# BOOST YOUR CONFIDENCE (INTERVALS) WITH SAS

Brought to you by:

Peter Langlois, PhD
Birth Defects Epidemiology & Surveillance Branch,
Texas Dept State Health Services

# Background

# Confidence Interval Definition

- DEFINITION:  An interval around a statistic that contains the true underlying value of the statistic (the population parameter) a certain amount of the time

# Confidence Interval Definition

- Example: a survey of 50 SAS programmers finds that the average IQ is 130 $\pm$ 10

- If we did 100 surveys, the average IQ should be between 120 and 140 in 95 of them

# Confidence Interval Definition

- 95% confidence interval bounded by the upper 95% confidence limit and the lower 95% confidence limit

- 95% just conventional.  Can have for e.g.:
  - 90% CIs (narrower)
  - 99% CIs (wider)
- CI for any level (95% etc) is narrower if based on more observations

# General Formula

- Can make CIs around almost any statistic you calculate, for example…

- Data summaries (1 var) such as:
  - Categorical variables: proportion
  - Continuous variables: mean

- Statistics resulting from hypothesis tests (2+ vars):
  - Correlation, regression slope
  - Relative risk, odds ratios

# Using Confidence Intervals

# Why Use CIs?

- Gives audience idea of impact of chance
- Gives reasonable bounds on your result(s)
- Can check if your data are compatible with a certain value
  - (From data summary): Does 95% confidence interval of IQ in SAS programmers include 100?
  - (From hypothesis testing): Is occurrence of schizophrenia higher in SAS programmers or SPSS users? (Does relative risk = 1.00?)

# Why Not Use CIs?

- Some organizations consider their figures to be a census, not a sample

- Increases statistical work for staff

- Some data users may find the concept confusing

# Overlapping CIs vs Hypothesis Tests

- Example: Want to compare prevalence of schizophrenia in SAS vs SPSS users

- You could:

(A) Calculate schizophrenia prevalence and 95% CI for each group and see if overlap OR

(B) Calculate prevalence ratio of one group vs another, and see if includes 1.00

# Overlapping CIs vs Hypothesis Tests

- Answer: (B) usually better
- Why? More statistical power

- Why even consider the first approach?
  - Easier to do if already have data summaries published
  - Can't anticipate all comparisons readers will want to make

# Confidence Intervals
# For Data Summaries

# CIs for Means

- Based on normal distribution

- Where your study sample is large (# of subjects > 30), the sampling distribution → normal, and you can use:

   CI = obsd mean ± 1.96 X standard error (of mean)

   = obsd mean ± 1.96 X standard dev'n / sqrt(n)

$$= \overline{x} \pm 1.96 \frac{s}{\sqrt{n}}$$

# CIs for Means: Using SAS

- Let SAS do it!

- Recall: Using SAS for simple desc stats:

```
proc means;
   var <variable name>;
```

# CIs for Means: Using SAS

- To get confidence limits, request "clm" as PROC MEANS option

- (Request "mean" to get the mean printed out too)

```
proc means mean clm;
    var <variable name>;
```

The MEANS Procedure

Analysis Variable : i_bwt_v

| Mean | Lower 95% CL for Mean | Upper 95% CL for Mean |
|---|---|---|
| 3305.77 | 3209.56 | 3401.98 |

# CIs for Proportions

- Based on binomial distribution

- Where your study sample is large (# with and without the characteristic > 5), the sampling distribution → normal, and you can use:

- = obsd proportion $\pm$ 1.96 X SE of prop'n

$$= p \pm 1.96 \sqrt{\frac{p(1-p)}{n}}$$

# CIs for Proportions: Using SAS

- To get confidence limits, request "binomial" as TABLE option in PROC FREQ:

```
proc freq;
   table m_edu2g_v / binomial;
```

File   Edit   View   Tools   Solutions   Window   Help

USING SAS FOR CONFIDENCE INTERVALS                    27
14:20 Monday, October 11, 2010

The FREQ Procedure

| m_edu2g_v | Frequency | Percent | Cumulative Frequency | Cumulative Percent |
|---|---|---|---|---|
| < High School | 35 | 35.00 | 35 | 35.00 |
| High School + | 65 | 65.00 | 100 | 100.00 |

Binomial Proportion for
m_edu2g_v = < High School

| Proportion | 0.3500 |
|---|---|
| ASE | 0.0477 |
| 95% Lower Conf Limit | 0.2565 |
| 95% Upper Conf Limit | 0.4435 |

Exact Conf Limits
| 95% Lower Conf Limit | 0.2573 |
| 95% Upper Conf Limit | 0.4518 |

Test of H0: Proportion = 0.5

| ASE under H0 | 0.0500 |
|---|---|
| Z | -3.0000 |
| One-sided Pr < Z | 0.0013 |
| Two-sided Pr > |Z| | 0.0027 |

Sample Size = 100

Results
   Results
      Freq: USING SAS FOR CONFIDENCE

Output - (Untitled)    Log - (Untitled)    CI talk for SAS user...

