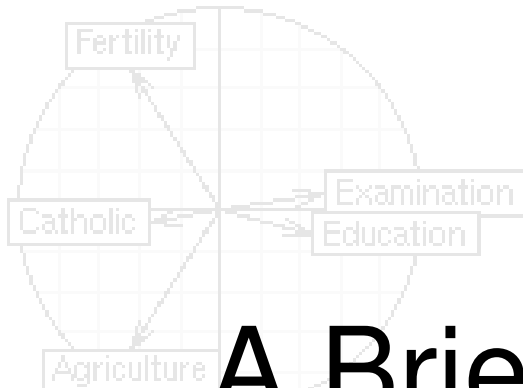
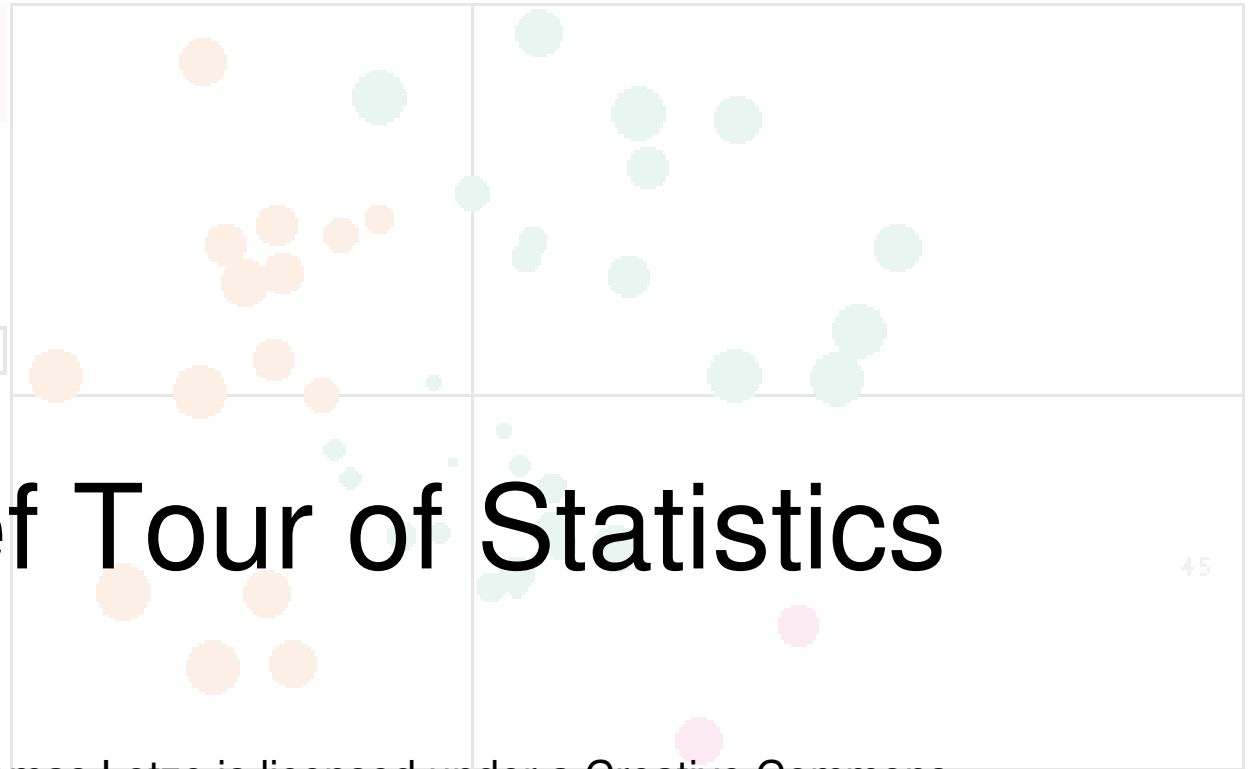


PCA 5 vars

`princomp(x = data, cor = cor)`

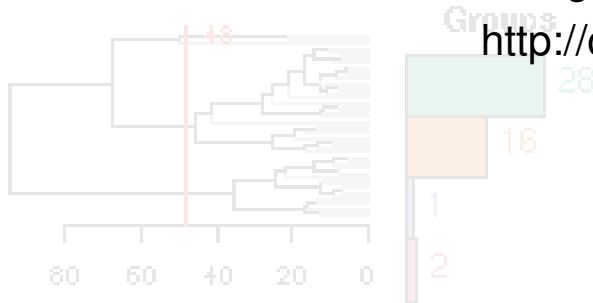


A Brief Tour of Statistics



This work by Thomas Lotze is licensed under a Creative Commons Attribution 3.0 United States License. You're free to share or change it, so long as you provide attribution to Thomas Lotze.

<http://creativecommons.org/licenses/by/3.0/us>



- What are Distributions?

- Models

- Binomial
- Poisson
- Uniform
- Gaussian (normal)

- Using Distributions to Answer Questions

- Distribution Parameters

- Estimating Parameters

- Confidence Intervals and P-values

- Why Gaussian?

- Regression

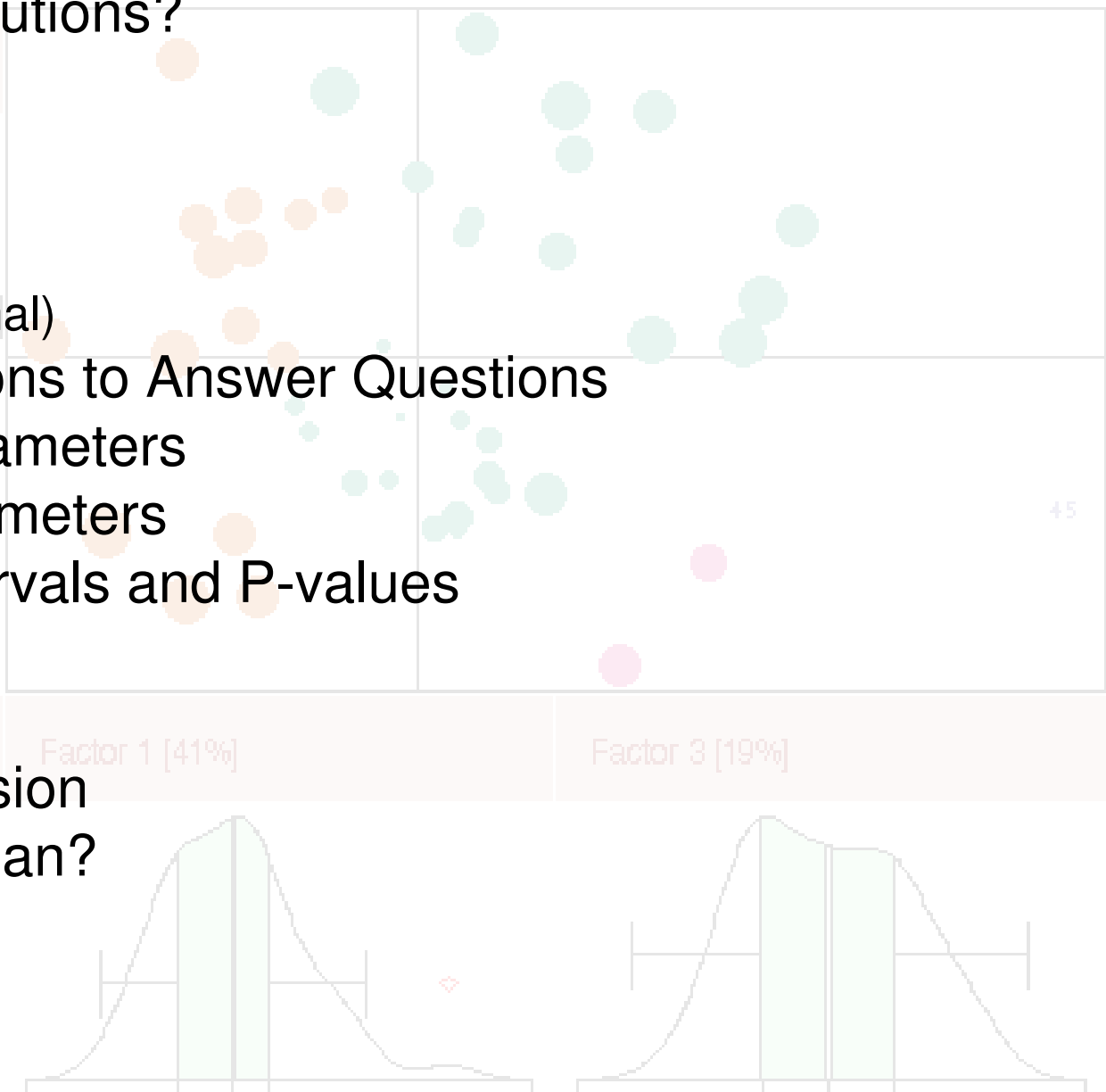
- Logistic Regression

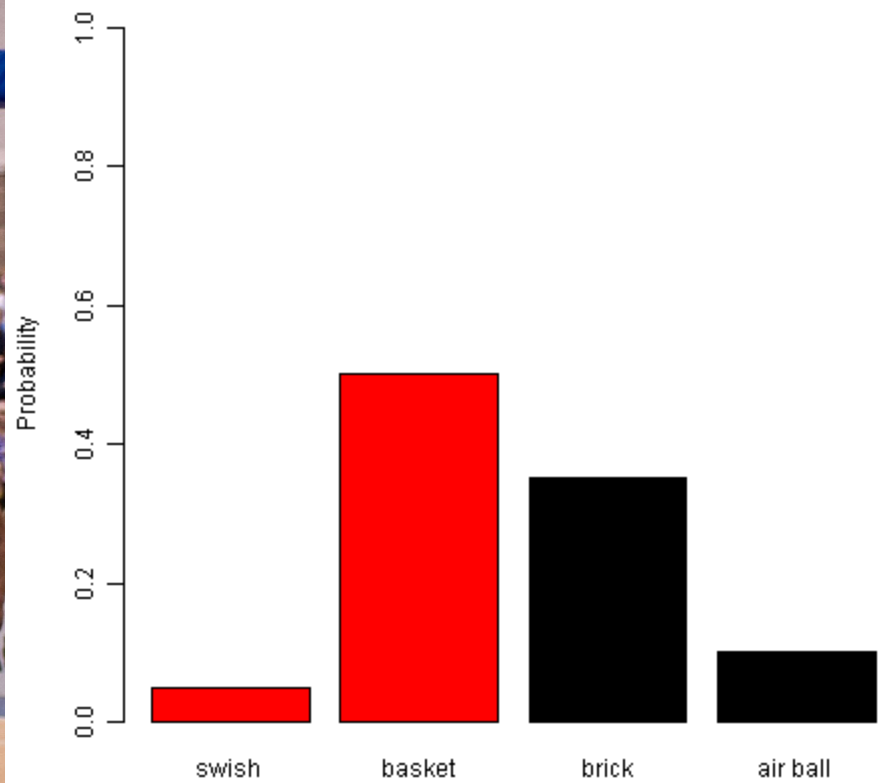
- Why Not Gaussian?

- Bootstrapping

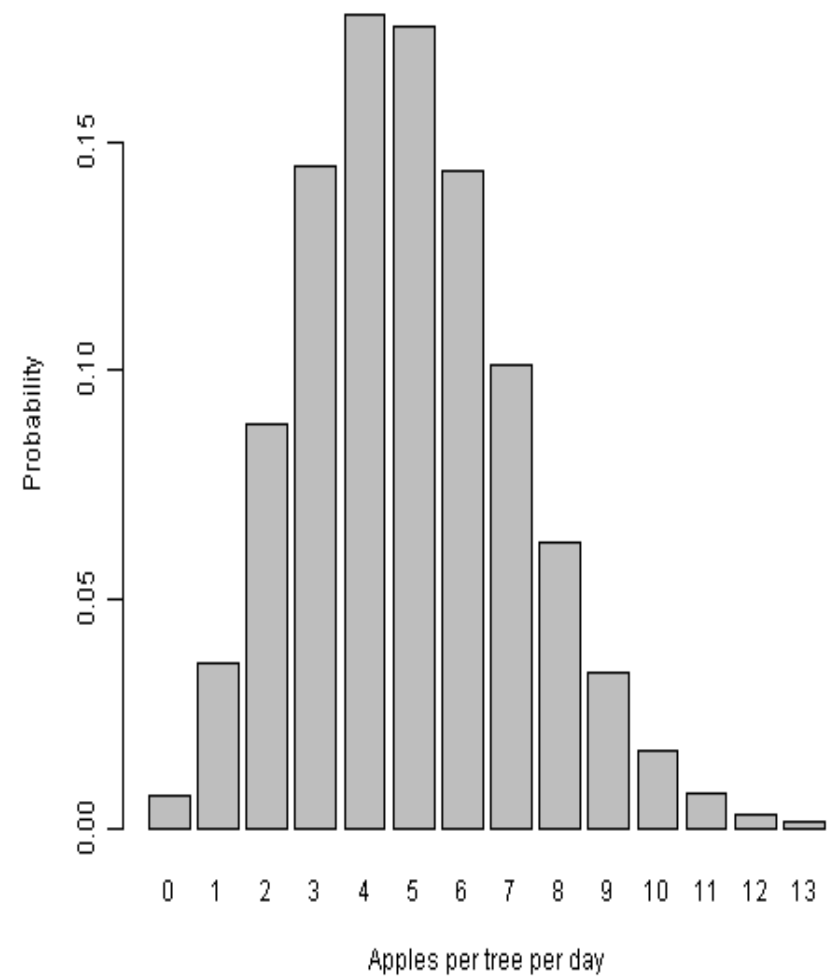
- Multiple Testing

- Useful Tools

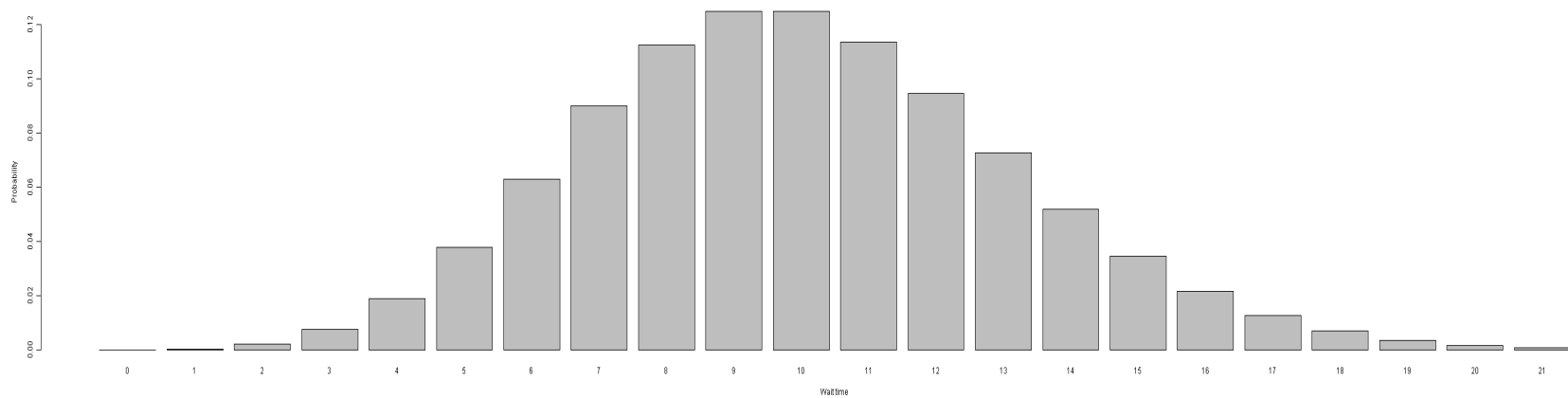




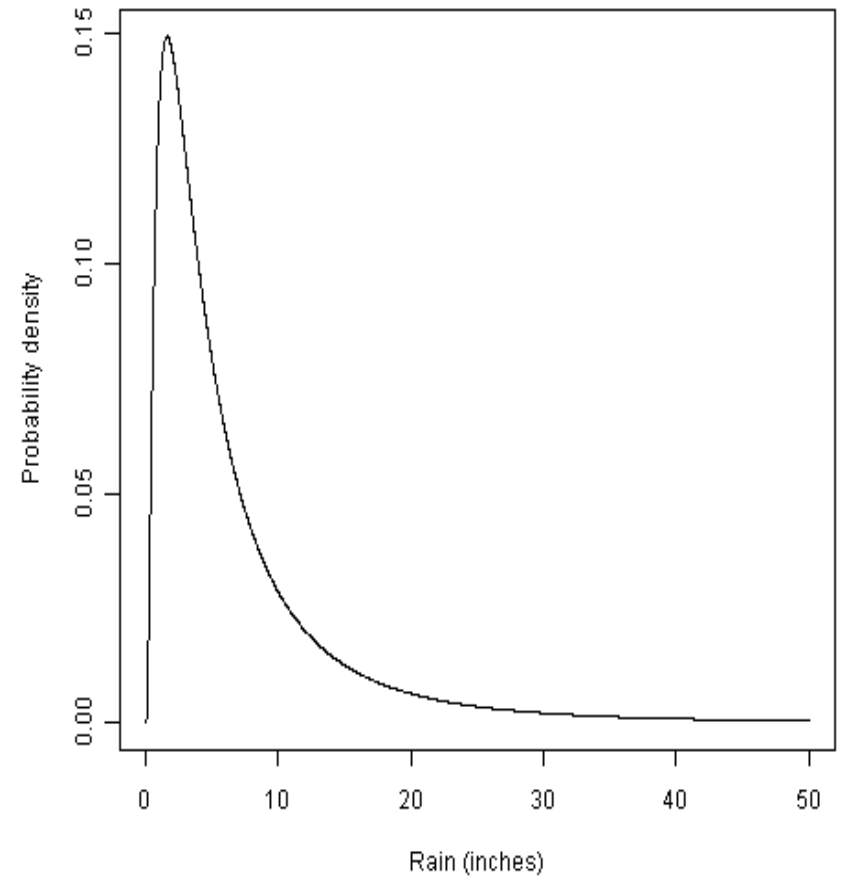
Creative Commons licensed, from gr0don on flickr



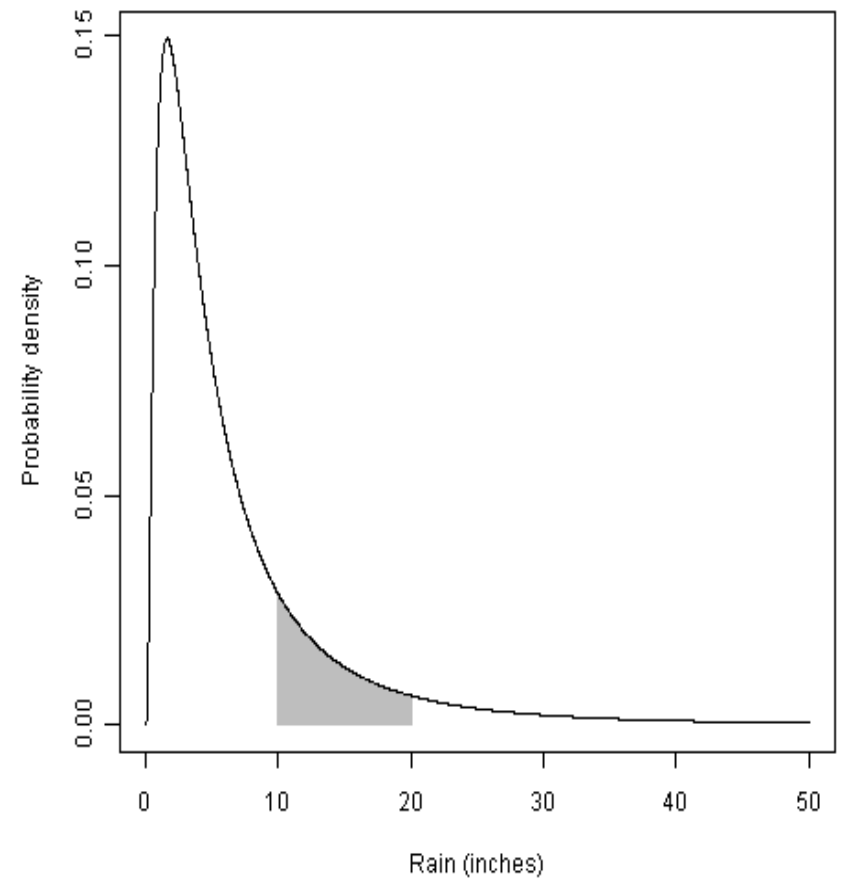
Creative Commons licensed, from WxMom on flickr



Creative Commons licensed, from stringberd on flickr



Creative Commons licensed, from benchilada on flickr



Creative Commons licensed, from benchilada on flickr

- What are Distributions?

- **Models**

- **Binomial**
- **Poisson**
- **Uniform**
- **Gaussian (normal)**

- Using Distributions to Answer Questions

- Distribution Parameters

- Estimating Parameters

- Confidence Intervals and P-values

- Why Gaussian?

