

AMS 572 Data Analysis I

Inference on one population mean μ

Pei-Fen Kuan

Applied Math and Stats, Stony Brook University

One population setting

Cross-sectional study: Collect data to test a hypothesis about the mean or median of X

- ▶ Take a random sample from the population.
- ▶ There are different types of sampling schemes. The simplest is the simple random sampling in which every subject in the population has equal chance to be selected.

Statistical inference on one population mean

Suppose we have a random sample of size n : X_1, X_2, \dots, X_n and we wish to draw inference about the population mean μ .

1. Point estimator

▶ $\bar{X} = \frac{X_1 + X_2 + \dots + X_n}{n} = \frac{\sum_{i=1}^n X_i}{n}$

▶ Other estimators: median, mode, trimmed mean

2. Confidence Interval (C.I.)

3. Hypothesis Test

▶ Example: μ is the height of adult US males

$$H_0 : \mu \leq 5'6'' \text{ vs } H_1 : \mu > 5'6''$$

Recap: Normal Distribution

- Probability Density Function (p.d.f.)
X follows normal distribution of mean μ and variance σ^2

$$X \sim N(\mu, \sigma^2)$$

$$f(x) = \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{(x-\mu)^2}{2\sigma^2}}, -\infty < x < \infty, x \in R$$

$P(a \leq X \leq b) = \int_a^b f(x)dx$ = area under the pdf curve bounded by a and b

- Cumulative Density Function (c.d.f.)

$$F(x) = P(X \leq x) = \int_{-\infty}^x f(t)dt$$

$$f(x) = [F(x)]' = \frac{d}{dx}F(x)$$

Recap: Normal Distribution

► Standard Normal Distribution

$$Z \sim N(0, 1)$$

$$X \sim N(\mu, \sigma^2)$$

$$Z = \frac{X - \mu}{\sigma} \sim N(0, 1)$$

$$\begin{aligned} P(a \leq X \leq b) &= P\left(\frac{a - \mu}{\sigma} \leq \frac{X - \mu}{\sigma} \leq \frac{b - \mu}{\sigma}\right) \\ &= P\left(\frac{a - \mu}{\sigma} \leq Z \leq \frac{b - \mu}{\sigma}\right) \end{aligned}$$

Also known as **Z-score**

Theorem: If two random variables have the same MGF, they have the same distribution.

Normal Distributions: $X \sim N(\mu, \sigma^2)$

$$M_X(t) = e^{\mu t + \frac{1}{2}\sigma^2 t^2}$$

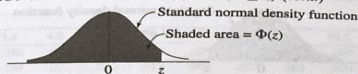
Example: If the mgf of W is $e^{-7t+10t^2}$, what is the distribution of W ?

Linear transformation of a normal random variable

- ▶ If $X \sim N(\mu, \sigma^2)$ and $Y = aX + b$, a, b are constants, then

$$Y \sim N(a\mu + b, a^2 \sigma^2)$$

Example: Suppose the height of students of AMS 572 is normally distributed. Ten percent of the students are over 6.5 feet tall, while the variance is 0.390625 (or 0.625^2). What is the probability that the height of a student is in between 6 and 7 feet?

Table A.3 Standard Normal Curve Areas $\Phi(z) = P(Z \leq z)$ (cont.)

