(x, y)}{\sigma_x \sigma_y} = \frac{E(XY) - E(X)E(Y)}{\sigma_x \sigma_y} \quad \begin{matrix} \sigma_x = SD(x) \\ \sigma_y = SD(y) \end{matrix}$$

Setw E1, 12

This is called coefficient of correlation b/w x and y .

$\rho = 1$	$\rho = -1$	$\rho > 0$	$\rho < 0$	$\rho = 0$
perfect +ve relationship	perfect -ve relationship	linear +ve relationship	linear -ve relationship	

"Measures only linear relationship b/w 2 rvs"

RESULT $\text{var}(ax+by) = a^2 \text{var}(x) + b^2 \text{var}(y) + 2ab \sigma_x \sigma_y \rho_{xy}$ correlation

* $\boxed{\text{var}(ax+by) = a^2 \text{var}(x) + b^2 \text{var}(y) + 2ab \text{covar}(x, y)}$

Proof:

$$\begin{aligned}
 \text{var}(ax+by) &= E[(ax+by) - (a\bar{x}+b\bar{y})]^2 \\
 &= E[a(x-M_x) + b(y-M_y)]^2 \\
 &= E[a^2(x-M_x)^2 + b^2(y-M_y)^2 + 2ab(x-M_x)(y-M_y)] \\
 &= a^2 E(x-M_x)^2 + b^2 E(y-M_y)^2 + 2ab E[(x-M_x)(y-M_y)] \\
 &= a^2 \text{var}(x) + b^2 \text{var}(y) + 2ab \text{covar}(x, y)
 \end{aligned}$$

$$a = \frac{1}{\sigma_x}, b = \frac{1}{\sigma_y}$$

$$\star \text{Var}\left(\frac{X}{\sigma_x} + \frac{Y}{\sigma_y}\right) = \frac{1}{\sigma_x^2} \text{Var}(X) + \frac{1}{\sigma_y^2} \text{Var}(Y) + 2ab\sigma_x\sigma_y r_{XY}$$

$$= 2 + 2r_{XY}$$

$$\star \text{Var}\left(\frac{X}{\sigma_x} - \frac{Y}{\sigma_y}\right) = 2 - 2r_{XY}$$

5

PROPERTIES:

$$\text{① } -1 \leq r_{XY} \leq 1$$

* Variance of any rv is non-negative

$$\text{Var}\left(\frac{X}{\sigma_x} + \frac{Y}{\sigma_y}\right) \geq 0 \quad \text{Var}\left(\frac{X}{\sigma_x} - \frac{Y}{\sigma_y}\right) \geq 0$$

10

$$2 + 2r_{XY} \geq 0$$

$$r_{XY} \geq -1$$

$$2 - 2r_{XY} \geq 0$$

$$r_{XY} \leq 1$$

$$\text{② } r_{XY} = 1 \text{ iff } Y = ax + b \text{ for constants } a, b \neq 0$$

$$r_{XY} = -1 \text{ iff } Y = ax + b \text{ for constants } a, b \neq 0$$

15 Proof: * Suppose $r_{XY} = 1$

$$\therefore 2 - 2r_{XY} \Rightarrow \text{Var}\left(\frac{X}{\sigma_x} - \frac{Y}{\sigma_y}\right) = 0 \quad \text{if } \frac{X}{\sigma_x} - \frac{Y}{\sigma_y} = c, \text{ constant}$$

$$20 \quad \frac{Y}{\sigma_y} = \frac{X}{\sigma_x} - c \quad Y = \left(\frac{\sigma_y}{\sigma_x}\right) X - \left(c\sigma_y\right)$$

$$* \text{Let } Y = ax + b$$

$$r_{XY} = \frac{\text{Cov}(X, Y)}{\sigma_x \sigma_y}$$

25

eg If x and y are r.v.s with joint pdf

$$f(x, y) = x + y \quad 0 < x, y < 1.$$

Show that x and y are not linearly correlated.

Ans:

First we find marginal densities of x .

$$f_x(x) = \int_0^1 x + y \, dy = xy + \frac{y^2}{2} \Big|_0^1$$

$$f_x(x) = x + 1/2 \quad 0 < x < 1$$

$$f_y(y) = y + 1/2 \quad 0 < y < 1$$

$$\begin{aligned} E(x) &= \int_0^1 x f_x(x) dx \\ &= \int_0^1 x^2 + x/2 dx = \frac{x^3}{3} + \frac{x^2}{4} \Big|_0^1 \end{aligned}$$

$$E(x) = 1/3 + 1/4 = 7/12$$

$$E(y) = 7/12$$

$$\begin{aligned} E(xy) &= \iint_{\mathbb{R}^2} xy f(x, y) dy dx \\ &= \int_0^1 \int_0^1 xy(x+y) dy dx \\ &= \int_0^1 \frac{x^2 y^2}{2} + \frac{xy^3}{3} \Big|_0^1 dx \\ &= \int_0^1 \frac{x^2}{2} + \frac{x}{3} dx = \frac{x^3}{6} + \frac{x^2}{6} \Big|_0^1 \end{aligned}$$

