

UCLA STAT 110A Applied Statistics

- **Instructor:** Ivo Dinov,
Asst. Prof. In Statistics and Neurology
 - **Teaching Assistant:** Helen Hu, UCLA Statistics
- University of California, Los Angeles, Spring 2002
<http://www.stat.ucla.edu/~dinov/>

STAT 110A, UCLA, Ivo Dinov

Slide 1

Inference & Estimation

- C + E model
- Types of Inference
- Sampling distributions
- CI's for μ & p
- Comparing 2 proportions
- How big should my study be?
- Paired vs. unpaired tests

STAT 110A, UCLA, Ivo Dinov

Slide 2

The C + E Model

- **Data = Center + Error : $Y = \mu + \epsilon$;**
- The response value Y is equal to unknown constant (μ), but because of normal variability we almost never observe μ exactly.
- Example **Speed of light (SOL)**, $\mu = 2.998 \times 10^9$ m/s. However, 100 measurements of the SOL are all going to be slightly different.
- **Model (population) parameter** – a quantity describing the model that can take on many values. Ex., μ .

Slide 3

STAT 110A, UCLA, Ivo Dinov

Types of inference

- **Estimation of model parameters:** Data-driven estimates of the model parameters. Also, includes how much uncertainty about those estimates is there.
- **Prediction of new (future) observations:** Uses past and current data to predict the value of new observations from the population.
- **Tolerance level:** a range of values that has user-specified probability of containing a particular proportion of the population.

Slide 4

STAT 110A, UCLA, Ivo Dinov

Estimation of model parameter(s) – μ

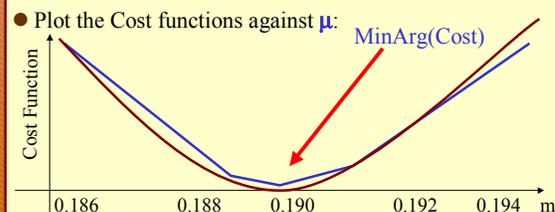
- **Least-Absolute-Error-Estimate(m)** – Suppose, $\mu = 3.5$ (unknown) and $Y = \{Y_1 = \mu + e_1, Y_2 = \mu + e_2, \dots, Y_{10} = \mu + e_{10}\}$ are our observed data. **Cost function** = Sum-of-Absolute-Errors = $SAE = \sum |Y_k - m| \rightarrow m = \text{MinArg}(SAE)$.
- **Least-Squares(m)** (in the same setting). Cost function = Sum-of-Squared-Errors = $SSE = \sum (Y_k - m)^2 \rightarrow m = \text{MinArg}(SSE)$, least-squares-estimate.
- **Solution (differentiate):**
 $d SSE(m) / d m = -2 \sum (Y_k - m) = 0$, **solve for m!**

Slide 5

STAT 110A, UCLA, Ivo Dinov

Estimation of model parameter(s) – μ (Example)

- **Data:** ball-bearing diameter: $\mu = ?$ (unknown) given the observed $Y = \{Y_1 = 0.1896, Y_2 = 0.1913, Y_{10} = 0.1900\}$.
 $SAE = \sum |Y_k - m|$ & $SSE = \sum (Y_k - m)^2$



Slide 6

STAT 110A, UCLA, Ivo Dinov

Parameters, Estimators, Estimates ...

- A **parameter** is a characteristic of the data – mean, 1st quartile, SD, etc.)
- An **estimator** is an abstract **rule** for calculating a quantity (or parameter) from the sample data.
- An **estimate** is the value obtained when real data are plugged-in the estimator rule.

Slide 7 STAT 1104, UCL, Lee Dinger

Parameters, Estimators, Estimates ...