z	0.00	0.01	0.02	0.03	0.04	0.05	0.06	0.07	0.08	0.09
0.0	0.5000	0.5040	0.5080	0.5120	0.5160	0.5199	0.5239	0.5279	0.5319	0.5359
0.1	0.5398	0.5438	0.5478	0.5517	0.5557	0.5596	0.5636	0.5675	0.5714	0.5753
0.2	0.5793	0.5832	0.5871	0.5910	0.5948	0.5987	0.6026	0.6064	0.6103	0.6141
0.3	0.6179	0.6217	0.6255	0.6293	0.6331	0.6368	0.6406	0.6443	0.6480	0.6517
0.4	0.6554	0.6591	0.6628	0.6664	0.6700	0.6736	0.6772	0.6808	0.6844	0.6879
0.5	0.6915	0.6950	0.6985	0.7019	0.7054	0.7088	0.7123	0.7157	0.7190	0.7224
0.6	0.7257	0.7291	0.7324	0.7357	0.7389	0.7422	0.7454	0.7486	0.7517	0.7549
0.7	0.7580	0.7611	0.7642	0.7673	0.7704	0.7734	0.7764	0.7794	0.7823	0.7852
0.8	0.7881	0.7910	0.7939	0.7967	0.7995	0.8023	0.8051	0.8078	0.8106	0.8133
0.9	0.8159	0.8186	0.8212	0.8238	0.8264	0.8289	0.8315	0.8340	0.8365	0.8389
1.0	0.8413	0.8438	0.8461	0.8485	0.8508	0.8531	0.8554	0.8577	0.8599	0.8621
1.1	0.8643	0.8665	0.8686	0.8708	0.8729	0.8749	0.8770	0.8790	0.8810	0.8830
1.2	0.8849	0.8869	0.8888	0.8907	0.8925	0.8944	0.8962	0.8980	0.8997	0.9015
1.3	0.9032	0.9049	0.9066	0.9082	0.9099	0.9115	0.9131	0.9147	0.9162	0.9177
1.4	0.9192	0.9207	0.9222	0.9236	0.9251	0.9265	0.9278	0.9292	0.9306	0.9319
1.5	0.9332	0.9345	0.9357	0.9370	0.9382	0.9394	0.9406	0.9418	0.9429	0.9441
1.6	0.9452	0.9463	0.9474	0.9484	0.9495	0.9505	0.9515	0.9525	0.9535	0.9545
1.7	0.9554	0.9564	0.9573	0.9582	0.9591	0.9599	0.9608	0.9616	0.9625	0.9633
1.8	0.9641	0.9649	0.9656	0.9664	0.9671	0.9678	0.9686	0.9693	0.9699	0.9706
1.9	0.9713	0.9719	0.9726	0.9732	0.9738	0.9744	0.9750	0.9756	0.9761	0.9767
2.0	0.9772	0.9778	0.9783	0.9788	0.9793	0.9798	0.9803	0.9808	0.9812	0.9817
2.1	0.9821	0.9826	0.9830	0.9834	0.9838	0.9842	0.9846	0.9850	0.9854	0.9857
2.2	0.9861	0.9864	0.9868	0.9871	0.9875	0.9878	0.9881	0.9884	0.9887	0.9890
2.3	0.9893	0.9896	0.9898	0.9901	0.9904	0.9906	0.9909	0.9911	0.9913	0.9916
2.4	0.9918	0.9920	0.9922	0.9925	0.9927	0.9929	0.9931	0.9932	0.9934	0.9936
2.5	0.9938	0.9940	0.9941	0.9943	0.9945	0.9946	0.9948	0.9949	0.9951	0.9952
2.6	0.9953	0.9955	0.9956	0.9957	0.9959	0.9960	0.9961	0.9962	0.9963	0.9964
2.7	0.9965	0.9966	0.9967	0.9968	0.9969	0.9970	0.9971	0.9972	0.9973	0.9974
2.8	0.9974	0.9975	0.9976	0.9977	0.9977	0.9978	0.9979	0.9979	0.9980	0.9981
2.9	0.9981	0.9982	0.9982	0.9983	0.9984	0.9984	0.9985	0.9985	0.9986	0.9986
3.0	0.9987	0.9987	0.9987	0.9988	0.9988	0.9989	0.9989	0.9989	0.9990	0.9990
3.1	0.9990	0.9991	0.9991	0.9991	0.9992	0.9992	0.9992	0.9992	0.9993	0.9993
3.2	0.9993	0.9993	0.9994	0.9994	0.9994	0.9994	0.9994	0.9995	0.9995	0.9995
3.3	0.9995	0.9995	0.9995	0.9996	0.9996	0.9996	0.9996	0.9996	0.9996	0.9997
3.4	0.9997	0.9997	0.9997	0.9997	0.9997	0.9997	0.9997	0.9997	0.9997	0.9998

Distribution of sample mean \bar{X}

Theorem: Let X_1, X_2, \dots, X_n be a random sample from a normal population with mean μ , variance σ^2 . Then, the distribution of \bar{X} is

$$\bar{X} \sim N\left(\mu, \frac{\sigma^2}{n}\right)$$

Example (cont'd): If 10 students are chosen, what is the probability that the average height is in between 6 and 7 feet?

