

Confidence Intervals II

18.05 Spring 2014

Agenda

- Polling: estimating θ in Bernoulli(θ).
- CLT \Rightarrow large sample confidence intervals for the mean.
- Three views of confidence intervals.
- Constructing a confidence interval without normality:
the exact binomial confidence interval for θ

Polling confidence interval

Also called a **binomial proportion confidence interval**

Polling means sampling from a Bernoulli(θ) distribution, i.e. data x_1, \dots, x_n Bernoulli(θ).

- **Consevative normal** confidence interval for θ :

$$\bar{x} \pm z_{\alpha/2} \cdot \frac{1}{2\sqrt{n}}$$

Proof uses the CLT and the observation $\sigma = \sqrt{\theta(1-\theta)} \leq 1/2$.

- **Rule-of-thumb 95%** confidence interval for θ :

$$\bar{x} \pm \frac{1}{\sqrt{n}}$$

(Reason: $z_{0.025} \approx 2$.)

Board question

For a poll to find the proportion θ of people supporting X we know that a $(1 - \alpha)$ confidence interval for θ is given by

$$\left[\bar{x} - \frac{z_{\alpha/2}}{2\sqrt{n}}, \bar{x} + \frac{z_{\alpha/2}}{2\sqrt{n}} \right].$$

1. How many people would you have to poll to have a margin of error of 0.01 with 95% confidence? (You can do this in your head.)
2. How many people would you have to poll to have a margin of error of 0.01 with 80% confidence. (You'll want R or other calculator here.)
3. If $n = 900$, compute the 95% and 80% confidence intervals for θ .

Concept question: overnight polling

During the presidential election season, pollsters often do 'overnight polls' and report a 'margin of error' of about $\pm 5\%$.

The number of people polled is in which of the following ranges?

- (a) 0 – 50
- (b) 50 – 100
- (c) 100 – 300
- (d) 300 – 600
- (e) 600 – 1000

National Council on Public Polls: Press Release, Sept 1992

“The National Council on Public Polls expressed concern today about the current spate of overnight Presidential polls. [...] Overnight polls do a disservice to both the media and the research industry because of the considerable potential for the results to be misleading. The overnight interviewing period may well mean some methodological compromises, the most serious of which is..”

...what?

“...the inability to make callbacks, resulting in samples that do not adequately represent such groups as single member households, younger people, and others who are apt to be out on any given night. As overnight polls often result in findings that are less reliable than those from more carefully conducted polls, if the media reports them, it should be with great caution.”

<http://www.ncpp.org/?q=node/42>

Large sample confidence interval

Data x_1, \dots, x_n independently drawn from a distribution that may not be normal but has finite mean and variance.

A version of the central limit theorem says that large n ,

$$\frac{\bar{x} - \mu}{s/\sqrt{n}} \approx N(0, 1)$$

i.e. the sampling distribution of the studentized mean is approximately standard normal:

So for large n the $(1 - \alpha)$ confidence interval for μ is approximately

$$\left[\bar{x} - \frac{s}{\sqrt{n}} \cdot z_{\alpha/2}, \bar{x} + \frac{s}{\sqrt{n}} \cdot z_{\alpha/2} \right]$$

This is called the **large sample confidence interval**.

Review: confidence intervals for normal data

Suppose the data x_1, \dots, x_n is drawn from $N(\mu, \sigma^2)$

Confidence level = $1 - \alpha$

- z confidence interval for the mean (σ known)

$$\left[\bar{x} - \frac{z_{\alpha/2} \cdot \sigma}{\sqrt{n}}, \bar{x} + \frac{z_{\alpha/2} \cdot \sigma}{\sqrt{n}} \right] \quad \text{or} \quad \bar{x} \pm \frac{z_{\alpha/2} \cdot \sigma}{\sqrt{n}}$$

- t confidence interval for the mean (σ unknown)

$$\left[\bar{x} - \frac{t_{\alpha/2} \cdot s}{\sqrt{n}}, \bar{x} + \frac{t_{\alpha/2} \cdot s}{\sqrt{n}} \right] \quad \text{or} \quad \bar{x} \pm \frac{t_{\alpha/2} \cdot s}{\sqrt{n}}$$

- χ^2 confidence interval for σ^2

$$\left[\frac{n-1}{c_{\alpha/2}} s^2, \frac{n-1}{c_{1-\alpha/2}} s^2 \right]$$

- t and χ^2 have $n - 1$ degrees of freedom.