- E.g., We are interested in the **population mean diameter (parameter)** of washers the **sample-average formula** represents an **estimator** we can use, where as the **value of the sample average** for a particular dataset is the **estimate** (for the **mean** parameter).

parameter = μ_y ; estimator = $\bar{Y} = \frac{1}{N} \sum_{k=1}^N Y_k$

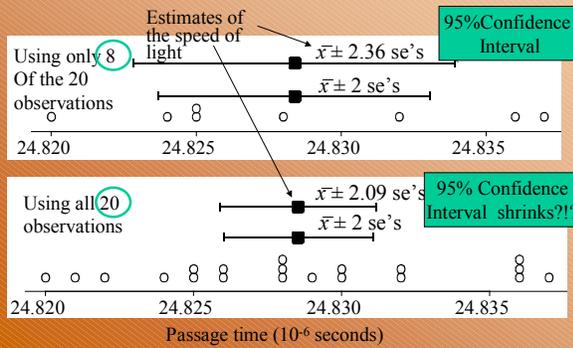
Data : $Y = \{0.1896, 0.1913, 0.1900\}$

estimate = $\bar{y} = \frac{1}{3} (0.1896 + 0.1913 + 0.1900)$

$\bar{y} = 0.1903$. How about $\bar{y} = \frac{2}{3} (0.1896 + 0.1913 + 0.1900)$

Slide 8 STAT 1104, UCL, Lee Dinger

20 replicated measurements to estimate the speed of light. Obtained by Simon Newcomb in 1882, by using distant (3.721 km) rotating mirrors.



Slide 9 STAT 1104, UCL, Lee Dinger

A 95% confidence interval

- A type of interval that contains the **true value of a parameter** for 95% of samples taken is called a **95% confidence interval** for that parameter, the ends of the CI are called **confidence limits**.
- (For the situations we deal with) a **confidence interval (CI)** for the true value of a **parameter** is given by **estimate $\pm t$ standard errors (SE)**

Value of the Multiplier, t , for a 95% CI

df :	7	8	9	10	11	12	13	14	15	16	17
t :	2.365	2.306	2.262	2.228	2.201	2.179	2.160	2.145	2.131	2.120	2.110
df :	18	19	20	25	30	35	40	45	50	60	∞
t :	2.101	2.093	2.086	2.060	2.042	2.030	2.021	2.014	2.009	2.000	1.960

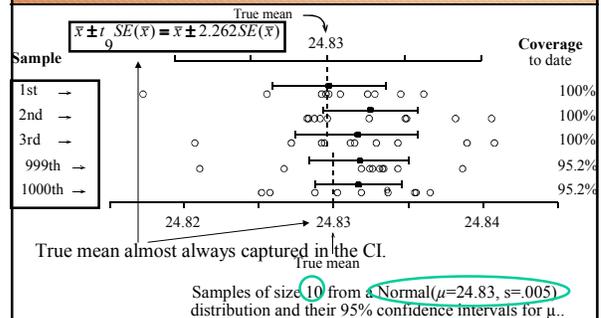
Slide 10 STAT 1104, UCL, Lee Dinger

(General) Confidence Interval (CI)

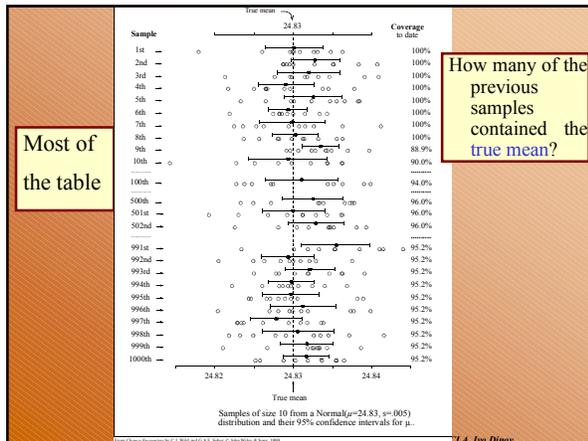
- A **level L confidence interval** for a parameter (θ), is an interval $(\theta_1^{\wedge}, \theta_2^{\wedge})$, where θ_1^{\wedge} & θ_2^{\wedge} , are estimators of θ , such that **$P(\theta_1^{\wedge} < \theta < \theta_2^{\wedge}) = L$** .
- E.g., **C+E model**: $Y = \mu + \epsilon$. Where $\epsilon \sim N(0, \sigma^2)$, then by CLT we have $\bar{Y} \sim N(\mu, \sigma^2/n)$
 $\rightarrow n^{1/2}(\bar{Y} - \mu) / \sigma \sim N(0, \sigma^2)$.
- **$L = P(z_{(1-L)/2} < n^{1/2}(\bar{Y} - \mu) / \sigma < z_{(1+L)/2})$** , where z_q is the q^{th} quartile.
- E.g., **$0.95 = P(z_{0.025} < n^{1/2}(\bar{Y} - \mu) / \sigma < z_{0.975})$** ,

Slide 11 STAT 1104, UCL, Lee Dinger

- CI are constructed using the sample \bar{x} and $s = SE$. But **different samples yield different estimates** and \rightarrow diff. CI's?!?
- Below is a computer simulation showing how the process of taking samples effects the estimates and the CI's.



