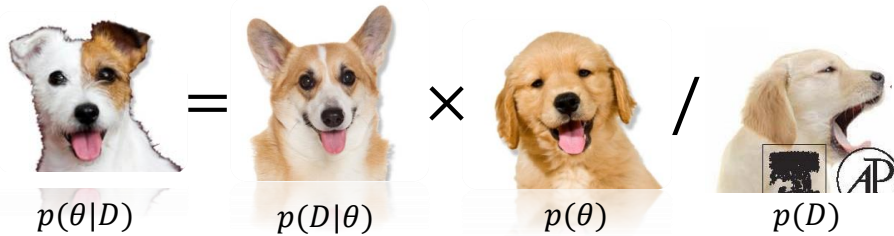


## *Doing Bayesian Data Analysis*



$$p(\theta|D) = p(D|\theta) \times p(\theta) / p(D)$$

*John K. Kruschke*

© John K. Kruschke, 2013

1

### Outline of Talk:

- Bayesian reasoning generally.
- Bayesian estimation applied to two groups. Rich information.
- The NHST  $t$  test: perfidious  $p$  values and the con game of confidence intervals.
- Conclusion: Bayesian estimation supersedes NHST.

© John K. Kruschke, 2013

2

## Bayesian Reasoning

The role of data is to re-allocate credibility:

**Prior Credibility** with **New Data**  
→ **Posterior Credibility**

*via Bayes' rule*

© John K. Kruschke, 2013

3

## Bayesian Reasoning

The role of data is to re-allocate credibility:

**Bayesian reasoning in everyday life is  
intuitive:**

© John K. Kruschke, 2013

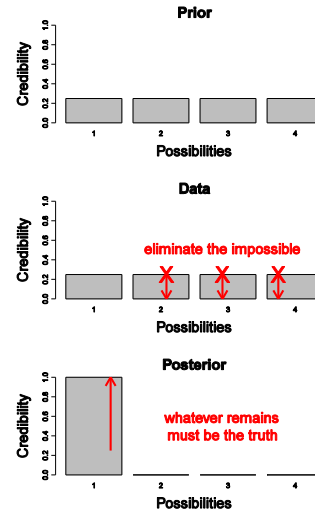
4

# Bayesian Reasoning

The role of data is to re-allocate credibility:

Bayesian reasoning in everyday life is intuitive:

**Sherlock Holmes:** "How often have I said to you that when you have eliminated the impossible, whatever remains, however improbable, must be the truth?" (Doyle, 1890)



© John K. Kruschke, 2013

5

# Bayesian Reasoning

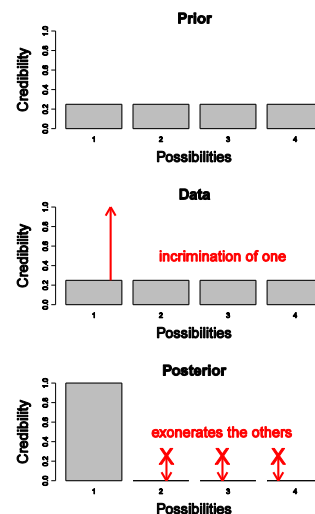
The role of data is to re-allocate credibility:

Bayesian reasoning in everyday life is intuitive:

**Sherlock Holmes:** "How often have I said to you that when you have eliminated the impossible, whatever remains, however improbable, must be the truth?" (Doyle, 1890)

**Judicial exoneration:** For unaffiliated suspects, the incrimination of one exonerates the others.

Credibility of the claim that the suspect committed the crime.



© John K. Kruschke, 2013

6

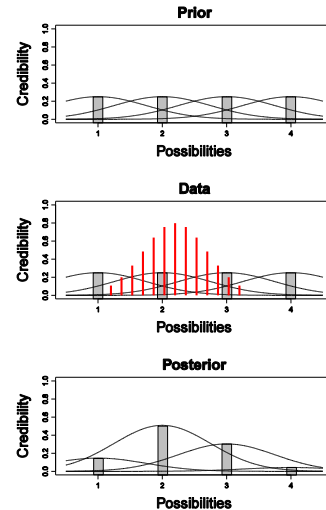
# Bayesian Data Analysis

The role of data is to re-allocate credibility:

**Bayesian reasoning in data analysis is intuitive:**

*Possibilities* are *parameter values* in a model, such as the *mean* of a normal distribution.

We reallocate credibility to parameter values that are consistent with the data.



© John K. Kruschke, 2013

7

# Bayesian Data Analysis

The role of data is to re-allocate credibility:

1. Define a meaningful descriptive model.
2. Establish prior credibility regarding parameter values in the model. The prior credibility must be acceptable to a skeptical scientific audience.
3. Collect data.
4. Use Bayes' rule to re-allocate credibility to parameter values that are most consistent with the data.

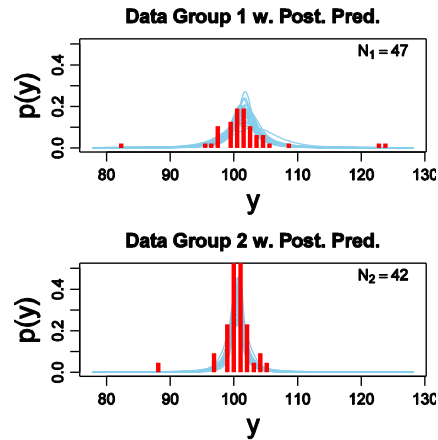
© John K. Kruschke, 2013

8

## Robust Bayesian estimation for comparing two groups

Consider two groups;  
e.g.,  
IQ of “smart drug” group  
and of control group.

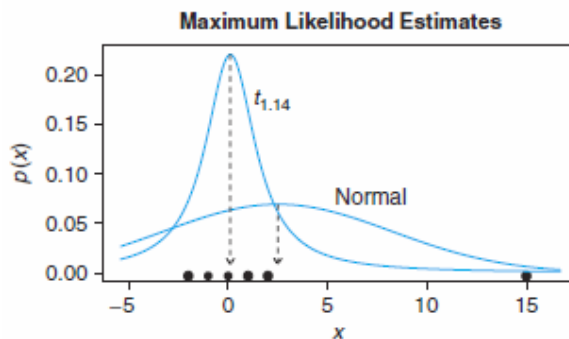
Step 1: Define a model  
for describing the data.