Three views of confidence intervals

View 1: Define/construct CI using a standardized point statistic.

View 2: Define/construct CI based on hypothesis tests.

View 3: Define CI as any interval statistic satisfying a formal mathematical property.

View 1: Using a standardized point statistic

Example. $x_1, \dots, x_n \sim N(\mu, \sigma^2)$, where σ is known.

The **standardized sample mean** follows a standard normal distribution.

$$z = \frac{\bar{x} - \mu}{\sigma/\sqrt{n}} \sim N(0, 1)$$

Therefore:

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

Pivot to:

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

This is the $(1 - \alpha)$ confidence interval:

$$\bar{x} \pm z_{\alpha/2} \cdot \frac{\sigma}{\sqrt{n}}$$

Think of it as $\bar{x} \pm \text{error}$

View 1: Other standardized statistics

The t and χ^2 statistics fit this paradigm as well:

$$t = \frac{\bar{x} - \mu}{s/\sqrt{n}} \sim t(n-1)$$

$$\chi^2 = \frac{(n-1)s^2}{\sigma^2} \sim \chi^2(n-1)$$

View 2: Using hypothesis tests

Set up: Unknown parameter θ . Test statistic x .

For any value θ_0 , we can run an NSHT with null hypothesis

$$H_0 : \theta = \theta_0$$

at significance level α .

Definition. Given x , the $(1 - \alpha)$ confidence interval contains all θ_0 which are not rejected when they are the null hypothesis.

Definition. A type 1 CI error occurs when the confidence interval does not contain the true value of θ .

For a $1 - \alpha$ confidence interval, the type 1 CI error rate is α .

Board question: exact binomial confidence interval

Use this table of binomial($8, \theta$) probabilities to:

- 1 find the (two-sided) rejection region with significance level 0.10 for each value of θ .
- 2 Given $x = 7$, find the 90% confidence interval for θ .
- 3 Repeat for $x = 4$.

θ/x	0	1	2	3	4	5	6	7	8
.1	0.430	0.383	0.149	0.033	0.005	0.000	0.000	0.000	0.000
.3	0.058	0.198	0.296	0.254	0.136	0.047	0.010	0.001	0.000
.5	0.004	0.031	0.109	0.219	0.273	0.219	0.109	0.031	0.004
.7	0.000	0.001	0.010	0.047	0.136	0.254	0.296	0.198	0.058
.9	0.000	0.000	0.000	0.000	0.005	0.033	0.149	0.383	0.430

Solution

For each θ , the non-rejection region is blue, the rejection region is red. In each row, the rejection region has probability at most $\alpha = 0.10$.

θ/x	0	1	2	3	4	5	6	7	8
.1	0.430	0.383	0.149	0.033	0.005	0.000	0.000	0.000	0.000
.3	0.058	0.198	0.296	0.254	0.136	0.047	0.010	0.001	0.000
.5	0.004	0.031	0.109	0.219	0.273	0.219	0.109	0.031	0.004
.7	0.000	0.001	0.010	0.047	0.136	0.254	0.296	0.198	0.058
.9	0.000	0.000	0.000	0.000	0.005	0.033	0.149	0.383	0.430

For $x = 7$ the 90% confidence interval for p is $[0.7, 0.9]$.

These are the values of θ we wouldn't reject as null hypotheses. They are the blue entries in the $x = 7$ column.

For $x = 4$ the 90% confidence interval for p is $[0.3, 0.7]$.

View 3: Formal

Recall: An interval statistic is an interval I_x computed from data x .

This is a random interval because x is random.

Suppose x is drawn from $f(x|\theta)$ with unknown parameter θ .

Definition:

A $(1 - \alpha)$ confidence interval for θ is an interval statistic I_x such that

$$P(I_x \text{ contains } \theta \mid \theta) = 1 - \alpha$$

for all possible values of θ (and hence for the true value of θ).

Note: equality in this equation is often relaxed to \geq or \approx .

$=$: z , t , χ^2

\geq : rule-of-thumb and exact binomial (polling)

\approx : large sample confidence interval

MIT OpenCourseWare
<https://ocw.mit.edu>

18.05 Introduction to Probability and Statistics

Spring 2014

For information about citing these materials or our Terms of Use, visit: <https://ocw.mit.edu/terms>.