- Regression

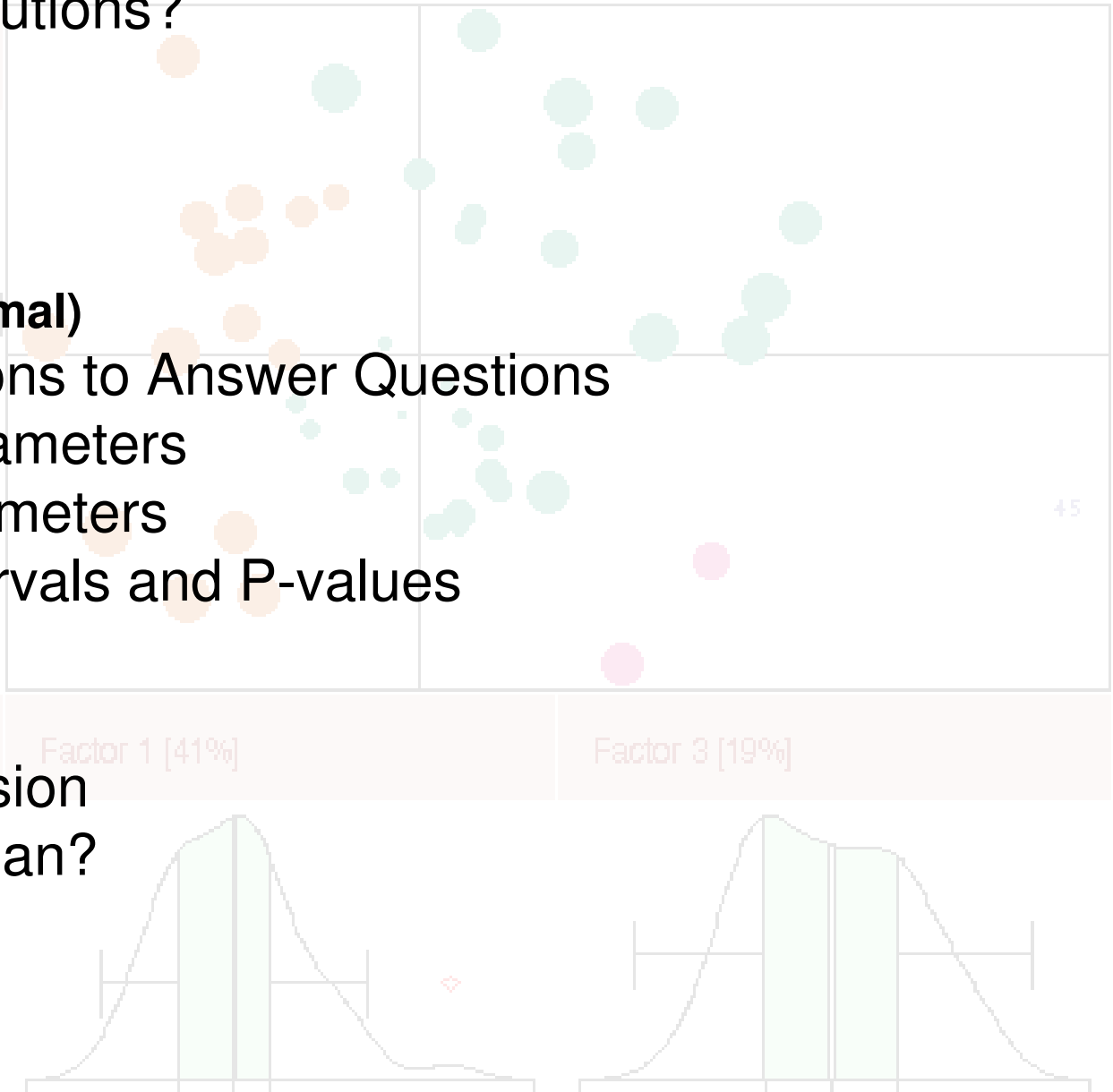
- Logistic Regression

- Why Not Gaussian?

- Bootstrapping

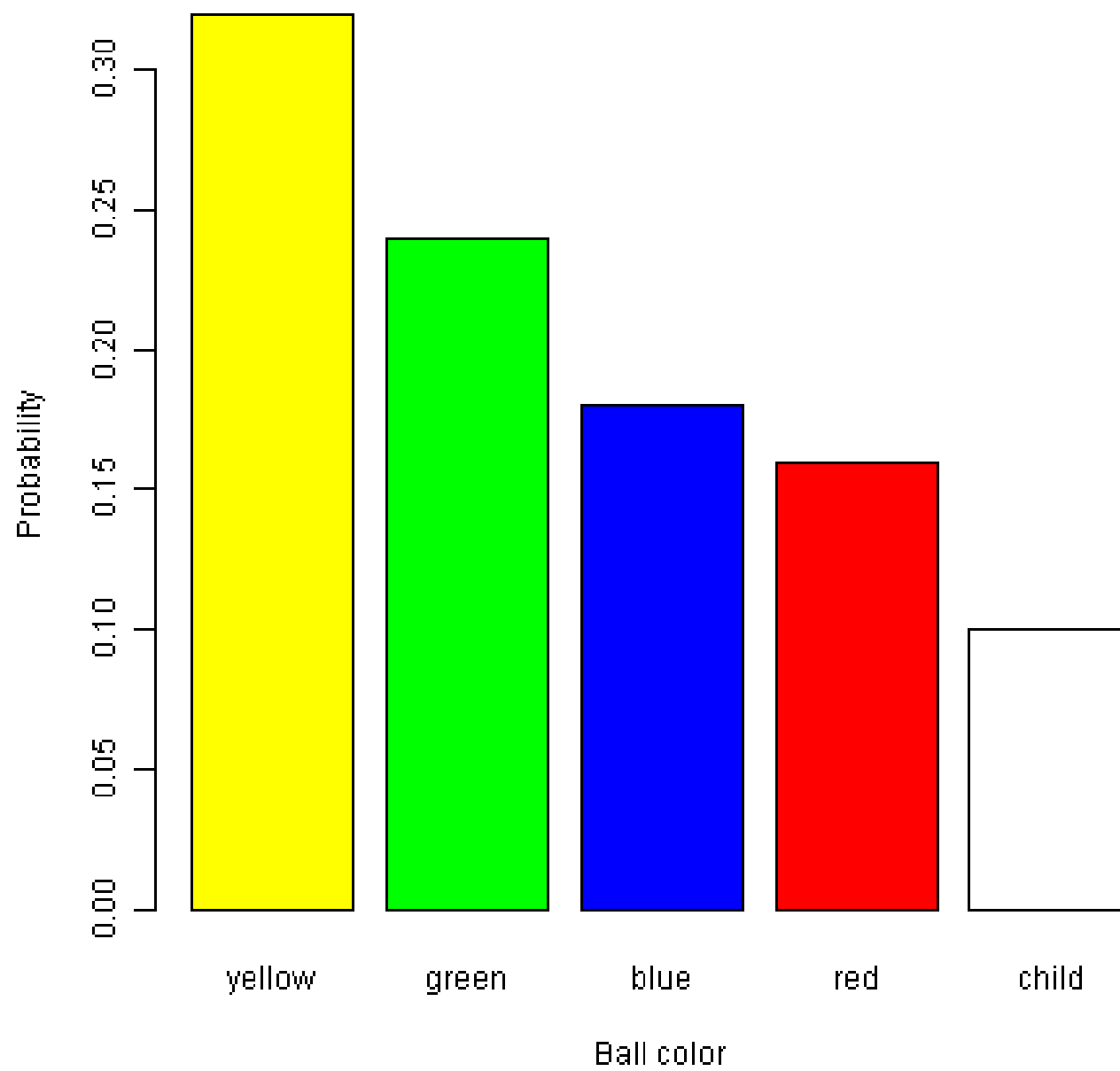
- Multiple Testing

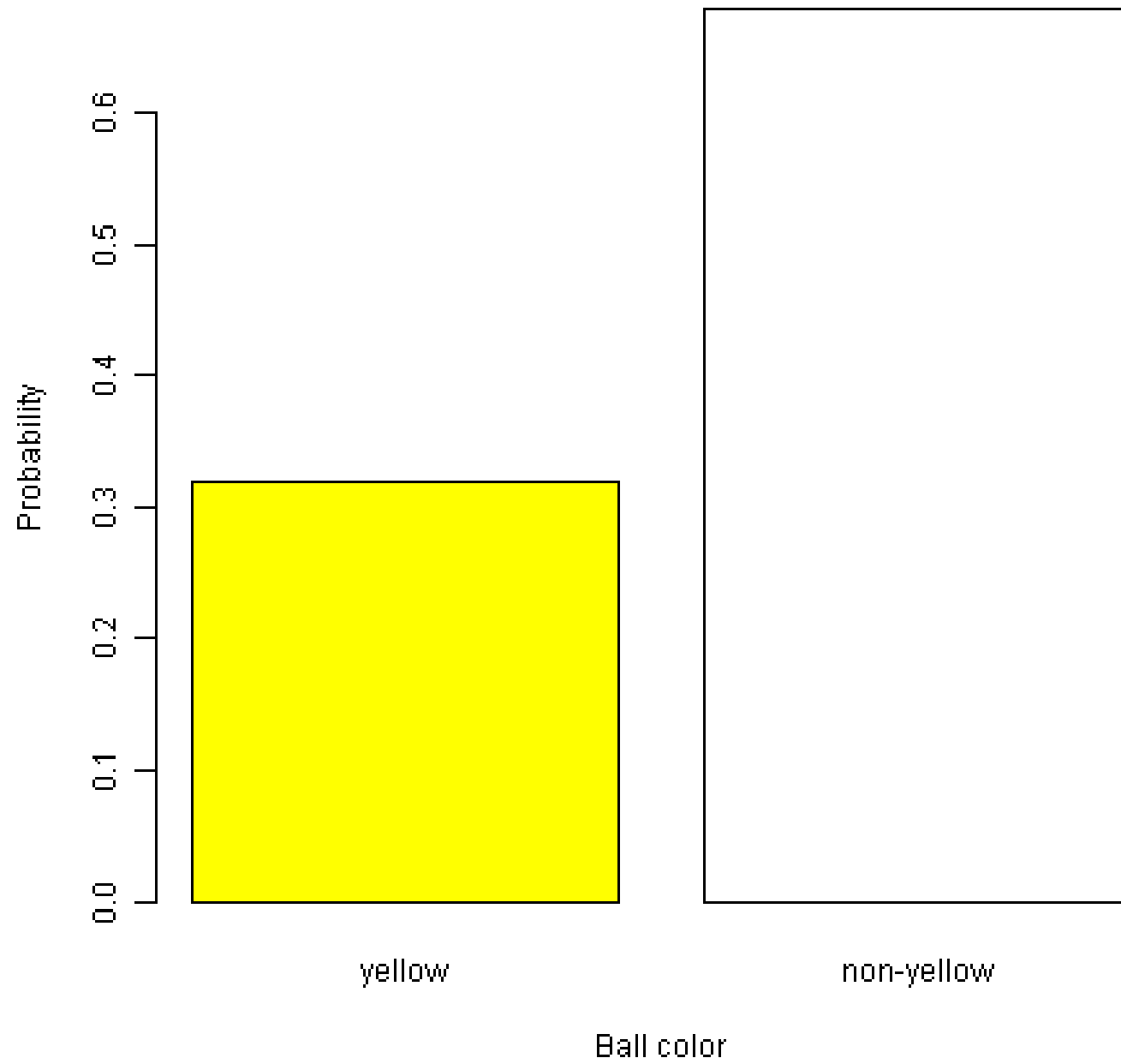
- Useful Tools

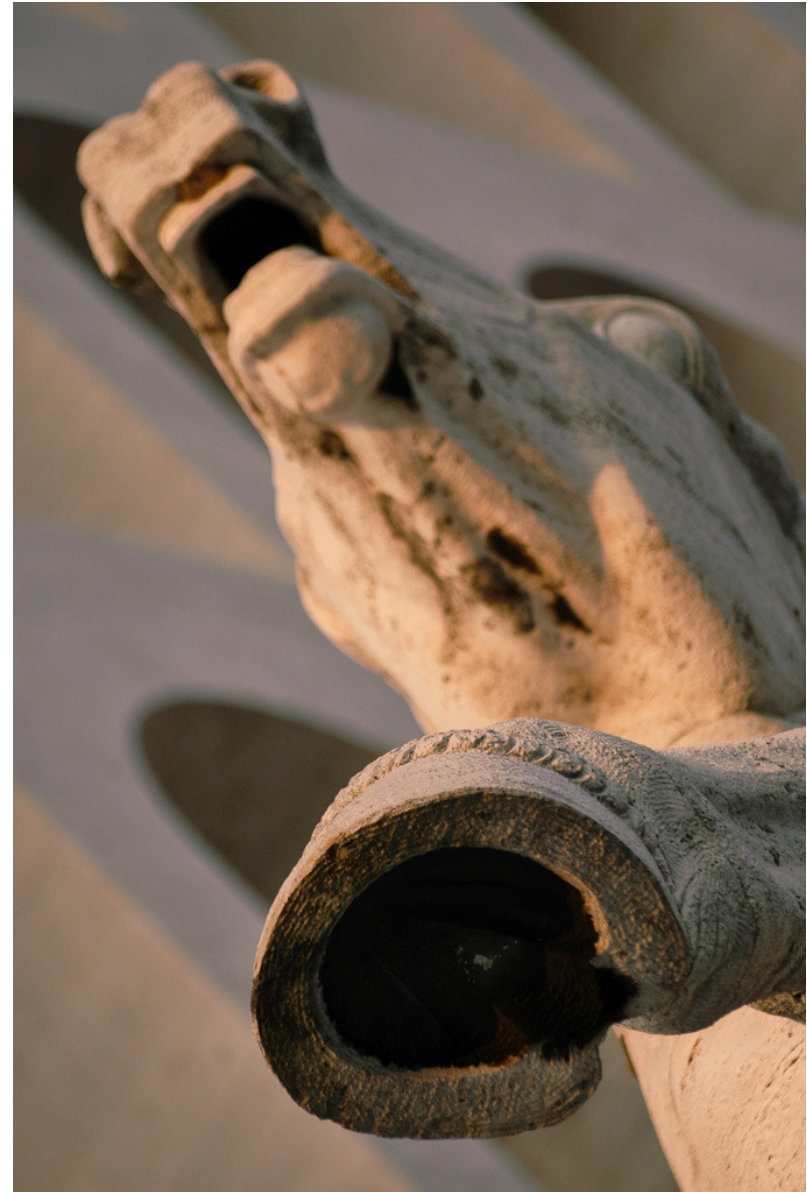
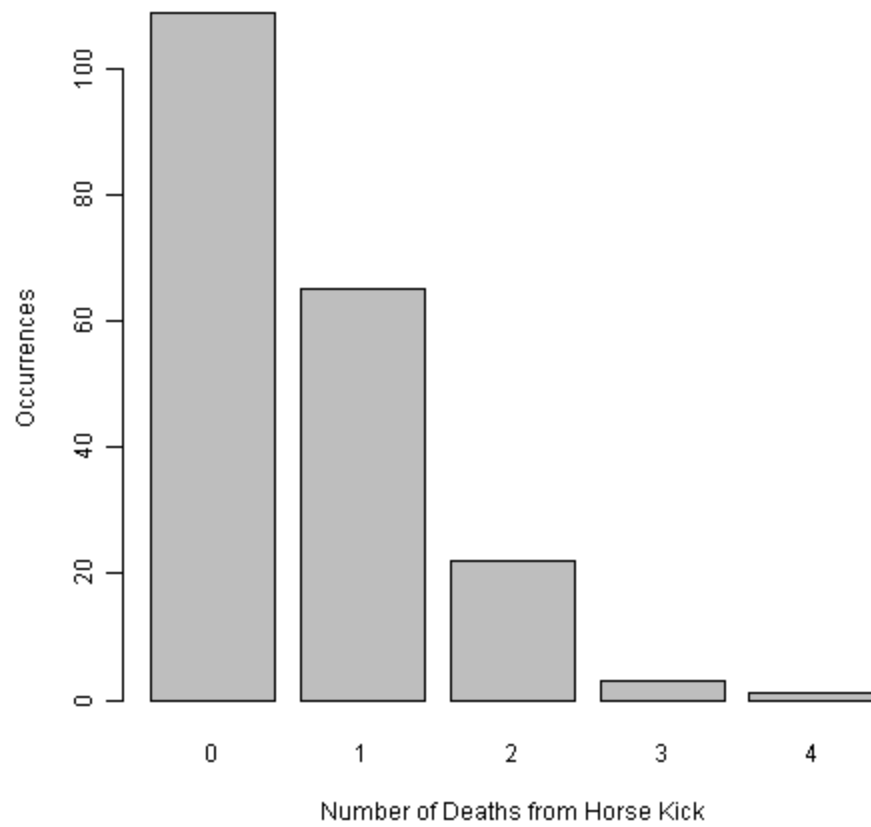




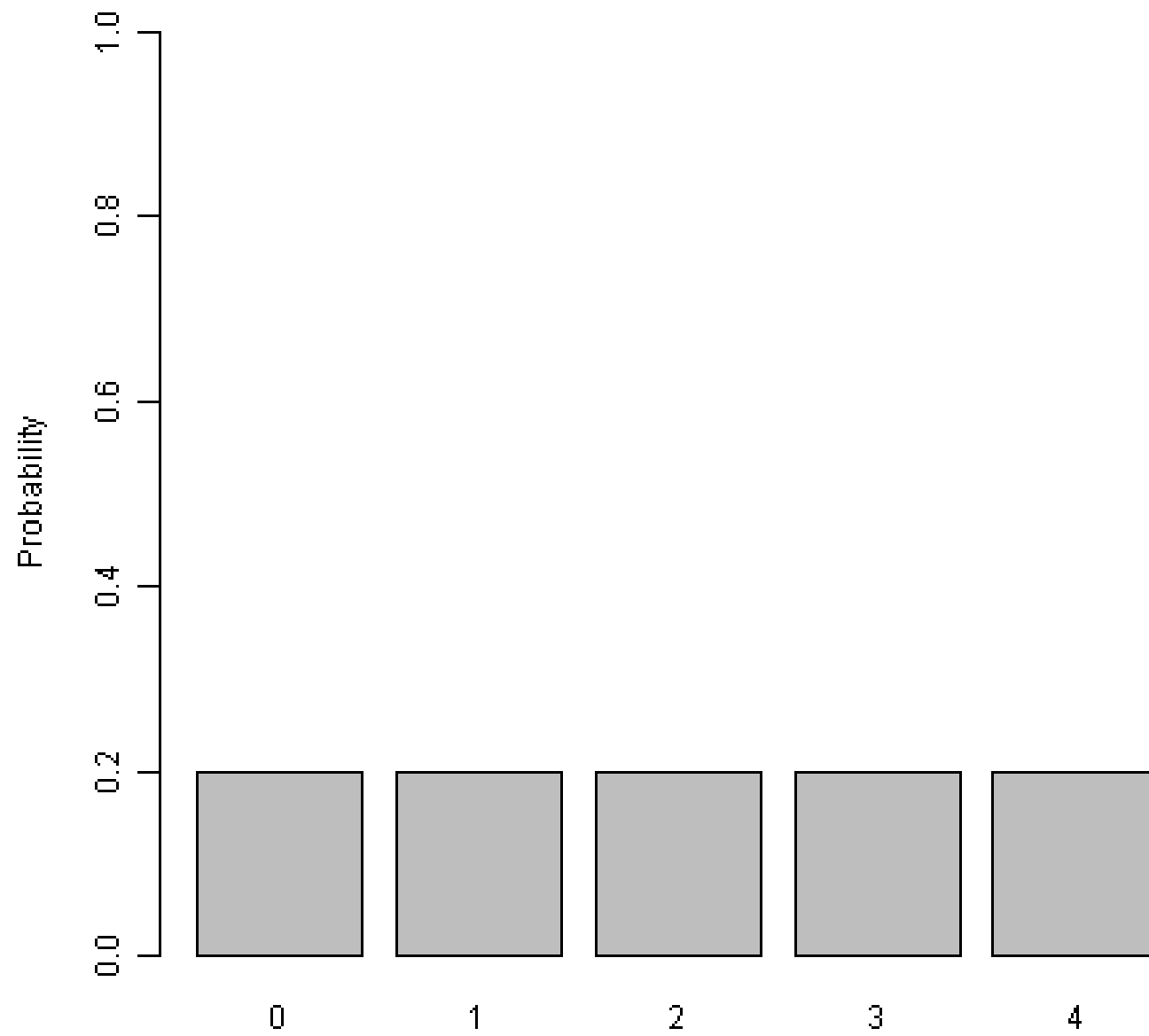
Creative Commons licensed, from Jason Tromm on flickr

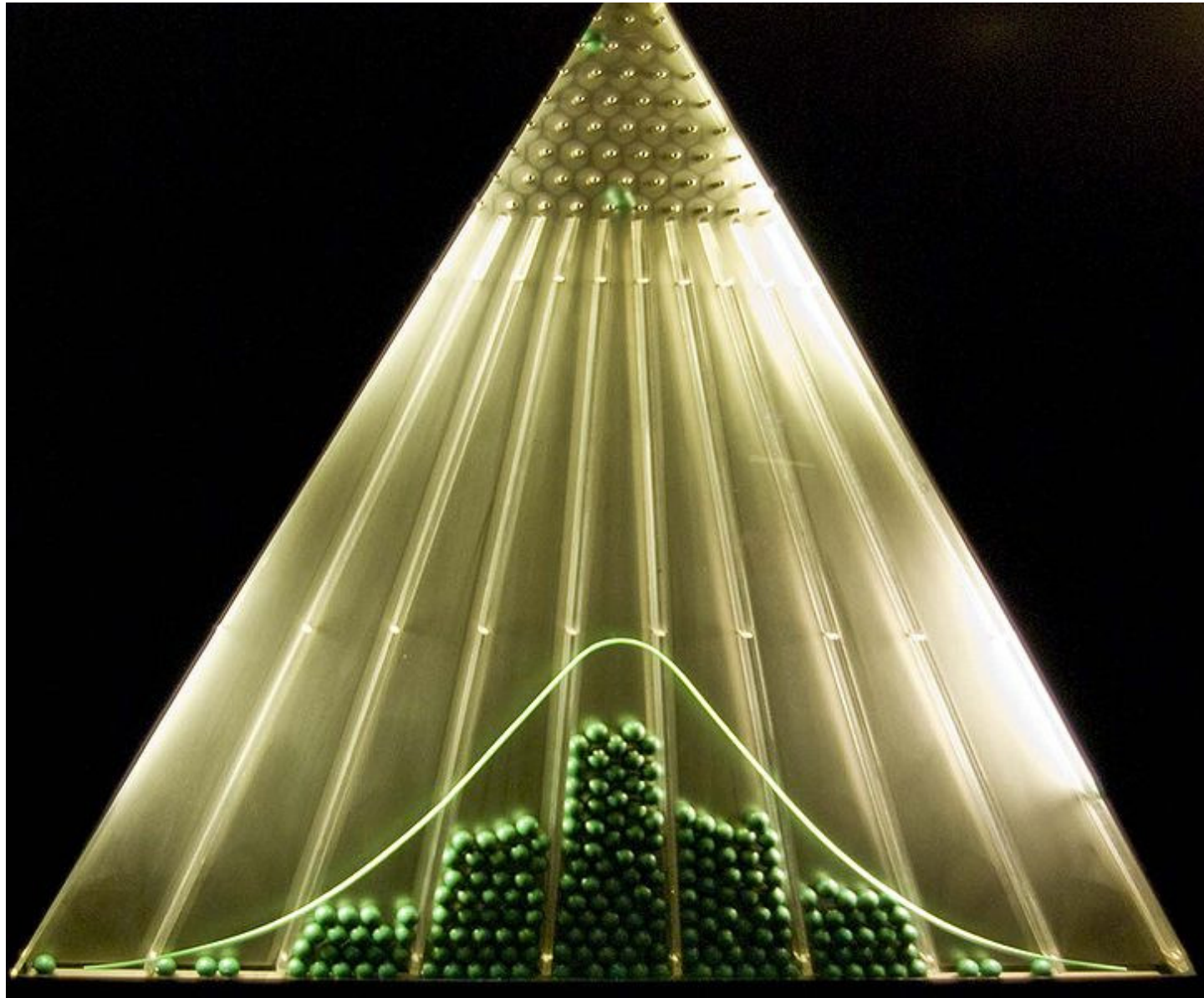






Creative Commons licensed, from Rickydavid on flickr





- What are Distributions?