© John K. Kruschke, 2013

10

## Descriptive distribution for data with outliers



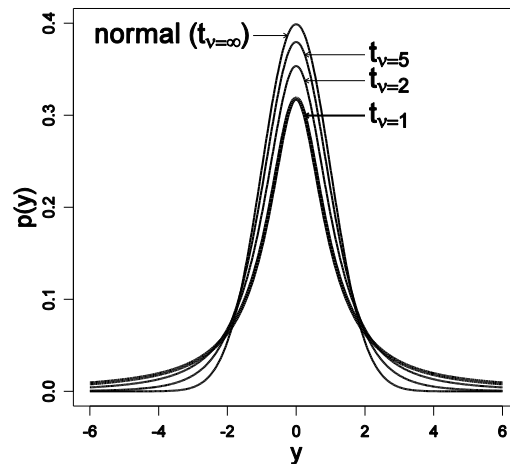
Normal is  
pulled by  
outliers, but  $t$   
distribution is  
not.

$t$  distribution is used here as a description of data,  
NOT as a sampling distribution for  $p$  values!

© John K. Kruschke, 2013

11

## Descriptive distribution for data with outliers

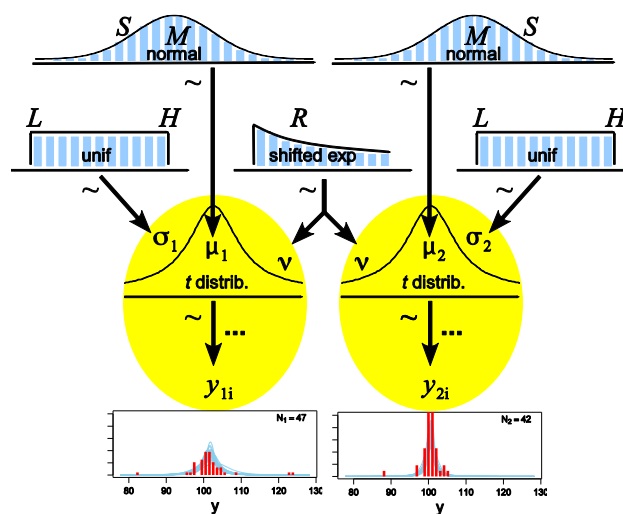


The  $t$  distribution has normality controlled by the parameter  $v$ .

© John K. Kruschke, 2013

13

## Robust Bayesian estimation for comparing two groups

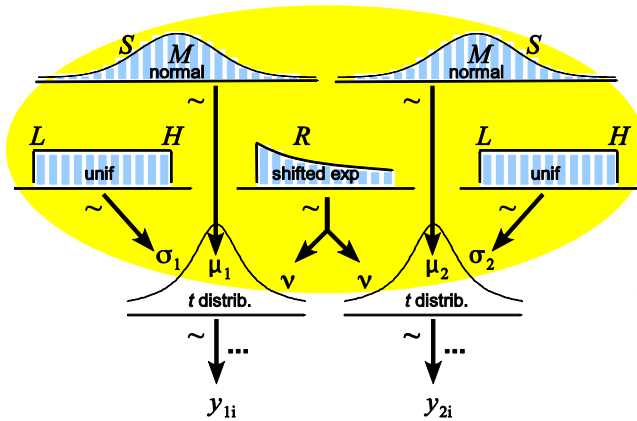


The data from each group are described by  $t$  distributions, using five parameters altogether.

© John K. Kruschke, 2013

14

## Robust Bayesian estimation for comparing two groups

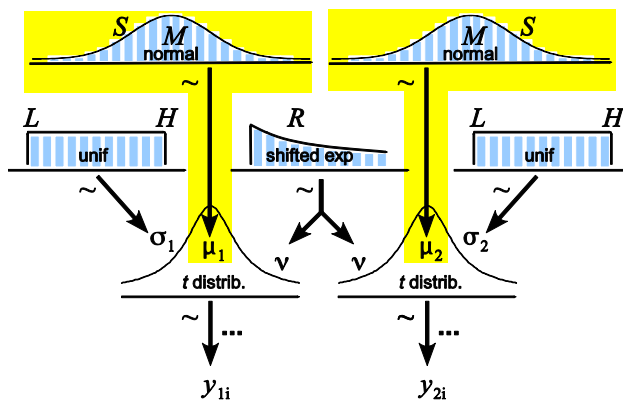


Step 2: Specify  
the prior.

© John K. Kruschke, 2013

15

## Robust Bayesian estimation for comparing two groups

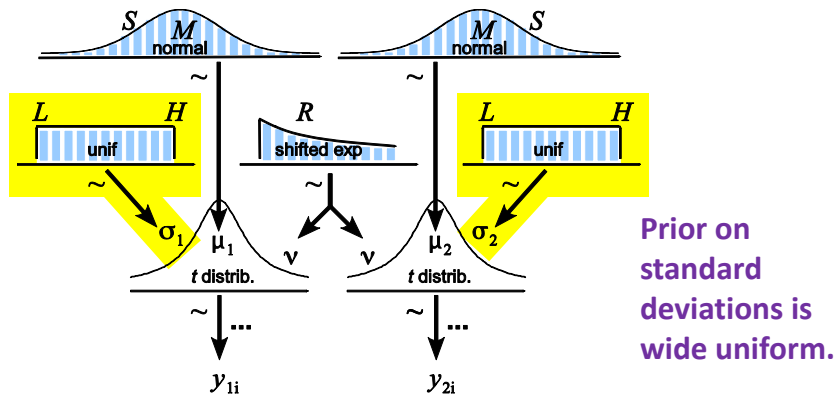


Prior on means  
is wide normal.

© John K. Kruschke, 2013

16

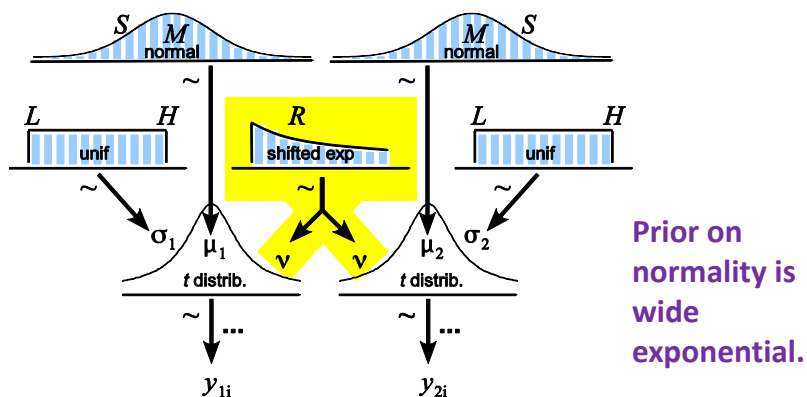
## Robust Bayesian estimation for comparing two groups