Slide 12 STAT 1104, UCL, Lee Dinger



CI for population mean

Confidence Interval for the true (population) mean μ :

sample mean \pm *t standard errors*

or $\bar{x} \pm t \text{ se}(\bar{x})$, where $\text{SE}(\bar{x}) = \frac{s_x}{\sqrt{n}}$ and $df = n - 1$

Value of the Multiplier, t , for a 95% CI											
df:	7	8	9	10	11	12	13	14	15	16	17
t:	2.365	2.306	2.262	2.228	2.201	2.179	2.160	2.145	2.131	2.120	2.110
df:	18	19	20	25	30	35	40	45	50	60	∞
t:	2.101	2.093	2.086	2.060	2.042	2.030	2.021	2.014	2.009	2.000	1.960

Slide 14 STAT 110A, UCLA, Jon Dinger

CI for population mean

- E.g., SYSTAT \rightarrow Data: **BirthDayDistribution_1978_systat.SYD**
- Statistics \rightarrow Descriptive Statistics \rightarrow Stem-&-Leaf-Plot
- Statistics \rightarrow Descriptive Statistics \rightarrow CI_for_mean

Slide 15 STAT 110A, UCLA, Jon Dinger

CI for population mean - Example

- E.g., Lab rats blood glucose levels: {266, 149, 161, 220}
- Estimate μ , the mean population blood sugar level.
- Assume the variance $\sigma^2=2958$, $\rightarrow \sigma=54.4$, from prior experience. Also assume data comes from $N(\mu, \sigma^2)$.
- Sample-avg=199, Compute the 95% CI, $L=0.95$.
- $(1-L)/2 = 0.025$, $(1+L)/2 = 0.975$,
- $Z_{(1-L)/2} = Z_{0.025} = -1.96$ & $Z_{(1+L)/2} = Z_{0.975} = 1.96$
- $L = P(z_{(1-L)/2} < n^{1/2}(Y_{\text{bar}} - \mu)/\sigma < z_{(1+L)/2})$,
- $CI(\mu) = (Y_{\text{bar}} - \sigma z_{(1+L)/2}/n^{1/2}; Y_{\text{bar}} - \sigma z_{(1-L)/2}/n^{1/2})$
- $CI(\mu) = (199 - 54.4 \times 1.96 / 4^{1/2}; 199 + 54.4 \times 1.96 / 4^{1/2})$
- $CI(\mu) = (145.7 : 252.3)$

Slide 16 STAT 110A, UCLA, Jon Dinger

CI - Interpretation

- Consider taking all possible samples from the population with parameter of interest (θ).

- Suppose we construct the **level L confidence interval** for a parameter (θ) for each sample. Then a proportion L of all constructed CI's will contain the value of θ .
- Note that this interpretation of CI's is in terms of **repeated sampling** from the same population ...

Slide 17 STAT 110A, UCLA, Jon Dinger

Effect of increasing the confidence level

99% CI, $\bar{x} \pm 2.576 \text{ se}(\bar{x})$

95% CI, $\bar{x} \pm 1.960 \text{ se}(\bar{x})$

90% CI, $\bar{x} \pm 1.645 \text{ se}(\bar{x})$

80% CI, $\bar{x} \pm 1.282 \text{ se}(\bar{x})$

Why?

The greater the confidence level, the wider the interval

from Chance Encounters by C.J. Wild and G.A.F. Seiber, © John Wiley & Sons, 2000.

Slide 18 STAT 110A, UCLA, Jon Dinger

Effect of increasing the sample size

$n = 10$ data points
 $n = 40$ data points
 $n = 90$ data points

Passage time

24.82 24.83 24.84

Three random samples from a Normal($\mu=24.83$, $s=.005$) distribution and their 95% confidence intervals for μ .

Increase Sample Size Decreases the size of the CI

from *Chance Encounters* by C.J. Wild and G.A.F. Seber, © John Wiley & Sons, 2000.

To **double the precision** we need **four times** as many observations.