- Models

- Binomial
- Poisson
- Uniform
- Gaussian (normal)

- **Using Distributions to Answer Questions**

- Distribution Parameters

- Estimating Parameters

- Confidence Intervals and P-values

- Why Gaussian?

- Regression

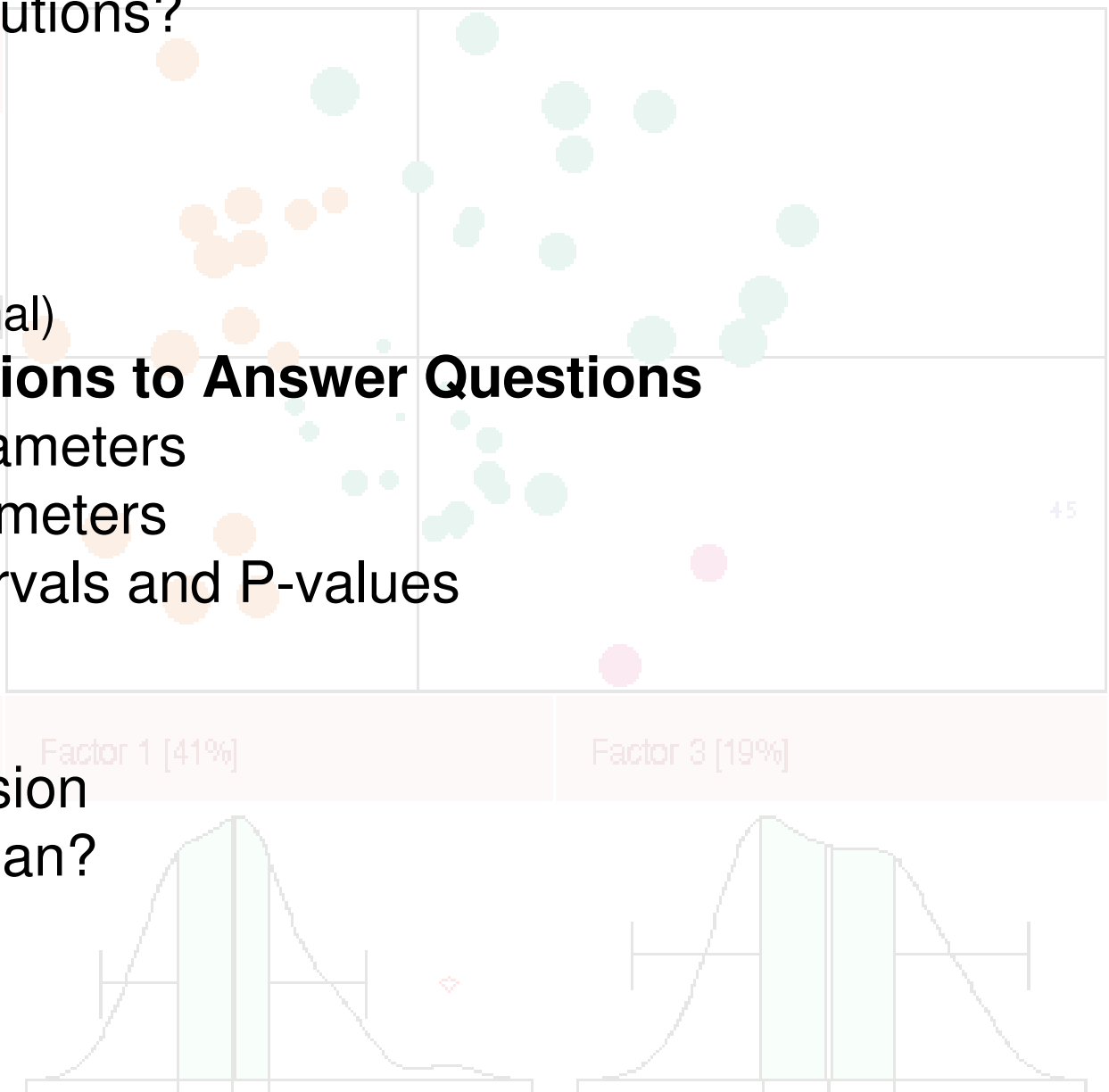
- Logistic Regression

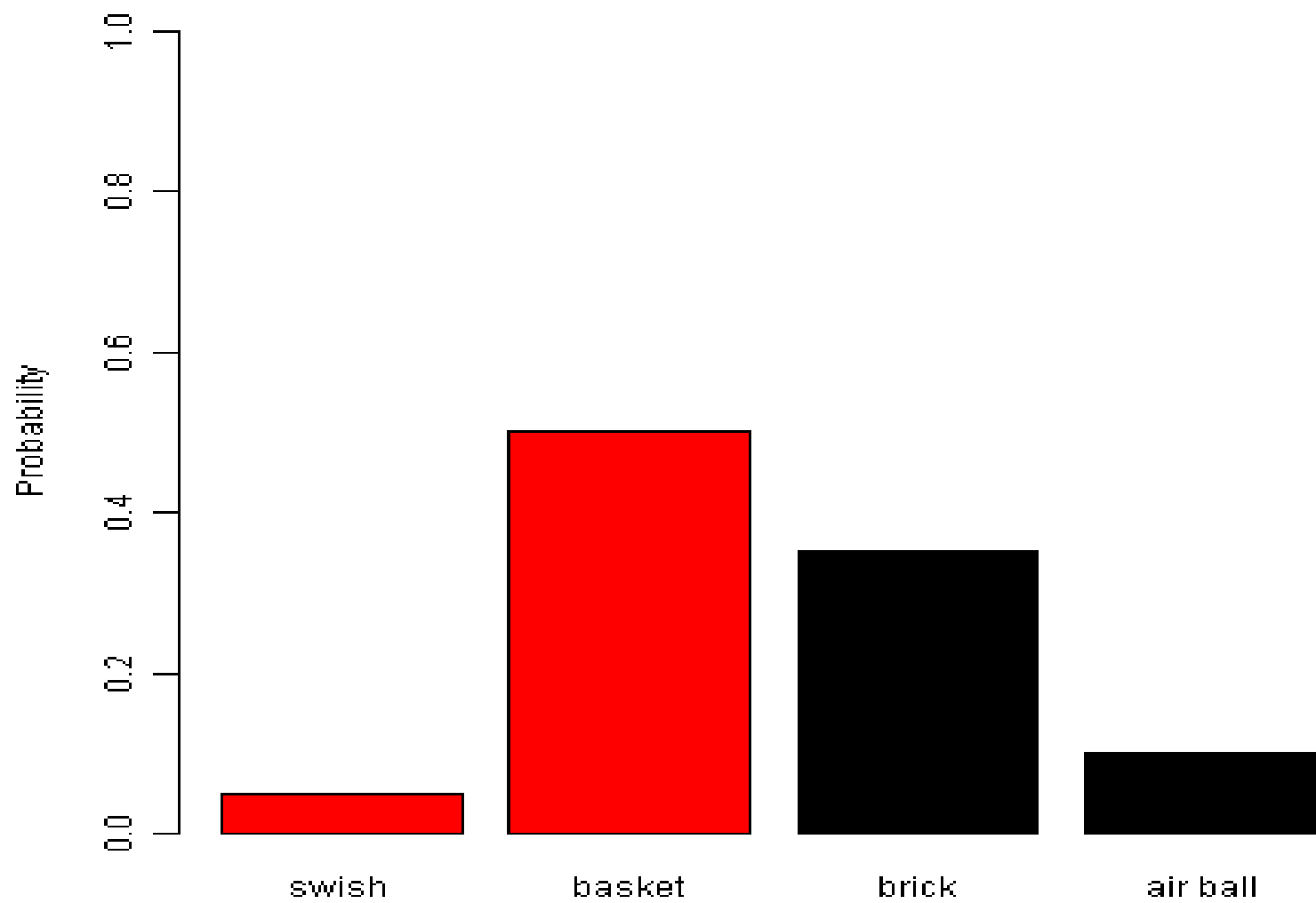
- Why Not Gaussian?

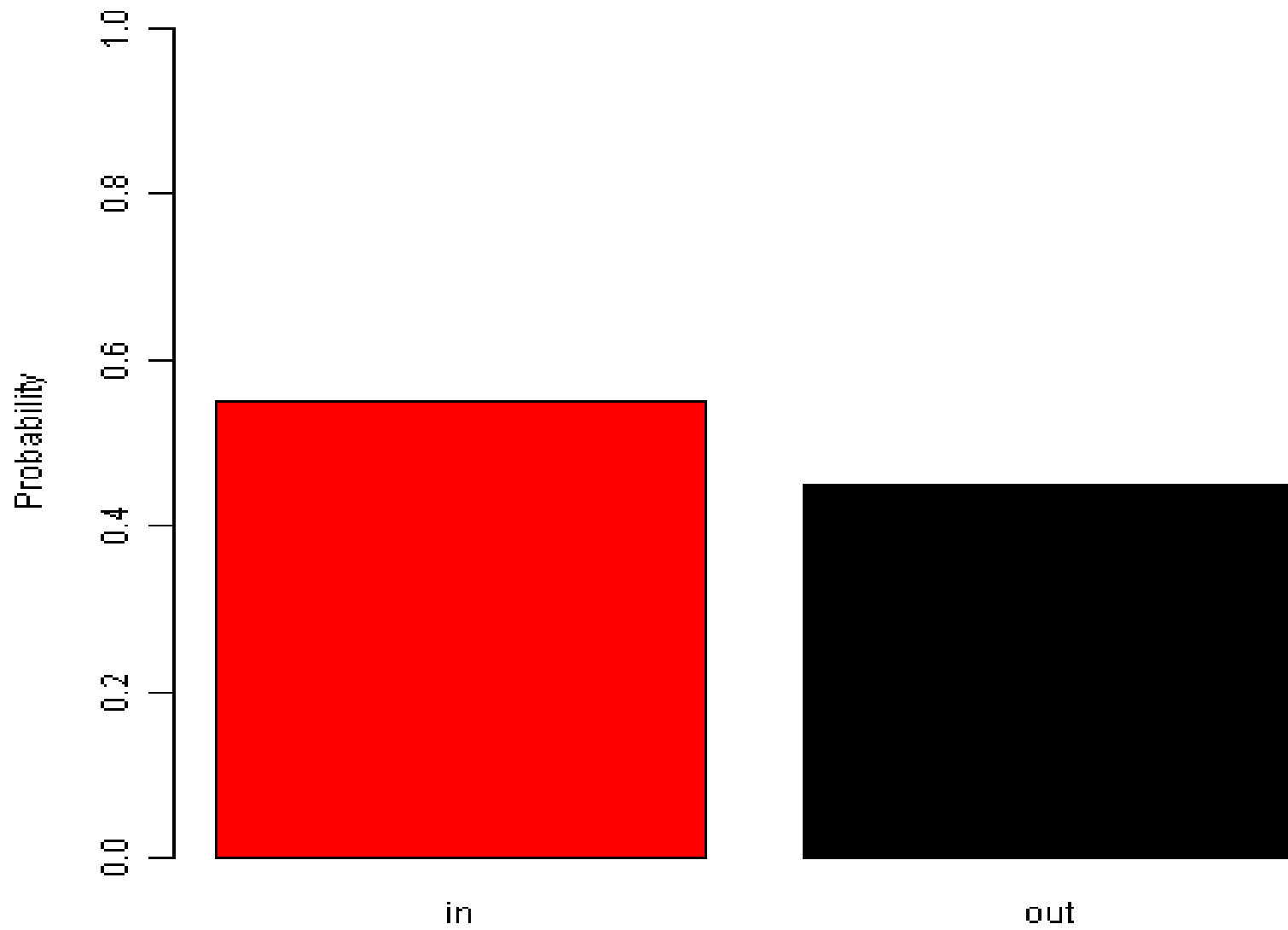
- Bootstrapping

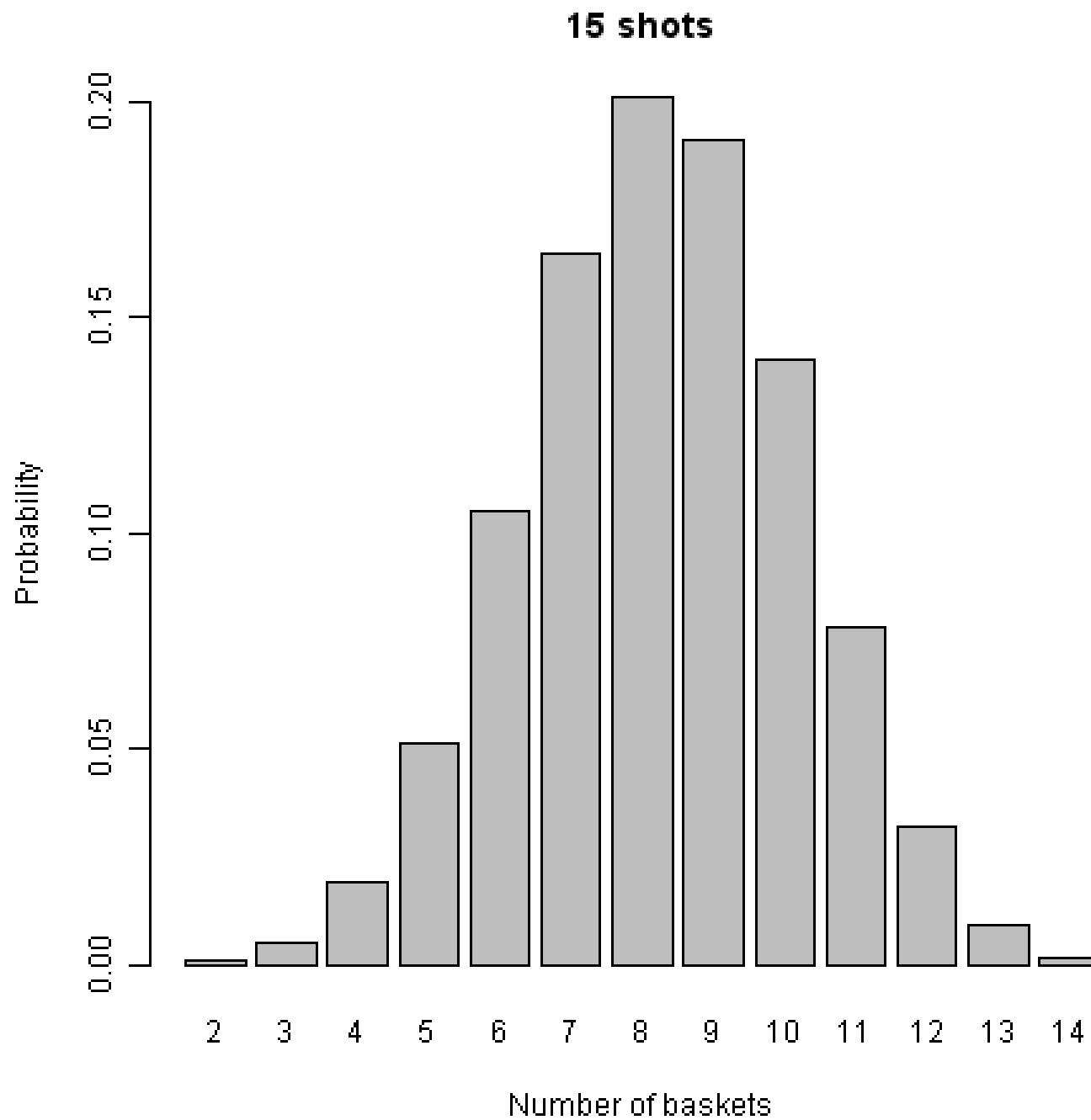
- Multiple Testing

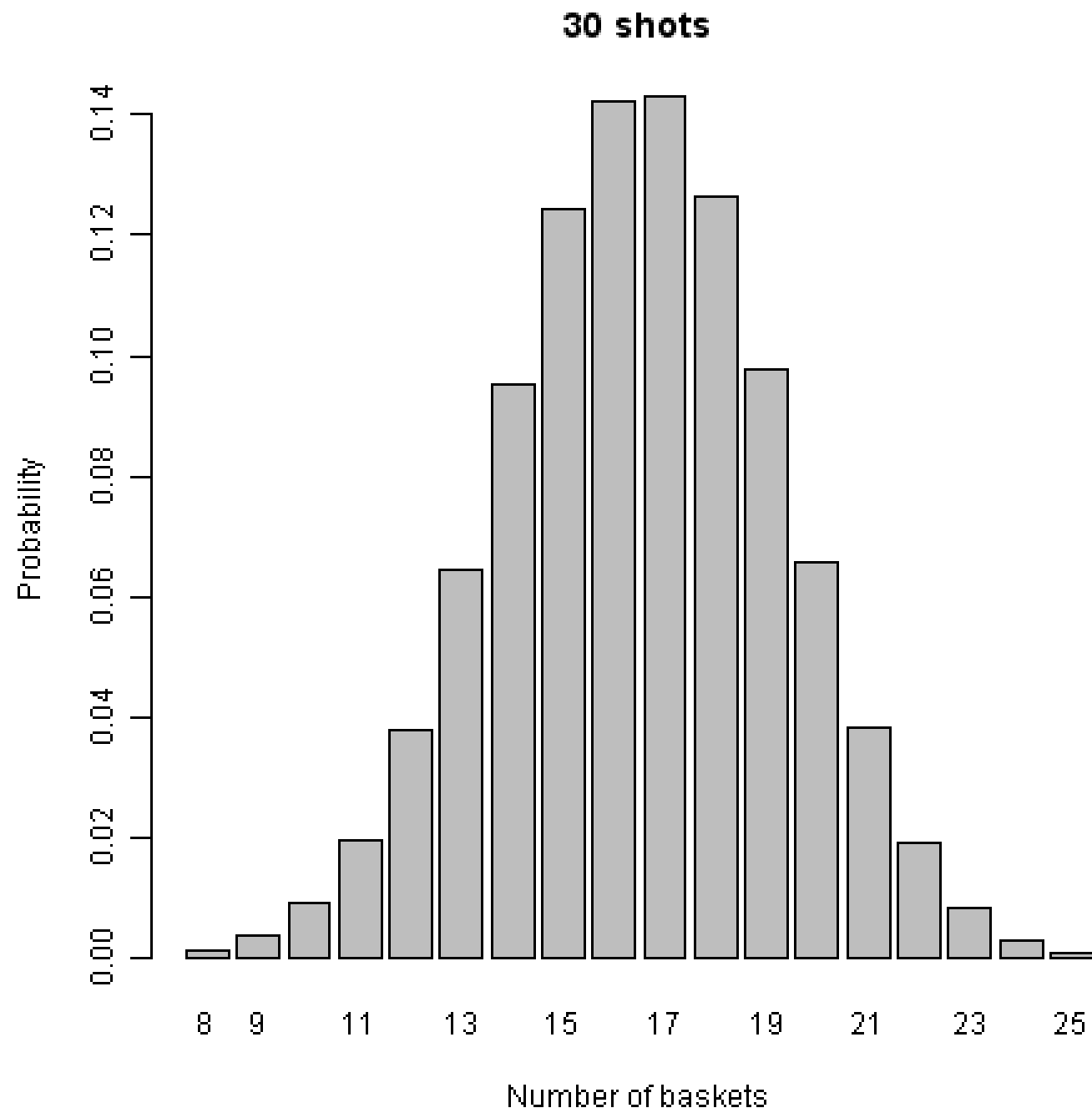
- Useful Tools



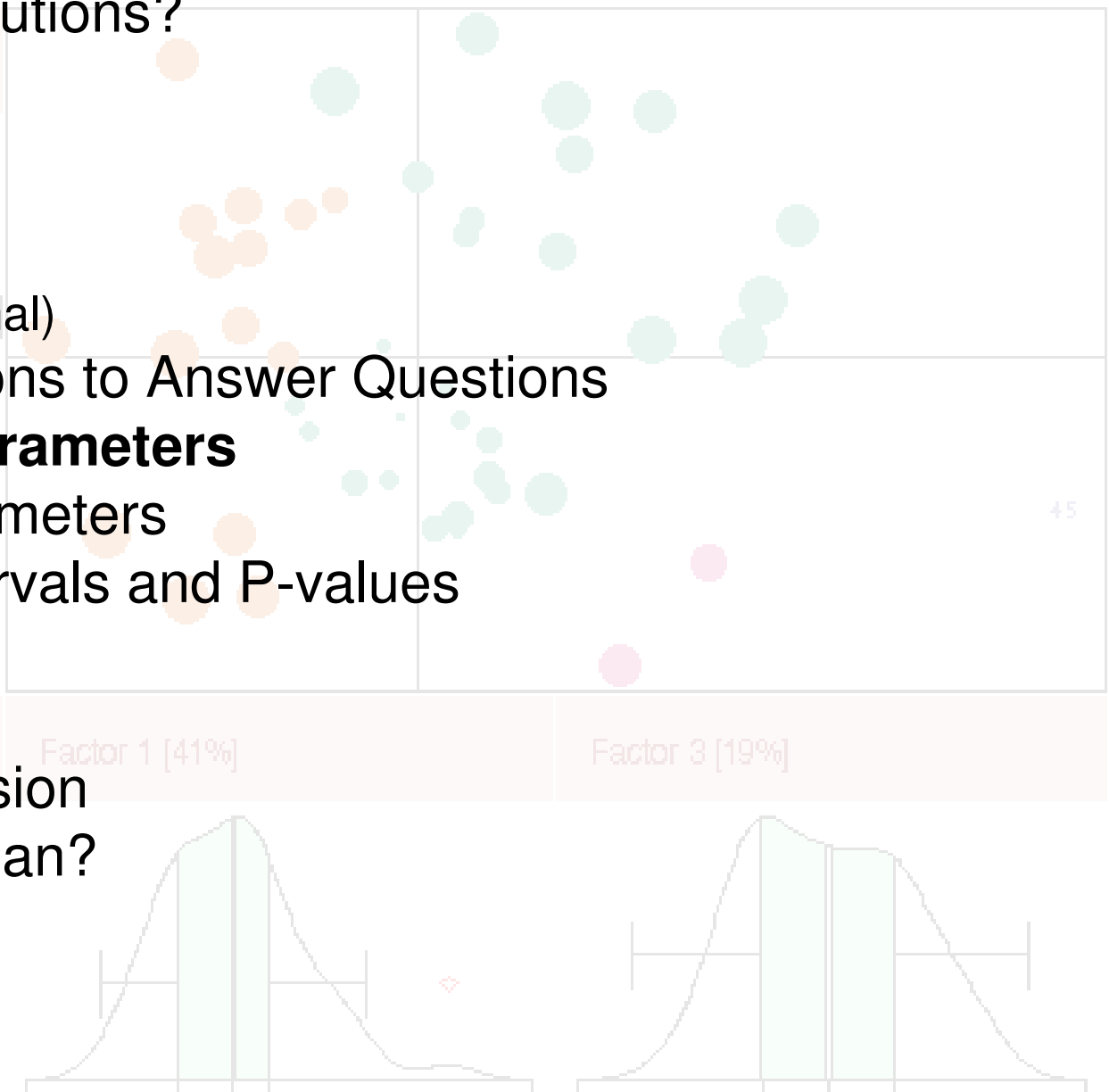








- What are Distributions?
- Models
 - Binomial
 - Poisson
 - Uniform
 - Gaussian (normal)
- Using Distributions to Answer Questions
- **Distribution Parameters**
- Estimating Parameters
- Confidence Intervals and P-values
- Why Gaussian?
- Regression
- Logistic Regression
- Why Not Gaussian?
- Bootstrapping
- Multiple Testing
- Useful Tools

