

Introduction Based on a Simple Example¹

STA442/2101 Fall 2013

¹See last slide for copyright information.

Background Reading

Optional

- ▶ Chapter 1 of *Data analysis with SAS*: What's going on and how would you say it to a client?
- ▶ Chapter 1 of Davison's *Statistical models: Data, and probability models for data*.

Goal of statistical analysis

The goal of statistical analysis is to draw reasonable conclusions from noisy numerical data.

Steps in the process of statistical analysis

One approach

- ▶ Consider a fairly realistic example or problem
- ▶ Decide on a statistical model
- ▶ Perhaps decide sample size
- ▶ Acquire data
- ▶ Examine and clean the data; generate displays and descriptive statistics
- ▶ Estimate model parameters, perhaps by maximum likelihood
- ▶ Carry out tests, compute confidence intervals, or both
- ▶ Perhaps re-consider the model and go back to estimation
- ▶ Based on the results of estimation and inference, draw conclusions about the example or problem

What is a statistical model?

You should always be able to state the model.

A *statistical model* is a set of assertions that directly or indirectly specify the probability distribution of the observable data.

- ▶ Let X_1, \dots, X_n be a random sample from a normal distribution with expected value μ and variance σ^2 .
- ▶ For $i = 1, \dots, n$, let $Y_i = \beta_0 + \beta_1 x_{i1} + \dots + \beta_k x_{ik} + \epsilon_i$, where

β_0, \dots, β_k are unknown constants.

x_{ij} are known constants.

$\epsilon_1, \dots, \epsilon_n$ are independent $N(0, \sigma^2)$ random variables.

σ^2 is an unknown constant.

Y_1, \dots, Y_n are observable random variables.

Is the model the same thing as the *truth*?

Parameter Space

The *parameter space* is the set of values that can be taken on by the parameter.

- ▶ Let X_1, \dots, X_n be a random sample from a normal distribution with expected value μ and variance σ^2 .
The parameter space is $\{(\mu, \sigma^2) : -\infty < \mu < \infty, \sigma^2 > 0\}$.
- ▶ For $i = 1, \dots, n$, let $Y_i = \beta_0 + \beta_1 x_{i1} + \dots + \beta_k x_{ik} + \epsilon_i$,
where
 - β_0, \dots, β_k are unknown constants.
 - x_{ij} are known constants.
 - $\epsilon_1, \dots, \epsilon_n$ are independent $N(0, \sigma^2)$ random variables.
 - σ^2 is an unknown constant.
 - Y_1, \dots, Y_n are observable random variables.

The parameter space is

$$\{(\beta_0, \dots, \beta_k, \sigma^2) : -\infty < \beta_j < \infty, \sigma^2 > 0\}.$$

Coffee taste test

A fast food chain is considering a change in the blend of coffee beans they use to make their coffee. To determine whether their customers prefer the new blend, the company plans to select a random sample of $n = 100$ coffee-drinking customers and ask them to taste coffee made with the new blend and with the old blend, in cups marked “ A ” and “ B .” Half the time the new blend will be in cup A , and half the time it will be in cup B . Management wants to know if there is a difference in preference for the two blends.

Statistical model

Letting θ denote the probability that a consumer will choose the new blend, treat the data Y_1, \dots, Y_n as a random sample from a Bernoulli distribution. That is, independently for $i = 1, \dots, n$,

$$P(y_i|\theta) = \theta^{y_i}(1 - \theta)^{1-y_i}$$

for $y_i = 0$ or $y_i = 1$, and zero otherwise.

Note that $Y = \sum_{i=1}^n Y_i$ is the number of consumers who choose the new blend. Because $Y \sim B(n, \theta)$, the whole experiment could also be treated as a single observation from a Binomial.

Find the MLE of θ

Show your work

Denoting the likelihood by $L(\theta)$ and the log likelihood by $\ell(\theta) = \log L(\theta)$, maximize the log likelihood.

$$\begin{aligned}\frac{\partial \ell}{\partial \theta} &= \frac{\partial}{\partial \theta} \log \left(\prod_{i=1}^n P(y_i | \theta) \right) \\ &= \frac{\partial}{\partial \theta} \log \left(\prod_{i=1}^n \theta^{y_i} (1 - \theta)^{1 - y_i} \right) \\ &= \frac{\partial}{\partial \theta} \log \left(\theta^{\sum_{i=1}^n y_i} (1 - \theta)^{n - \sum_{i=1}^n y_i} \right) \\ &= \frac{\partial}{\partial \theta} \left(\left(\sum_{i=1}^n y_i \right) \log \theta + \left(n - \sum_{i=1}^n y_i \right) \log(1 - \theta) \right) \\ &= \frac{\sum_{i=1}^n y_i}{\theta} - \frac{n - \sum_{i=1}^n y_i}{1 - \theta}\end{aligned}$$

Setting the derivative to zero and solving

- ▶ $\theta = \frac{\sum_{i=1}^n y_i}{n} = \bar{y}$
- ▶ Second derivative test: $\frac{\partial^2 \log \ell}{\partial \theta^2} = -n \left(\frac{1-\bar{y}}{(1-\theta)^2} + \frac{\bar{y}}{\theta^2} \right) < 0$
- ▶ Concave down, maximum, and the MLE is the sample proportion.