Slide 19 STAT 110A, UCLA, Jon Dinger

Why \uparrow in sample-size \downarrow CI?

Confidence Interval for the true (population) mean μ :

sample mean $\pm t$ standard errors

or $\bar{x} \pm t \text{ se}(\bar{x})$, where $\text{se}(\bar{x}) = \frac{s}{\sqrt{n}}$ and $df = n - 1$

Slide 20 STAT 110A, UCLA, Jon Dinger

Student's t -distribution

- For random samples from a Normal distribution,

$$T = \frac{(\bar{X} - \mu)}{SE(\bar{X})}$$

Recall that for samples from $N(\mu, \sigma)$

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

is exactly distributed as Student($df = n - 1$) ← Approx/Exact Distributions

- but methods we shall base upon this distribution for T work well even for small samples sampled from distributions which are **quite non-Normal**.
- df is number of observations $- 1$, **degrees of freedom**.

Slide 21 STAT 110A, UCLA, Jon Dinger

Density curves for Student's t

Student(df) density curves for various df .

Slide 22 STAT 110A, UCLA, Jon Dinger

Notation

- By $t_{df}(prob)$, we mean the number t such that when $T \sim \text{Student}(df)$, $P(T \geq t_{df}) = prob$; that is, the **tail area above t** (that is to the right of t on the graph) is $prob$.

Normal(0,1) density

$z(prob)$

Student(df) density

$t_{df}(prob)$

The $z(prob)$ and $t(prob)$ notations.

from *Chance Encounters* by C.J. Wild and G.A.F. Seber, © John Wiley & Sons, 2000.

Slide 23 STAT 110A, UCLA, Jon Dinger

The central 90% of the Student(df) distribution.

from *Chance Encounters* by C.J. Wild and G.A.F. Seber, © John Wiley & Sons, 2000.

Slide 24 STAT 110A, UCLA, Jon Dinger

Reading Student's t table

Desired upper-tail prob

Extracts from the Student's t-Distribution Table

df	.20	.15	.10	.05	.025	.01	.005	.001	.0005	.0001
6	0.906	1.134	1.440	1.943	2.447	3.143	3.707	5.208	5.959	8.025
7	0.896	1.119	1.415	1.895	2.365	2.998	3.499	4.785	5.408	7.063
8	0.889	1.108	1.397	1.860	2.306	2.896	3.355	4.501	5.041	6.442
...
10	0.879	1.093	1.372	1.812	2.228	2.764	3.169	4.144	4.587	5.694
...
15	0.866	1.074	1.341	1.753	2.131	2.602	2.947	3.733	4.073	4.880
...
∞	0.842	1.036	1.282	1.645	1.960	2.326	2.576	3.090	3.291	3.719

Desired df

t-value

Slide 25 STAT 1104, UCL4, Jon Dinger

Comparison of the CI using T (unknown σ) & Z (known σ) distributions

- For the old data: glucose levels: $\{266, 149, 161, 220\}$ $\hat{\sigma} = \sqrt{\frac{1}{N-1} \sum_{k=1}^N (y_k - \bar{y})^2}$
- CI(μ), when σ is unknown (T-distr.), small-sample-size, and data comes from (approx.) Normal distribution. $\bar{x} = 199$
 $\hat{\sigma} = 54.39$

$L = P(t_{N-1, (1-L)/2} < n^{1/2}(\bar{Y}_{bar} - \mu) / \hat{\sigma} < t_{N-1, (1+L)/2})$

CI(μ) = $(\bar{Y}_{bar} - \hat{\sigma} \cdot t_{N-1, (1+L)/2} / n^{1/2}; \bar{Y}_{bar} + \hat{\sigma} \cdot t_{N-1, (1+L)/2} / n^{1/2})$

95% CI(μ) = $(199 - 54.39 \times 3.18 / 4^{1/2}; 199 + 54.39 \times 3.18 / 4^{1/2})$

$t_{N-1, (1+L)/2} = t_{3, 0.975} = 3.18$ & $t_{N-1, (1-L)/2} = t_{3, 0.025} = -3.18 \rightarrow CI_T(\mu) = (112.4; 285.6)$