© John K. Kruschke, 2013

17

## Robust Bayesian estimation for comparing two groups

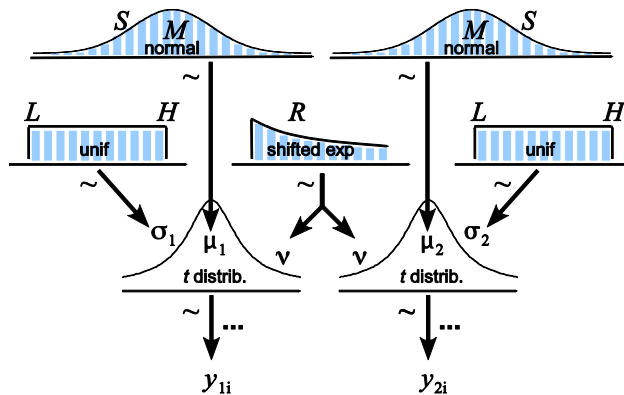


© John K. Kruschke, 2013

18



## Robust Bayesian estimation for comparing two groups

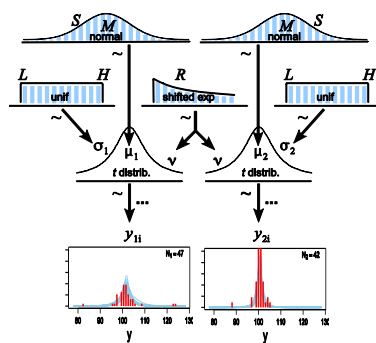


Parameter distributions will be represented by histograms: A huge number of representative parameter values.

© John K. Kruschke, 2013

19

## Step 3: Collect Data.

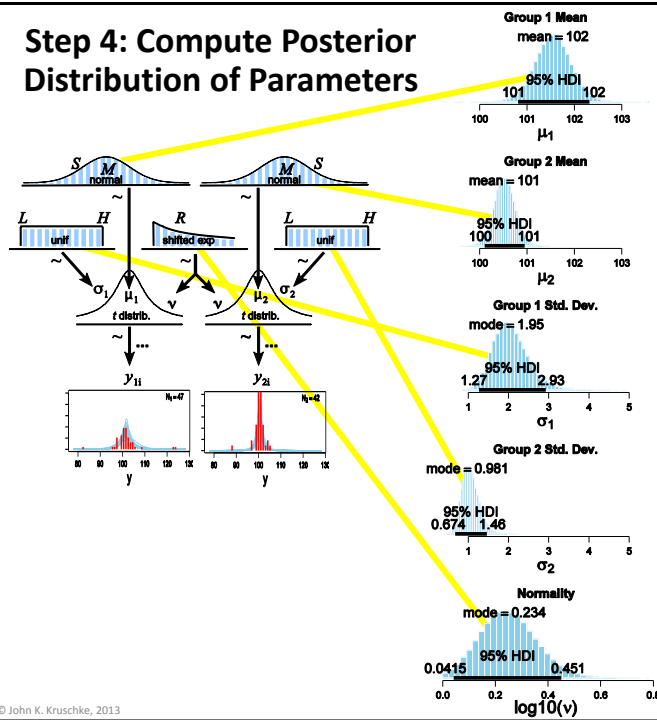


One fixed data set,  
shown as red  
histograms.

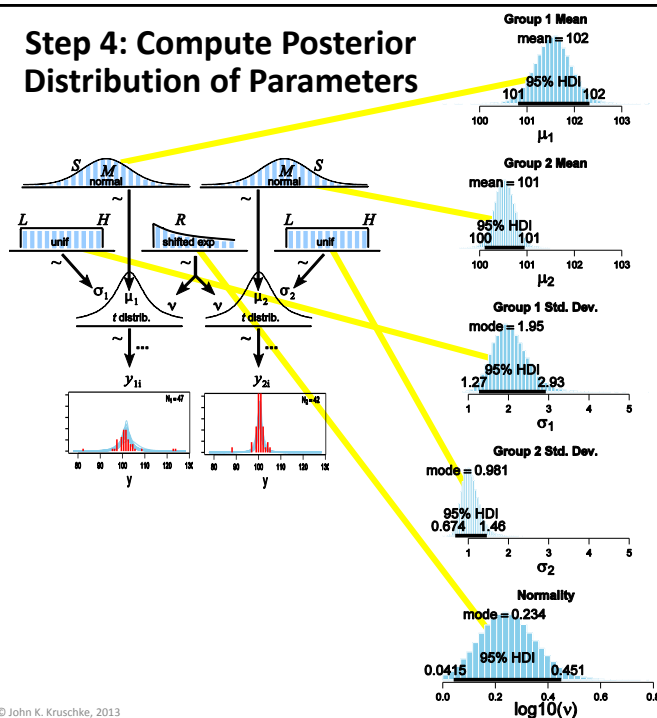
© John K. Kruschke, 2013

20

## Step 4: Compute Posterior Distribution of Parameters



## Step 4: Compute Posterior Distribution of Parameters



**Important:**  
These are histograms  
of parameter values  
from the posterior  
distribution:  
A huge number of  
*combinations* of  
 $\mu_1, \mu_2, \sigma_1, \sigma_2, v$   
that are jointly  
credible given the  
data.

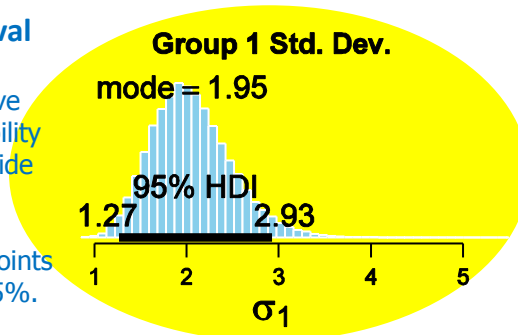
These are *not* data  
distributions, and  
*not* sampling  
distributions from a  
null hypothesis.

**95% HDI:****Highest density interval**

Points within the HDI have higher credibility (probability density) than points outside the HDI.

The total probability of points within the 95% HDI is 95%.

Points outside the HDI may be deemed not credible.



© John K. Kruschke, 2013

23

**Robust Bayesian estimation for comparing two groups**

Differences between groups?