$$E(xy) = 1/6 + 1/6 = 2/6 = 1/3$$

$$\begin{aligned} E(x^2) &= \int_0^1 x^2 f_x(x) dx \\ &= \int_0^1 x^2 (x+1/2) dx \end{aligned}$$

$$E(x^2) = \int_0^1 x^4/4 + 1/2 x^3/3 y \Big|_0^1 = 1/4 + 1/6 = \frac{10}{24} = \frac{5}{12}$$

$$\text{Var}(x) = E(x^2) - E(x)^2 = \frac{5}{12} - \frac{49}{144} = \frac{60-49}{144} = \frac{11}{144}$$

$$E(y^2) - E(y)^2 = \text{Var}(y) = 11/144$$

$$\text{correlation} = \frac{\text{covar}(x, y)}{\sigma_x \sigma_y} = \frac{E(xy) - E(x)E(y)}{\sqrt{\text{var}(x)} \sqrt{\text{var}(y)}}$$

$$= \frac{\frac{1}{3} - \frac{49}{144}}{\frac{\sqrt{11}}{12} \cdot \frac{\sqrt{11}}{12}} = \frac{\frac{-1}{144}}{\frac{11}{144}} = \frac{-1}{144} \times \frac{144}{11} = \frac{-1}{11} \neq 1$$

\therefore They are not linearly related.

Mean & variance of sum of rvs

10 For any list of rvs x_1, x_2, \dots, x_n

$$\text{Let } w_n = x_1 + x_2 + \dots + x_n$$

$$E(w_n) = \sum E(x_i)$$

$$15 \quad \text{var}(w_n) = \sum_i \text{var}(x_i) + 2 \sum_{i < j} \text{covar}(x_i, x_j)$$

If x_1, x_2, \dots, x_n are pairwise independent,
or pairwise correlated then

$$20 \quad \text{var}(w_n) = \sum_i \text{var}(x_i) \quad (\text{Ans})$$

RESULT: Sum of iid rvs

If $x_1 = x_2 = \dots = x_n$ are iid rvs each with
mean "M" and SD " σ " then

$$E(w_n) = n\mu$$

$$25 \quad \text{var}(w_n) = n\sigma^2$$

eg: x_1, x_2, \dots be a sequence of rvs with $E(x_i) = 0$
and $\text{covar}(x_i, x_j) = 0.8^{|i-j|}$. Find the mean &
variance of $y_i = x_i + x_{i+1} + x_{i+2}$

30 Soln:

$$E(y_i) = 0$$

$$\text{Var}(y_i) = \text{var}(x_i + x_{i+1} + x_{i+2})$$

$$= \text{var}(x_i) + \text{var}(x_{i+1}) + \text{var}(x_{i+2}) +$$

$$2\text{covar}(x_i, x_{i+1}) + 2\text{covar}(x_{i+1}, x_{i+2}) +$$

$$2\text{covar}(x_i, x_{i+2})$$

$$\text{covar}(x, y) = E(xy) - E(x)E(y)$$

$$\text{covar}(x, x) = E(x^2) - E(x)^2$$

$$= 1 + 1 + 1 + 2(0.8) + 2(0.8) + 2(0.64)$$

$$\text{covar}(Y_1) = 7.48$$

Eg: At a party of $n \geq 2$ people, each person has his/her hat in a common box. They are shuffled and each person draws a hat, without replacement. If a person draws his/her hat it is success. What are mean & variance of the # no. of successes?

Let $X_i = \begin{cases} 1 & \text{if person } i \text{ draws his/her own hat} \\ 0 & \text{o/w} \end{cases}$

$$P(X_i = 1) = \frac{1}{n} \quad P(X_i = 0) = 1 - \frac{1}{n}$$

No. of successes $S_n = X_1 + X_2 + \dots + X_n$

$$E(X_i) = \frac{1}{n} \quad \text{var}(X_i) = \frac{1}{n}(1 - \frac{1}{n})^{pq}$$

$$E(S_n) = n \times \frac{1}{n} = 1$$

$$\text{covar}(X_i, X_j) = E(X_i X_j) - E(X_i)E(X_j)$$

To find $E(X_i X_j)$

$$X_i X_j = \begin{cases} 0 & \text{o/w} \\ 1 & \text{iff } X_i = X_j = 1 \end{cases}$$

$$E(X_i X_j) = 1 \cdot P_{X_i X_j}(1, 1)$$

$$f_{(X_i Y)}(x|y) = \frac{f(x, y)}{f_y(y)}$$

$$= 1 \cdot P_{X_i X_j}(1|1) P_{X_j}(1)$$

Given $X_j = 1$, the j^{th} person draws his hat
 then $X_i = 1$ iff i^{th} person draws his hat from
 $n-1$ other hats $P(1|1) = \frac{1}{n-1}$