Inference for μ when σ^2 is known under normal distribution

Inference for μ when σ^2 is known under normal distribution

Setup:

- ▶ Assume that the distribution is normal
- ▶ Let X_1, X_2, \dots, X_n be a random sample for a normal distribution with mean μ and variance σ^2 . That is,
$$X \stackrel{iid.}{\sim} N(\mu, \sigma^2), i = 1, \dots, n.$$
- ▶ Assume that σ^2 is known.

Inference for μ when σ^2 is known under normal distribution

1. Point estimator for μ

$$\hat{\mu} = \bar{X} \sim N\left(\mu, \frac{\sigma^2}{n}\right)$$

Note that $E(\hat{\mu}) = E(\bar{X}) = \mu \Rightarrow \hat{\mu} = \bar{X}$ is an unbiased estimator of μ

- ▶ \bar{X} is also a maximum likelihood estimator (MLE) and method of moment estimator (MME) of μ .

Inference for μ when σ^2 is known under normal distribution

2 Confidence Interval for μ :

Intuitive approach : $P(c_1 \leq \mu \leq c_2) = 0.95$

$$P(\bar{X} - c_1 \geq \bar{X} - \mu \geq \bar{X} - c_2) = 0.95$$

$$\bar{X} - \mu \sim N(0, \frac{\sigma^2}{n})$$

$$P(\frac{\bar{X} - c_1}{\sigma/\sqrt{n}} \geq \frac{\bar{X} - \mu}{\sigma/\sqrt{n}} \geq \frac{\bar{X} - c_2}{\sigma/\sqrt{n}}) = 0.95$$

There ARE many ways to choose the c 's.

Pivotal quantity with symmetric pdf, the symmetric CIs are optimal, i.e., they have the shortest lengths for a given confidence level.

Pivotal Quantity (P.Q.) approach for deriving confidence interval

Definition: A pivotal quantity is a function of the sample and parameter of interest whose probability distribution does not depend on the unknown parameters.

Pivotal Quantity (P.Q.) approach for deriving confidence interval

Consider $\bar{X} \sim N(\mu, \frac{\sigma^2}{n})$, where σ^2 is known

- ▶ Is \bar{X} a pivotal quantity for μ ?
- ▶ Function of \bar{X} and μ : $\bar{X} - \mu \sim N(0, \frac{\sigma^2}{n})$
- ▶ Another function of \bar{X} and μ : $Z = \frac{\bar{X} - \mu}{\sigma/\sqrt{n}}$

Derivation for the symmetrical CI's for μ based on pivotal quantity Z

CI for μ , $0 < \alpha < 1$ (e.g. $\alpha = 0.05 \Rightarrow 95\%$ C.I.)

$$P(-z_{\alpha/2} \leq Z \leq z_{\alpha/2}) = 1 - \alpha$$

$$P(-z_{\alpha/2} \leq \frac{\bar{X} - \mu}{\sigma/\sqrt{n}} \leq z_{\alpha/2}) = 1 - \alpha$$

$$P(-z_{\alpha/2} \cdot \frac{\sigma}{\sqrt{n}} \leq \bar{X} - \mu \leq z_{\alpha/2} \cdot \frac{\sigma}{\sqrt{n}}) = 1 - \alpha$$

$$P(\bar{X} - z_{\alpha/2} \cdot \frac{\sigma}{\sqrt{n}} \leq \mu \leq \bar{X} + z_{\alpha/2} \cdot \frac{\sigma}{\sqrt{n}}) = 1 - \alpha$$

Thus the $100(1 - \alpha)\%$ C.I. for μ is $[\bar{X} - z_{\alpha/2} \cdot \frac{\sigma}{\sqrt{n}}, \bar{X} + z_{\alpha/2} \cdot \frac{\sigma}{\sqrt{n}}]$

Confidence Interval Interpretation

- If we draw 100 different random samples, on average $100(1 - \alpha)\%$ of them will contain μ

Sample $n \longleftrightarrow$ the corresponding $100(1-\alpha)\%$ C.I.