Slide 26 STAT 1104, UCL4, Jon Dinger

Comparison of the CI using T (unknown σ) & Z (known σ) distributions

- CI(μ), when $\sigma = 54.4$ is known (Normal distr.)
CI(μ) = $(\bar{Y}_{bar} - \sigma \cdot z_{(1+L)/2} / n^{1/2}; \bar{Y}_{bar} + \sigma \cdot z_{(1+L)/2} / n^{1/2})$
 $z_{(1+L)/2} = 1.96$

95% CI(μ) = $(199 - 54.4 \times 1.96 / 4^{1/2}; 199 + 54.4 \times 1.96 / 4^{1/2})$
CI_Z(μ) = (145.7 : 252.3)

- Comparison:
CI_T(μ) = (112.4; 285.6) \leftarrow compare \rightarrow
CI_Z(μ) = (145.7; 252.3)
Which one is better?!? More appropriate?!?

Slide 27 STAT 1104, UCL4, Jon Dinger

Prediction vs. Confidence intervals

- Confidence Intervals (for the population mean μ):
 $(\bar{Y} - \frac{\hat{\sigma} \times t_{n-1, (1+L)/2}}{\sqrt{n}}; \bar{Y} + \frac{\hat{\sigma} \times t_{n-1, (1+L)/2}}{\sqrt{n}})$
- Prediction Intervals: L-level prediction interval (PI) for a new value of the process Y is defined by:
 $(\hat{Y}_{new} - \hat{\sigma} \times t_{n-1, (1+L)/2}; \hat{Y}_{new} + \hat{\sigma} \times t_{n-1, (1+L)/2})$
where the predicted value $\hat{Y}_{new} = \bar{Y}$, is obtained as an estimator of the unknown process mean μ .

Slide 29 STAT 1104, UCL4, Jon Dinger

Prediction vs. Confidence intervals – Differences?

- Confidence Intervals (for the population mean μ):
 $(\bar{Y} - \frac{\hat{\sigma} \times t_{n-1, (1+L)/2}}{\sqrt{n}}; \bar{Y} + \frac{\hat{\sigma} \times t_{n-1, (1+L)/2}}{\sqrt{n}})$
 $\hat{\sigma} = \hat{\sigma}(\bar{Y}) = \sqrt{\frac{1}{n-1} \sum_{k=1}^n (y_k - \bar{y})^2}$ **Which SD is bigger?!?**
- Prediction Intervals:
 $(\hat{Y}_{new} - \hat{\sigma} \times t_{n-1, (1+L)/2}; \hat{Y}_{new} + \hat{\sigma} \times t_{n-1, (1+L)/2})$ where $\hat{Y}_{new} = \bar{Y}$
 $\hat{\sigma} = \hat{\sigma}(Y_{new} - \hat{Y}_{new}) = \sqrt{\frac{1}{n-1} \sum_{k=1}^n (y_k - \bar{y})^2} \times \sqrt{1 + \frac{1}{n}}$

Slide 30 STAT 1104, UCL4, Jon Dinger

Classical Prediction for the C+E model

- $Y = C + E$. When why, how to use prediction?
- When: $E \sim N(0, \sigma^2) \leftarrow \rightarrow Y \sim N(\mu, \sigma^2)$, there are more general situations, of course. Here we only consider this case.
- Why: Future predictions are of paramount importance in any area of science/engineering/medicine.
- How: μ is mostly unknown, so we estimate it by: m^{\wedge} , (the sample average).

If population proportion, p, is unknown we estimate it by the sample-proportion, p^{\wedge} , etc.

Slide 31 STAT 1104, UCL4, Jon Dinger

Classical Prediction for the C+E model

- **How:** μ is mostly unknown, so we estimate it by: m^\wedge ,
 - Let Y^\wedge_{new} be the predicted value
 - Error made by using Y^\wedge_{new} , instead of observing a new value, Y_{new} is:

$$(1) Y_{new} - Y^\wedge_{new} = (\mu - \epsilon_{new}) - Y^\wedge_{new} = (\mu - Y^\wedge_{new}) + \epsilon_{new}$$
 - But if we use μ^\wedge to predict a new value for Y , $Y^\wedge_{new} = \mu^\wedge$.
 - $\text{Var}(\mu - Y^\wedge_{new}) = \text{Var}(Y^\wedge_{new}) = \text{Var}(\mu^\wedge) = \text{Var}(\text{SampleAvg}) = \sigma^2/n$.
- The variance of the second term is just σ^2 .
- Since the first-term in (1) is obtained from $\{Y_1, Y_2, \dots, Y_n\}$, and $\epsilon_{new} = \epsilon_{n+1}$, we have two independent terms \rightarrow Variances add up!
- $\text{Var}(Y_{new} - Y^\wedge_{new}) = \text{Var}(\mu - Y^\wedge_{new}) + \text{Var}(\epsilon_{new}) = \sigma^2/n + \sigma^2$.