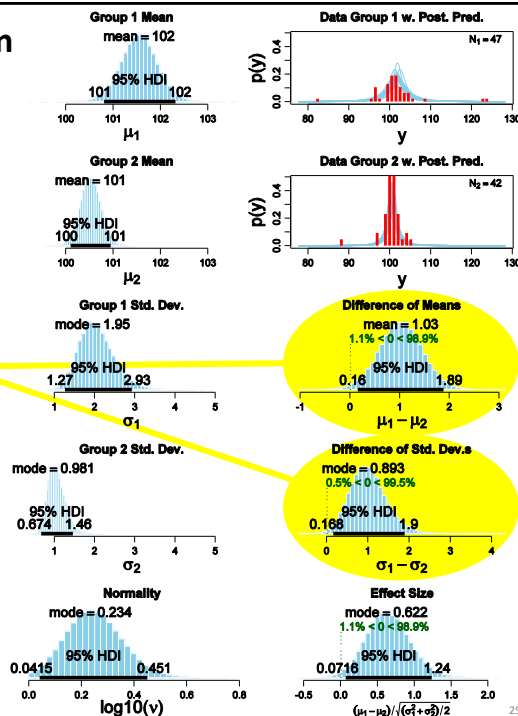
Compute  $\mu_1 - \mu_2$

and  $\sigma_1 - \sigma_2$

at each of the many credible combinations.

Here, both differences are credibly non-zero.

(NHST would require two tests...)



© John K. Kruschke, 2013

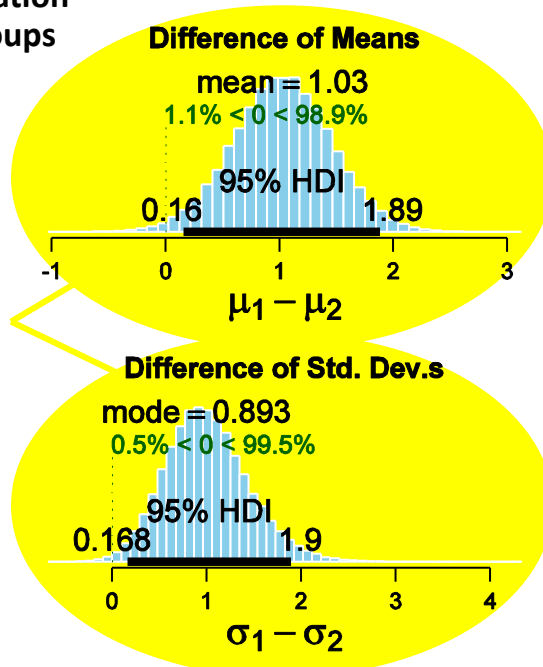
25

## Robust Bayesian estimation for comparing two groups

Differences between  
groups?  
Compute  $\mu_1 - \mu_2$   
and  $\sigma_1 - \sigma_2$   
at each of the many  
credible combinations.

Here, both differences  
are credibly non-zero.

(NHST would require  
two tests...)

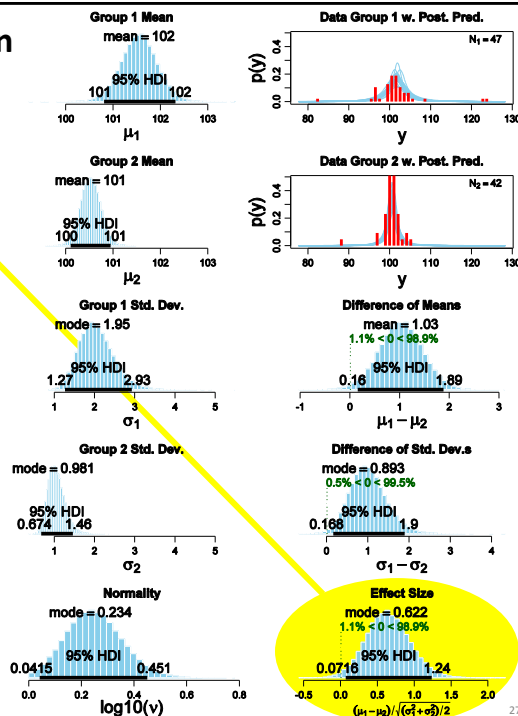


© John K. Kruschke, 2013

26

## Robust Bayesian estimation for comparing two groups

Complete distribution  
on effect size!

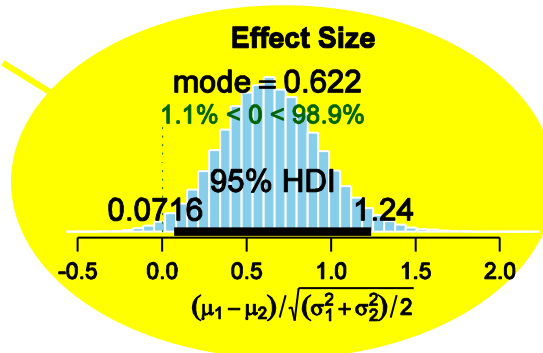


© John K. Kruschke, 2013

27

## Robust Bayesian estimation for comparing two groups

Complete distribution  
on effect size!



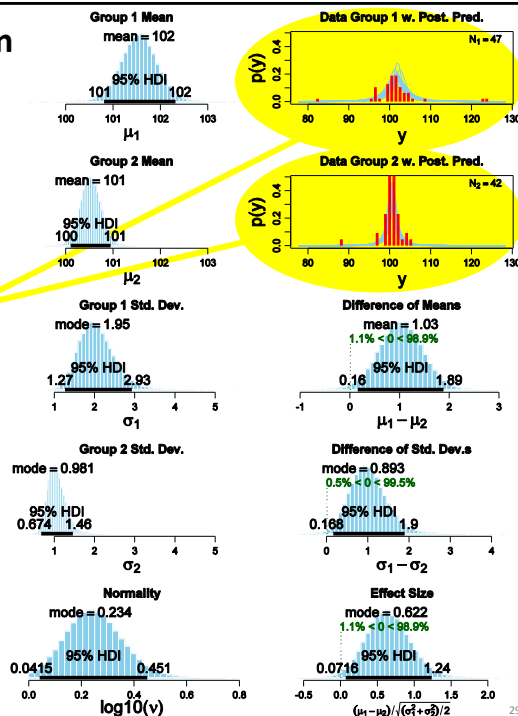
© John K. Kruschke, 2013

28

## Robust Bayesian estimation for comparing two groups

Are the data described  
well by the model?