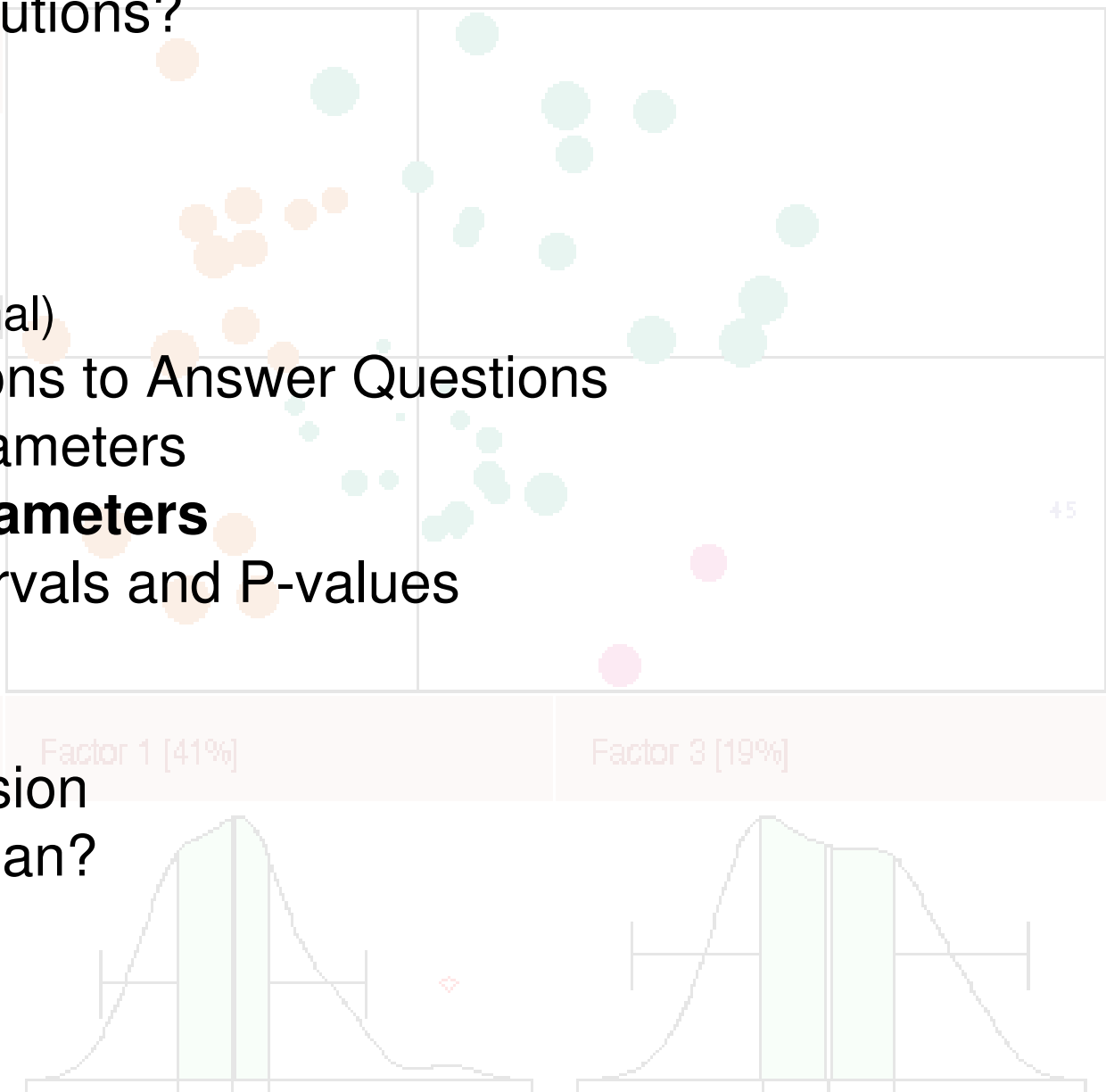




Creative Commons licensed, from gr0don on flickr

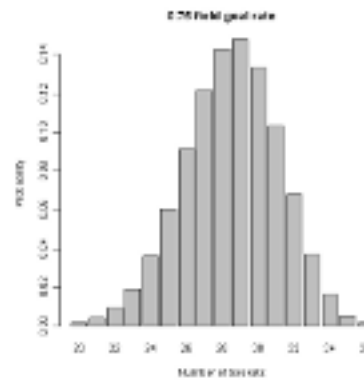
21/61

- What are Distributions?
- Models
 - Binomial
 - Poisson
 - Uniform
 - Gaussian (normal)
- Using Distributions to Answer Questions
- Distribution Parameters
- **Estimating Parameters**
- Confidence Intervals and P-values
- Why Gaussian?
- Regression
- Logistic Regression
- Why Not Gaussian?
- Bootstrapping
- Multiple Testing
- Useful Tools