Numerical estimate

Suppose 60 of the 100 consumers prefer the new blend. Give a point estimate of the parameter θ . Your answer is a number.

```
> p = 60/100; p  
[1] 0.6
```

Carry out a test to answer the question

Is there a difference in preference for the two blends?

Start by stating the null hypothesis.

- ▶ $H_0 : \theta = 0.50$
- ▶ $H_1 : \theta \neq 0.50$
- ▶ A case could be made for a one-sided test, but we'll stick with two-sided.
- ▶ $\alpha = 0.05$ as usual.
- ▶ Central Limit Theorem says $\hat{\theta} = \bar{Y}$ is approximately normal with mean θ and variance $\frac{\theta(1-\theta)}{n}$.

Several valid test statistics for $H_0 : \theta = \theta_0$ are available

Two of them are

$$Z_1 = \frac{\sqrt{n}(\bar{Y} - \theta_0)}{\sqrt{\theta_0(1 - \theta_0)}}$$

and

$$Z_2 = \frac{\sqrt{n}(\bar{Y} - \theta_0)}{\sqrt{\bar{Y}(1 - \bar{Y})}}$$

What is the critical value? Your answer is a number.

```
> alpha = 0.05
> qnorm(1-alpha/2)
[1] 1.959964
```

Calculate the test statistic and the p -value for each test

```
> theta0 = .5; ybar = .6; n = 100
> Z1 = sqrt(n)*(ybar-theta0)/sqrt(theta0*(1-theta0)); Z1
[1] 2
> pval1 = 2 * (1-pnorm(Z1)); pval1
[1] 0.04550026

> Z2 = sqrt(n)*(ybar-theta0)/sqrt(ybar*(1-ybar)); Z2
[1] 2.041241
> pval2 = 2 * (1-pnorm(Z2)); pval2
[1] 0.04122683
```

Conclusions

- ▶ Do you reject H_0 ? *Yes, just barely.*
- ▶ Isn't the $\alpha = 0.05$ significance level pretty arbitrary? *Yes, but if people insist on a Yes or No answer, this is what you give them.*
- ▶ What do you conclude, in symbols? $\theta \neq 0.50$. *Specifically, $\theta > 0.50$.*
- ▶ What do you conclude, in plain language? Your answer is a statement about coffee. *More consumers prefer the new blend of coffee beans.*
- ▶ Can you really draw directional conclusions when all you did was reject a non-directional null hypothesis? *Yes. Decompose the two-sided size α test into two one-sided tests of size $\alpha/2$. This approach works in general.*

It is very important to state directional conclusions, and state them clearly in terms of the subject matter. **Say what happened!** If you are asked state the conclusion in plain language, your answer *must* be free of statistical mumbo-jumbo.

What about negative conclusions?

What would you say if $Z = 1.84$?

Here are two possibilities, in plain language.

- ▶ “This study does not provide clear evidence that consumers prefer one blend of coffee beans over the other.”
- ▶ “The results are consistent with no difference in preference for the two coffee bean blends.”

In this course, we will not just casually *accept* the null hypothesis.

Confidence Intervals

Approximately for large n ,

$$\begin{aligned}1 - \alpha &= Pr\{-z_{\alpha/2} < Z_2 < z_{\alpha/2}\} \\ &\approx Pr\left\{-z_{\alpha/2} < \frac{\sqrt{n}(\bar{Y} - \theta)}{\sqrt{\bar{Y}(1 - \bar{Y})}} < z_{\alpha/2}\right\} \\ &= Pr\left\{\bar{Y} - z_{\alpha/2}\sqrt{\frac{\bar{Y}(1 - \bar{Y})}{n}} < \theta < \bar{Y} + z_{\alpha/2}\sqrt{\frac{\bar{Y}(1 - \bar{Y})}{n}}\right\}\end{aligned}$$

- ▶ Could express this as $\bar{Y} \pm z_{\alpha/2}\sqrt{\frac{\bar{Y}(1 - \bar{Y})}{n}}$
- ▶ $z_{\alpha/2}\sqrt{\frac{\bar{Y}(1 - \bar{Y})}{n}}$ is sometimes called the *margin of error*.
- ▶ If $\alpha = 0.05$, it's the 95% margin of error.