Superimpose a  
smattering of credible  
descriptive distributions  
on data.  
= “posterior predictive  
check”



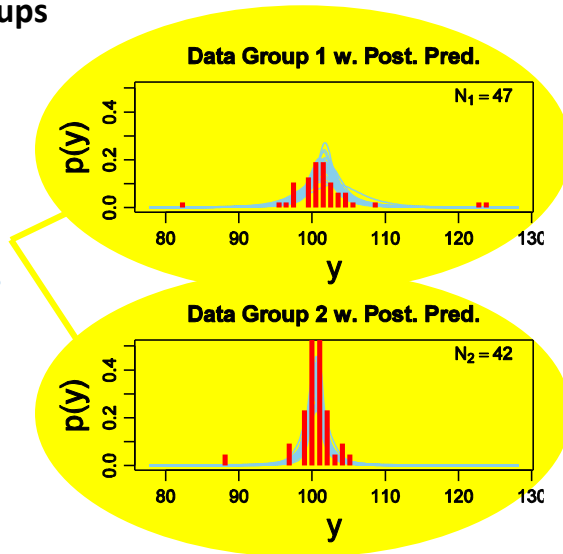
© John K. Kruschke, 2013

29

## Robust Bayesian estimation for comparing two groups

Are the data described  
well by the model?

Superimpose a  
smattering of credible  
descriptive distributions  
on data.  
= “posterior predictive  
check”



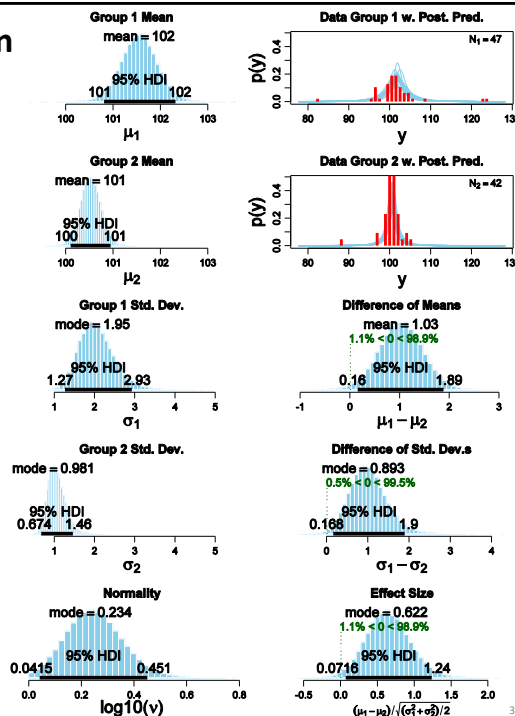
© John K. Kruschke, 2013

30

## Robust Bayesian estimation for comparing two groups

Summary:

- Complete distribution of credible parameter values (not merely point estimate with ends of confidence interval).
- Decisions about multiple aspects of parameters (without reference to  $p$  values).
- Flexible descriptive model, robust to outliers (unlike NHST  $t$  test).



© John K. Kruschke, 2013

31

## Computer Software:

Packaged for easy use!  
Underlying program is never seen.

```
source("BEST.R") # load the program

# Specify data as vectors (replace with your own data):
y1 = c(101,100,102,104,102,97,105,105,98,...,101)
y2 = c(99,101,100,101,102,100,97,101,104,...,99)

# Run the Bayesian analysis:
mcmcChain = BESTmcmc( y1 , y2 )

# Plot the results of the Bayesian analysis:
BESTplot( y1 , y2 , mcmcChain )
```

© John K. Kruschke, 2013

32

## An example of a $t$ test:

### Data:

Group 1: 5.70 5.40 5.75 5.25 4.25 4.74;  $M1 = 5.18$

Group 2: 4.55 4.98 4.70 4.78 3.26 3.67;  $M2 = 4.32$

$t = 2.33$

Show of hands please:

Who bets that  $p < .05$  ?      Who bets that  $p > .05$  ?

© John K. Kruschke, 2013

33

## An example of a $t$ test:

### Data:

Group 1: 5.70 5.40 5.75 5.25 4.25 4.74;  $M1 = 5.18$

Group 2: 4.55 4.98 4.70 4.78 3.26 3.67;  $M2 = 4.32$

$t = 2.33$

Show of hands please:

Who bets that  $p < .05$  ?      Who bets that  $p > .05$  ?

You're right!

You're right!

© John K. Kruschke, 2013

34

## Null Hypothesis Significance Testing (NHST)

Consider how we draw conclusions from data:

- Collect data, *carefully insulated from our intentions*.
  - Double blind clinical designs.
  - No datum is influenced by any other datum before or after.
- Compute a summary statistic, e.g., for a difference between groups, the  $t$  statistic.
- Compute  $p$  value of  $t$ . If  $p < .05$ , declare the result to be "significant."

© John K. Kruschke, 2013

35



## Null Hypothesis Significance Testing (NHST)

Consider how we draw conclusions from data:

- Collect data, *carefully insulated from our intentions*.
  - Double blind clinical d
  - No datum is influence
- Compute a summary s between groups, the  $t$  statistic.

Value of  $p$  depends on the intention of the experimenter!

- Compute  $p$  value of  $t$ . If  $p < .05$ , declare the result to be "significant."

© John K. Kruschke, 2013

36

## The road to NHST is paved with good intentions.

The  $p$  value is the probability that the actual sample statistic, or a result more extreme, would be obtained from the null hypothesis, *if the **intended** experiment were repeated *ad infinitum*.*

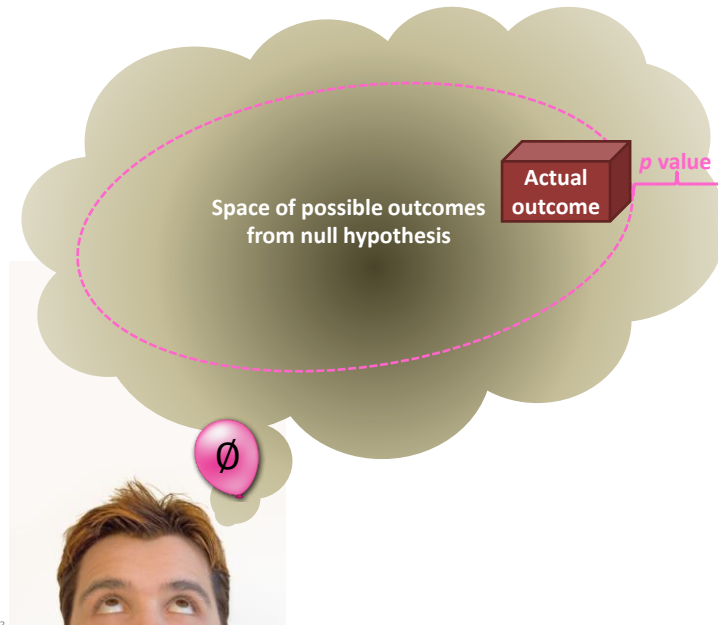
$$p \text{ value} = p(|t_{\text{null}}| > |t_{\text{act}}|)$$