Slide 32 STAT 110A, UCLA, Jon Dinger

Classical Prediction for the C+E model

- **How:** Let Y^\wedge_{new} be the predicted value
 - Error $Y_{new} - Y^\wedge_{new} = (\mu - \epsilon_{new}) - Y^\wedge_{new} = (\mu - Y^\wedge_{new}) + \epsilon_{new}$
 - $\text{Var}(Y_{new} - Y^\wedge_{new}) = \text{Var}(\mu - Y^\wedge_{new}) + \text{Var}(\epsilon_{new}) = \sigma^2/n + \sigma^2$.
 - Often σ is unknown, and we estimate it by the sample SD, $S \rightarrow$
 - $\text{SD}(Y_{new} - Y^\wedge_{new}) = [S^2(1+1/n)]^{1/2}$
- We can show that
$$T = \frac{Y_{new} - \hat{Y}_{new} - 0}{\hat{\sigma} \sqrt{(Y_{new} - \hat{Y}_{new})^2}} \sim t_{n-1}$$
- \rightarrow The L-level prediction interval ($\text{PI}(Y_{new})$) is:

$$L = P(t_{n-1, (1-L)/2} < T < t_{n-1, (1+L)/2}) \rightarrow$$

$\left(\hat{Y}_{new} + \hat{\sigma} \times t_{n-1, (1-L)/2} \right)$

$\left(\hat{Y}_{new} + \hat{\sigma} \times t_{n-1, (1+L)/2} \right)$

Solve for T

$\left(\hat{Y}_{new} - \hat{\sigma} \times t_{n-1, (1+L)/2} \right)$

$\left(\hat{Y}_{new} - \hat{\sigma} \times t_{n-1, (1-L)/2} \right)$

By symmetry of t_{n-1} .

Slide 33 STAT 110A, UCLA, Jon Dinger

CI for a population proportion

Confidence Interval for the true (population) proportion p :

sample proportion $\pm z$ standard errors

or $\hat{p} \pm z \text{se}(\hat{p})$, where $\text{se}(\hat{p}) = \sqrt{\frac{\hat{p}(1-\hat{p})}{n}}$

Slide 34 STAT 110A, UCLA, Jon Dinger

Example – higher blood thiol concentrations associated with rheumatoid arthritis?!

Thiol Concentration (mmol)		
	Normal	Rheumatoid
Research question:	1.84	2.81
Is the change in the Thiol status in the lysate of packed blood cells substantial to be indicative of a non trivial relationship between Thiol-levels and rheumatoid arthritis?	1.92	4.06
	1.94	3.62
	1.92	3.27
	1.85	3.27
	1.91	3.76
	2.07	
Sample size	7	6
Sample mean	1.92143	3.46500
Sample standard deviation	0.07559	0.44049

Slide 35 STAT 110A, UCLA, Jon Dinger

Example – higher blood thiol concentrations with rheumatoid arthritis

Dot plot of Thiol concentration data.

Two groups of subjects are studied: 1. NC (normal controls)
2. RA (rheumatoid arthritis).

Observations: 1. The avg. levels of thiol seem diff. in NC & RA
2. NC and RA groups are separated completely.

Question: Is there **statistical evidence** that thiol-level correlates with the disease?

