# *Doing Bayesian Data Analysis*



$p(\theta|D) \qquad p(D|\theta) \qquad p(\theta) \qquad p(D)$

*John K. Kruschke*

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.

2

# Bayesian Reasoning

The role of data is to re-allocate credibility:

**Prior Credibility**   with   **New Data**
$\rightarrow$ **Posterior Credibility**

*via Bayes' rule*

3

# Bayesian Reasoning
The role of data is to re-allocate credibility:

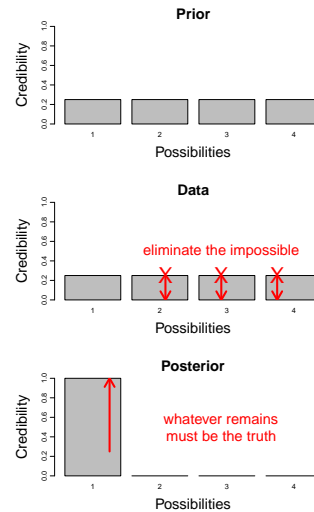**Bayesian reasoning in everyday life is intuitive:**

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)



**Prior**

Credibility — Possibilities

**Data**

eliminate the impossible

**Posterior**

whatever remains must be the truth
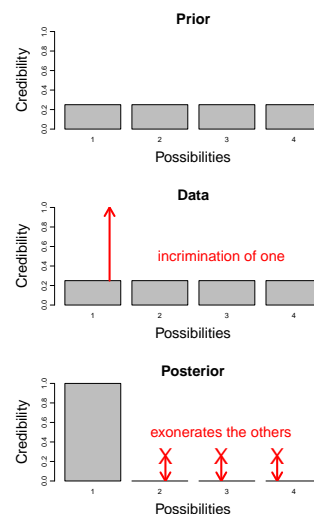
5

# 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.



**Prior**

Credibility — Possibilities

**Data**

incrimination of one
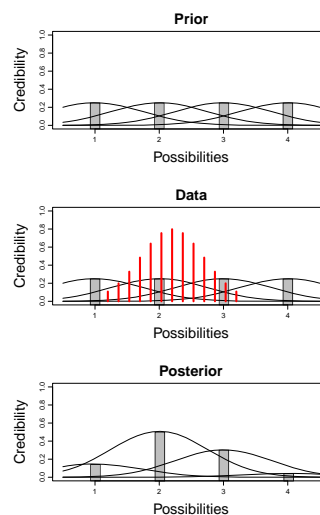
**Posterior**

exonerates the others

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.

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.**
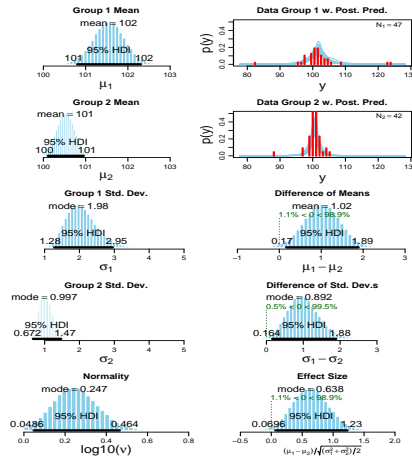4. **Use Bayes' rule to re-allocate credibility to parameter values that are most consistent with the data.**

8

**Robust Bayesian estimation
for comparing two groups**

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

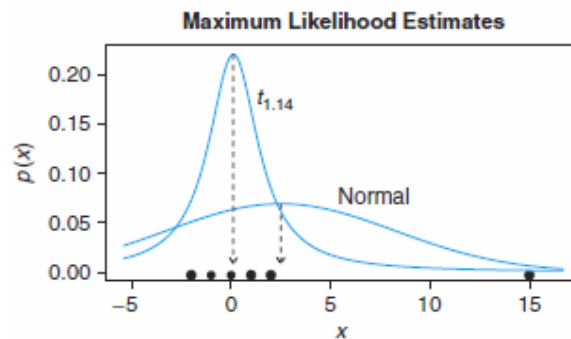Step 1: Define a model
for describing the data.

10

---

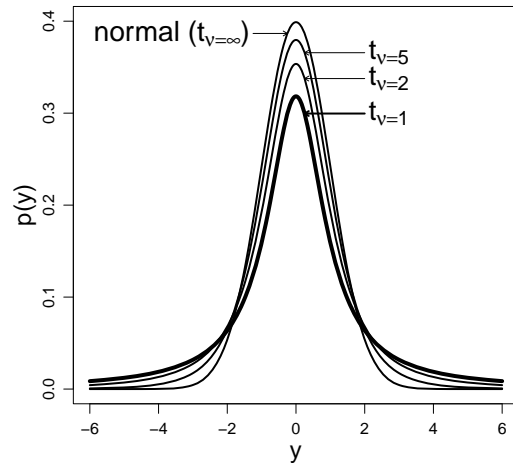# Descriptive distribution for data with outliers



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!