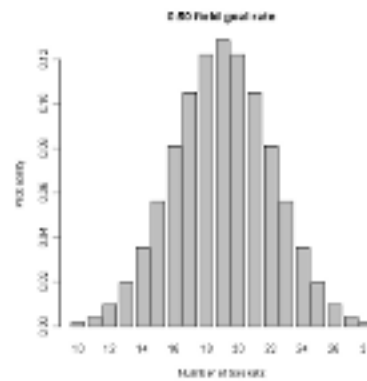
21/38

$p=0.75$



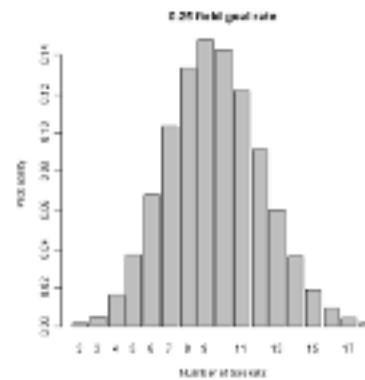
20,28,31,25,29/38

$p=0.5$



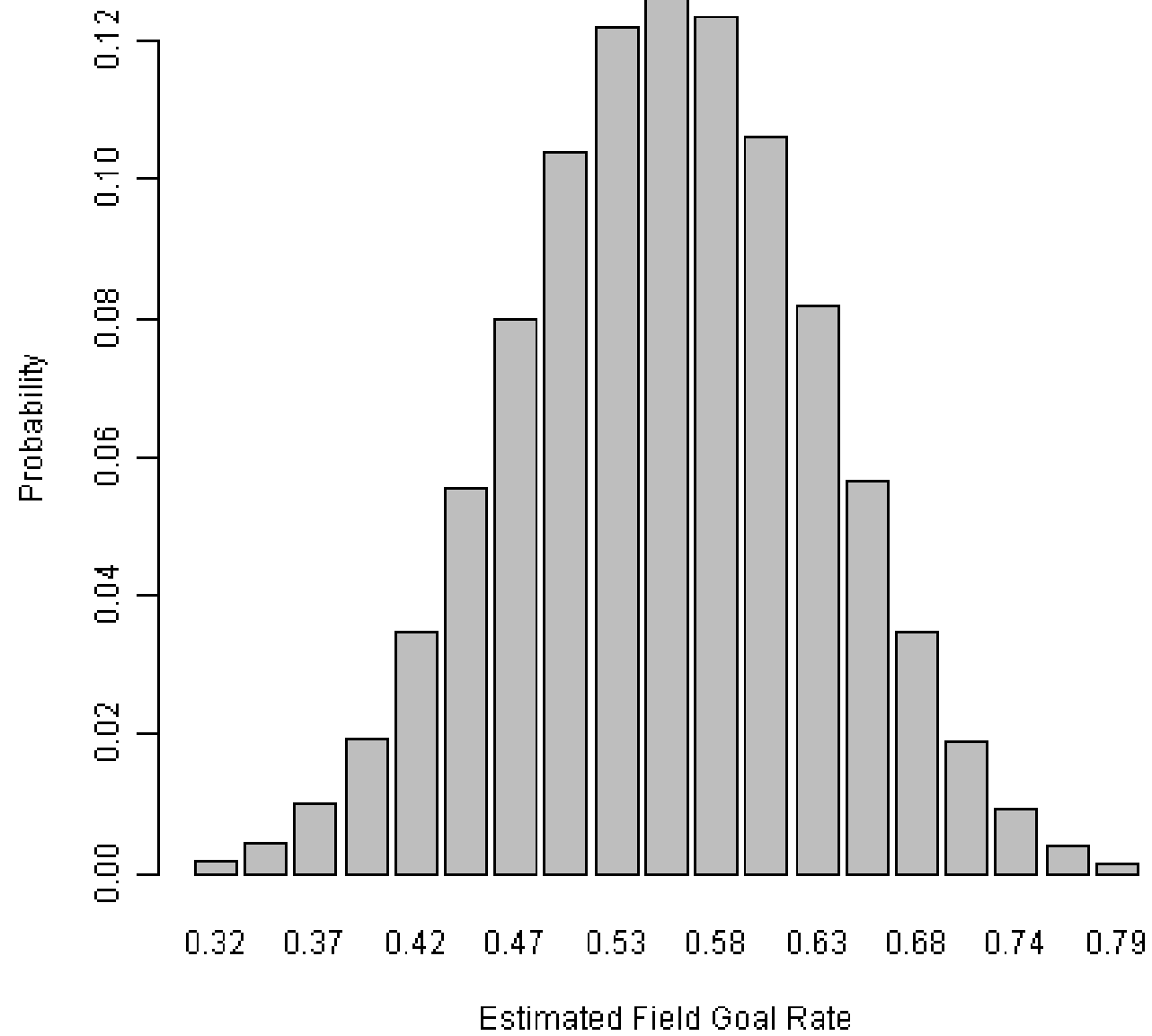
22,15,21,18,18/38

$p=0.25$

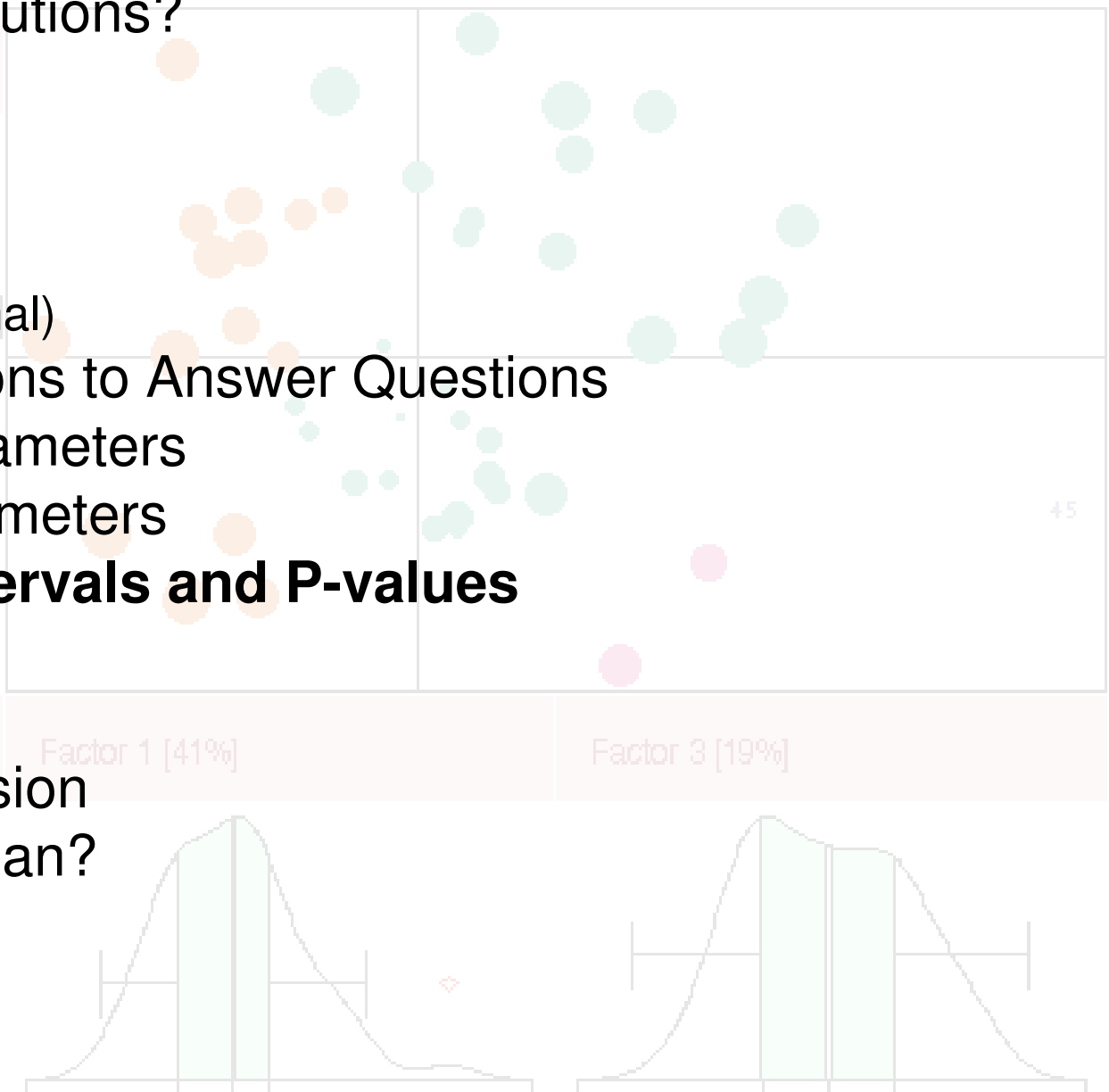


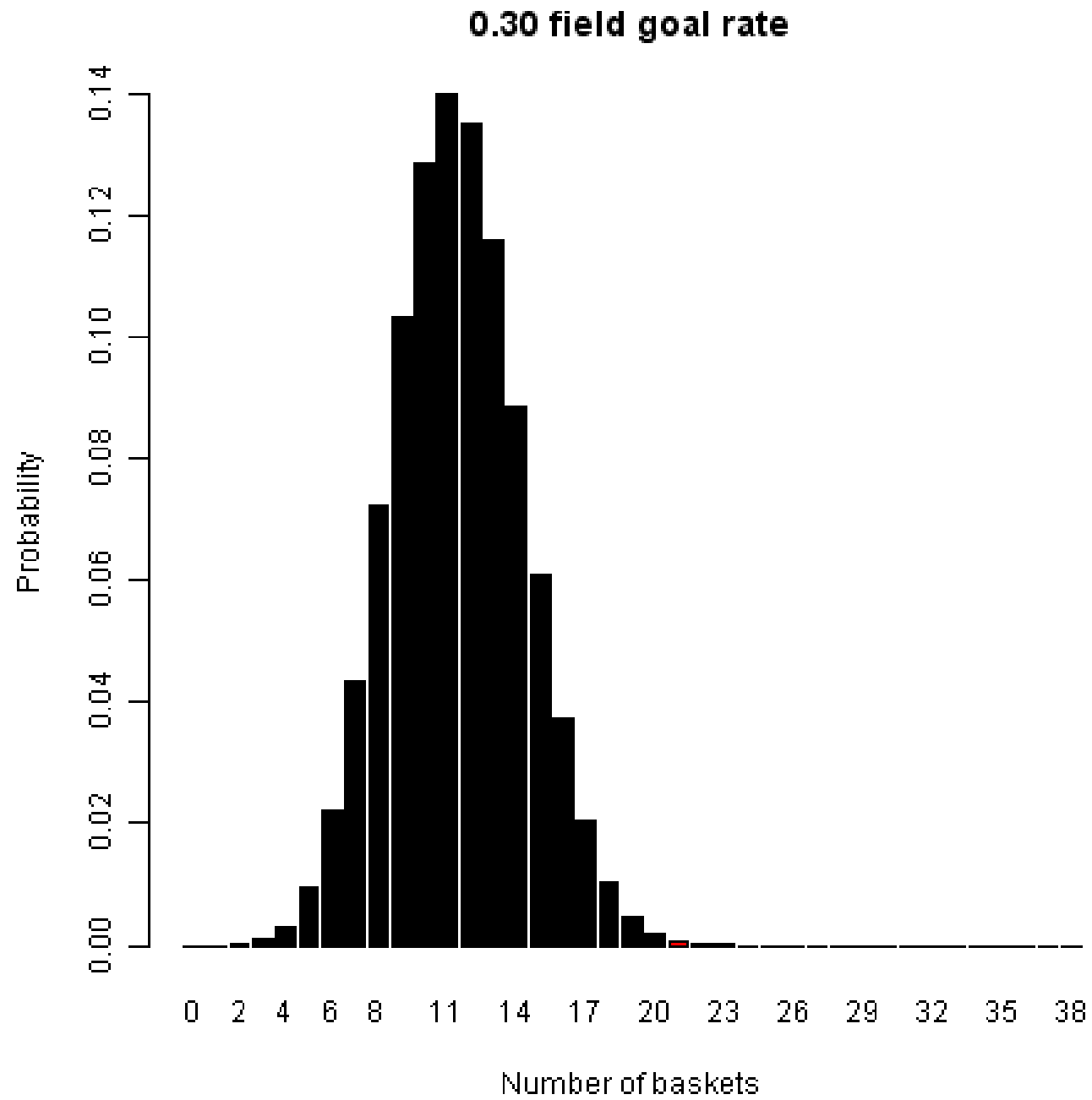
12,10,11,11,7/38

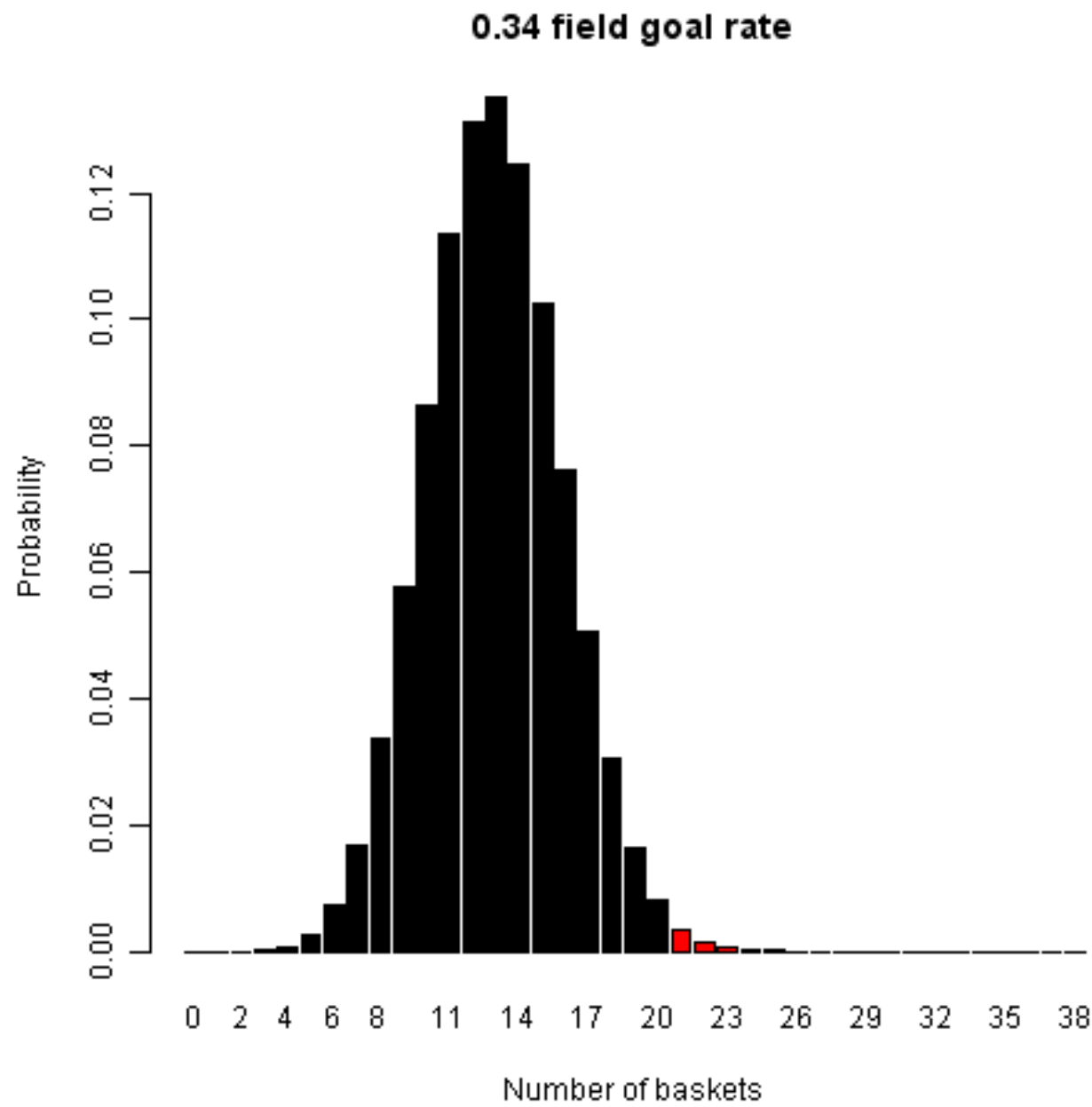
0.55 Field Goal rate

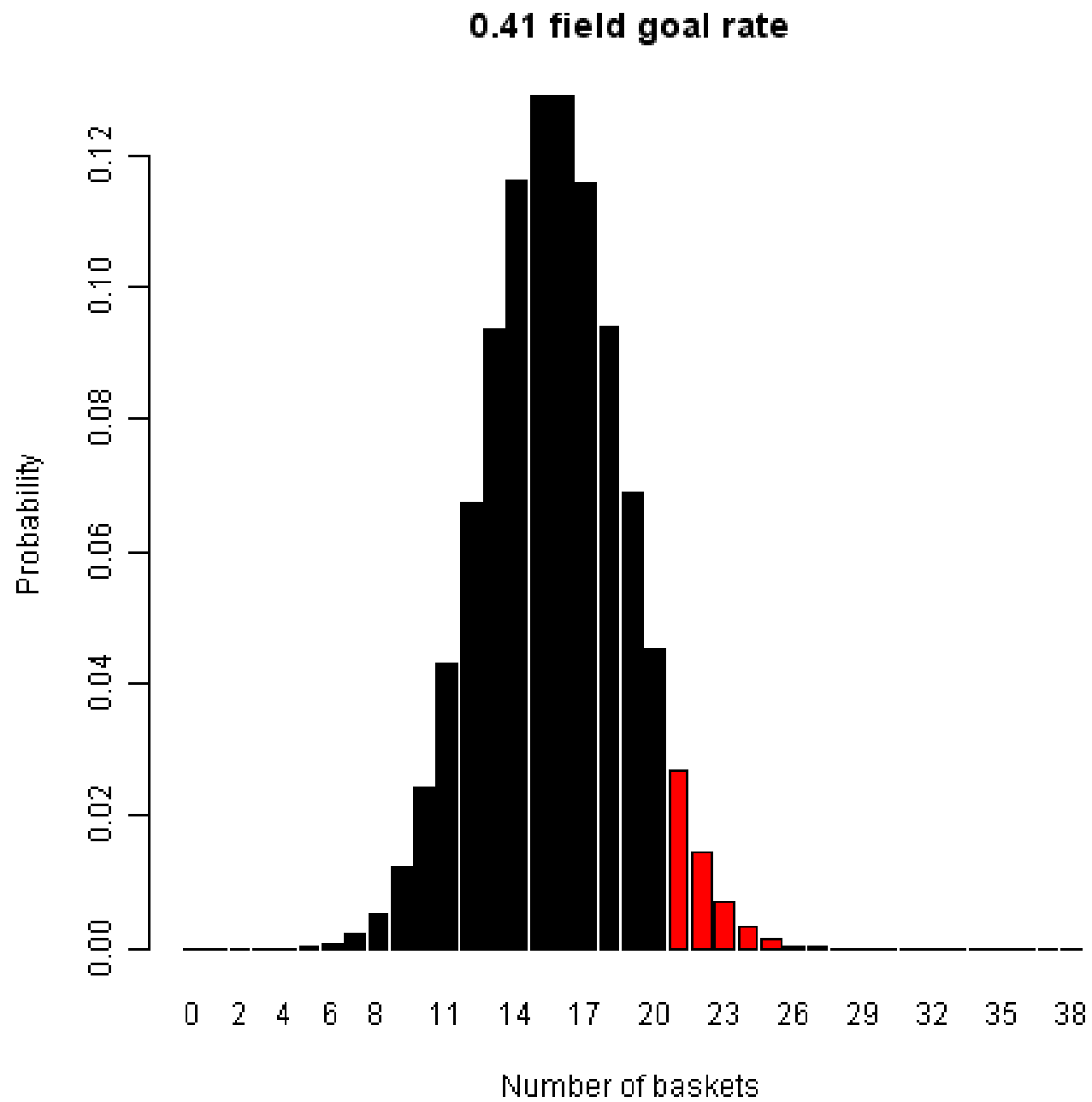


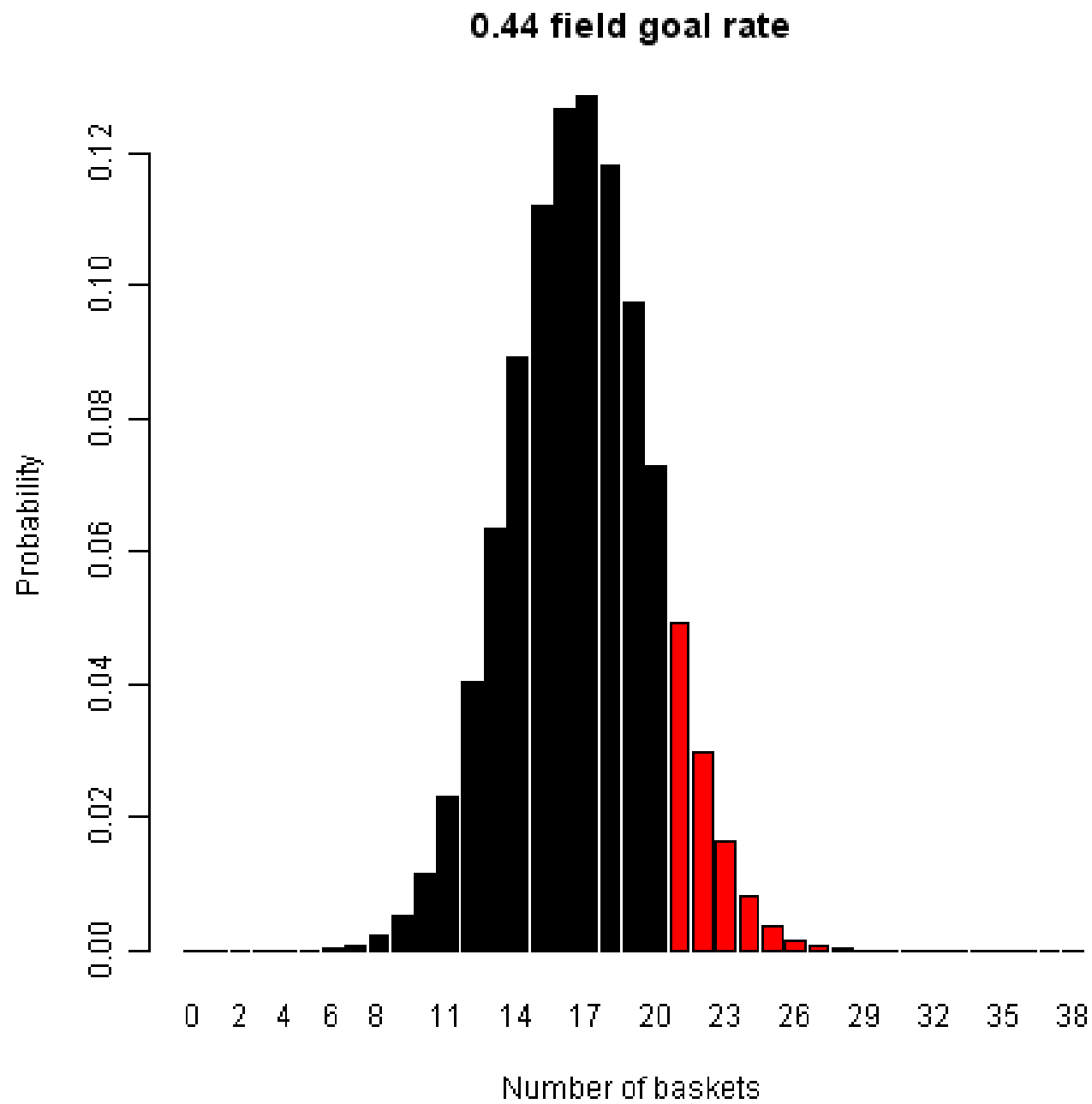
- What are Distributions?
- Models
 - Binomial
 - Poisson
 - Uniform
 - Gaussian (normal)
- Using Distributions to Answer Questions
- Distribution Parameters
- Estimating Parameters
- **Confidence Intervals and P-values**
- Why Gaussian?
- Regression
- Logistic Regression
- Why Not Gaussian?
- Bootstrapping
- Multiple Testing
- Useful Tools











| Actual Field Goal Rate | Probability of 21 or more successes |
|------------------------------|--|
| 0.30 | 0.001 |
| 0.34 | 0.01 |
| 0.41 | 0.05 |
| 0.44 | 0.1 |

- What are Distributions?

- Models

- Binomial
- Poisson
- Uniform
- Gaussian (normal)

- Using Distributions to Answer Questions

- Distribution Parameters

- Estimating Parameters

- Confidence Intervals and P-values

- **Why Gaussian?**

- Regression

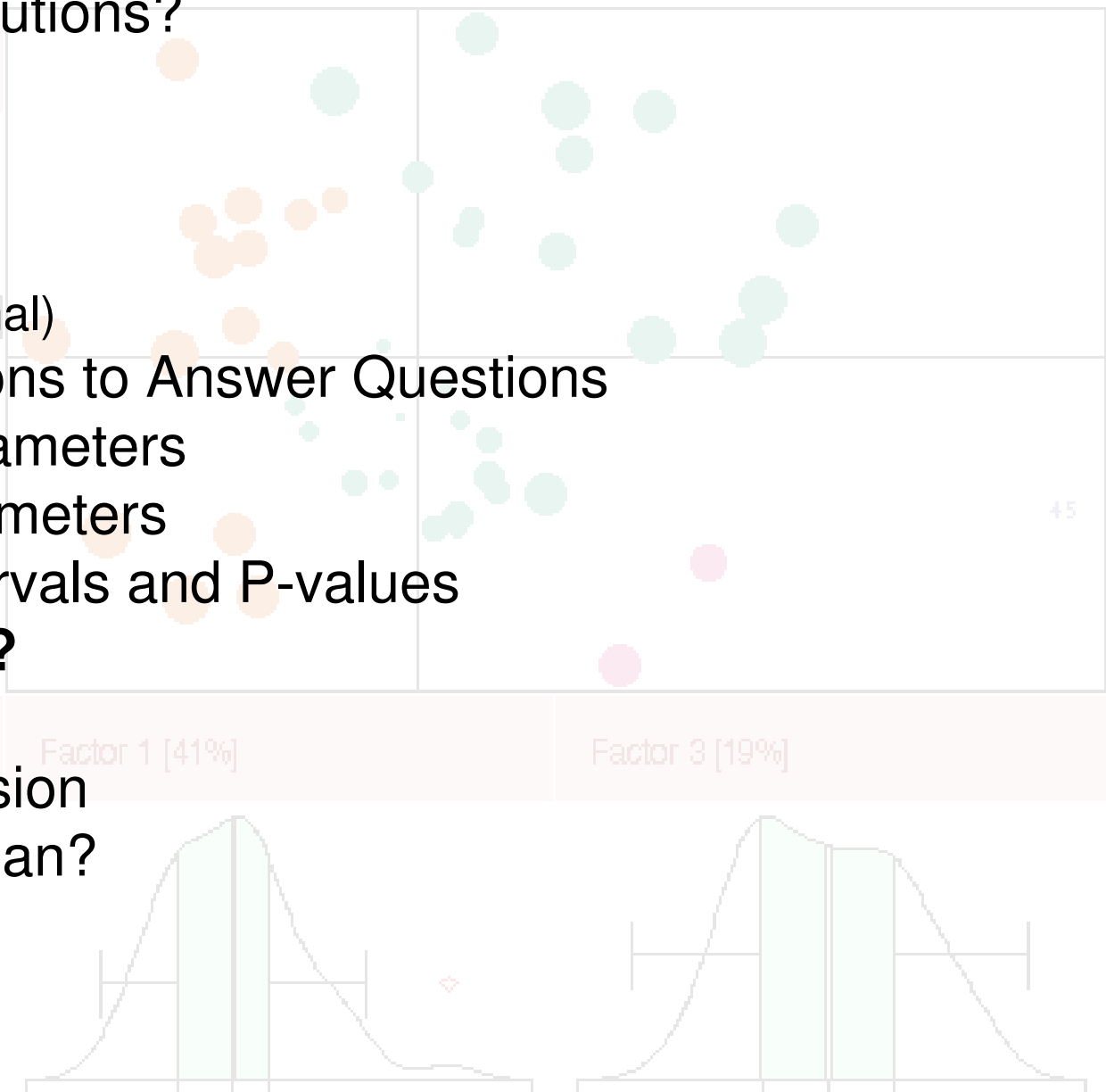
- Logistic Regression

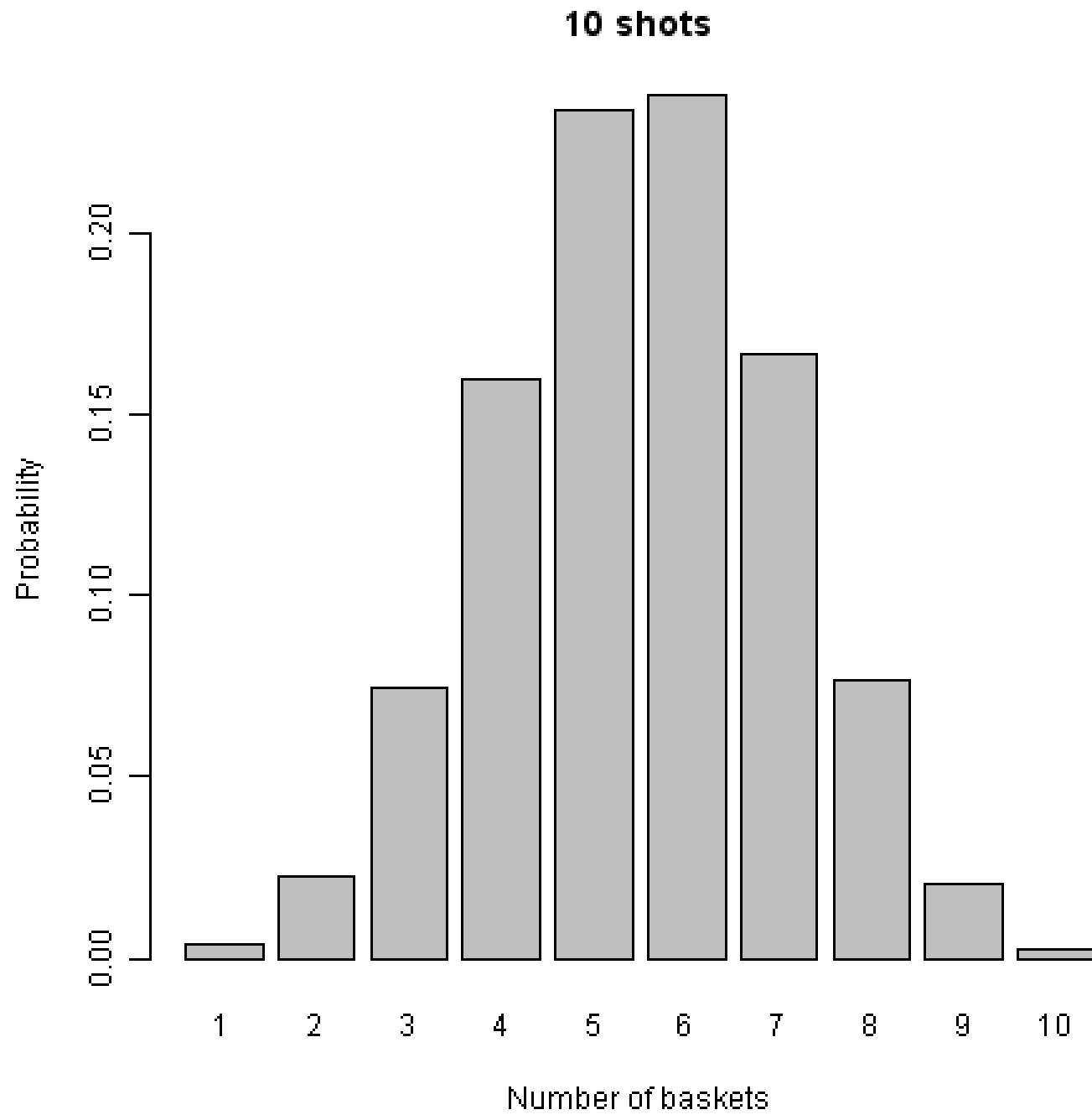
- Why Not Gaussian?

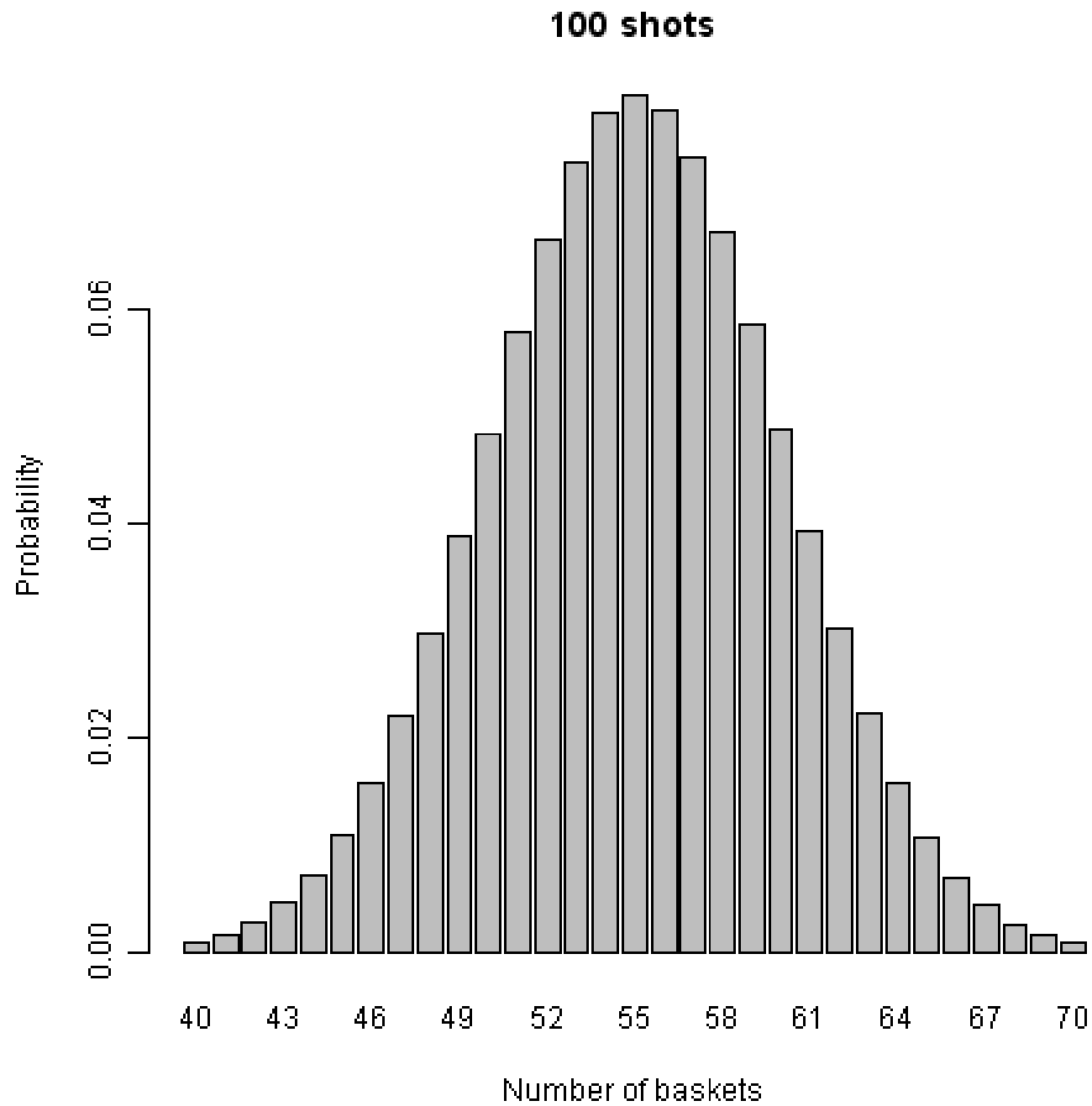
- Bootstrapping

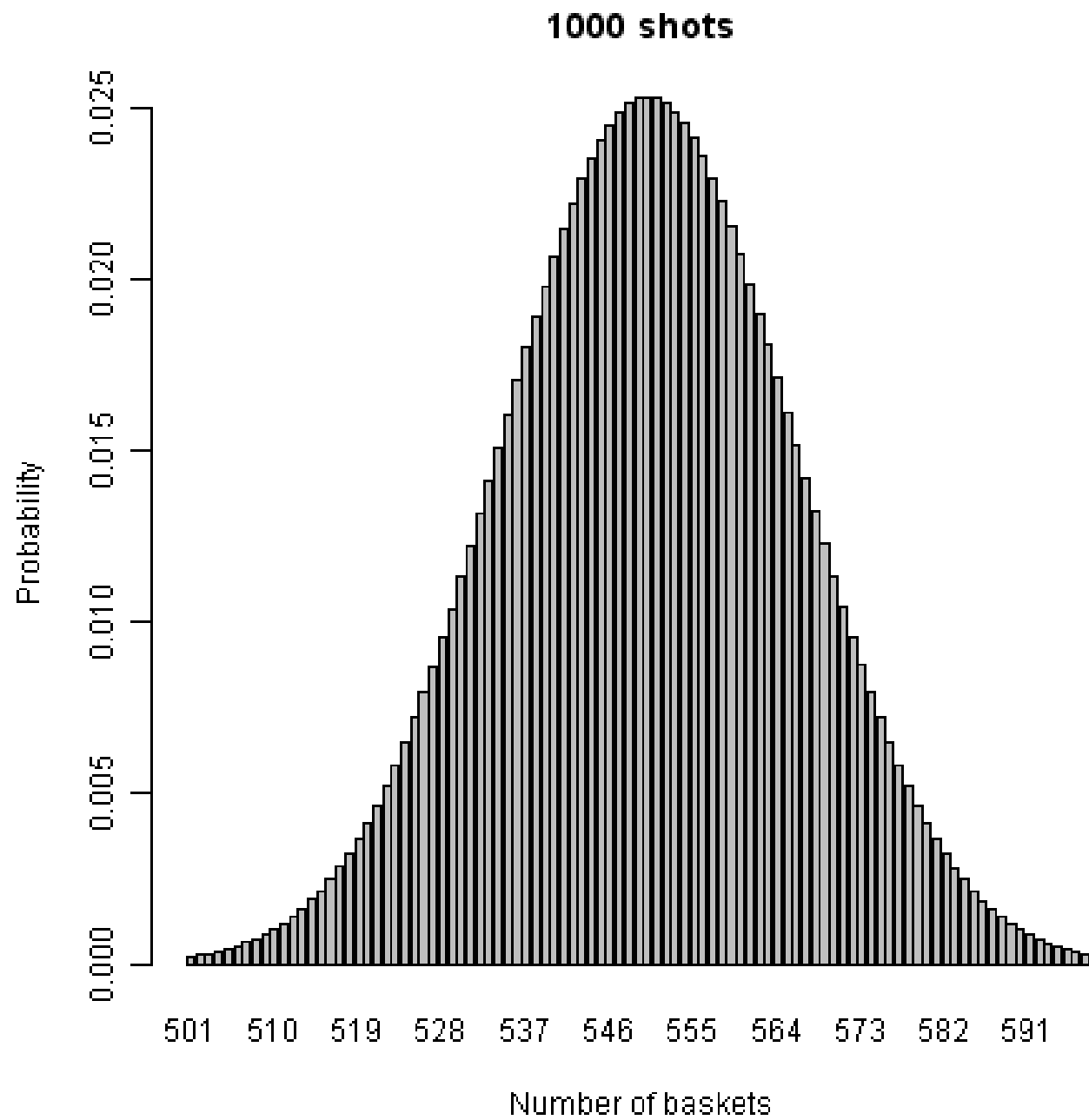
- Multiple Testing

- Useful Tools









- What are Distributions?