Slide 36 STAT 110A, UCLA, Jon Dinger

Difference between means

Confidence Interval for a difference between population means ($\mu_1 - \mu_2$):

Difference between sample means
 $\pm t$ standard errors of the difference

or
$$\bar{x}_1 - \bar{x}_2 \pm t \text{se}(\bar{x}_1 - \bar{x}_2)$$

Slide 37 STAT 110A, UCLA, Jon Dinger

Difference between proportions

Confidence Interval for a difference between population proportions ($p_1 - p_2$):

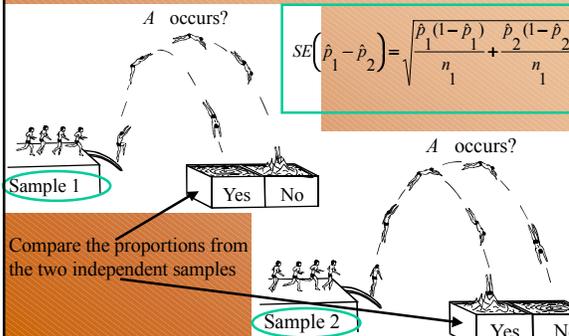
Difference between sample proportions
 $\pm z$ standard errors of the difference

$$\hat{p}_1 - \hat{p}_2 \pm z \text{ se}(\hat{p}_1 - \hat{p}_2)$$

How do we compute the $SE(\hat{p}_1 - \hat{p}_2)$ for different cases?

Slide 38 STAT 110A, UCLA, Jon Dinger

Proportions from 2 independent samples

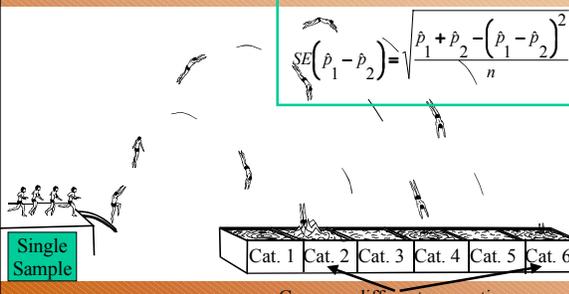


$$SE(\hat{p}_1 - \hat{p}_2) = \sqrt{\frac{\hat{p}_1(1-\hat{p}_1)}{n_1} + \frac{\hat{p}_2(1-\hat{p}_2)}{n_2}}$$

Compare the proportions from the two independent samples

Slide 39 STAT 110A, UCLA, Jon Dinger

Single sample, several response categories



$$SE(\hat{p}_1 - \hat{p}_2) = \sqrt{\frac{\hat{p}_1 + \hat{p}_2 - (\hat{p}_1 - \hat{p}_2)^2}{n}}$$

Compare different proportions from the same sample

Slide 40 STAT 110A, UCLA, Jon Dinger

Example – 1996 US Presidential Election

State	n	Pre-election Polls				Election Results		
		Clinton	Dole	Perot	Other/Undecided	Clinton	Dole	Perot
New Jersey	1,000	51	33	8	8	53	36	9
New York	1,000	59	25	7	9	59	31	8
Connecticut	1,000	51	29	11	9	52	35	10

Compare proportions of NJ and NY voters supporting Clinton and Dole, pre- and post election

$$\hat{p}_1 - \hat{p}_2 \pm z \text{ se}(\hat{p}_1 - \hat{p}_2)$$

Note the **independence-case SE formula** is only applicable for the cases when the samples are independent. In this case, the **pre-election poll** and the **election results** are **not independent** (obviously these are highly correlated observations).

Slide 41 STAT 110A, UCLA, Jon Dinger

Example – 1996 US Presidential Election

State	n	Pre-election Polls				Election Results		
		Clinton	Dole	Perot	Other/Undecided	Clinton	Dole	Perot
New Jersey	1,000	51	33	8	8	53	36	9
New York	1,000	59	25	7	9	59	31	8
Connecticut	1,000	51	29	11	9	52	35	10

Proportions from 2 independent samples

How far is Clinton ahead in NY compared to NJ? Diff. proportions = 59-51% = 8%
 CI: [4% : 12%]
 Actual diff 59-53=6

$$\hat{p}_1 - \hat{p}_2 \pm z \text{ se}(\hat{p}_1 - \hat{p}_2)$$

estimate $\pm z \times SE = \hat{p}_1 - \hat{p}_2 \pm 1.96 \times SE(\hat{p}_1 - \hat{p}_2) =$
 $\hat{p}_1 - \hat{p}_2 \pm 1.96 \times \sqrt{\frac{\hat{p}_1(1-\hat{p}_1)}{n_1} + \frac{\hat{p}_2(1-\hat{p}_2)}{n_2}} =$
 $0.08 \pm 1.96 \times 0.02842 = [4\% : 12\%]$