for  $t_{\text{null}}$  sampled according to the intended experiment

© John K. Kruschke, 2013

37

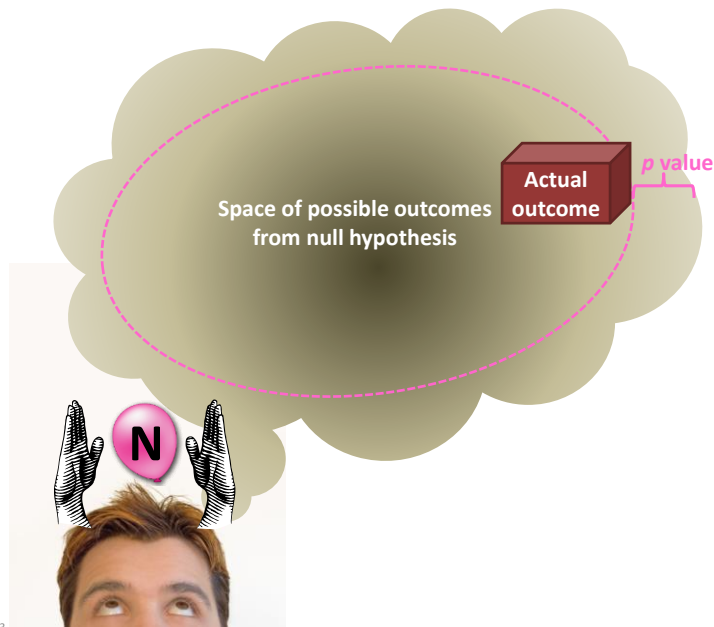
“The”  $p$  value...



© John K. Kruschke, 2013

38

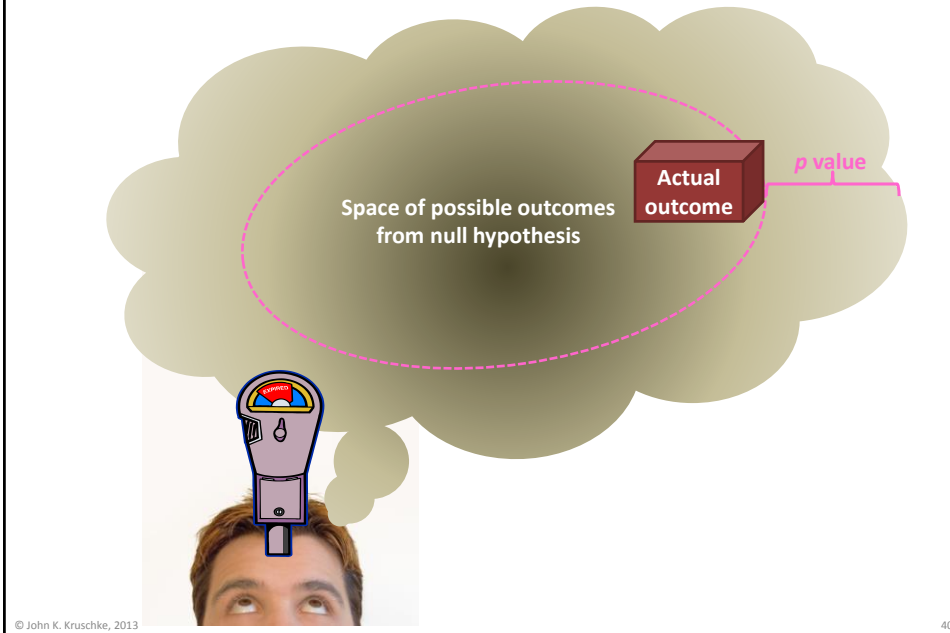
$p$  value for intention to sample until  $N$



© John K. Kruschke, 2013

39

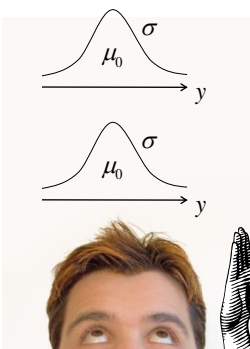
## $p$ value for intention to sample until Time



40

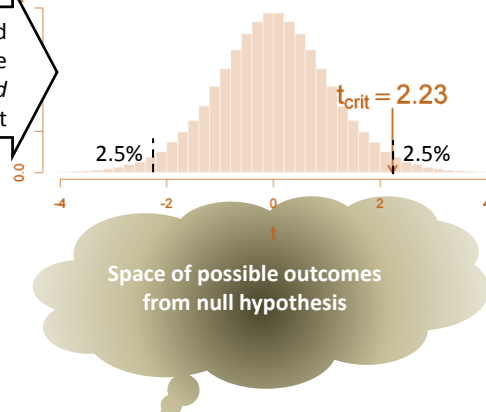
## The distribution of $t$ when the intended experiment is repeated many times

**Null Hypothesis:**  
Groups are identical



Many simulated  
repetitions of the  
*intended*  
experiment

Fixed  $N=6$  per group (x2 groups)



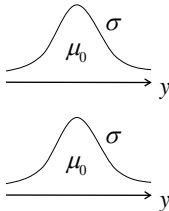
© John K. Kruschke, 2013

42

## The distribution of $t$

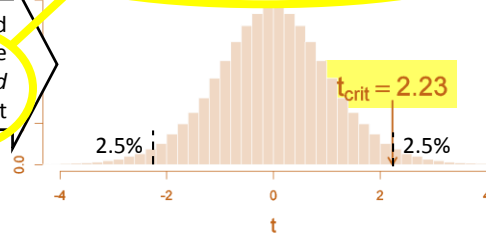
when the intended experiment is repeated many times

**Null Hypothesis:**  
Groups are identical



Many simulated repetitions of the intended experiment

Fixed  $N=6$  per group (x2 groups)

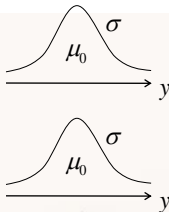


© John K. Kruschke, 2013

43

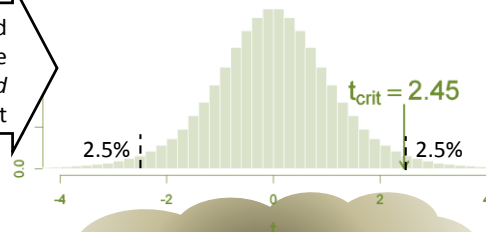
## The intention to collect data until the end of the week