$$E(X_i X_j) = \frac{1}{n-1} \cdot \frac{1}{n}$$

$$\text{covar}(x_i, x_j) = \frac{1}{n(n-1)} - \frac{1}{n^2}$$

$$\text{var}(S_n) = \text{var}(x_1) + \text{var}(x_2) + \dots + \text{var}(x_n) + \\ 2nC_2 \text{cov}(x_i, x_j)$$

$$= n\left(\frac{1}{n} - \frac{1}{n^2}\right) + 2 \cdot \frac{n(n-1)}{2} \left[\frac{1}{n(n-1)} - \frac{1}{n^2}\right]$$

$$\text{var}(S_n) = \left(1 - \frac{1}{n}\right) + \left(1 - \frac{n-1}{n}\right) = 1$$

Note: Suppose each person immediately returns the hat into iron. What are mean & variance?

$$E(x_i) = E(S_n) = 1 \quad \text{var}(S_n) = n\left(\frac{1}{n} - \frac{1}{n^2}\right) = 1 - \frac{1}{n}$$

LIMIT THEOREMS

fundamental issues related to the asymptotic behaviour of sequence of r.v.s

1. Markov inequality

2. Chebyshev's inequality

3.

4.

① Markov Inequality

If X is a r.v. that takes non-negative value

then for any $a > 0$

$$P[X \geq a] \leq \frac{E(X)}{a}$$

upper bound for $P[X \geq a]$

Proof:

$$E(x) = \int_0^\infty x f(x) dx$$

$$= \int_0^a x f(x) dx + \int_a^\infty x f(x) dx$$

$$\leq \int_a^\infty a f(x) dx$$

> removing this will become \leq

taking min value for x in interval a to ∞

$$\leq \int_a^{\infty} xf(x)dx$$

$$\leq a \int_a^{\infty} f(x)dx$$

$$E(x) \leq a P[x \geq a]$$

② Chebyshev's inequality.

If X is a rv with μ and variance σ^2 , then
for any values of $K > 0$

$$P[|X-\mu| \geq K] \leq \frac{\sigma^2}{K^2}$$

Proof: $(X-\mu)^2$ is a non-negative rv

By Markov's inequality,

$$P[(X-\mu)^2 \geq K^2] \leq \frac{E(X-\mu)^2}{K^2} = \frac{\sigma^2}{K^2}$$

$$P[|X-\mu| \geq K] \leq \frac{\sigma^2}{K^2}$$

If a rv has small var then the probability
that it takes values far away from mean
is small.

Note: If $K=n\sigma$, Chebyshev's inequality between

$$P[|X-\mu| \geq n\sigma] \leq \frac{1}{n^2}$$

The probability that X deviates from mean
by atleast n SD is $< 1/n^2$.

eg Suppose that no. of items produced in a factory
in a week is a rv with a mean of 500

30. i. What is it can be said about the prob that
this week's production will be atleast 1000?
ii. If the variance is known to be 100, what can
be said about the probability that the
week's production will be betw 400 & 600?

Soln:

- i. X no. of items produced in a week
 $X \geq 0$ & $E(X) = 500$

By Markov's inequality,

$$P[X \geq 1000] \leq \frac{500}{1000} = 0.5$$

- ii. Given $\sigma^2 = 100$

By Chebychev's inequality,

$$P[|X - \mu| \geq k] \leq \frac{\sigma^2}{k^2}$$

$$\begin{cases} k = 100 \\ 400 - 500 = -100 \\ 600 - 500 = 100 \end{cases}$$

$$P[|X - 500| \geq 100] \leq \frac{100}{100^2} = 0.01$$

$$P[|X - 500| \leq 100] \geq 0.99$$

$$P[400 \leq X \leq 600] \geq 0.99$$

eg A person post office handles an average of 10000 letters per day. What can be said about the prob that it will handle i. atleast 15000 letters tomorrow?

ii. fewer than 15000 letters tomorrow?

g.e. the variance is 2000, what can be said about the probability that it will handle between 8000 and 12000 letters tomorrow?