11

## Descriptive distribution for data with outliers



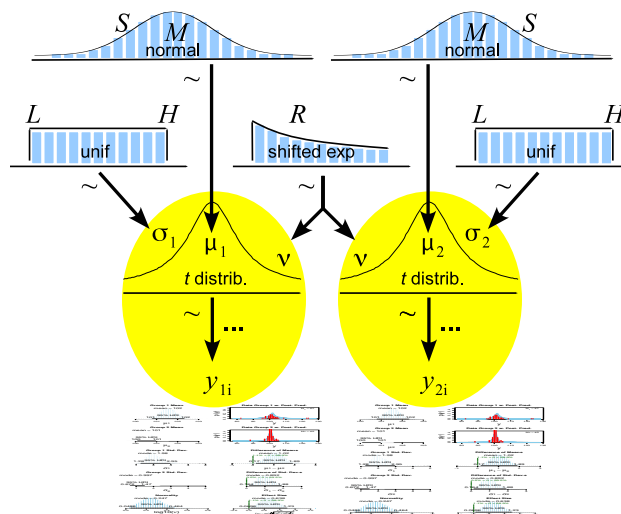The *t* distribution has normality controlled by the parameter ν.
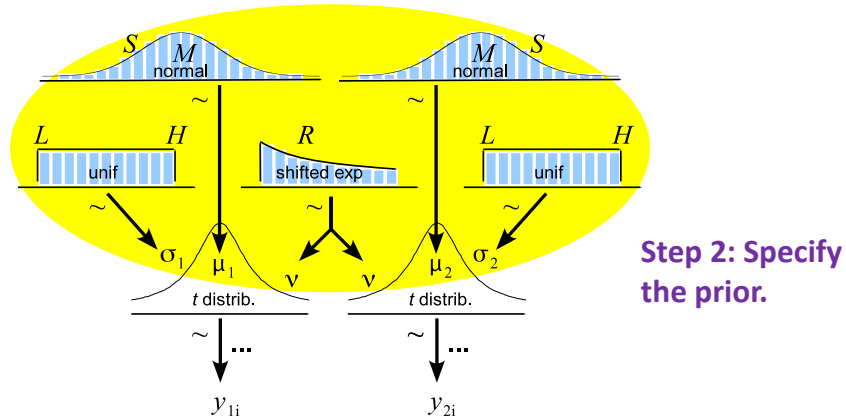
**Robust Bayesian estimation for comparing two groups**



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

**Robust Bayesian estimation
for comparing two groups**



**Step 2: Specify
the prior.**

15

**Robust Bayesian estimation
for comparing two groups**



**Prior on means
is wide normal.**

16

7

## Robust Bayesian estimation
## for comparing two groups



**Prior on standard deviations is wide uniform.**

## Robust Bayesian estimation
## for comparing two groups



**Prior on normality is wide exponential.**

## Robust Bayesian estimation for comparing two groups



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

19

# Step 3: Collect Data.



**One fixed data set, shown as red histograms.**

20

## Slide 23

*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.

## Slide 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…)

## 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, Oct. 2012

26

## Robust Bayesian estimation for comparing two groups

**Complete distribution on effect size!**



© John K. Kruschke, Oct. 2012

27

## Robust Bayesian estimation for comparing two groups

**Complete distribution on effect size!**

## 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"

# 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"

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)**.**

31

14

## Computer Software:

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

# Specify data as vectors (replace with your own data):
y1 = c(101,100,102,104,102,97,105,105,98,101,100,123,105,
       109,102,82,102,100,102,102,101,102,102,103,103,97,
       96,103,124,101,101,100,101,101,104,100,101)
y2 = c(99,101,100,101,102,100,97,101,104,101,102,102,100,
       104,100,100,100,101,102,103,97,101,101,100,101,99,
       101,100,99,101,100,102,99,100,99)

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

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

32

## Robust Bayesian estimation for comparing two groups

**Download the programs from**
**http://www.indiana.edu/~kruschke/BEST/BEST.zip**

Now for a look
under the hood



http://www.autonationconnect.com/2010/07/backseat-mechanic-under-the-hoo/

33

# Doing it with JAGS

"JAGS" = Just Another Gibbs Sampler
but other sampling methods are incorporated.

| R programming language | rjags commands | ⇄ | JAGS executables |
|---|---|---|---|

JAGS makes it easy. You specify only the
• prior function
• likelihood function
*and JAGS does the rest! You do no math, no selection of sampling methods.*

34

# JAGS and BUGS



R

rjags
runs on MacOS, Linux, & Windows  →  JAGS

BRugs
requires 32-bit Windows  →  OpenBUGS

R2WinBUGS  →  WinBUGS
requires Windows

35

## Installation: See Blog Entry

http://doingbayesiandataanalysis.blogspot.com/2012/01/complete-steps-for-installing-software.html



36

**Robust Bayesian estimation for comparing two groups**

Program BEST.R:  JAGS model specification.

```
model {
  for ( i in 1:Ntotal ) {
    y[i] ~ dt( mu[x[i]] , tau[x[i]] , nu )
  }
  for ( j in 1:2 ) {
    mu[j] ~ dnorm( muM , muP )
    tau[j] <- 1/pow( sigma[j] , 2 )
    sigma[j] ~ dunif( sigmaLow , sigmaHigh )
  }
  nu <- nuMinusOne+1
  nuMinusOne ~ dexp(1/29)
}
```