**Null Hypothesis:**  
Groups are identical



Many simulated repetitions of the intended experiment

Fixed Duration=2 weeks (x6/week)



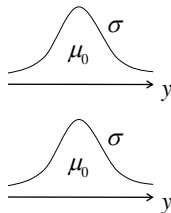
Space of possible outcomes from null hypothesis

© John K. Kruschke, 2013

44

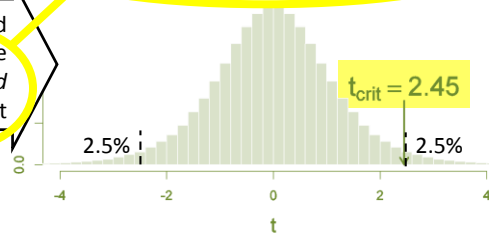
# The intention to collect data until the end of the week

**Null Hypothesis:**  
Groups are identical



Many simulated repetitions of the intended experiment

Fixed Duration=2 weeks (x6/week)



© John K. Kruschke, 2013

45

## Two labs collect the same data:

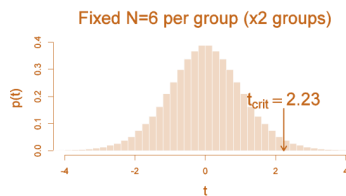
**Lab A: Collect data until N=6 per group.**

**Data:**

Group 1: 5.70 5.40 5.75 5.25 4.25 4.74; M1 = 5.18

Group 2: 4.55 4.98 4.70 4.78 3.26 3.67; M2 = 4.32

**t = 2.33**



**Lab A: Reject the null.**

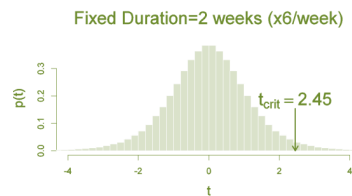
**Lab B: Collect data for two weeks.**

**Data:**

Group 1: 5.70 5.40 5.75 5.25 4.25 4.74; M1 = 5.18

Group 2: 4.55 4.98 4.70 4.78 3.26 3.67; M2 = 4.32

**t = 2.33**



**Lab B: Do not reject the null.**

© John K. Kruschke, 2013

47

## The *real* use of the Neuralyzer:

You *meant* to collect data until  $N=12$  !

Now *that's* significant!



© John K. Kruschke, 2013

48

## Problem is not solved by “fixing” the intention

- All we need to do is decide in advance exactly what our intention is (or use a Neuralyzer after the fact), and have everybody chant a mantra to keep that intention fixed in their minds while the experiment is being conducted. Right?
- Wrong. The data don't know our intention, and the same data could have been collected under many other intentions.

© John K. Kruschke, 2013

49

## The intention to examine data thoroughly

Many experiments involve multiple groups, and **multiple comparisons** of means.

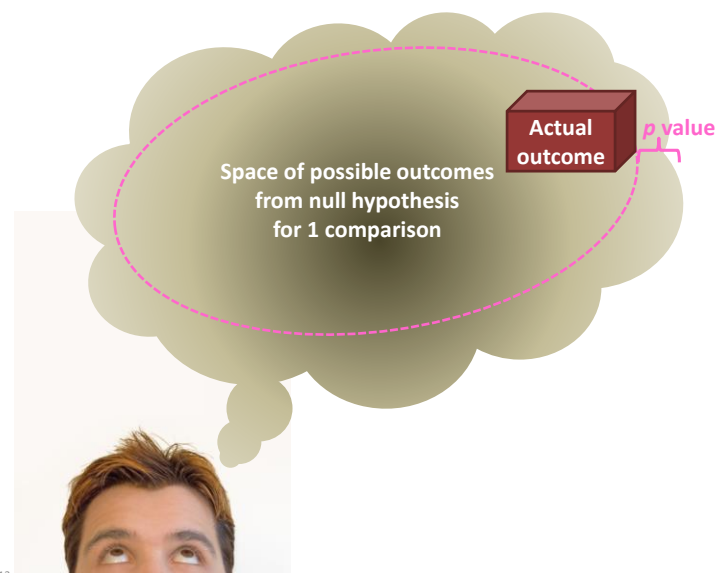
Example: Consider 2 different drugs from chemical family A, 2 different drugs from chemical family B, and a placebo group. Lots of possible comparisons...

Problem: With every test, there is possibility of false alarm!  
False alarms are bad; therefore, keep the experimentwise false alarm rate down to 5%.

© John K. Kruschke, 2013

50

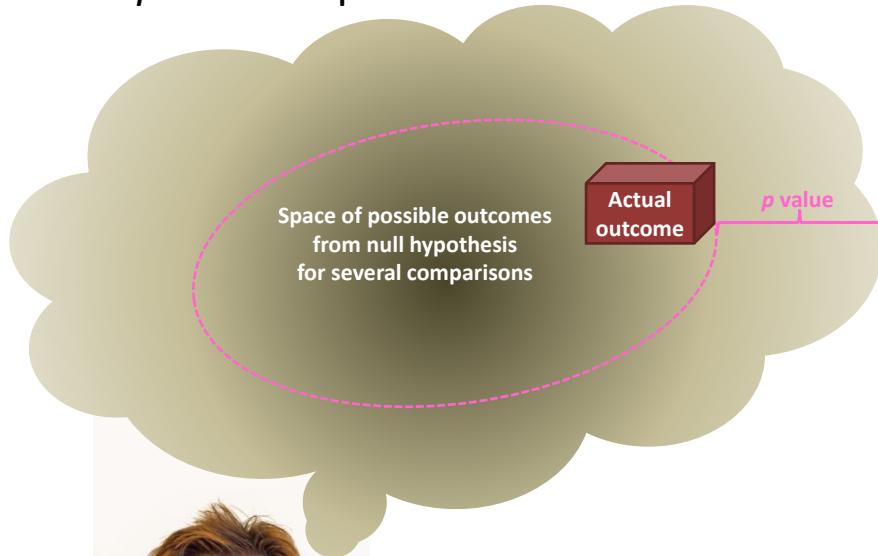
“The”  $p$  value depends on intended tests:



© John K. Kruschke, 2013

51

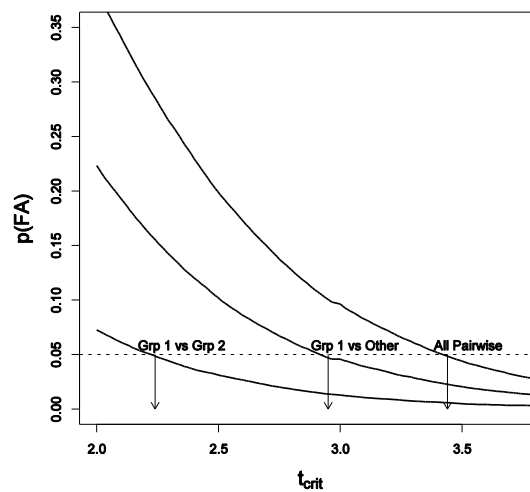
“The”  $p$  value depends on intended tests:



© John K. Kruschke, 2013

52

Experimentwise false alarm rate



© John K. Kruschke, 2013

53



## Multiple Corrections for Multiple Comparisons

Begin: **Is goal to identify the best treatment?**

Yes: Use **Hsu's method**.