Soln: X: no. of letters handled

$$E(X) = 10000 \quad \rightarrow E(x)/\sigma$$

Markov's inequality

$$P[X \geq 15000] \leq \frac{10000}{15000} = \frac{2}{3}$$

$$i. P[X \leq 15000] > \frac{1}{3}$$

$$\sigma^2 = 2000$$

Chebychev's inequality

$$P[8000 \leq X \leq 12000] = P[|X - \mu| \leq 2000] = ?$$

$$P[|X - \mu| \geq 2000] \leq \frac{2000}{2000^2} = \frac{1}{2000}$$

$$\therefore P[8000 \leq X \leq 12000] = 1 - \frac{1}{2000} = \frac{1999}{2000}$$

Note: The bounds obtained for the year
in Markov's and Chebyshev's inequalities
are not close to actual probabilities.

e.g. roll a dice X - outcome of the die.

$$E(X) = \frac{1}{6} (1+2+3+4+5+6) = \frac{7}{2} \quad (2\frac{1}{2})$$

$$E(X^2) = \frac{1}{6} (1^2 + 2^2 + 3^2 + 4^2 + 5^2 + 6^2) = \frac{91}{6}$$

$$\text{var}(X) = E(X^2) - E(X)^2 = \frac{91}{6} - \left(\frac{21}{6}\right)^2 = \frac{35}{12}$$

By Markov's inequality:

$$P[X \geq 6] \leq \frac{\frac{21}{6}}{6} = \frac{21}{36} = 0.583$$

By Chebyshev's inequality:

$$P[|X - \frac{21}{6}| \geq \frac{3}{2}] \leq \frac{\frac{35}{12}}{\frac{91}{4}} = 1.296$$

But the trivial bound $P[|X - \frac{21}{6}| \geq \frac{3}{2}] \leq 1$

Exact probabilities:

$$P[X \geq 6] = \frac{1}{6} = 0.167 \neq 0.583 \quad (\text{too far})$$

$$P[|X - \frac{21}{6}| \geq \frac{3}{2}]$$

$$|X - \frac{21}{6}| \leq \frac{3}{2}$$

$$P[X \leq 2 \text{ or } X \geq 5] P[2 \leq X \leq 5] = \frac{4}{6} = \frac{2}{3}$$

$$-\frac{3}{2} \leq X - \frac{21}{6} \leq \frac{3}{2}$$

$$= 0.667 \leq 1$$

$$\frac{21}{6} - \frac{3}{2} \leq X \leq \frac{3}{2} + \frac{21}{6}$$

(too far)

$$\frac{12}{6} \leq X \leq \frac{30}{6}$$

$$2 \leq X \leq 5$$

Note: If $\text{var}(X)=0$, then X is constant with probability 1.

We can prove this \uparrow using Chebyshev's inequality.

Laws of Large numbers

weak law

describes how a seq of prob converges

strong law

describes how a seq of n behave in the limit

* Weak law of large numbers (WLL)

Let x_1, x_2, \dots, x_n be a sequence of iid rvs each having limit mean $E(x_i) = \mu$ and $\text{var}(x_i) = \sigma^2$

Then for any $\epsilon > 0$

$$P\left\{\left|\frac{x_1 + x_2 + \dots + x_n}{n} - \mu\right| > \epsilon\right\} \rightarrow 0 \text{ as } n \rightarrow \infty$$

$$\lim_{n \rightarrow \infty} P\left\{\left|\frac{x_1 + x_2 + \dots + x_n}{n} - \mu\right| > \epsilon\right\} = 0$$

15

(or)

$$P\left\{\left|\frac{x_1 + x_2 + \dots + x_n}{n} - \mu\right| \leq \epsilon\right\} \rightarrow 1 \text{ as } n \rightarrow \infty$$

sample mean

true mean

20 Proof:

$$\text{Let } M_n = \frac{x_1 + x_2 + \dots + x_n}{n}$$

$$E(M_n) = \frac{1}{n} E(x_1) + E(x_2) + \dots + E(x_n) = \frac{1}{n} n\mu = \mu$$

$$25 \quad \text{var}(M_n) = \frac{1}{n^2} [\text{var}(x_1) + \text{var}(x_2) + \dots + \text{var}(x_n)] = \frac{1}{n^2} n\sigma^2$$

$$\text{var}(M_n) = \frac{\sigma^2}{n}$$

$$\text{var}(ax) = a^2 \text{var}(x)$$

$$\text{var}(a+x) = \text{var}(a) + \text{var}(x)$$

By Chebyshev's inequality,

25

$$P\{|M_n - \mu| \geq \epsilon\} \leq \frac{\sigma^2}{n\epsilon^2} \text{ for any } \epsilon > 0$$

for any fixed $\epsilon > 0$

$$P\left\{\left|\frac{x_1 + x_2 + \dots + x_n}{n} - \mu\right| \geq \epsilon\right\} \rightarrow 0 \text{ as } n \rightarrow \infty$$

WLL and relative frequency interpretation of probability:

In a repeatable experiment with sample space S , let A be an event. Suppose the experiment is repeated independently.