Give a 95% confidence interval for the taste test data.

The answer is a pair of numbers. Show some work.

$$\begin{aligned} & \left(\bar{y} - z_{\alpha/2} \sqrt{\frac{\bar{y}(1-\bar{y})}{n}} , \bar{y} + z_{\alpha/2} \sqrt{\frac{\bar{y}(1-\bar{y})}{n}} \right) \\ &= \left(0.60 - 1.96 \sqrt{\frac{0.6 \times 0.4}{100}} , 0.60 + 1.96 \sqrt{\frac{0.6 \times 0.4}{100}} \right) \\ &= (0.504, 0.696) \end{aligned}$$

In a report, you could say

- ▶ The estimated proportion preferring the new coffee bean blend is 0.60 ± 0.096 , or
- ▶ “Sixty percent of consumers preferred the new blend. These results are expected to be accurate within 10 percentage points, 19 times out of 20.”

Meaning of the confidence interval

- ▶ We calculated a 95% confidence interval of $(0.504, 0.696)$ for θ .
- ▶ Does this mean $Pr\{0.504 < \theta < 0.696\} = 0.95$?
- ▶ No! The quantities 0.504, 0.696 and θ are all constants, so $Pr\{0.504 < \theta < 0.696\}$ is either zero or one.
- ▶ The endpoints of the confidence interval are random variables, and the numbers 0.504 and 0.696 are *realizations* of those random variables, arising from a particular random sample.
- ▶ Meaning of the probability statement: If we were to calculate an interval in this manner for a large number of random samples, the interval would contain the true parameter around 95% of the time.
- ▶ So we sometimes say that we are “95% confident” that $0.504 < \theta < 0.696$.

Confidence intervals (regions) correspond to tests

Recall $Z_1 = \frac{\sqrt{n}(\bar{Y} - \theta_0)}{\sqrt{\theta_0(1 - \theta_0)}}$ and $Z_2 = \frac{\sqrt{n}(\bar{Y} - \theta_0)}{\sqrt{\bar{Y}(1 - \bar{Y})}}$.

From the derivation of the confidence interval,

$$-z_{\alpha/2} < Z_2 < z_{\alpha/2}$$

if and only if

$$\bar{Y} - z_{\alpha/2} \sqrt{\frac{\bar{Y}(1 - \bar{Y})}{n}} < \theta_0 < \bar{Y} + z_{\alpha/2} \sqrt{\frac{\bar{Y}(1 - \bar{Y})}{n}}$$

- ▶ So the confidence interval consists of those parameter values θ_0 for which $H_0 : \theta = \theta_0$ is *not* rejected.
- ▶ That is, the null hypothesis is rejected at significance level α if and only if the value given by the null hypothesis is outside the $(1 - \alpha) \times 100\%$ confidence interval.
- ▶ There is a confidence interval corresponding to Z_1 too.
- ▶ In general, any test can be inverted to obtain a confidence region.

Selecting sample size

- ▶ Where did that $n = 100$ come from?
- ▶ Probably off the top of someone's head.
- ▶ We can (and should) be more systematic.
- ▶ Sample size can be selected
 - ▶ To achieve a desired margin of error
 - ▶ To achieve a desired statistical power
 - ▶ In other reasonable ways

Power

The power of a test is the probability of rejecting H_0 when H_0 is false.

- ▶ More power is good.
- ▶ Power is not just one number. It is a *function* of the parameter(s).
- ▶ Usually,
 - ▶ For any n , the more incorrect H_0 is, the greater the power.
 - ▶ For any parameter value satisfying the alternative hypothesis, the larger n is, the greater the power.

Statistical power analysis

To select sample size

- ▶ Pick an effect you'd like to be able to detect – a parameter value such that H_0 is false. It should be just over the boundary of interesting and meaningful.
- ▶ Pick a desired power, a probability with which you'd like to be able to detect the effect by rejecting the null hypothesis.
- ▶ Start with a fairly small n and calculate the power. Increase the sample size until the desired power is reached.

There are two main issues.

- ▶ What is an “interesting” or “meaningful” parameter value?
- ▶ How do you calculate the probability of rejecting H_0 ?