- Models

- Binomial
- Poisson
- Uniform
- Gaussian (normal)

- Using Distributions to Answer Questions

- Distribution Parameters

- Estimating Parameters

- Confidence Intervals and P-values

- Why Gaussian?

- **Regression**

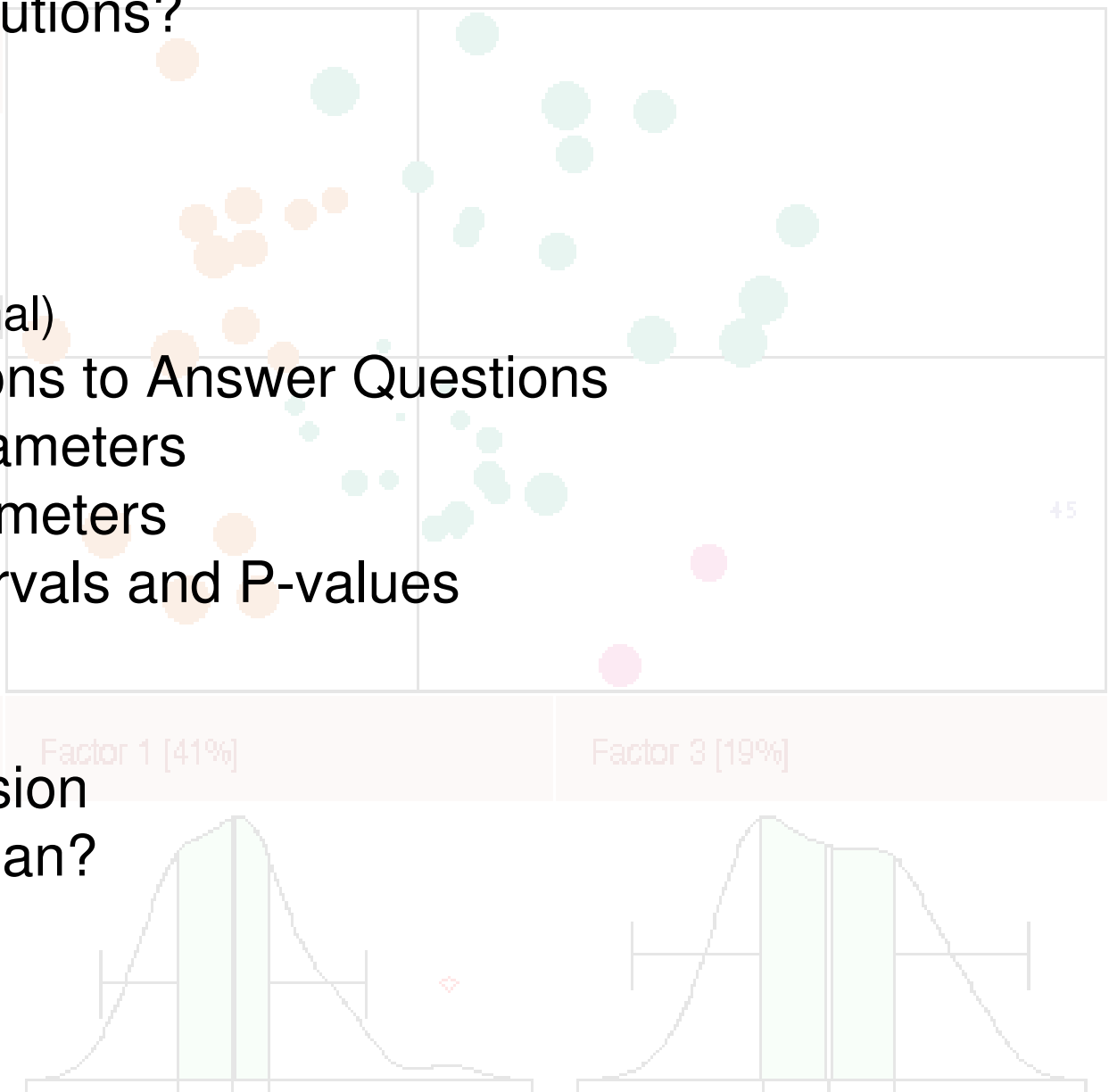
- Logistic Regression

- Why Not Gaussian?

- Bootstrapping

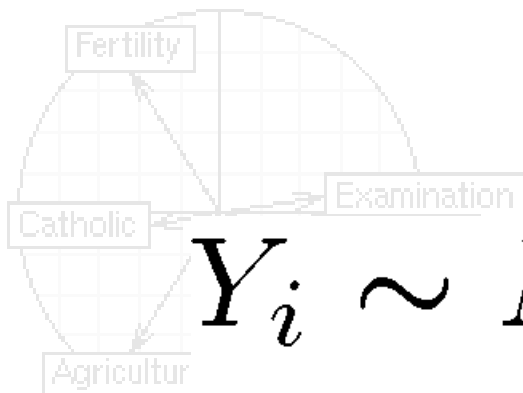
- Multiple Testing

- Useful Tools



PCA 5 vars

`princomp(x = data, cor = cor)`



$$Y_i \sim N(\beta_0 + \beta_1 x_i, \sigma^2)$$

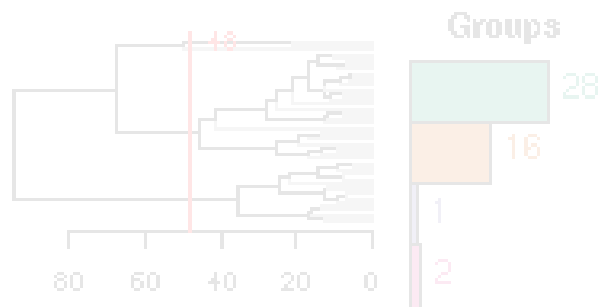
$$Y_i = \beta_0 + \beta_1 x_i + \epsilon_i$$

45

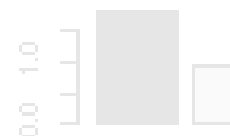
Clustering 4 groups

Factor 1 [41%]

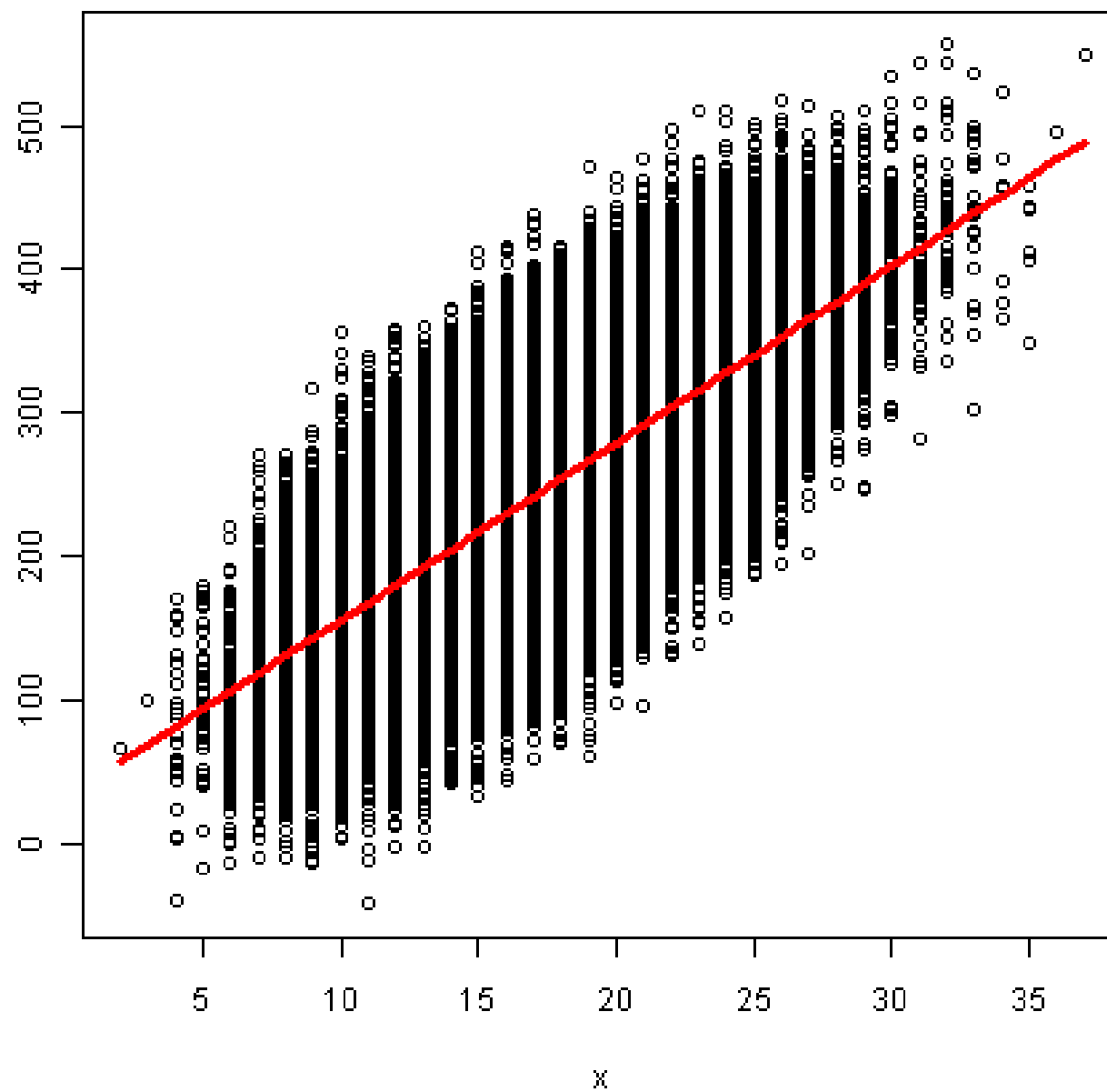
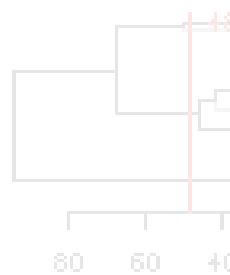
Factor 3 [19%]



PCA 5 vars
princomp(x = dat



Clustering 4

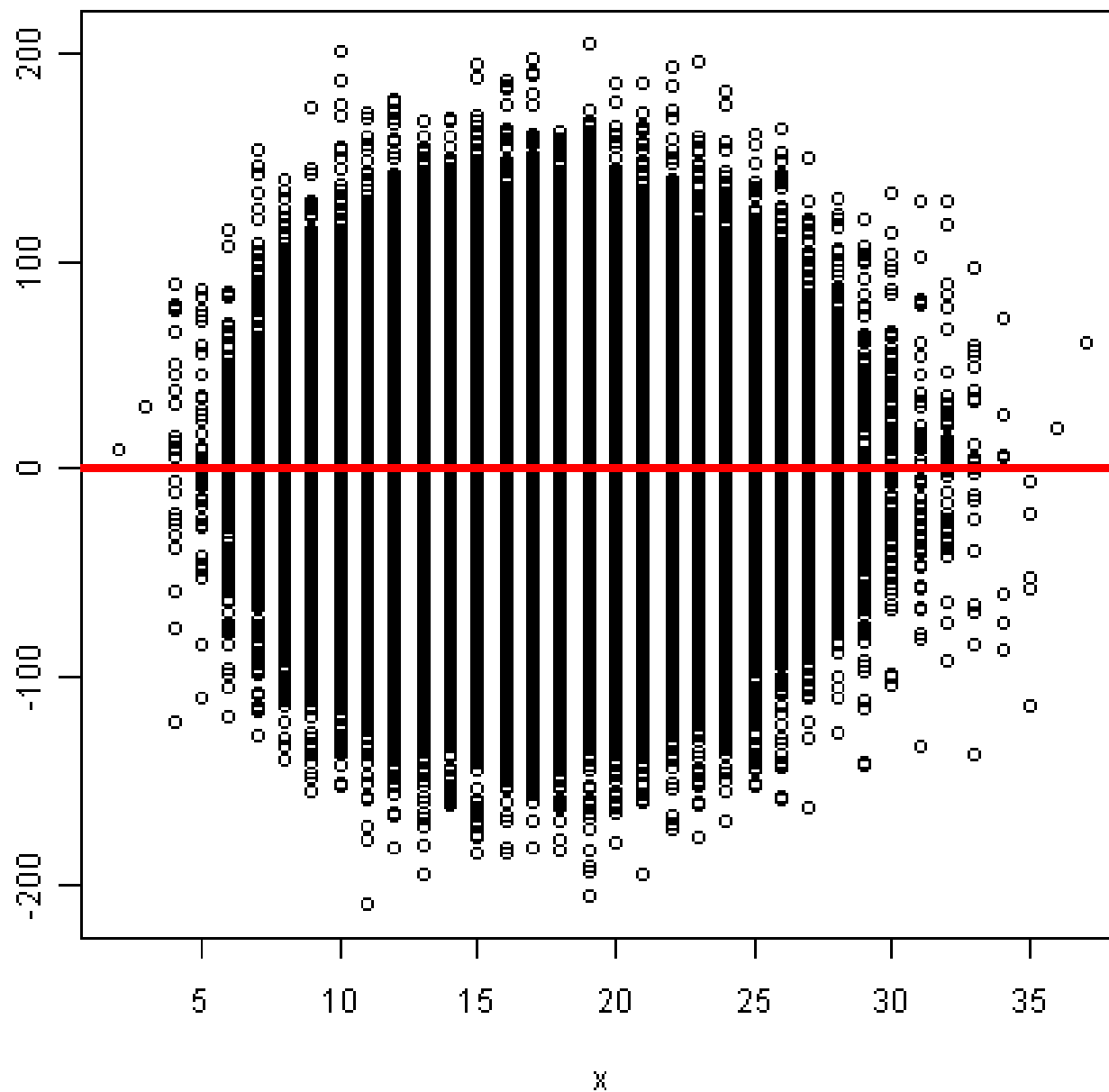
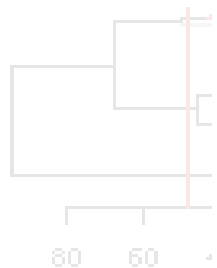


PCA 5 vars

princomp(x = d)

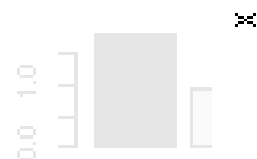


Clustering

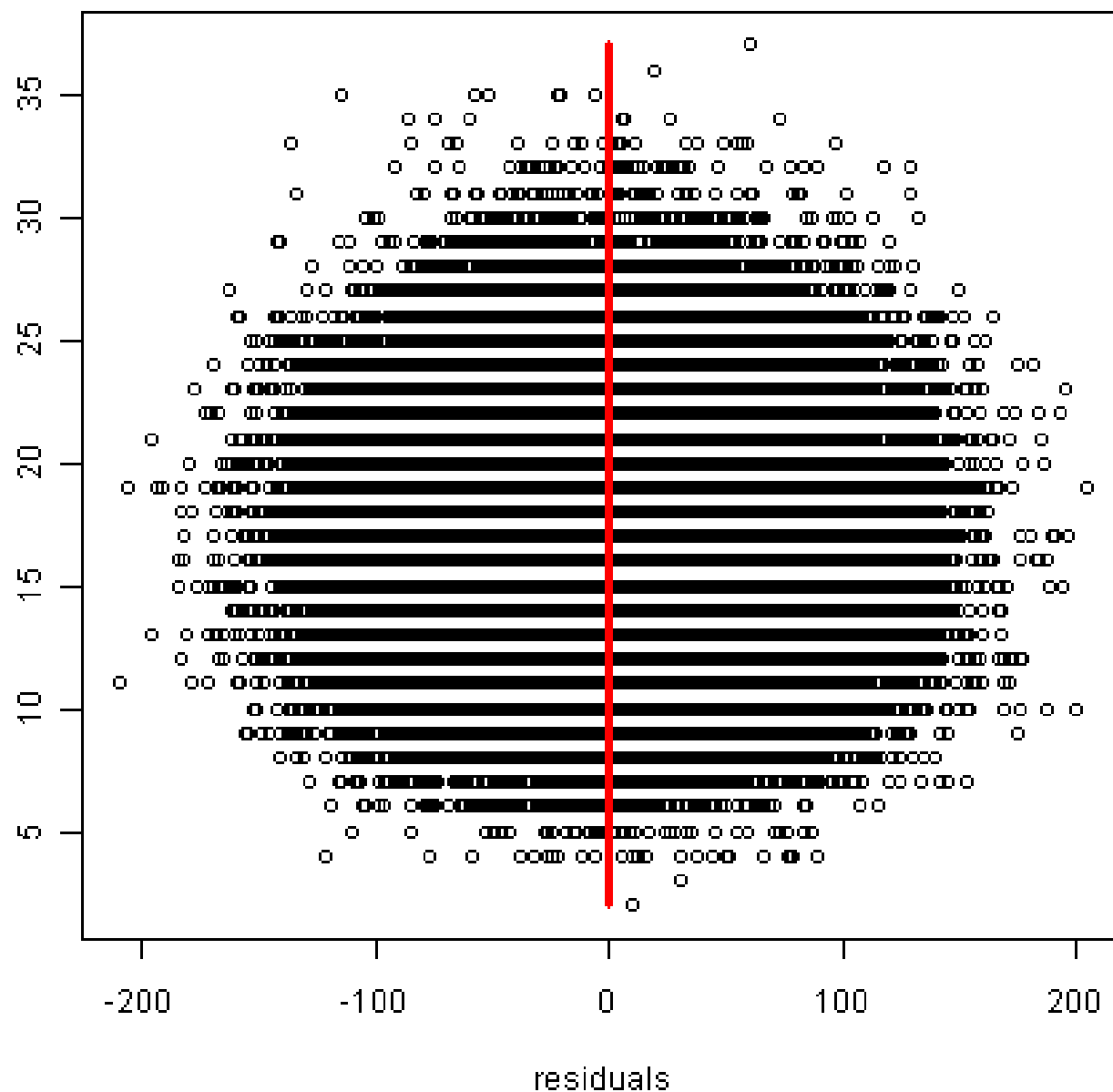
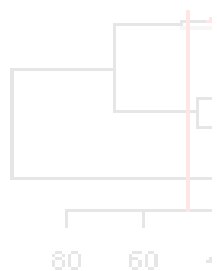


PCA 5 vars

princomp(x = d)