For $i = 1, \dots, n$,

Let $X_i = \begin{cases} 1 & \text{if } A \text{ occurs in } i^{\text{th}} \text{ trial} \\ 0 & \text{if } A \text{ does not occur} \end{cases}$

10) $E(X_i) = 1P(A) + 0P(\bar{A}) = P(A)$

and $S_n = X_1 + X_2 + \dots + X_n = \#(A) \begin{bmatrix} \text{no. of times } A \text{ occurs} \\ \text{in } n \text{ trials} \end{bmatrix}$

Then by WLL

15)

$$\lim_{n \rightarrow \infty} P\left\{ \left| \frac{X_1 + X_2 + \dots + X_n}{n} - E(X_i) \right| \geq \epsilon \right\} = 0$$

$$\lim_{n \rightarrow \infty} P\left\{ \left| \frac{\#(A)}{n} - P(A) \right| \geq \epsilon \right\} = 0$$

20) When n is large, no. of times A occurs out of n times is close to $P(A)$.

* Strong Law of Large Numbers (SLL)

25) $P\left\{ \lim_{n \rightarrow \infty} \frac{X_1 + X_2 + \dots + X_n}{n} = \mu \right\} = 1$

With probability 1, the sequence of sample means converges to the true mean.

30) Central Limit theorem (CLT)

Sum of large number of iid rvs have a distribution that is approximately normal.

Let x_1, x_2, \dots be a sequence of iid rvs each with mean $E(x_i) = \mu$ and $\text{var}(x_i) = \sigma^2$. Then distribution of $\frac{x_1 + x_2 + \dots + x_n - n\mu}{\sigma/\sqrt{n}}$ tends to

⁵ follow SND as $n \rightarrow \infty$

$$\text{ie.) } \frac{x_1 + x_2 + \dots + x_n - n\mu}{\sigma/\sqrt{n}} = \frac{\bar{x} - \mu}{\sigma/\sqrt{n}} \text{ is close to SND as } n \rightarrow \infty$$

$$\bar{x} \sim N(\mu, \sigma/\sqrt{n}) \quad z = \frac{\bar{x} - \mu}{\sigma/\sqrt{n}}$$

¹⁰ eg The lifetime in hrs, of a component is a rv with mean 500 and SD 35. What is the prob that a rv random sample of 49 bulbs will have mean life betw 488 and 505 hrs?

Soln:

\bar{x} - mean life of $n=49$ bulbs

By CLT $\mu = 500, \sigma_{\bar{x}} = 35/\sqrt{49}$

$$\bar{x} \sim N(500, \frac{35}{\sqrt{49}}) \sim N(500, 5)$$

$$z = \frac{\bar{x} - \mu}{\sigma/\sqrt{n}} = \frac{\bar{x} - 500}{5} = \frac{\bar{x} - 500}{5}$$

20

$$P[488 < \bar{x} < 505] = P[-2.4 < z < 1]$$

$$= \Phi(1) - \Phi(-2.4)$$

$$= 0.8413 - 0.0082$$

$$= 0.8331$$

²⁵ eg: The distribution of annual earnings of all bank managers has a negatively skewed distribution with mean \$19000 and \$2000.

³⁰ In a random sample of 30 managers, what is the probability that mean annual earnings is more than \$19750?

Soln: \bar{x} -mean annual earnings of 30 mamas

$$M_{\bar{x}} = M \text{ and } \sigma_{\bar{x}} = \frac{\sigma}{\sqrt{n}}$$

$$\text{By CLT } M_{\bar{x}} = 19000 \quad \sigma_{\bar{x}} = \frac{2000}{\sqrt{30}} =$$

$$\bar{x} \sim N\left(19000, \frac{2000}{\sqrt{30}}\right) \quad z = \frac{\bar{x} - M}{\sigma/\sqrt{n}}$$

$$P[\bar{x} > 19750] = P\left[z > \frac{19750 - 19000}{\frac{2000}{\sqrt{30}}}\right]$$

365.148

$$= P[z > 2.054]$$

$$= 1 - \Phi(2.054)$$

$$= 1 - 0.9798$$

$$= 0.0202$$

eg: The time it takes for a student to finish an aptitude test has a pdf $f(x) = \begin{cases} 6(x-1)(2-x) & 1 < x < 2 \\ 0 & \text{elsewhere} \end{cases}$

Approximate the prob that the average length of time it takes for a random sample of 15 students to complete the test is less than 1 hr 25 mins?

Soln:

For $i=1$ to 15, let x_i be the time for i^{th} student to finish the test

$$\bar{x} = \frac{x_1 + x_2 + \dots + x_{15}}{15}$$

$$M = E(x_i) = \int_1^2 x f(x) dx = \int_1^2 6(x-1)(2-x) dx = \frac{1+2}{2} - \frac{3}{2} = 1.5$$

$$E(x_i^2) = \int_1^2 x^2 f(x) dx = \frac{23}{10} = 2.3$$

$$\sigma^2 = \text{var}(x) = E(x_i^2) - E(x_i)^2 = 2.3 - 1.5^2 = 0.05$$

$$\bar{x} \sim N\left(1.5, \frac{\sqrt{0.05}}{\sqrt{15}}\right) \quad P\left(\bar{x} < \frac{85}{60}\right) = 0.0749$$

$$P(z < -1.44)$$

If 20 random numbers are selected from $(0,1)$ approximate the prob that sum of them is atleast 8.