37

17

**Robust Bayesian estimation for comparing two groups**

Program BEST.R: JAGS model specification.

```
model {
  for ( i in 1:Ntotal ) {
    y[i] ~ dt( mu[x[i]] , tau[x[i]] , nu )
  }
  for ( j in 1:2 ) {
    mu[j] ~ dnorm( muM , muP )
    tau[j] <- 1/pow( sigma[j] , 2 )
    sigma[j] ~ dunif( sigmaLow , sigmaHigh )
  }
  nu <- nuMinusOne+1
  nuMinusOne ~ dexp(1/29)
}
```

© John K. Kruschke, Sept. 2012                                    38



**Robust Bayesian estimation for comparing two groups**

Program BEST.R: JAGS model specification.

```
model {
  for ( i in 1:Ntotal ) {
    y[i] ~ dt( mu[x[i]] , tau[x[i]] , nu )
  }
  for ( j in 1:2 ) {
    mu[j] ~ dnorm( muM , muP )
    tau[j] <- 1/pow( sigma[j] , 2 )
    sigma[j] ~ dunif( sigmaLow , sigmaHigh )
  }
  nu <- nuMinusOne+1
  nuMinusOne ~ dexp(1/29)
}
```

Nested indexing:
x[i] is the group (1 or 2)
of the i[th] score.

© John K. Kruschke, Sept. 2012                                    39

18

**Robust Bayesian estimation for comparing two groups**

Program BEST.R: JAGS model specification.

```
model {
  for ( i in 1:Ntotal ) {
    y[i] ~ dt( mu[x[i]] , tau[x[i]] , nu )
  }
  for ( j in 1:2 ) {
    mu[j] ~ dnorm( muM , muP )
    tau[j] <- 1/pow( sigma[j] , 2 )
    sigma[j] ~ dunif( sigmaLow , sigmaHigh )
  }
  nu <- nuMinusOne+1
  nuMinusOne ~ dexp(1/29)
}
```

© John K. Kruschke, Sept. 2012                                          40



**Robust Bayesian estimation for comparing two groups**

Program BEST.R: JAGS model specification.

```
model {
  for ( i in 1:Ntotal ) {
    y[i] ~ dt( mu[x[i]] , tau[x[i]] , nu )
  }
  for ( j in 1:2 ) {
    mu[j] ~ dnorm( muM , muP )
    tau[j] <- 1/pow( sigma[j] , 2 )
    sigma[j] ~ dunif( sigmaLow , sigmaHigh )
  }
  nu <- nuMinusOne+1
  nuMinusOne ~ dexp(1/29)
}
```

© John K. Kruschke, Sept. 2012                                          41

19

**Robust Bayesian estimation for comparing two groups**

Program BEST.R: JAGS model specification.

```
model {
  for ( i in 1:Ntotal ) {
    y[i] ~ dt( mu[x[i]] , tau[x[i]] , nu )
  }
  for ( j in 1:2 ) {
    mu[j] ~ dnorm( muM , muP )
    tau[j] <- 1/pow( sigma[j] , 2 )
    sigma[j] ~ dunif( sigmaLow , sigmaHigh )
  }
  nu <- nuMinusOne+1
  nuMinusOne ~ dexp(1/29)
}
```

© John K. Kruschke, Sept. 2012                                        42

---

# Programs in R + rjags + JAGS:

**Five main sections in all programs:**

1. **Specify model** (we just did this)**.**

2. **Load data.**

3. **Initialize the MCMC chain.**

4. **Run the MCMC chain.**

5. **Examine the results.**

© John K. Kruschke, Sept. 2012                                        43

# BEST.R

```
BESTmcmc = function( y1, y2, numSavedSteps=100000, thinSteps=1, showMCMC=FALSE) {
  # This function generates an MCMC sample from the posterior distribution.
  # Description of arguments:
  # showMCMC is a flag for displaying diagnostic graphs of the chains.
  #    If F (the default), no chain graphs are displayed. If T, they are.

  require(rjags)

  #------------------------------------------------------------------------------
  # THE MODEL.
  modelString = "
  model {
    for ( i in 1:Ntotal ) {
      y[i] ~ dt( mu[x[i]] , tau[x[i]] , nu )
    }
    for ( j in 1:2 ) {
      mu[j] ~ dnorm( muM , muP )
      tau[j] <- 1/pow( sigma[j] , 2 )
      sigma[j] ~ dunif( sigmaLow , sigmaHigh )
    }
    nu <- nuMinusOne+1
    nuMinusOne ~ dexp(1/29)
  }
  " # close quote for modelString
  # Write out modelString to a text file
  writeLines( modelString , con="BESTmodel.txt" )

  #------------------------------------------------------------------------------
  # THE DATA.
  # Load the data:
  y = c( y1 , y2 ) # combine data into one vector
  x = c( rep(1,length(y1)) , rep(2,length(y2)) ) # create group membership code
  Ntotal = length(y)
  # Specify the data in a list, for later shipment to JAGS:
  dataList = list(
    y = y ,
    x = x ,
    Ntotal = Ntotal ,
```