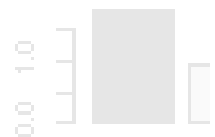


Clustering 2

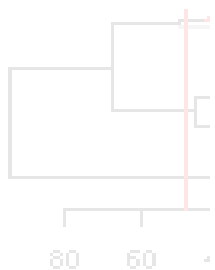


PCA 5 vars

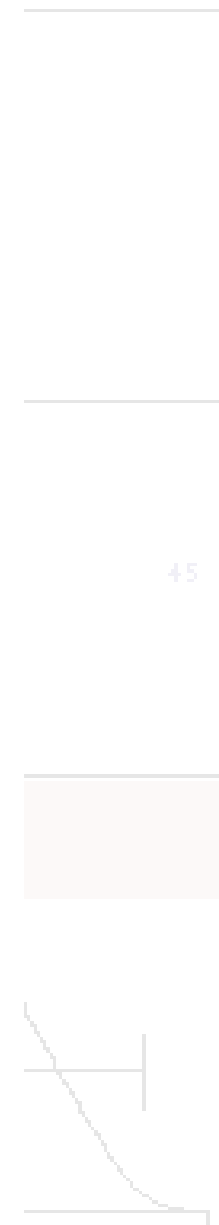
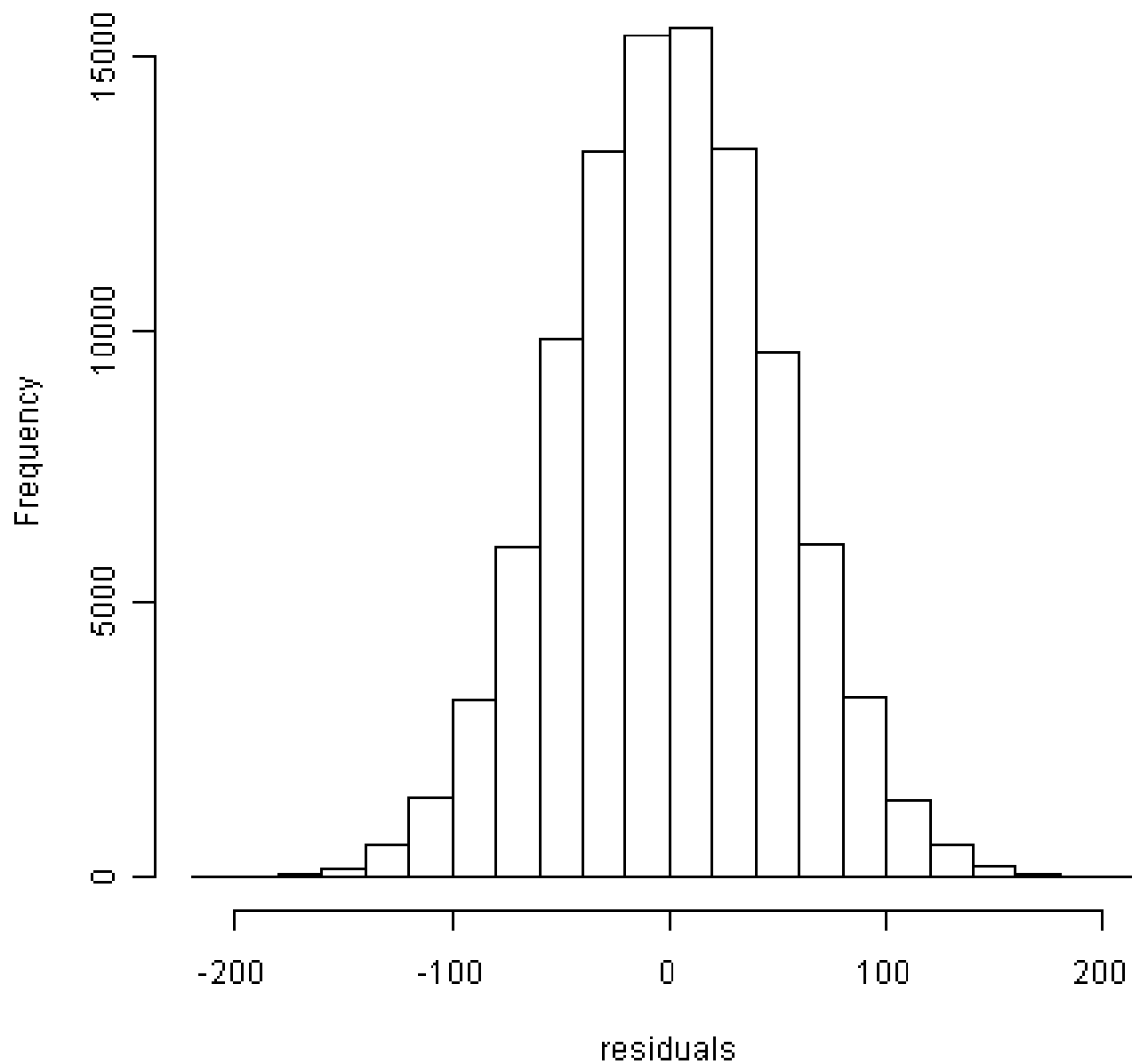
```
princomp(x = d[
```



Clustering 4

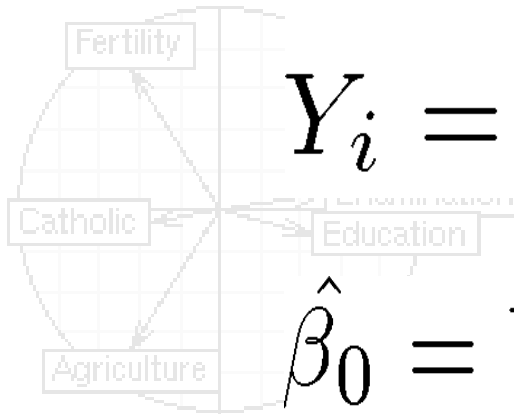


Histogram of residuals



PCA 5 vars

princomp(x = data, cor = cor)



$$Y_i \sim N(\beta_0 + \beta_1 x_i, \sigma^2)$$

$$Y_i = \beta_0 + \beta_1 x_i + \epsilon_i$$

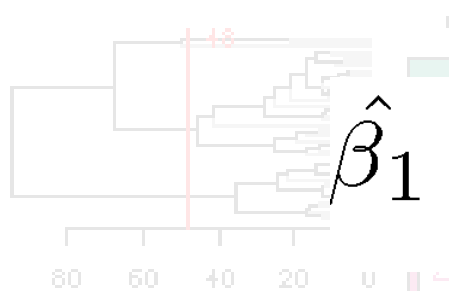
$$\hat{\beta}_0 = \bar{Y} - \hat{\beta}_1 \bar{x}$$

$$\hat{\beta}_1 = \frac{\sum_{i=1}^n (x_i - \bar{x})(Y_i - \bar{Y})}{\sum_{i=1}^n (x_i - \bar{x})^2}$$

Clustering 4 groups

Factor 1 [41%]

Factor 3 [19%]



$$\hat{\beta}_1 \sim N\left(\beta_1, \frac{\sigma^2}{\sum_{i=1}^n (x_i - \bar{x})^2}\right)$$

- What are Distributions?

- Models

- Binomial
- Poisson
- Uniform
- Gaussian (normal)

- Using Distributions to Answer Questions

- Distribution Parameters

- Estimating Parameters

- Confidence Intervals and P-values

- Why Gaussian?

- Regression

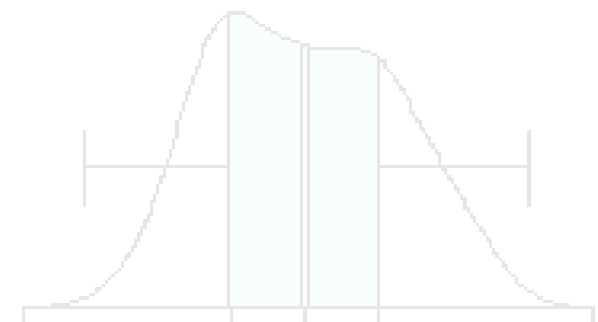
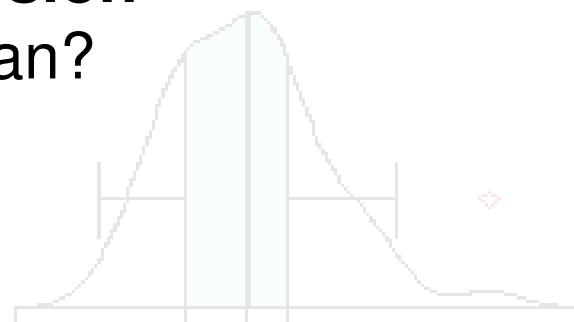
- **Logistic Regression**

- Why Not Gaussian?

- Bootstrapping

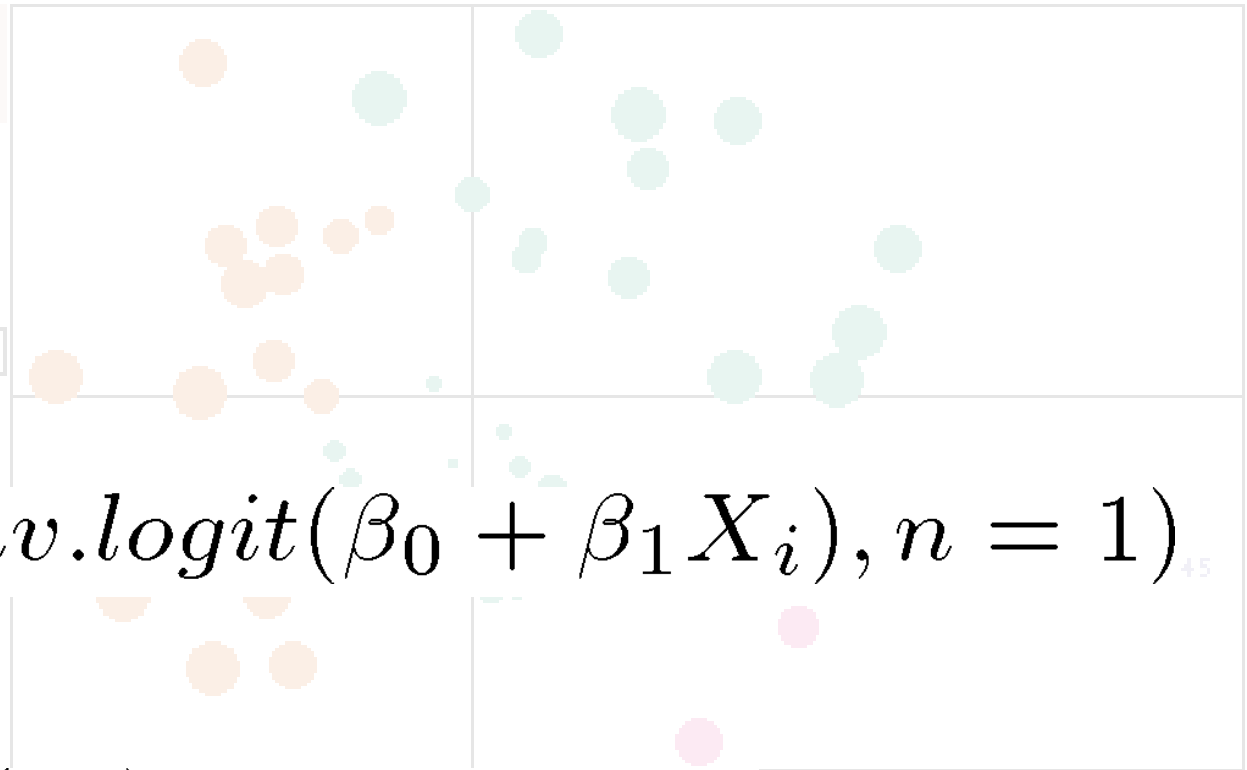
- Multiple Testing

- Useful Tools



PCA 5 vars

`princomp(x = data, cor = cor)`

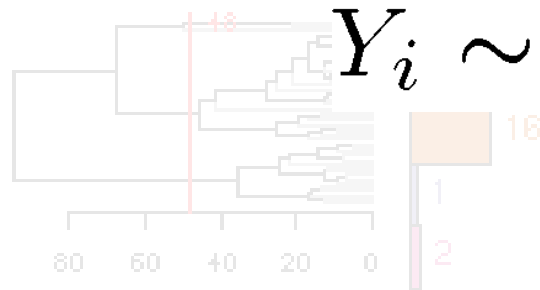


$$Y_i \sim \text{Bin}(\text{inv.logit}(\beta_0 + \beta_1 X_i), n = 1)$$

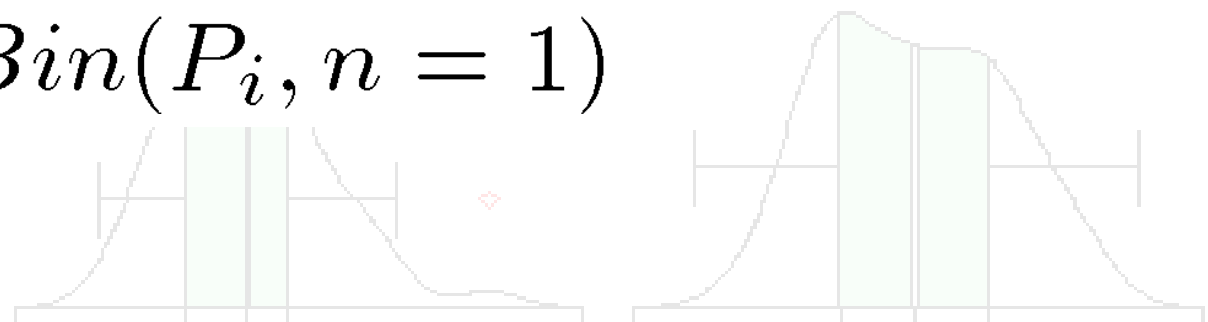


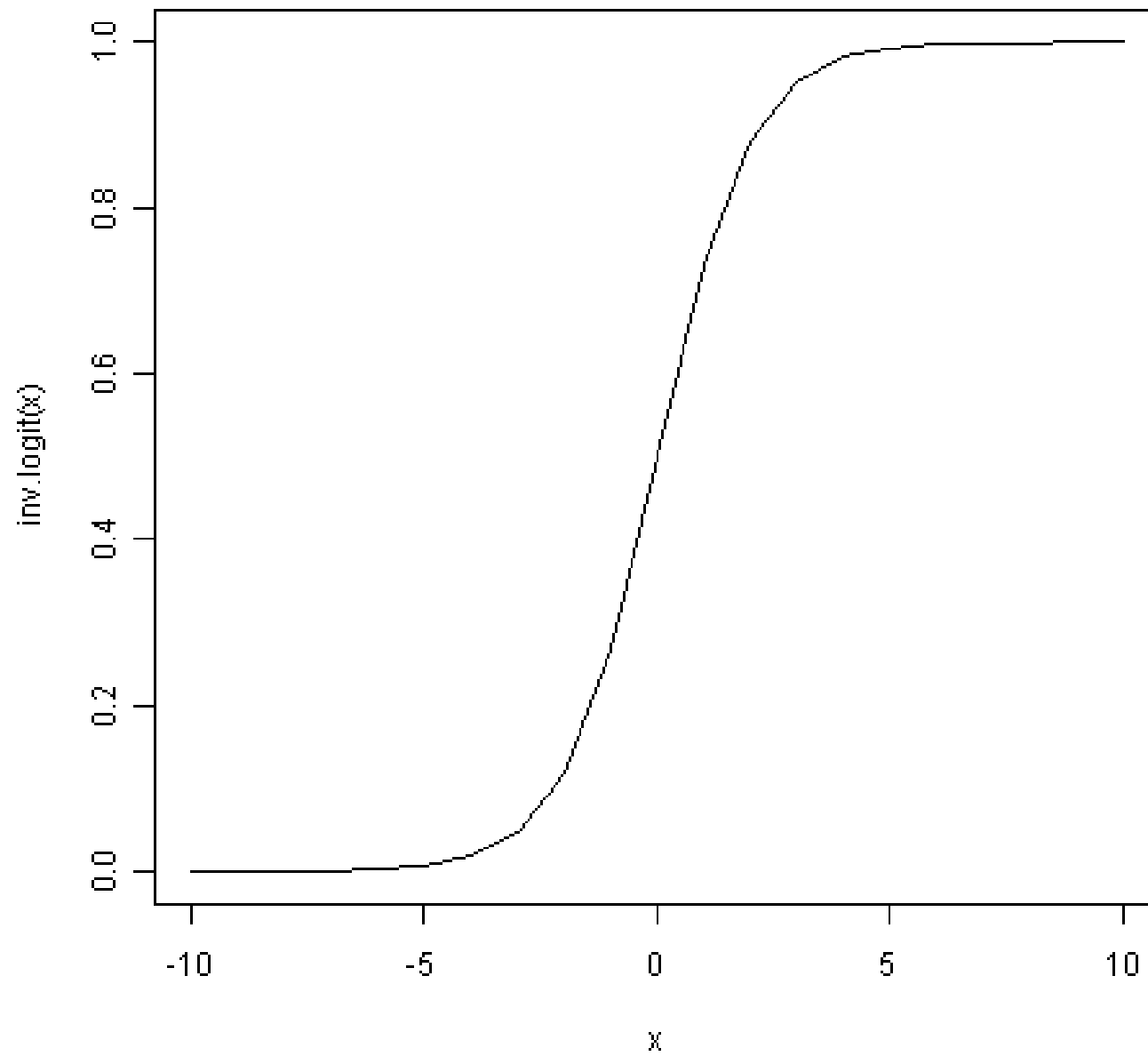
$$\text{Logit}(P_i) = \beta_0 + \beta_1 X_i$$

Clustering 4 groups



$$Y_i \sim \text{Bin}(P_i, n = 1)$$





- What are Distributions?

- Models

- Binomial
- Poisson
- Uniform
- Gaussian (normal)

- Using Distributions to Answer Questions

- Distribution Parameters

- Estimating Parameters

- Confidence Intervals and P-values

- Why Gaussian?