Ans:

For $i=1$ to 20, let x_i be the i th no. selected

$$P\left[\sum_{i=1}^{20} x_i \geq 8\right] = P\left[\bar{x} \geq \frac{8}{20}\right]$$

$$x \sim \text{uniform}(0,1) * \mu = \frac{a+b}{2} = \frac{0+1}{2} = 0.5$$

$$* \sigma^2 = \frac{(a-b)^2}{12} = \frac{(0-1)^2}{12} = \frac{1}{12}$$

$$\bar{x} \sim N(\mu, \sigma/\sqrt{n}) = N(0.5, \frac{1}{\sqrt{240}})$$

$$P\left[\bar{x} \geq \frac{8}{20}\right] = P\left[z > \frac{\frac{8}{20} - 0.5}{\frac{1}{\sqrt{240}}}\right] \rightarrow 0.0645$$

$$= P\left[z \geq -15.504\right] = 1 - \varphi(-1.55)$$

$$= 1 - 0.0606 = 0.9394$$

STATISTICS METHODS

descriptive

inferential

estimation

hypothesis

$$CI = \bar{X} \pm Z_{\alpha/2} \frac{\sigma}{\sqrt{n}}$$

Confidence interval

 $n \geq 30$ Large (becomes normal dist) $n < 30$ (follows t-distribution)degrees of freedom $\rightarrow n-1$

$$CI = \bar{X} \pm t_{\alpha/2} \frac{s}{\sqrt{n}}$$

eg: $n=70 \geq 30$, σ unknown $\bar{X} \sim \text{normal}(0,1)$

$$\bar{X} = 1759 \quad s = 380$$

$$99\% \rightarrow 2.58 \quad 95\% \rightarrow 1.96$$

$$CI = 1759 \pm \frac{2.58 \cdot 380}{\sqrt{70}}$$

$$90\% \rightarrow 1.645$$

$$= 1759 \pm 117.180$$

$$CI = 1641.82 / 1876.18$$

$$\sigma_p = \sqrt{\frac{p(1-p)}{n}}$$

population proportion p

$$CI = \hat{p} \pm Z_{\alpha/2} \sqrt{\frac{\hat{p}(1-\hat{p})}{n}}$$

$$eg: \hat{p} = \frac{25}{100} = 0.25 \quad 95\%$$

$$CI = 0.25 \pm 1.96 \sqrt{\frac{0.25(0.75)}{100}}$$

$$= 0.25 \pm 0.00735 \quad 0.08487$$

$$CI = 0.25735 / 0.24265 \quad 0.3348 / 0.1651$$

Sampling error

parameter = statistic + error

$$\bar{X} + \left(\frac{Z \sigma}{\sqrt{n}} \right) \rightarrow \text{margin of error}$$

$$e = \frac{Z \sigma}{\sqrt{n}} \quad n = \frac{Z^2 \sigma^2}{e^2} = \left(\frac{Z \sigma}{e} \right)^2$$

"Population mean" check

eg: $\sigma = 45$ $e = \pm 5$ 90% n (sample size) = ?

$$n = \left(\frac{z\sigma}{e}\right)^2 = \left(\frac{1.645 \times 45}{5}\right)^2 = 219.188$$

≈ 220 (always round up)

eg: 95% $e = \pm 5$ $\sigma = 15$

$$n = \left(\frac{z\sigma}{e}\right)^2 = \left(\frac{1.96 \times 15}{5}\right)^2 = 34.5744$$

≈ 35

Population proportion

$$e = z \sqrt{\frac{\hat{p}(1-\hat{p})}{n}} \quad n = \frac{z^2 \hat{p}(1-\hat{p})}{e^2}$$

eg: $\hat{p} = 0.12$, $e = 3$, 95%

$$n = \frac{(1.96)^2 0.12 (0.88)}{9/100^2} = 0.0450 \times 100^2 = 450$$

- * Characteristics for selecting estimators (unbiased)
- * unbiased (how much the calculated value differs from the true value)
- * efficient (σ^2 small more efficient)
- * sufficient
- * consistent

- * unbiased

$$E(\hat{\theta}) - \theta = 0 \quad \text{bias} = 0 \quad [\text{if } \hat{\theta} \text{ is an unbiased estimation of } \theta]$$

sample mean & sample median \rightarrow unbiased

Hypothesis:

is a statement or claim about the parameters of one or more populations

Two sided hypothesis

$H_0: \mu = 8.0$ (null hypothesis)

$H_1: \mu \neq 8.0$ (opp to null hypothesis)

→ alternate hypothesis / research

One sided hypothesis

$H_0: \mu \leq 8.0$ $\mu \geq 8.0$

$H_1: \mu > 8.0$ $\mu < 8.0$

* Null hypothesis is a hypothesis of "no diff".