C:\Documents and Settings\planglois029

start    Inbox - Microsoft Out...    SAS - [Output - (Untit...    SAS Help and Docum...    Microsoft PowerPoint ...    3:57 PM

# CIs for Rates of Rare Outcomes

- Follows Poisson distribution

- Standard error based on number of cases

# CIs for Rates of Rare Outcomes
## PREVALENCE BASED ON < 30 CASES

- Find upper and lower 95% conf limits for # cases:
  - In table (e.g. poisson.xls)
  - Using PEPI program POISSON

- Multiply both by (10000) and divide by # population or live births to express as birth prevalence

| | COUNT | LFACTOR | UFACTOR |
|---|---|---|---|
| 1 | COUNT | LFACTOR | UFACTOR |
| 2 | 0 | 0.0000 | 3.6889 |
| 3 | 1 | 0.0253 | 5.5716 |
| 4 | 2 | 0.2422 | 7.2247 |
| 5 | 3 | 0.6187 | 8.7673 |
| 6 | 4 | 1.0899 | 10.2416 |
| 7 | 5 | 1.6235 | 11.6683 |
| 8 | 6 | 2.2019 | 13.0595 |
| 9 | 7 | 2.8144 | 14.4227 |
| 10 | 8 | 3.4538 | 15.7632 |
| 11 | 9 | 4.1154 | 17.0848 |
| 12 | 10 | 4.7954 | 18.3904 |
| 13 | 11 | 5.4912 | 19.6820 |
| 14 | 12 | 6.2006 | 20.9616 |
| 15 | 13 | 6.9220 | 22.2304 |
| 16 | 14 | 7.6539 | 23.4896 |
| 17 | 15 | 8.3954 | 24.7402 |
| 18 | 16 | 9.1454 | 25.9830 |
| 19 | 17 | 9.9031 | 27.2186 |
| 20 | 18 | 10.6679 | 28.4478 |
| 21 | 19 | 11.4392 | 29.6709 |
| 22 | 20 | 12.2165 | 30.8884 |
| 23 | 21 | 12.9993 | 32.1007 |
| 24 | 22 | 13.7873 | 33.3083 |
| 25 | 23 | 14.5800 | 34.5113 |
| 26 | 24 | 15.3773 | 35.7101 |
| 27 | 25 | 16.1787 | 36.9049 |
| 28 | 26 | 16.9841 | 38.0960 |
| 29 | 27 | 17.7932 | 39.2836 |
| 30 | 28 | 18.6058 | 40.4678 |
| 31 | 29 | 19.4218 | 41.6488 |
| 32 | 30 | 20.2409 | 42.8269 |
| 33 | 31 | 21.0630 | 44.0020 |

POISSON

# CIs for Rates of Rare Outcomes

EXAMPLE

- # cases of anophthalmia in CA 1983-1986 = 18
- Prevalence = (18 * 10,000 / 452,287) = 0.40
- Looking in table, lower 95% CL is 10.67 and upper 95% CL is 28.45 for # cases
- To express as prevalence CLs, multiply both by 10,000 and divide by 452,287 live births
- Lower 95% CL = 0.24, upper 95% CL = 0.63
- Thus we say prev = 0.40 cases per 10,000 live births, 95% CI = 0.24 - 0.63

# CIs for Rates of Rare Outcomes

**PREVALENCE BASED ON 30+ CASES**

- Considered large # cases (more or less)

- Poisson → normal distribution
- Can use normal approximation in SAS code
- Several equations for doing this, yielding similar results

# CIs for Rates of Rare Outcomes

- SAS code for obs with few cases: Combine with Poisson lower and upper limits for cases (get lfactor and ufactor for the observed # cases):

```
proc sort data=b1;  by count;

proc sort data=lib2.poisson out=poisson;
  by count;

data c1 prob2;
  merge b1(in=b) poisson(in=p);
  by count;
  if b;
```