Sample 1, $\bar{x}=5'7'' \longleftrightarrow [5'5'', 5'9'']$

Sample 2, $\bar{x}=5'5'' \longleftrightarrow [5'3'', 5'7'']$

Sample 3, $\bar{x}=5'8'' \longleftrightarrow [5'6'', 5'10'']$

... e.g.) 95% C.I., $\mu=5'7''$, 95% of all these CI's will cover μ

Confidence Interval Interpretation

- ▶ For the $100(1-\alpha)\%$ C.I. for μ is $[\bar{X} - z_{\alpha/2} \cdot \frac{\sigma}{\sqrt{n}}, \bar{X} + z_{\alpha/2} \cdot \frac{\sigma}{\sqrt{n}}]$, the length of this CI is:

$$L_{sy} = 2 \cdot z_{\alpha/2} \cdot \frac{\sigma}{\sqrt{n}}$$

- ▶ To decrease the length of confidence interval:
 - ▶ increase α , i.e., decrease confidence
 - ▶ increase sample size

Non-symmetric confidence interval

- ▶ Note that $P(-z_{\alpha/3} \leq Z \leq z_{2\alpha/3}) = 1 - \alpha$
- ▶ 100(1- α)% C.I. for μ

$$[\bar{X} - z_{2\alpha/3} \cdot \frac{\sigma}{\sqrt{n}}, \bar{X} + z_{\alpha/3} \cdot \frac{\sigma}{\sqrt{n}}]$$

The length of this CI is:

$$L_{nsy} = (z_{\alpha/3} + z_{2\alpha/3}) \cdot \frac{\sigma}{\sqrt{n}}$$

- ▶ It can be shown that: $L_{sy} \leq L_{nsy}$
- ▶ Try a few numerical values for α , and see for yourself.

Hypothesis testing

- ▶ Null hypothesis H_0
- ▶ Example of null hypothesis for folic acid study:
The incidence of stroke in the current study cohort will be the same as the reported population incidence last year.

Null and Alternative

- ▶ In a test of a hypothesis, we are testing whether some population parameter has a particular value or range
- ▶ For example,

$$H_0 : \mu \leq \mu_0$$

where μ_0 is a known constant

- ▶ The **alternative hypothesis** is complement of null hypothesis

$$H_a : \mu > \mu_0$$

Null and Alternative

- ▶ A law suit analogy:
Example: The OJ Simpson trial
 H_0 : OJ is innocent
 H_a : OJ is guilty

		The truth	
		H_0 : OJ innocent	H_a : OJ guilty
Jury's Decision	H_0	Right decision	Type II error
	H_a	Type I error	Right decision (Power)

Tests of Hypotheses: Seven Steps

1. Design study
2. Establish null hypothesis
3. Determine test statistic to be employed
4. Choose significance level α and establish critical region C_α .
 α is also $P(\text{Type I error})$.
5. Carry out study and collect data
6. Compute statistic from data
7. If statistic is in C_α , reject H_0

Types of hypothesis

One-sided (or one-tailed) test:

- ▶ $H_0 : \mu = \mu_0$ vs $H_a : \mu > \mu_0$
 $H_0 : \mu \leq \mu_0$ vs $H_a : \mu > \mu_0$
- ▶ $H_0 : \mu = \mu_0$ vs $H_a : \mu < \mu_0$
 $H_0 : \mu \geq \mu_0$ vs $H_a : \mu < \mu_0$

Two-sided (or two-tailed) test:

- ▶ $H_0 : \mu = \mu_0$ vs $H_a : \mu \neq \mu_0$

Hypothesis test for $H_0 : \mu = \mu_0$ vs $H_a : \mu > \mu_0$

Setting

- ▶ Data : $X_1, X_2, \dots, X_n \stackrel{iid}{\sim} N(\mu, \sigma^2)$, σ^2 is known
- ▶ The significance level α (e.g., $\alpha = 0.05$) is given.

Hypothesis test for $H_0 : \mu = \mu_0$ vs $H_a : \mu > \mu_0$

Pivotal Quantity Method

1. Identify a pivotal quantity: Recall that

$$Z = \frac{\bar{X} - \mu}{\sigma / \sqrt{n}} \sim N(0, 1) \text{ is a pivotal quantity.}$$

2. The **test statistic** is the pivotal quantity with the value of the parameter of interest under the null hypothesis

Hypothesis test for $H_0 : \mu = \mu_0$ vs $H_a : \mu > \mu_0$

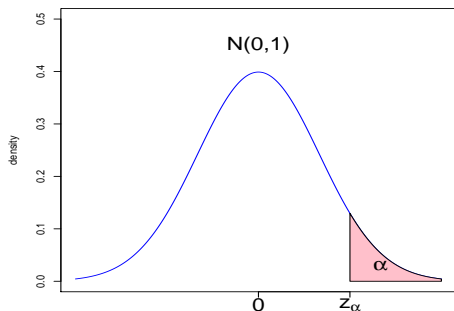
- 3 Derive the decision threshold for your test based on the Type I error rate, that is, the significance level α

For the pair of hypotheses: $H_0 : \mu = \mu_0$ versus $H_a : \mu > \mu_0$

It is intuitive that one should reject H_0 , in support of the H_a , when the $\bar{X} > \mu_0$. Equivalently, this means to reject H_0 when the test statistic Z_0 is larger than certain positive value c ($Z_0 > c$).