44

# BEST.R

```
BESTmcmc = function( y1, y2, numSavedSteps=100000, thinSteps=1, showMCMC=FALSE) {
  # This function generates an MCMC sample from the posterior distribution.
  # Description of arguments:
  # showMCMC is a flag for displaying diagnostic graphs of the chains.
  #    If F (the default), no chain graphs are displayed. If T, they are.

  require(rjags)

  #------------------------------------------------------------------------------
  # THE MODEL.
  modelString = "
  model {
    for ( i in 1:Ntotal ) {
      y[i] ~ dt( mu[x[i]] , tau[x[i]] , nu )
    }
    for ( j in 1:2 ) {
      mu[j] ~ dnorm( muM , muP )
      tau[j] <- 1/pow( sigma[j] , 2 )
      sigma[j] ~ dunif( sigmaLow , sigmaHigh )
    }
    nu <- nuMinusOne+1
    nuMinusOne ~ dexp(1/29)
  }
  " # close quote for modelString
  # Write out modelString to a text file
  writeLines( modelString , con="BESTmodel.txt" )

  #------------------------------------------------------------------------------
  # THE DATA.
  # Load the data:
  y = c( y1 , y2 ) # combine data into one vector
  x = c( rep(1,length(y1)) , rep(2,length(y2)) ) # create group membership code
  Ntotal = length(y)
  # Specify the data in a list, for later shipment to JAGS:
  dataList = list(
    y = y ,
    x = x ,
    Ntotal = Ntotal ,
```
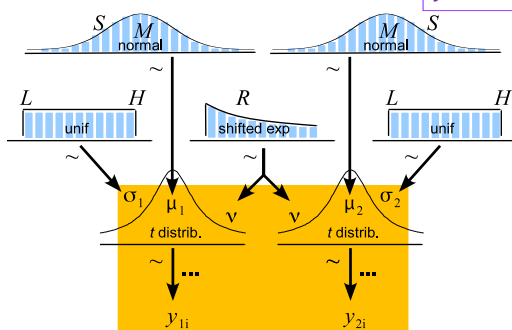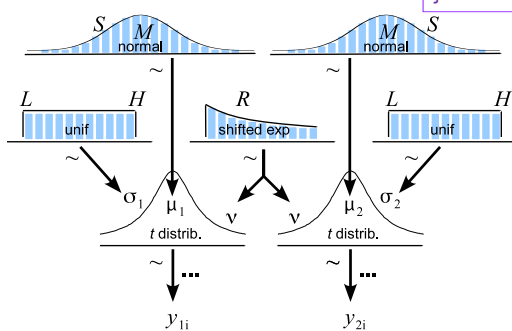
45

```
      }
      nu <- nuMinusOne+1
      nuMinusOne ~ dexp(1/29)                BEST.R
    }
    " # close quote for modelString
    # Write out modelString to a text file
    writeLines( modelString , con="BESTmodel.txt" )

    #------------------------------------------------------------------------
    # THE DATA.
    # Load the data:
    y = c( y1 , y2 ) # combine data into one vector
    x = c( rep(1,length(y1)) , rep(2,length(y2)) ) # create group membership code
    Ntotal = length(y)
    # Specify the data in a list, for later shipment to JAGS:
    dataList = list(
      y = y ,
      x = x ,
      Ntotal = Ntotal ,
      muM = mean(y) ,
      muP = 0.000001 * 1/sd(y)^2 ,
      sigmaLow = sd(y) / 1000 ,
      sigmaHigh = sd(y) * 1000
    )

    #------------------------------------------------------------------------
    # INTIALIZE THE CHAINS.
    # Initial values of MCMC chains based on data:
    mu = c( mean(y1) , mean(y2) )
    sigma = c( sd(y1) , sd(y2) )
    # Regarding initial values in next line: (1) sigma will tend to be too big if
    # the data have outliers, and (2) nu starts at 5 as a moderate value. These
    # initial values keep the burn-in period moderate.
    initsList = list( mu = mu , sigma = sigma , nuMinusOne = 4 )

    #------------------------------------------------------------------------
    # RUN THE CHAINS

    parameters = c( "mu" , "sigma" , "nu" )     # The parameters to be monitored
    adaptSteps = 500                  # Number of steps to "tune" the samplers
    burnInSteps = 1000
    nChains = 3
```

© John K. Kruschke, Sept. 2012                                                    46

**BEST.R**

```
                # initial values keep the burn-in period moderate.
                initsList = list( mu = mu , sigma = sigma , nuMinusOne = 4 )
        #------------------------------------------------------------------------------
        # RUN THE CHAINS

        parameters = c( "mu" , "sigma" , "nu" )     # The parameters to be monitored
        adaptSteps = 500                 # Number of steps to "tune" the samplers
        burnInSteps = 1000
        nChains = 3
        nIter = ceiling( ( numSavedSteps * thinSteps ) / nChains )
        # Create, initialize, and adapt the model:
        jagsModel = jags.model( "BESTmodel.txt" , data=dataList , inits=initsList ,
                               n.chains=nChains , n.adapt=adaptSteps )
        # Burn-in:
        cat( "Burning in the MCMC chain...\n" )
        update( jagsModel , n.iter=burnInSteps )
        # The saved MCMC chain:
        cat( "Sampling final MCMC chain...\n" )
        codaSamples = coda.samples( jagsModel , variable.names=parameters ,
                                    n.iter=nIter , thin=thinSteps )
        # resulting codaSamples object has these indices:
        #   codaSamples[[ chainIdx ]][ stepIdx , paramIdx ]

        #------------------------------------------------------------------------------
        # EXAMINE THE RESULTS
        if ( showMCMC ) {
          windows()
          autocorr.plot( codaSamples[[1]] , ask=FALSE )
        }

        # Convert coda-object codaSamples to matrix object for easier handling.
        # But note that this concatenates the different chains into one long chain.
        # Result is mcmcChain[ stepIdx , paramIdx ]
        mcmcChain = as.matrix( codaSamples )
        return( mcmcChain )

      } # end function BESTmcmc
```