Calculating power for the Z_1 test of a single proportion

True parameter value is θ

$$\begin{aligned}\text{Power} &= 1 - Pr\{-z_{\alpha/2} < Z_1 < z_{\alpha/2}\} \\ &= 1 - Pr\left\{-z_{\alpha/2} < \frac{\sqrt{n}(\bar{Y} - \theta_0)}{\sqrt{\theta_0(1 - \theta_0)}} < z_{\alpha/2}\right\} \\ &= \dots \\ &= 1 - Pr\left\{\frac{\sqrt{n}(\theta_0 - \theta)}{\sqrt{\theta(1 - \theta)}} - z_{\alpha/2}\sqrt{\frac{\theta_0(1 - \theta_0)}{\theta(1 - \theta)}} < \frac{\sqrt{n}(\bar{Y} - \theta)}{\sqrt{\theta(1 - \theta)}}\right. \\ &\quad \left.< \frac{\sqrt{n}(\theta_0 - \theta)}{\sqrt{\theta(1 - \theta)}} + z_{\alpha/2}\sqrt{\frac{\theta_0(1 - \theta_0)}{\theta(1 - \theta)}}\right\} \\ &\approx 1 - \Phi\left(\frac{\sqrt{n}(\theta_0 - \theta)}{\sqrt{\theta(1 - \theta)}} + z_{\alpha/2}\sqrt{\frac{\theta_0(1 - \theta_0)}{\theta(1 - \theta)}}\right) + \Phi\left(\frac{\sqrt{n}(\theta_0 - \theta)}{\sqrt{\theta(1 - \theta)}} - z_{\alpha/2}\sqrt{\frac{\theta_0(1 - \theta_0)}{\theta(1 - \theta)}}\right),\end{aligned}$$

where $\Phi(\cdot)$ is the cumulative distribution function of the standard normal.

An R function to calculate approximate power

For the Z_1 test of a single proportion

$$\text{Power} = 1 - \Phi\left(\frac{\sqrt{n}(\theta_0 - \theta)}{\sqrt{\theta(1 - \theta)}} + z_{\alpha/2}\sqrt{\frac{\theta_0(1 - \theta_0)}{\theta(1 - \theta)}}\right) + \Phi\left(\frac{\sqrt{n}(\theta_0 - \theta)}{\sqrt{\theta(1 - \theta)}} - z_{\alpha/2}\sqrt{\frac{\theta_0(1 - \theta_0)}{\theta(1 - \theta)}}\right)$$

```
Z1power = function(theta,n,theta0=0.50,alpha=0.05)
{
  a = sqrt(n)*(theta0-theta)/sqrt(theta*(1-theta))
  b = qnorm(1-alpha/2) * sqrt(theta0*(1-theta0)/(theta*(1-theta)))
  Z1power = 1 - pnorm(a+b) + pnorm(a-b)
  Z1power
} # End of function Z1power
```

```
> Z1power(theta=0.50,n=100) # Should be alpha = 0.05
[1] 0.05
```

Some numerical examples

```
> Z1power(0.55,100)
```

```
[1] 0.168788
```

```
> Z1power(0.60,100)
```

```
[1] 0.5163234
```

```
> Z1power(0.65,100)
```

```
[1] 0.8621995
```

```
> Z1power(0.40,100)
```

```
[1] 0.5163234
```

```
> Z1power(0.55,500)
```

```
[1] 0.6093123
```

```
> Z1power(0.55,1000)
```

```
[1] 0.8865478
```

Find smallest sample size needed to detect $\theta = 0.60$ as different from $\theta_0 = 0.50$ with probability at least 0.80

```
> samplesize = 1; power=0
> while(power < 0.80)
+ {
+ samplesize = samplesize+1
+ power = Z1power(theta=0.60,n=samplesize)
+ }
>
> samplesize
[1] 194
> power
[1] 0.8003138
```

Conclusions from the power analysis

Assuming true $\theta = 0.60$

- ▶ Power of 0.52 with $n = 100$ is pathetic.
- ▶ As Fisher said, “To call in the statistician after the experiment is done may be no more than asking him to perform a postmortem examination: he may be able to say what the experiment died of.”
- ▶ $n = 200$ is better.
> `Z1power(theta=0.60,n=200)`
[1] 0.8122918
- ▶ What sample size is required for power of 90%?

What sample size is needed for a power of 90%?

```
> # Find sample size needed for power = 0.90
> samplesize = 1; power=0
> while(power < 0.90)
+ {
+ samplesize = samplesize+1
+ power = Z1power(theta=0.60,n=samplesize)
+ }
>
> samplesize
[1] 259
> power
[1] 0.9005493
```

What is required of the scientist

Who wants to select sample size by power analysis

The scientist must specify

- ▶ Parameter values that he or she wants to be able to detect as different from H_0 value.
- ▶ Desired power (probability of detection)

It's not always easy for a scientist to think in terms of the parameters of a statistical model.

Copyright Information

This slide show was prepared by **Jerry Brunner**, Department of Statistics, University of Toronto. It is licensed under a **Creative Commons Attribution - ShareAlike 3.0 Unported License**. Use any part of it as you like and share the result freely. The \LaTeX source code is available from the course website:
<http://www.utstat.toronto.edu/~brunner/oldclass/appliedf13>