* It always contain equality sign. It may / may not be rejected.

Hypothesis testing begins with the assumption that null hyp is true. Either one hypothesis will be the result / true / claim.

CRITICAL VALUE

if $\bar{x} \sim N(\mu, \sigma/\sqrt{n})$

test statistics $z = \frac{\bar{x} - \mu}{\sigma/\sqrt{n}}$ NSND

$$z = \frac{t - E(t)}{SE}$$

t - statistics

SE - standard error

* P-value method: compare p-value & α

* if p-value $< \alpha$, reject H_0

* critical value method: compare $|z|$ & $|z_{\alpha}|$

$|z| > |z_{\alpha}|$ reject

$|z| < |z_{\alpha}|$ do not reject

σ known case

$$H_0: \mu \leq 12 \quad H_1: \mu > 12 \rightarrow \text{right-tailed}$$

$n = 40 \quad \bar{x} = 13.25 \quad \sigma = 3.2 \quad \alpha = 0.05 \quad z_\alpha = 1.96$

Test statistics: $Z = \frac{\bar{x} - \mu_0}{\sigma/\sqrt{n}} = \frac{13.25 - 12}{3.2/\sqrt{40}} = 0.39 \times \sqrt{40} = 2.466$

By critical value approach,

$$|z| > z_\alpha \quad 2.466 > 1.96$$

reject H_0

\therefore Service goal of 12 mins or less is not achieved

* Left tailed $P[Z \leq z]$

* Right tailed $1 - P[Z \leq z]$

* Two-tailed $\frac{1}{2}[P[Z \leq z]] + \frac{1}{2}[1 - P[Z \leq z]]$

σ unknown case

use t-distribution

rejection rule for critical value approach:

$H_0: \mu \geq \mu_0$ rej H_0 if $t \leq -t_\alpha$ left

$H_0: \mu \leq \mu_0$ rej H_0 if $t \geq t_\alpha$ right

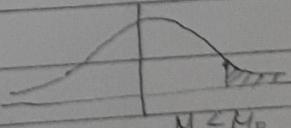
$H_0: \mu = \mu_0$ rej H_0 if $|t| > t_\alpha$ 2 tailed

known

eg. 25 $H_0: \mu \leq 130$

$H_1: \mu > 130$

$n = 9 \quad \bar{x} = 131.08 \quad \sigma = 1.5 \quad \alpha = 0.01$



$\mu \leq \mu_0$

Test statistics: $\frac{\bar{x} - \mu_0}{\sigma/\sqrt{n}} = \frac{131.08 - 130}{1.5/\sqrt{9}} = \frac{2.16}{1.5} = 2$

critical value approach

For $\alpha = 0.01$, $z_\alpha = 2.58$

$$P[Z \leq z_\alpha] = 0.01$$

$$z_\alpha = -2.33$$

$|z| < |z_\alpha| \quad \therefore \text{do not reject } H_0$

p-value approach:

$$1 - P[Z \leq 2.16] = 1 - 0.9846 = 0.0154$$

$$0.0154 > 0.01$$

$p\text{-value} > \alpha$ Hence do not reject H_0

\therefore Manufacturer's claim is correct

eg:

$$H_0: \mu = 1000$$

$$H_1: \mu > 1000$$

$$n=64 \quad \bar{X} = 1038 \quad s = 146 \quad \alpha = 0.05$$

$$\text{Test statistics: } \frac{\bar{X} - \mu_0}{s/\sqrt{n}} = \frac{1038 - 1000}{146/\sqrt{64}} = 2.0821$$

critical value approach

$$\text{For } \alpha = 0.05 \quad z_{\alpha} = 1.645$$

$$z > z_{\alpha} \quad \text{reject } H_0$$

 σ -unknown

eg:

$$H_0: \mu \geq 8$$

$$H_1: \mu < 8$$

$$n=25 \quad \bar{X} = 7.73 \quad s = 0.77 \quad \alpha = 0.05$$

$$\begin{aligned} \text{Test statistics} &= \frac{\bar{X} - \mu_0}{s/\sqrt{n}} \quad \text{with dof } n-1 = 24 \\ &= \frac{7.73 - 8}{0.77/\sqrt{25}} \end{aligned}$$

$$t = -1.75$$

Critical value approach:

$$t = \bar{X} - \mu_0$$

$$\text{for } \alpha = 0.05, \text{ one tailed dof} = 24 \quad t_{\alpha} = 1.71 \quad \text{Refer table (t-dist)}$$

$$-1.75 > -1.71$$

$$|t| > |t_{\alpha}| \quad \therefore \text{reject } H_0$$

P-value approach:

$$P[Z \leq -1.75] = P[Z \geq 1.75] \quad (\because \text{They are symmetric})$$

eg:

$$H_0: M \geq 5$$

$$H_1: M < 5$$

$$n=20 \quad \bar{x}=4.6 \quad s=2.2 \quad \alpha=0.05$$

$$t = \frac{\bar{x} - \mu_0}{s/\sqrt{n}} = \frac{4.6 - 5}{2.2/\sqrt{20}} = -0.813$$