**BEST.R**

```
                # initial values keep the burn-in period moderate.
                initsList = list( mu = mu , sigma = sigma , nuMinusOne = 4 )
        #------------------------------------------------------------------------------
        # RUN THE CHAINS

        parameters = c( "mu" , "sigma" , "nu" )     # The parameters to be monitored
        adaptSteps = 500                 # Number of steps to "tune" the samplers
        burnInSteps = 1000
        nChains = 3
        nIter = ceiling( ( numSavedSteps * thinSteps ) / nChains )
        # Create, initialize, and adapt the model:
        jagsModel = jags.model( "BESTmodel.txt" , data=dataList , inits=initsList ,
                               n.chains=nChains , n.adapt=adaptSteps )
        # Burn-in:
        cat( "Burning in the MCMC chain...\n" )
        update( jagsModel , n.iter=burnInSteps )
        # The saved MCMC chain:
        cat( "Sampling final MCMC chain...\n" )
        codaSamples = coda.samples( jagsModel , variable.names=parameters ,
                                    n.iter=nIter , thin=thinSteps )
        # resulting codaSamples object has these indices:
        #   codaSamples[[ chainIdx ]][ stepIdx , paramIdx ]

        #------------------------------------------------------------------------------
        # EXAMINE THE RESULTS
        if ( showMCMC ) {
          windows()
          autocorr.plot( codaSamples[[1]] , ask=FALSE )
        }

        # Convert coda-object codaSamples to matrix object for easier handling.
        # But note that this concatenates the different chains into one long chain.
        # Result is mcmcChain[ stepIdx , paramIdx ]
        mcmcChain = as.matrix( codaSamples )
        return( mcmcChain )

      } # end function BESTmcmc
```

## 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(101,100,102,104,102,97,105,105,98,101,100,123,105,
       109,102,82,102,100,102,102,101,102,102,103,103,97,
       96,103,124,101,101,100,101,101,104,100,101)
y2 = c(99,101,100,101,102,100,97,101,104,101,102,102,100,
       104,100,100,100,101,102,103,97,101,101,100,101,99,
       101,100,99,101,100,102,99,100,99)

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

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

© John K. Kruschke, Sept. 2012                                    50

### *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, Oct. 2012                                    51

24

**Recall Bayesian estimation:**

**What does NHST say?**

*t* test of means:
*t*(87)=1.62, *p*=0.110 (>.05)
95%CI: -0.361 to 3.477.

*F* test of variances:
F(46,41)=5.72, p < .001.
95%CI on *difference*: ?

*But,* must apply corrections for multiple tests.

© John K. Kruschke, Oct. 2012



**Recall Bayesian estimation:**

**What does NHST say?**
Oops! Data are not normal, so do **resampling** instead.

Resampling test of means:
*p*=0.116 (>.05)

Resampling test of difference of standard deviations:
*p* = 0.072 (>.05)

*And,* still must apply corrections for multiple tests. And there are no CI's.

© John K. Kruschke, Oct. 2012

## Example with outliers: BESTexample.R

Bayesian estimation:
- Credible differences between means and standard deviations.
- Complete distributional information on effect size and everything else.
- Non-normality indicated.

NHST *t* test:
- Outliers invalidate classic test.
- Resampling shows p>.05 for difference of means, p>.05 for difference of standard deviations.
- Need correction for multiple tests.
- No CI's. (And CI's would have no distributional info and fickle end points linked to fickle *p* values.)

55

## Example with small N

Bayesian estimation:
- **Zero is among credible differences between means** and standard deviations, and for effect size.
- Complete distributional information on effect size and everything else.
- Normality is credible.

NHST *t* test:
- $t(14)=2.33$, ***p=0.035***, 95% CI: 0.099, 2.399. ($F(7,7)=1.00$, $p=.999$, CI on ratio: 0.20, 5.00.)
- Need correction for multiple tests, if intended.
- CI's have no distributional info and fickle end points linked to fickle *p* values.
- ***t*** test fails to reveal true uncertainty in parameter estimates when simultaneously estimating SD's and normality.

56

26

## Region of Practical Equivalence (ROPE)

**Marginal Posterior**

mean = 108

0.5% < 100 < 99.5%

1% in ROPE

95% HDI
102          114

95   100   105   110   115   120

$\mu$

Consider a landmark value. Values that are equivalent to that landmark for all practical purposes define the ROPE around that value.