# CIs for Rates of Rare Outcomes

- SAS code: Calculate CIs for obs with many cases too

```
data c2;
  set c1;
  calcrate = count * 10000 / births;
  if count le 30 then do;
    calclcl = lfactor * 10000 / births;
    calcucl = ufactor * 10000 / births;
    end;
  else if count > 30 then do;
    calclcl = ((count / births) - (1.96 * sqrt(count) / births)) * 10000;
    calcucl = ((count / births) + (1.96 * sqrt(count) / births)) * 10000;
    end;
  rate = round(calcrate,.01);
  lcl = round(calclcl,.01);
  ucl = round(calcucl,.01);
  rename count=cases;
```

# Reminder

- To compare groups (e.g. whether rates are statistically different), can calculate 95% confidence intervals and see if they overlap

- Better to do hypothesis testing

# Confidence Intervals
# For Statistics From
# Hypothesis Tests / Measures
# of Association

# Types of Hypothesis Tests / Measures of Association

| Indep Var | Dep Var | Approach (SAS PROC) |
|---|---|---|
| Categorical | Categorical | Contingency tables (FREQ) |
| Cat. / Cont. | Categorical | Logistic reg. (LOGISTIC) |
| Categorical | Cont. | ANOVA (GLM); if 2 levels then T-test (TTEST or GLM) |
| Cont. | Cont. | Linear reg. (REG or GLM) |
| Cat. / Cont. | Rates | Poisson reg. (GENMOD) |

# Contingency Tables Using SAS

- Recall: Using SAS to produce the basic 2x2 table:

```
proc freq;
   tables <indep var> * <outcome var>;
```

- To get odds ratios and their CIs, request measures of association:

```
proc freq;
   tables <indep var> * <outcome var>
   / measures;
```

# Contingency Tables Using SAS

QUESTION: Is mother's education associated with % low birth weight babies?

- Independent variable = m_edu2g_v

- Outcome variable = lbw

- If no statistically significant association, 95% CI will include 1.00

```
proc freq;
    tables m_edu2g_v * lbw / measures;
```

The FREQ Procedure

Table of m_edu2g_v by lbw

m_edu2g_v          lbw

| Frequency<br>Percent<br>Row Pct<br>Col Pct | Low bwt | Normal b<br>wt | Total |
|---|---|---|---|
| < High School | 10769<br>2.65<br>8.64<br>31.40 | 113937<br>28.05<br>91.36<br>30.63 | 124706<br>30.70 |
| High School + | 23525<br>5.79<br>8.36<br>68.60 | 258010<br>63.51<br>91.64<br>69.37 | 281535<br>69.30 |
| Total | 34294<br>8.44 | 371947<br>91.56 | 406241<br>100.00 |

Frequency Missing = 539

Statistics for Table of m_edu2g_v by lbw

| Statistic | Value | ASE |
|---|---|---|
| Gamma | 0.0180 | 0.0061 |
| Kendall's Tau-b | 0.0046 | 0.0016 |
| Stuart's Tau-c | 0.0024 | 0.0008 |
| Somers' D C\|R | 0.0028 | 0.0010 |
| Somers' D R\|C | 0.0077 | 0.0026 |
| Pearson Correlation | 0.0046 | 0.0016 |
| Spearman Correlation | 0.0046 | 0.0016 |

# Logistic Regression: Using SAS

- Can use PROC LOGISTIC in SAS; nice since it will exponentiate the slope (b) and its 95% confidence interval

- Basic syntax:

```
proc logistic;
    model <outcome var> = <independent var>;
```

- For low birthweight example:

```
proc logistic;
    model lbw = m_edu2g_v;
```

# Logistic Regression: Using SAS

- (One way) to get correct odds ratio: declare independent var to be a classification (categorical) var:

```
proc logistic;
   class m_edu2g_v;
   model lbw = m_edu2g_v;
```

# Logistic Regression: Using SAS

- Comparing results from PROC FREQ and PROC LOGISTIC:

| SAS Proc | OR | 95% CI |
|----------|--------|-----------------|
| FREQ | 1.0366 | 1.0122 – 1.0616 |
| LOGISTIC | 1.037 | 1.012 – 1.062 |

# Logistic Regression
# With Multiple Predictor Variables

- Like other regression, the slope (b) is adjusted for all other independent variables in the model

- SAS takes both cont and categorical vars
  - SAS assumes ind vars are continuous
  - If categorical, list in CLASS statement and SAS creates dummy vars automatically

```
proc logistic;
   class <categorical independent vars>;
   model <dependent var> = <independent vars>;
```

peter.langlois@dshs.state.tx.us

Phone: 512-458-7111 x6183

Thanks