Slide 42 STAT 110A, UCLA, Jon Dinger

Example – 1996 US Presidential Election

State	n	Pre-election Polls				Election Results		
		Clinton	Dole	Perot	Other/Undecided	Clinton	Dole	Perot
New Jersey	1,000	51	33	8	8	53	36	9
New York	1,000	59	25	7	9	59	31	8
Connecticut	1,000	51	29	11	9	52	35	10

Single sample, several response categories

How far is Clinton ahead of Dole in NJ? Diff. proportions = 18%
 CI: [12% : 24%]
 Actual diff 53-36=17

$$\hat{p}_1 - \hat{p}_2 \pm z \text{ se}(\hat{p}_1 - \hat{p}_2)$$

estimate $\pm z \times SE = \hat{p}_1 - \hat{p}_2 \pm 1.96 \times SE(\hat{p}_1 - \hat{p}_2) =$
 $\hat{p}_1 - \hat{p}_2 \pm 1.96 \times \sqrt{\frac{\hat{p}_1 + \hat{p}_2 - (\hat{p}_1 - \hat{p}_2)^2}{n}} =$
 $0.18 \pm 1.96 \times 0.02842 = [12\% : 24\%]$

Slide 43 STAT 110A, UCLA, Jon Dinger

SE's for the 2 cases of differences in proportion

(a) Proportions from two independent samples of sizes n_1 and n_2 , respectively

$$se(\hat{p}_1 - \hat{p}_2) = \sqrt{\frac{\hat{p}_1(1-\hat{p}_1)}{n_1} + \frac{\hat{p}_2(1-\hat{p}_2)}{n_2}}$$

(b) One sample of size n , several response categories

$$se(\hat{p}_1 - \hat{p}_2) = \sqrt{\frac{\hat{p}_1 + \hat{p}_2 - (\hat{p}_1 - \hat{p}_2)^2}{n}}$$

Slide 44 STAT 110A, UCLA, Jon Dinger

Sample size - proportion

- For a 95% CI, margin = $1.96 \times \sqrt{\hat{p}(1-\hat{p})/n}$

- Sample size for a desired margin of error:

For a margin of error no greater than m , use a sample size of approximately

$$n = \left(\frac{z}{m}\right)^2 \times p^*(1-p^*)$$

- p^* is a guess at the value of the proportion -- err on the side of being too close to 0.5

- z is the multiplier appropriate for the confidence level

- m is expressed as a proportion (between 0 and 1), not a percentage (basically, What's n , so that $m \gg \text{margin?}$)

Slide 45 STAT 110A, UCLA, Jon Dinger

Sample size -- mean

- Sample size for a desired margin of error:

For a margin of error no greater than m , use a sample size of approximately

$$n = \left(\frac{z\sigma^*}{m}\right)^2$$

- σ^* is an estimate of the variability of individual observations

- z is the multiplier appropriate for the confidence level

Slide 46 STAT 110A, UCLA, Jon Dinger

Paired vs. Unpaired comparisons

- We will discuss these later, when we get to the hypothesis testing (ch6_HT_Paired_Indep_Tests.ppt)

Slide 47 STAT 110A, UCLA, Jon Dinger

Confidence intervals

- We construct an interval estimate of a parameter to summarize our level of uncertainty about its true value.

- The uncertainty is a consequence of the sampling variation in point estimates.

- If we use a method that produces intervals which contain the true value of a parameter for 95% of samples taken, the interval we have calculated from our data is called a 95% confidence interval for the parameter.

- Our confidence in the particular interval comes from the fact that the method works 95% of the time (for 95% CI's).

Slide 48 STAT 110A, UCLA, Jon Dinger

Summary cont.

- For a great many situations,

an (approximate) confidence interval is given by

$$\text{estimate} \pm t \text{ standard errors}$$

The size of the multiplier, t , depends both on the desired confidence level and the degrees of freedom (df).

[With proportions, we use the Normal distribution (i.e., $df = \infty$) and it is conventional to use z rather than t to denote the multiplier.]

- The *margin of error* is the quantity added to and subtracted from the estimate to construct the interval (i.e. t standard errors).

Slide 49 STAT 110A, UCLA, Jon Dinger

Summary cont.

- If we want greater confidence that an interval calculated from our data will contain the true value, we have to use a wider interval.
- To double the precision of a 95% confidence interval (i.e. halve the width of the confidence interval), we need to take 4 times as many observations.