For example, the landmark value is 100, and the ROPE is 99 to 101.

57

---

## Region of Practical Equivalence (ROPE)

**Marginal Posterior**

mean = 108

0.5% < 100 < 99.5%

1% in ROPE

95% HDI
102          114

95   100   105   110   115   120

$\mu$

A parameter value is declared to be not credible, or rejected, if its entire ROPE lies outside the 95% HDI of the posterior distribution of that parameter.

A parameter value is declared to be accepted for practical purposes if that value's ROPE completely contains the 95% HDI of the posterior of that parameter.

58

27

## Example of accepting null value

Bayesian estimation:
- 95% HDI for difference on means falls within ROPE; same for SD's **(enlarged in next slide)**.
- Complete distributional information on effect size and everything else.
- Normality is credible.

NHST *t* test:
- *p* is large for both *t* and *F* tests, but NHST cannot accept null hypothesis.
- Need correction for multiple tests, if intended.
- CI's have no distributional info and fickle end points linked to fickle *p* values, and CI does not indicate probability of parameter value. Hence, **cannot use ROPE method in NHST.**

© John K. Kruschke, Oct. 2012



## Example of accepting null value

Bayesian estimation:
- 95% HDI for difference on means falls within ROPE; same for SD's.
- Complete distributional information on effect size and everything else.
- Normality is credible.

NHST *t* test:
- *p* is large for both *t* and *F* tests, but NHST cannot accept null hypothesis.
- Need correction for multiple tests, if intended.
- CI's have no distributional info and fickle end points linked to fickle *p* values, and CI does not indicate probability of parameter value. Hence, **cannot use ROPE method in NHST.**

© John K. Kruschke, Oct. 2012

## Sequential Testing

For simulated data *from the null hypothesis:*



61

## Sequential Testing

For simulated data *from the null hypothesis:*



62

29

## Many other topics are in the book, e.g.

❖ Bayesian hierarchical **ANOVA, oneway and twoway with interaction contrasts**.
❖The **generalized linear model**.
❖ Many types of **regression**, including multiple linear regression, logistic regression, ordinal regression.
❖ **Log-linear models vs chi-square test**.
❖ **Power**: Probability of achieving the goals of research.
❖ All preceded by **extensive introductory chapters** covering notions of probability, Bayes' rule, MCMC, model comparison, etc.

John K. Kruschke

# Doing Bayesian Data Analysis

A Tutorial with R and BUGS

---

# 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 ?**

en

# 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, Oct. 2012                                                                                       65

---

# 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, Oct. 2012                                                                                       66

31

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

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

67

---

## 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*.

$$p \text{ value} = p\left( |t_{\text{null}}| > |t_{\text{act}}| \right)$$
for $t_{\text{null}}$ sampled according to the intended experiment

68

## "The" *p* value…

**Space of possible outcomes from null hypothesis**

**Actual outcome**

*p* value

Ø

69

## *p* value for intention to sample until N

**Space of possible outcomes from null hypothesis**

**Actual outcome**

*p* value

N

70

*p* value for intention to sample until Time

Space of possible outcomes from null hypothesis

Actual outcome

*p* value

© John K. Kruschke, Oct. 2012    71



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 each group

$t_{act} = 2.33$

$p = 0.042$

Space of possible outcomes from null hypothesis

© John K. Kruschke, Oct. 2012    73

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

**Null Hypothesis:**
Groups are identical

$\sigma$

$\mu_0$

$y$

$\sigma$

$\mu_0$

$y$

Many simulated repetitions of the *intended* experiment

Fixed N=6 each group

0.4

$t_{act} = 2.33$

$p = 0.042$

-4   -2   0   2   4

$t_{null}$

© John K. Kruschke, Oct. 2012

74

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

**Null Hypothesis:**
Groups are identical

$\sigma$

$\mu_0$

$y$

$\sigma$

$\mu_0$

$y$

Many simulated repetitions of the *intended* experiment

**T**

Fixed time, modal&max N=6 each group

0.4

$t_{act} = 2.33$

$p = 0.056$

0.0

-4   -2   0   2   4

$t_{null}$

**Space of possible outcomes from null hypothesis**

© John K. Kruschke, Oct. 2012

75

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

**Null Hypothesis:**
Groups are identical



Many simulated repetitions of the *intended* experiment

Fixed time, modal&max N=6 each group

$t_{act} = 2.33$

$p = 0.056$

76

# 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**

***Can the null hypothesis be rejected?*** *To answer, we must know the intention of the data collector.*
• We ask the research assistant who collected the data. The assistant says, "I just collected data for two weeks. It's my job. I happened to get 6 subjects in each group."
• We ask the graduate student who oversaw the assistant. The student says, "I knew we needed 6 subjects per group, so I told the assistant to run for two weeks, because we usually get about 6 subjects per week."
• We ask the lab director, who says, "I told my graduate student to collect 6 subjects per group."
• Therefore, for the lab director, t = 2.33 **rejects the null hypothesis** (because p < .05), but for the research assistant who actually collected the data, t = 2.33 **fails to reject** the null hypothesis (because p **>** .05).