Critical value approach, $\alpha=0.05$, $df=19$, $t_{\alpha}=1.73$

$$|t| < |t_{\alpha}| \quad t < -t_{\alpha} \text{ is false}$$

\therefore Don't reject H_0

95.1. confidence interval for μ :

$$CI = \bar{x} + t_{\alpha/2} \frac{s}{\sqrt{n}}$$

For 95% CI, $df=19$,
 $t_{\alpha/2}=2.09$ (always 2-tail)

$$= 4.6 + (2.09) \frac{2.2}{\sqrt{20}} = 4.6 \pm 1.028$$

$$= 5.628 \pm 3.572$$

eg: For the sample, $\bar{x}=14.4$, $s=0.158$, $n=25$

$$H_0: M=14.0$$

$$H_1: M \neq 14.0$$

$$\alpha=0.05$$

test statistics: $\frac{\bar{x}-\mu_0}{s/\sqrt{n}} = \frac{14.4 - 14}{0.158/\sqrt{25}} = 5.66$

$$t_{\alpha}=2.13$$

$$p\text{-value} = 0.004$$

\therefore claim is wrong.

independant samples

eg: $n_1 = 20 \quad n_2 = 25$

$\bar{x}_1 = 29.8$

$\sigma_1 = 4.0$

$\bar{x}_2 = 34.7$

$\sigma_2 = 5.0$

 σ_1, σ_2 known \rightarrow Nunknown $\rightarrow t$ -dist

$H_0: \mu_1 = \mu_2$

$\alpha = 0.01$

$H_1: \mu_1 \neq \mu_2$

$$\text{test statistic } Z = \frac{\bar{x}_1 - \bar{x}_2}{\sqrt{\frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2}}} = \frac{29.8 - 34.7}{\sqrt{\frac{16}{20} + \frac{25}{25}}} = \frac{-4.9}{\sqrt{0.04 + 0.04}} = -3.6522$$

independant small samples

assumption: $\sigma_1 = \sigma_2 = \sigma$ (var are equal but unknown)

$T = \frac{\bar{x}_1 - \bar{x}_2}{S_p \sqrt{\frac{1}{n_1} + \frac{1}{n_2}}}$

$S_p^2 = \frac{(n_1-1)s_1^2 + (n_2-1)s_2^2}{(n_1+n_2-2)}$

where $S_p^2 = \frac{(n_1-1)s_1^2 + (n_2-1)s_2^2}{(n_1+n_2-2)}$ $S_p \Rightarrow \text{pooled variance}$

eg: $\bar{x}_1 = 85 \quad s_1 = 4 \quad \bar{x}_2 = 81 \quad s_2 = 5 \quad \alpha = 0.05$
 $n_1 = 12 \quad n_2 = 10$

$H_0: \mu_1 - \mu_2 \leq 2$

$H_1: \mu_1 - \mu_2 > 2$

$d_0 = \mu_1 - \mu_2 = 2$

test statistic $t = \frac{(\bar{x}_1 - \bar{x}_2) - d_0}{S_p \sqrt{\frac{1}{n_1} + \frac{1}{n_2}}}$

with df = $n_1 + n_2 - 2$

= 20

$S_p^2 = \frac{11(16) + 9(25)}{20} = \sqrt{20.05} = 4.478$

$t = \frac{(85-81)-2}{4.478 \sqrt{\frac{1}{12} + \frac{1}{10}}} = \frac{2}{4.478 \sqrt{0.183}} = 1.044$

Critical value approach:

at $1-\alpha = 5\%$, one tailed, $df = 20$, $t_{\alpha} = 1.72$

$|t| < t_{\alpha}$ \therefore do not reject H_0

5. Matched pair comparisons

e.g. $\bar{x}_i \times_i$ = before

y_i = after

$$d_i = x_i - y_i, \bar{d} = \text{mean}(d_i)$$

$$T = \frac{\bar{d} - 0}{S_d / \sqrt{n}} \quad df = n-1$$

e.g. $d_i = 9 \ 13 \ 2 \ 5 \ -2 \ 6 \ 6 \ 5 \ 2 \ 6$

$$\bar{d} = 5.2 \quad S_d = 7.08 \quad H_0: M_d = 0$$

$$\alpha = 0.05$$

$$H_1: M_d < 0$$

$$t = \frac{\bar{d} - 0}{S_d / \sqrt{n}} - \frac{5.2}{7.08 / \sqrt{10}} = 4.0303$$

$$df = 9 \quad \alpha = 0.05 \quad t_{\alpha} = 1.83$$

$$t > t_{\alpha}$$