No: **Contrasts between control group and all other groups?**

Yes: Use **Dunnett's method**.

No: **Testing all pairwise and no complex comparisons (either planned or post hoc) and choosing to test only some pairwise comparisons post hoc?**

Yes: Use **Tukey's method**.

No: **Are all comparisons planned?**

Yes: Use **Scheffe's method**.

No: Is Bonferroni critical value less than Scheffe critical value?

Yes: Use **Bonferroni's method**.

No: Use Scheffe's method (or, prior to collecting the data, reduce the number of contrasts to be tested).

Adapted from Maxwell & Delaney (2004). Designing experiments and analyzing data: A model comparison perspective. Erlbaum.

© John K. Kruschke, 2013

54

## Multiple Corrections for Multiple Comparisons

Begin: **Is goal to identify the best treatment?**

Yes: Use **Hsu's method**.

No: **Contrasts between control group and all other groups?**

Yes: Use **Dunnett's method**.

No: **Testing all pairwise and no complex comparisons (either planned or post hoc) and choosing to test only some pairwise comparisons post hoc?**

Yes: Use **Tukey's method**.

No: **Are all comparisons planned?**

Yes: Use **Scheffe's method**.

No: Is Bonferroni critical value less than Scheffe critical value?

Yes: Use ~~**Bonferroni's method**~~.

! **No: Use Scheffe's method (or, prior to collecting the data, reduce the number of contrasts to be tested).**

Adapted from Maxwell & Delaney (2004). Designing experiments and analyzing data: A model comparison perspective. Erlbaum.

© John K. Kruschke, 2013

55

## Good intentions make any result *insignificant*

- Consider an experiment with two groups.
- Collect data; compute  $t$  test on difference of means. Suppose it yields  $p < .05$
- Now, think thoroughly about all the other comparison groups and other experiment groups you should and could meaningfully run.
- Earnestly intend to run them eventually, and to compare your current results with those results.
- *Poof! Your current data are no longer significantly different.*

© John K. Kruschke, 2013

56



© John K. Kruschke, 2013

57

## Good intentions make many results *significant*

- Consider an experiment with two groups.
- Collect data; compute  $t$  test on difference of means, using df corresponding to actual  $N$ . Suppose  $p > .05$ , but not by much.
- *You had intended to collect a much larger sample size, but you were unexpectedly interrupted.*
- Use the larger intended  $N$  for df in the  $t$  test.
- *Poof! Your current data are now significantly different!*

© John K. Kruschke, 2013

58

## ? Confidence Intervals provide no confidence ?

### General definition of CI:

95% CI is the range of parameter values (e.g.,  $\mu_1 - \mu_2$ ) that would not be rejected by  $p < .05$

Hence, *the 95% CI is as ill-defined as the  $p$  value.*

We see this dramatically in confidence intervals corrected for multiple comparisons.

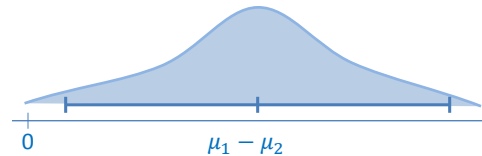
© John K. Kruschke, 2013

60

## ? Confidence Intervals provide no confidence ?

*Confidence intervals provide no distributional information:*

We have no idea whether a point at the limit of the confidence interval is any less credible than a point in the middle of the interval.



© John K. Kruschke, 2013

61

## ? Confidence Intervals provide no confidence ?

*Confidence intervals provide no distributional information:*

We have no idea whether a point at the limit of the confidence interval is any less credible than a point in the middle of the interval.

Implies

vast range for predictions of new data, and  
“virtually unknowable” power.

© John K. Kruschke, 2013

62

## NHST autopsy

- $p$  values are ill-defined: depend on sampling intentions of data collector. Any set of data has many different  $p$  values.
- Confidence intervals are as ill-defined as  $p$  values because they are defined in terms of  $p$  values.
- Confidence intervals carry no distributional information.

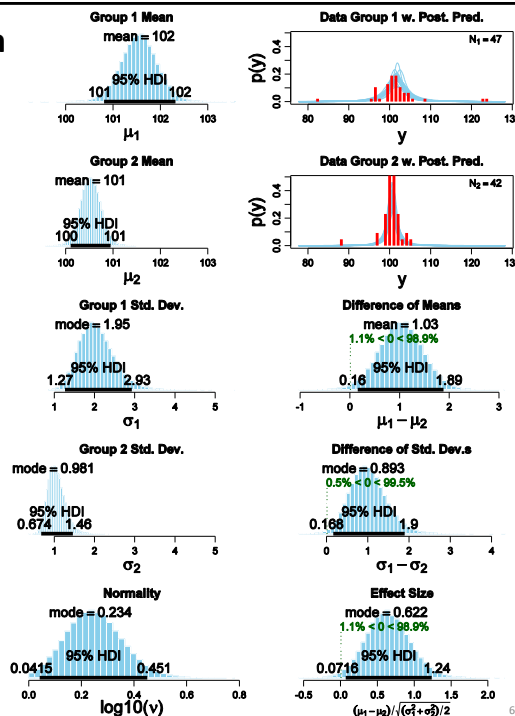
© John K. Kruschke, 2013

63

### Recall Bayesian estimation for comparing two groups

#### Summary:

- Complete distribution of credible parameter values (not merely point estimate with ends of confidence interval).
- Decisions about multiple aspects of parameters (without reference to  $p$  values).
- Flexible descriptive model, robust to outliers (unlike NHST  $t$  test).



© John K. Kruschke, 2013

64