Fixed N=6 each group

$t_{act} = 2.33$

$p = 0.042$

Fixed time, modal&max N=6 each group

$t_{act} = 2.33$

$p = 0.056$

77

## Two labs collect data with same *t* and N:

**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**

Fixed N=6 each group

$t_{act} = 2.33$

p = 0.042

**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**

Fixed time, modal&max N=6 each group

$t_{act} = 2.33$

p = 0.056

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

78

## The *real* use of the Neuralyzer:



79

## 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.

80

# 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%.

81

"The" *p* value depends on intended tests:

Actual outcome

*p* value

Space of possible outcomes from null hypothesis for 1 comparison

© John K. Kruschke, Oct. 2012        82



"The" *p* value depends on intended tests:

Actual outcome

*p* value

Space of possible outcomes from null hypothesis for several comparisons

© John K. Kruschke, Oct. 2012        83

10/4/2012

# Experimentwise false alarm rate



© John K. Kruschke, Oct. 2012

84

# 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, Oct. 2012

85

40

## 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.

86

---

# Good intentions make any result *in*significant

- 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.*

87

88

# 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!

89

## ? Confidence Intervals ?
### provide no confidence

**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

**Under assumption of fixed N:**
$$CI = (M_1 - M_2) \pm t_{crit} \times se$$
$$= (5.18 - 4.32) \pm 2.23 \times 0.370$$
$$= [\, \mathbf{0.036} \, , \mathbf{1.68} \,]$$
which *ex*cludes zero.

95% CI constructed with fixed-N $t_{crit}$ will span true difference *less* than 95% of time if data are sampled according to fixed duration.

**Under assumption of fixed duration:**
$$CI = (M_1 - M_2) \pm t_{crit} \times se$$
$$= (5.18 - 4.32) \pm 2.45 \times 0.370$$
$$= [\, \mathbf{-0.046} \, , \mathbf{1.77} \,]$$
which *in*cludes zero.

95% CI constructed with fixed-duration $t_{crit}$ will span true difference *more* than 95% of the time if data are sampled according to fixed N.

90

## ? 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.

91

? **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.

92

---

# 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.

93

## Bayesian Estimation or NHST?

When Bayesian estimation and NHST *agree*, which should be used?

Bayesian estimation gives the most complete and informative answer. Answer from NHST is not informative and is fickle.

When Bayesian estimation and NHST *disagree*, which should be used?

Bayesian estimation gives the most complete and informative answer. Answer from NHST is not informative and is fickle.

94

# Conclusion

- p values are not well defined, nor are the limits of confidence intervals, and confidence intervals have no distributional info.

- Bayesian data analysis is the most complete and normatively correct way to estimate parameters in any model, for all your data.

- Bayesian data analysis is taking hold in 21$^{st}$ century science, from astronomy to zoology. *Don't be left behind.*

- And, for more info, …

95

**The blog: http://doingbayesiandataanalysis.blogspot.com/**

Kruschke, J. K. (2012). **Bayesian estimation supersedes the *t* test.** *Journal of Experimental Psychology: General*.
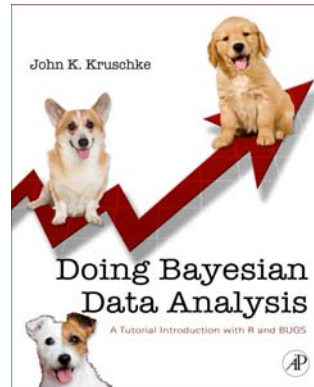
Kruschke, J. K. (2011). **Bayesian assessment of null values via parameter estimation and model comparison.** *Perspectives on Psychological Science*, 6(3), 299-312.

Kruschke, J. K. (2010). **What to believe: Bayesian methods for data analysis.** *Trends in Cognitive Sciences*, 14(7), 293-300.

Kruschke, J. K. (2010). **Bayesian data analysis.** *Wiley Interdisciplinary Reviews: Cognitive Science*, 1(5), 658-676.
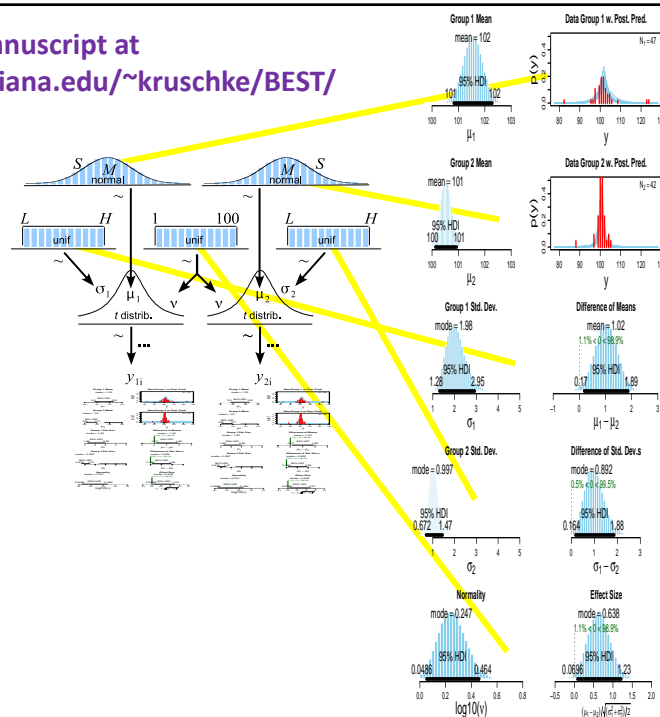
Kruschke, J. K. (2011). *Doing Bayesian Data Analysis: A Tutorial with R and BUGS.* Academic Press / Elsevier.