The question is what is the exact value of c , which can be determined based on the significance level α , i.e., how much Type I error we would allow ourselves to commit.

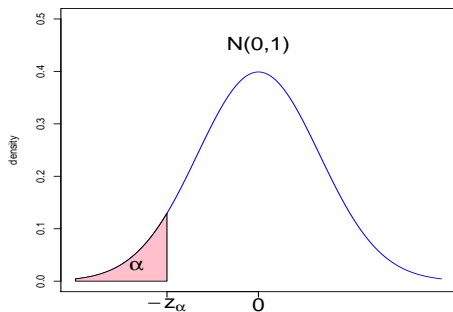
Hypothesis test for $H_0 : \mu = \mu_0$ vs $H_a : \mu > \mu_0$



\therefore At the significance level α , we will reject H_0 in favor of H_a if $Z_0 \geq z_\alpha$

Hypothesis test for $H_0 : \mu = \mu_0$ versus $H_a : \mu < \mu_0$

Test statistic: $Z_0 = \frac{\bar{X} - \mu_0}{\sigma/\sqrt{n}} \sim N(0, 1)$

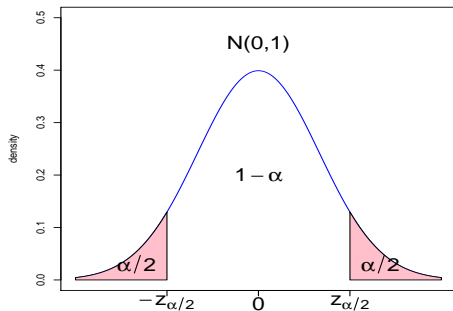


\therefore At the significance level α , we will reject H_0 in favor of H_a if $Z_0 \leq -z_\alpha$

Hypothesis test for $H_0 : \mu = \mu_0$ versus $H_a : \mu \neq \mu_0$

Test statistic : $Z_0 = \frac{\bar{X} - \mu_0}{\sigma/\sqrt{n}} \sim N(0, 1)$

Hypothesis test for $H_0 : \mu = \mu_0$ versus $H_a : \mu \neq \mu_0$



P-value

- ▶ We have just discussed the “**rejection/critical region**” approach for decision making.
- ▶ There is another approach for decision making, it is the “**p-value**” approach.

Definition: p-value is the probability that of observing a test statistic value that is as extreme, or more extreme, than the one we observed.

- ▶ Reject H_0 in favor of H_a if $\text{p-value} \leq \alpha$.
For example, for $H_0 : \mu = \mu_0$ versus $H_a : \mu > \mu_0$, the p-value is $P(Z_0 \geq z_0 | H_0)$.

Summary

$H_0 : \mu = \mu_0$ $H_a : \mu > \mu_0$	$H_0 : \mu = \mu_0$ $H_a : \mu < \mu_0$	$H_0 : \mu = \mu_0$ $H_a : \mu \neq \mu_0$
Observed value of test statistic $Z_0 = \frac{\bar{X} - \mu_0}{\sigma/\sqrt{n}} \stackrel{H_0}{\sim} N(0, 1)$		
Rejection region : we reject H_0 in favor of H_a at the significance level α if		
$Z_0 \geq z_\alpha$	$Z_0 \leq -z_\alpha$	$ Z_0 \geq z_{\alpha/2}$
p-value = $P(Z_0 \geq z_0 H_0)$	p-value = $P(Z_0 \leq z_0 H_0)$	p-value $= P(Z_0 \geq z_0 H_0)$ $= 2 * P(Z_0 \geq z_0 H_0)$
the area under $N(0, 1)$ pdf to the right of z_0	the area under $N(0, 1)$ pdf to the left of z_0	twice the area to the right of $ z_0 $