© John K. Kruschke, Oct. 2012    96



**Program and manuscript at http://www.indiana.edu/~kruschke/BEST/**

© John K. Kruschke, Oct. 2012    97

# Priors are not capricious

1. Priors are explicitly specified and must be acceptable to a skeptical scientific audience.
2. Typically, priors are set to be noncommittal and have very little influence on the posterior.
3. Priors can be informed by well-established data and theory, thereby giving inferential leverage to small samples.
4. When there is disagreement about the prior, then the influence of the prior on the posterior can be, and is, directly investigated. Different theoretically-informed priors can be checked.
5. Not using priors can be a serious blunder! E.g., drug/disease testing without incorporating prior knowledge of base rates.
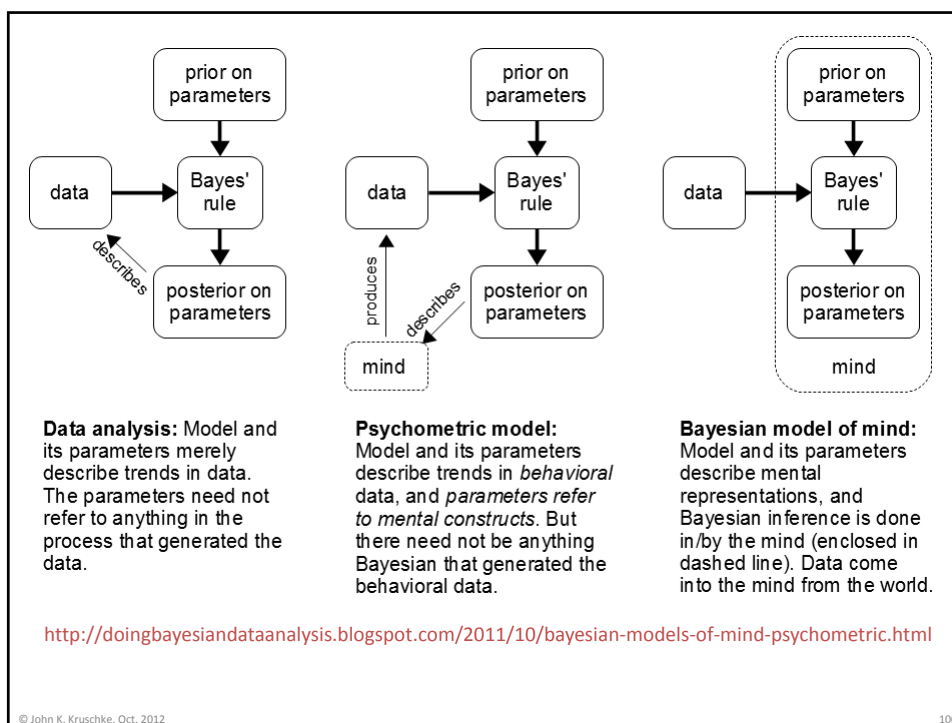
98

# Prior credibility is *not* intentions

| Bayesian Prior | NHST Intention (e.g., stopping rule, number of comparisons) |
|---|---|
| Explicit and supported by previous data. | Unknowable |
| Should influence interpretation of data. | Should *not* influence interpretation of data |

99

## Slide 100



**Data analysis:** Model and its parameters merely describe trends in data. The parameters need not refer to anything in the process that generated the data.

**Psychometric model:** Model and its parameters describe trends in *behavioral* data, and *parameters refer to mental constructs*. But there need not be anything Bayesian that generated the behavioral data.

**Bayesian model of mind:** Model and its parameters describe mental representations, and Bayesian inference is done in/by the mind (enclosed in dashed line). Data come into the mind from the world.

http://doingbayesiandataanalysis.blogspot.com/2011/10/bayesian-models-of-mind-psychometric.html
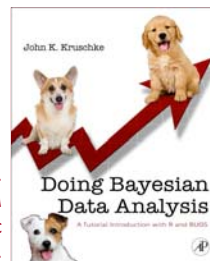
100

## Slide 101

# Bayesian estimation or Bayesian model comparison?

Bayesian estimation is also better than the "Bayesian *t* test," which uses the "Bayes factor" from Bayesian model comparison…

Kruschke, J. K. (2011). **Bayesian assessment of null values via parameter estimation and model comparison.** *Perspectives on Psychological Science*, 6(3), 299-312.

*Chapter 12* of Kruschke, J. K. (2011). ***Doing Bayesian Data Analysis: A Tutorial with R and BUGS.*** Academic Press / Elsevier.

Kruschke, J. K. (in press). **Bayesian estimation supersedes the *t* test.** *Journal of Experimental Psychology: General*. Appendix D.

101