- Regression

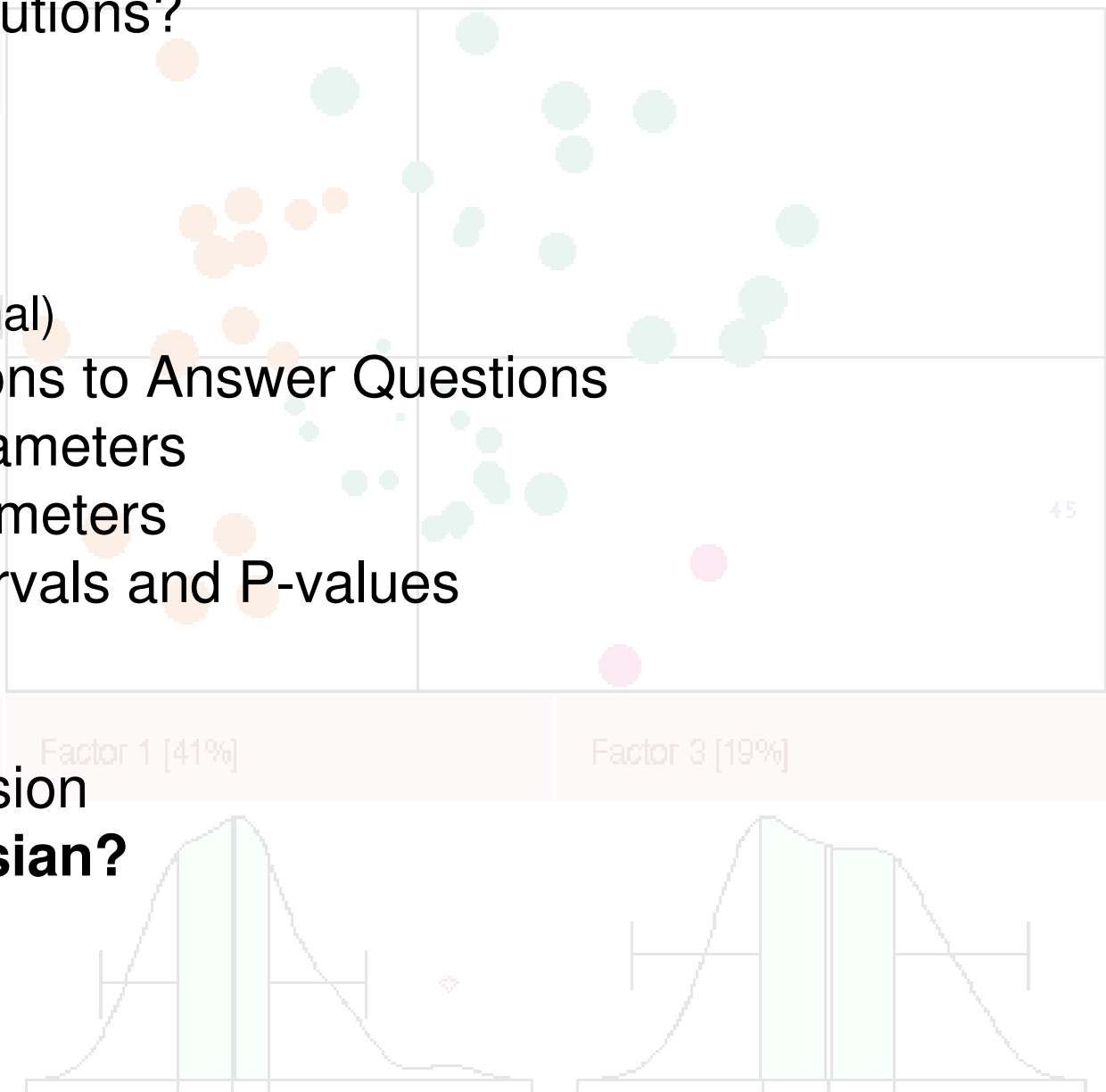
- Logistic Regression

- **Why Not Gaussian?**

- Bootstrapping

- Multiple Testing

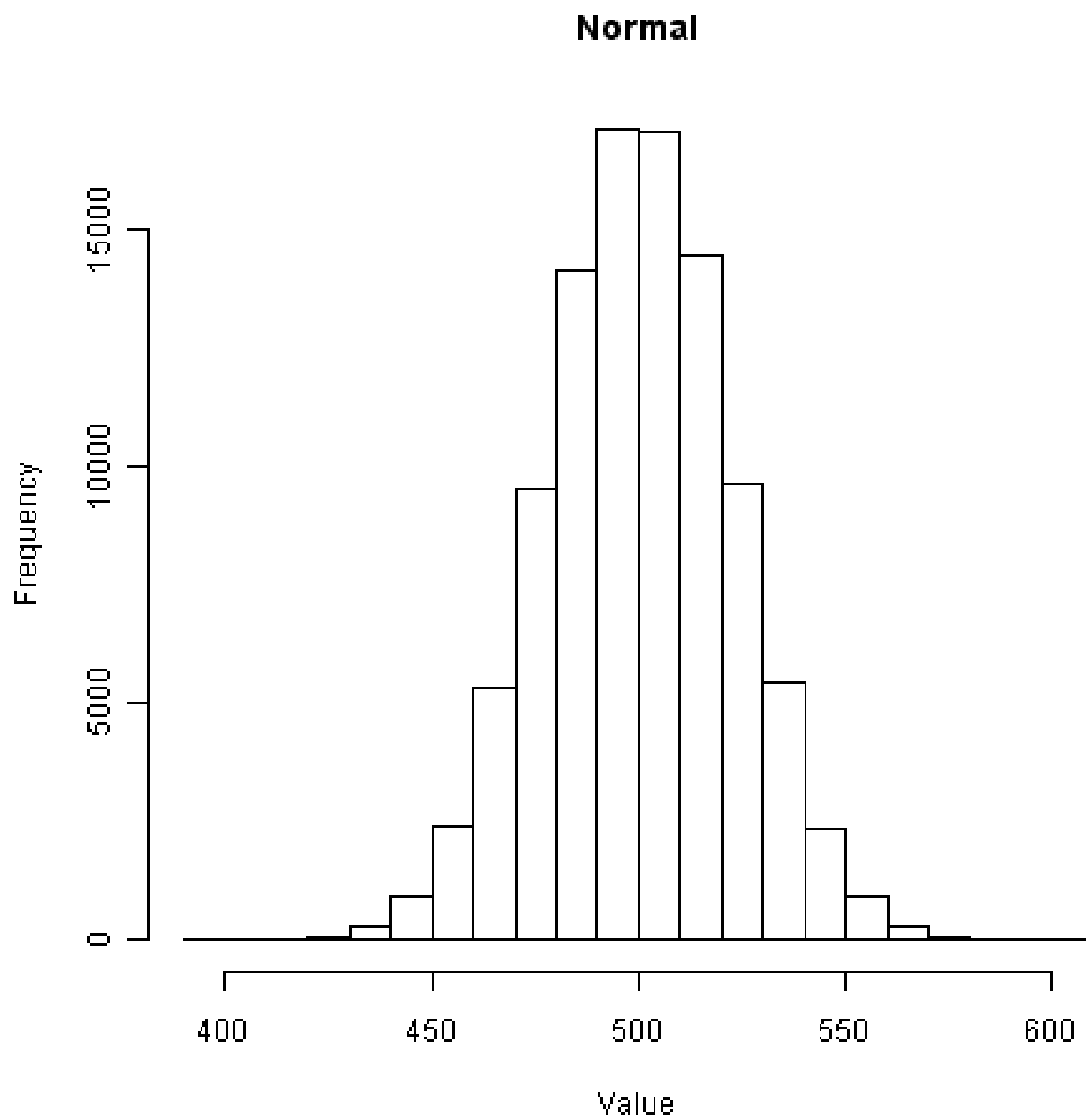
- Useful Tools



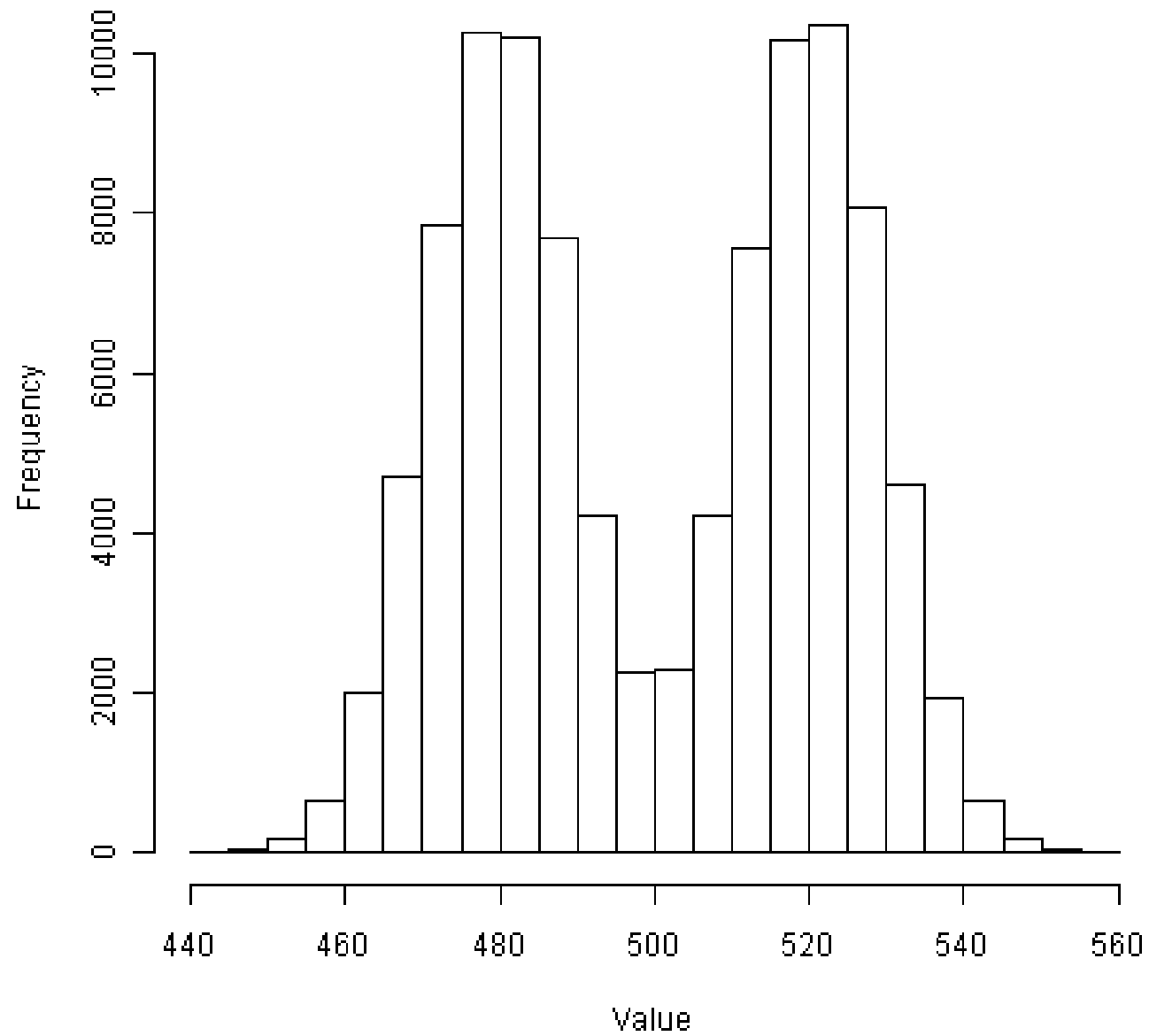


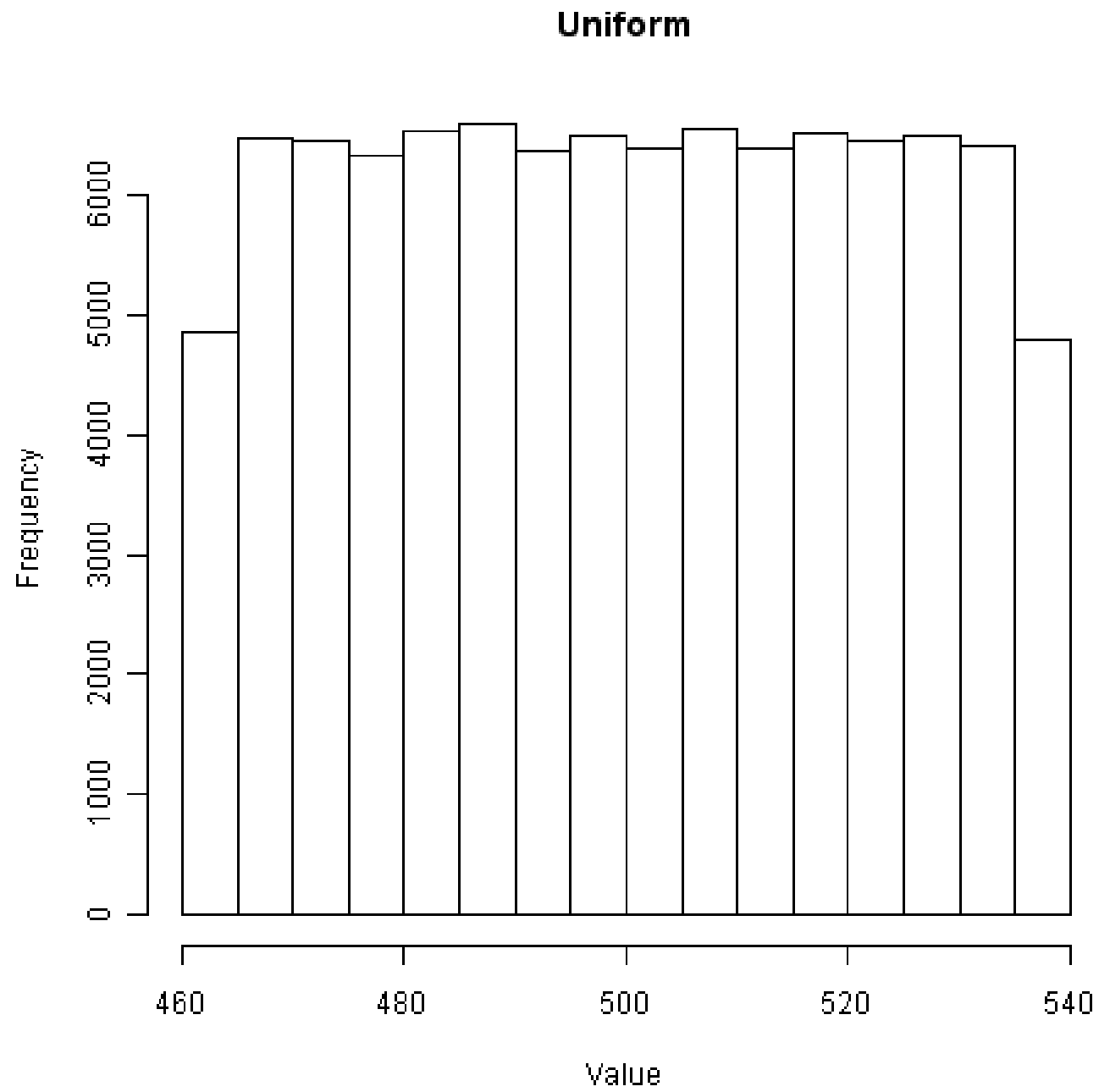
Creative Commons licensed, from Alan Vernon on flickr

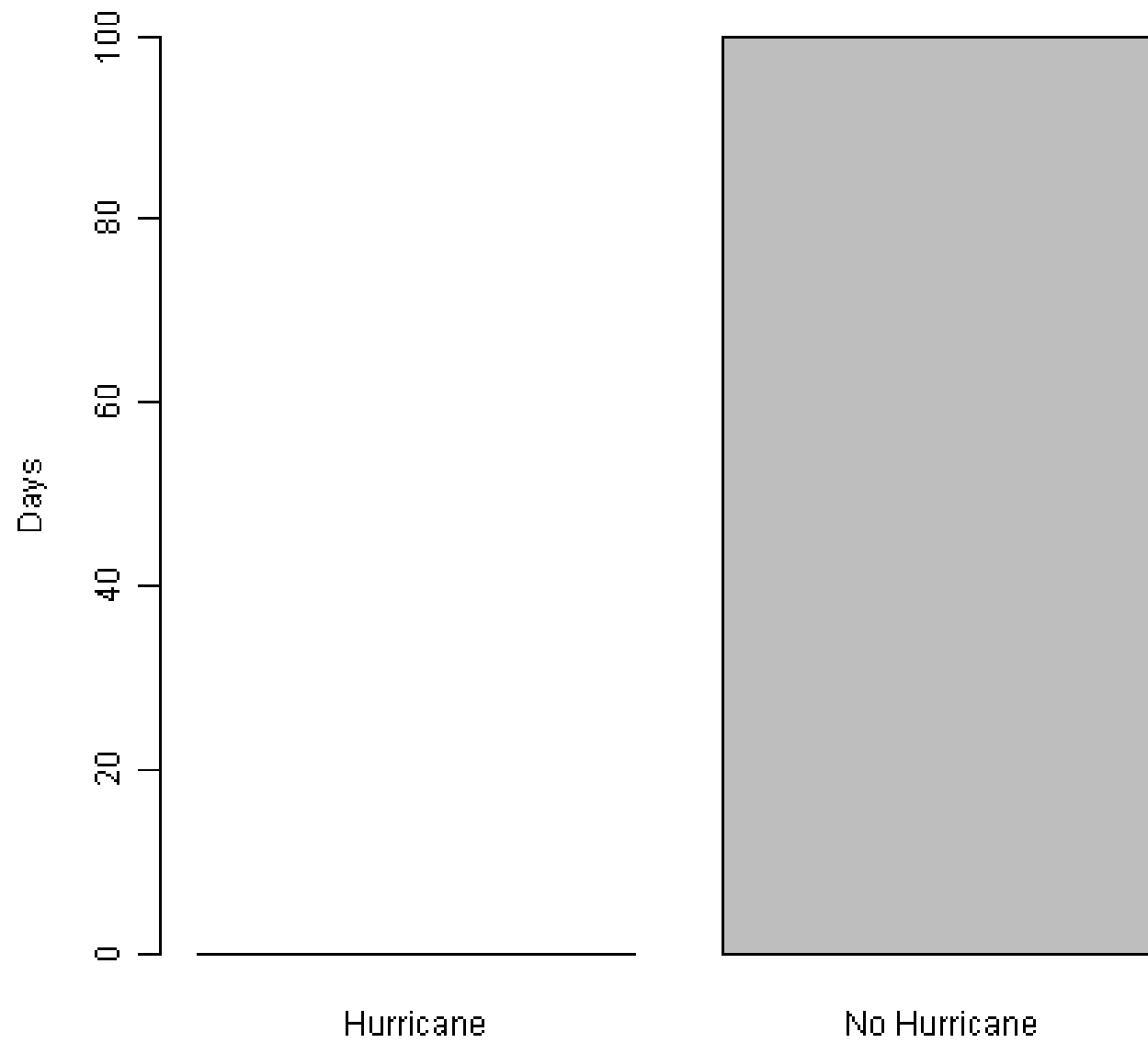
47/61

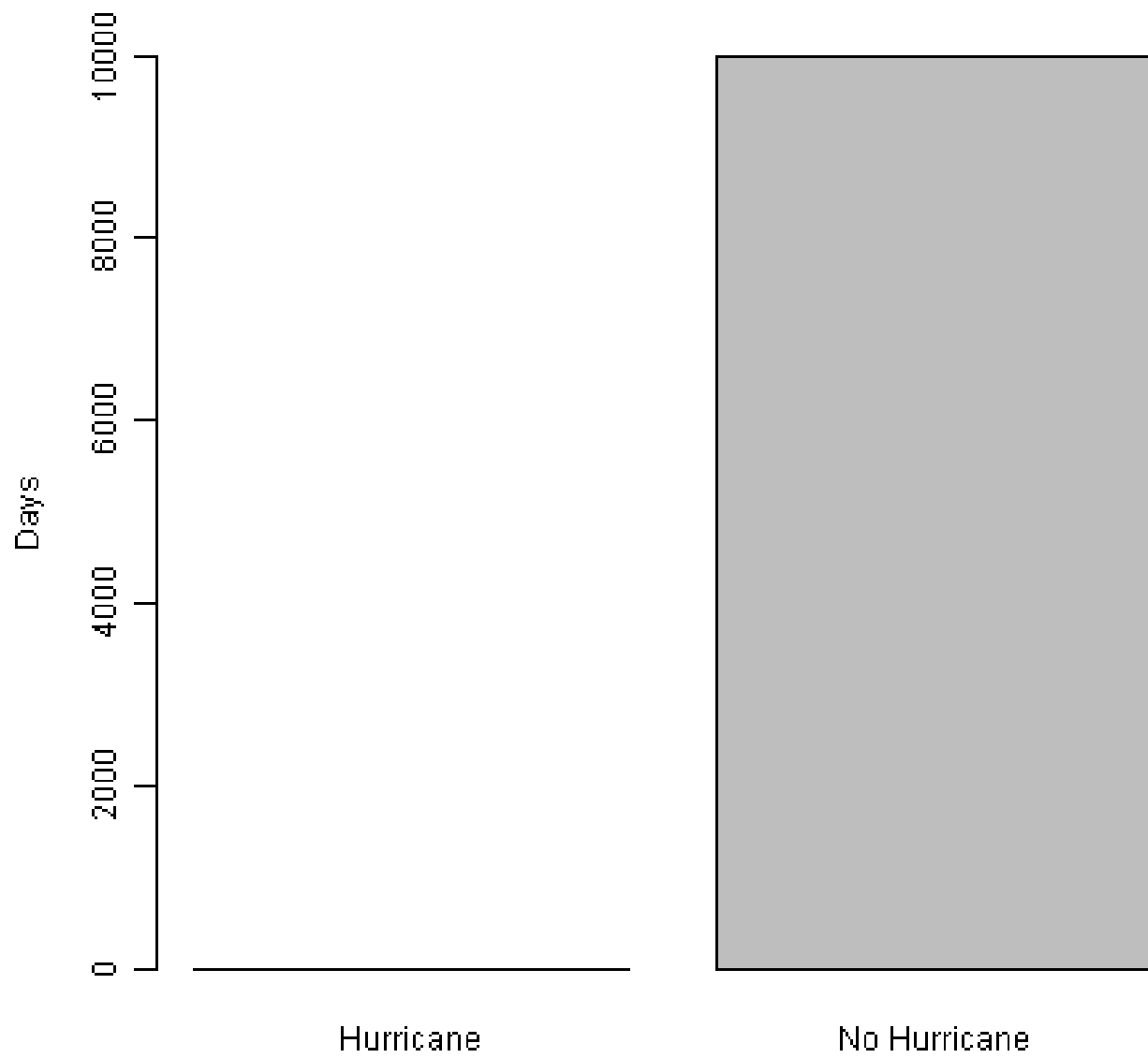


Bimodal









- What are Distributions?

- Models

- Binomial
- Poisson
- Uniform
- Gaussian (normal)

- Using Distributions to Answer Questions

- Distribution Parameters

- Estimating Parameters

- Confidence Intervals and P-values

- Why Gaussian?

- Regression

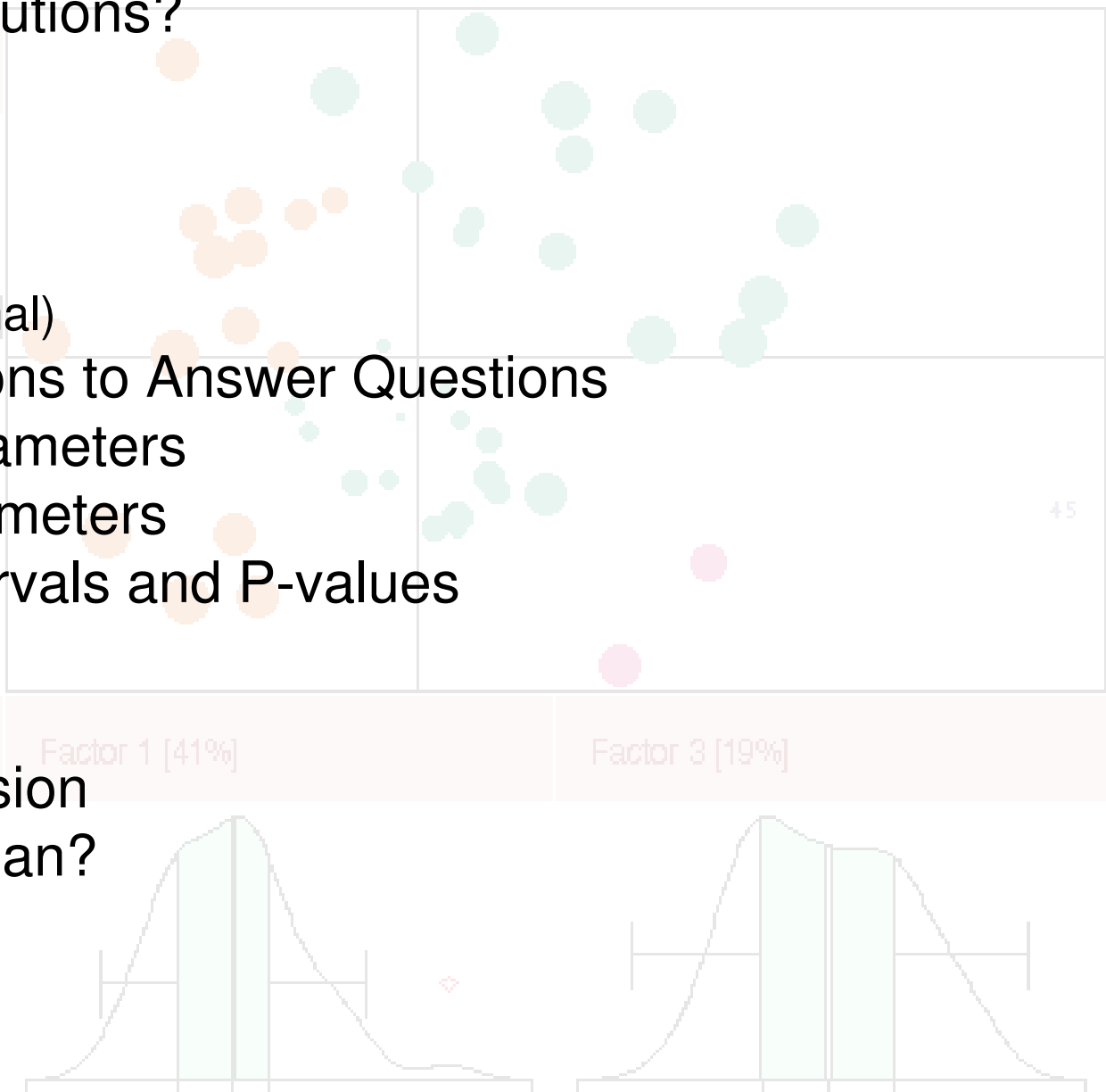
- Logistic Regression

- Why Not Gaussian?

- **Bootstrapping**

- Multiple Testing

- Useful Tools



- What are Distributions?

- Models

- Binomial
- Poisson
- Uniform
- Gaussian (normal)

- Using Distributions to Answer Questions

- Distribution Parameters

- Estimating Parameters

- Confidence Intervals and P-values

- Why Gaussian?

- Regression

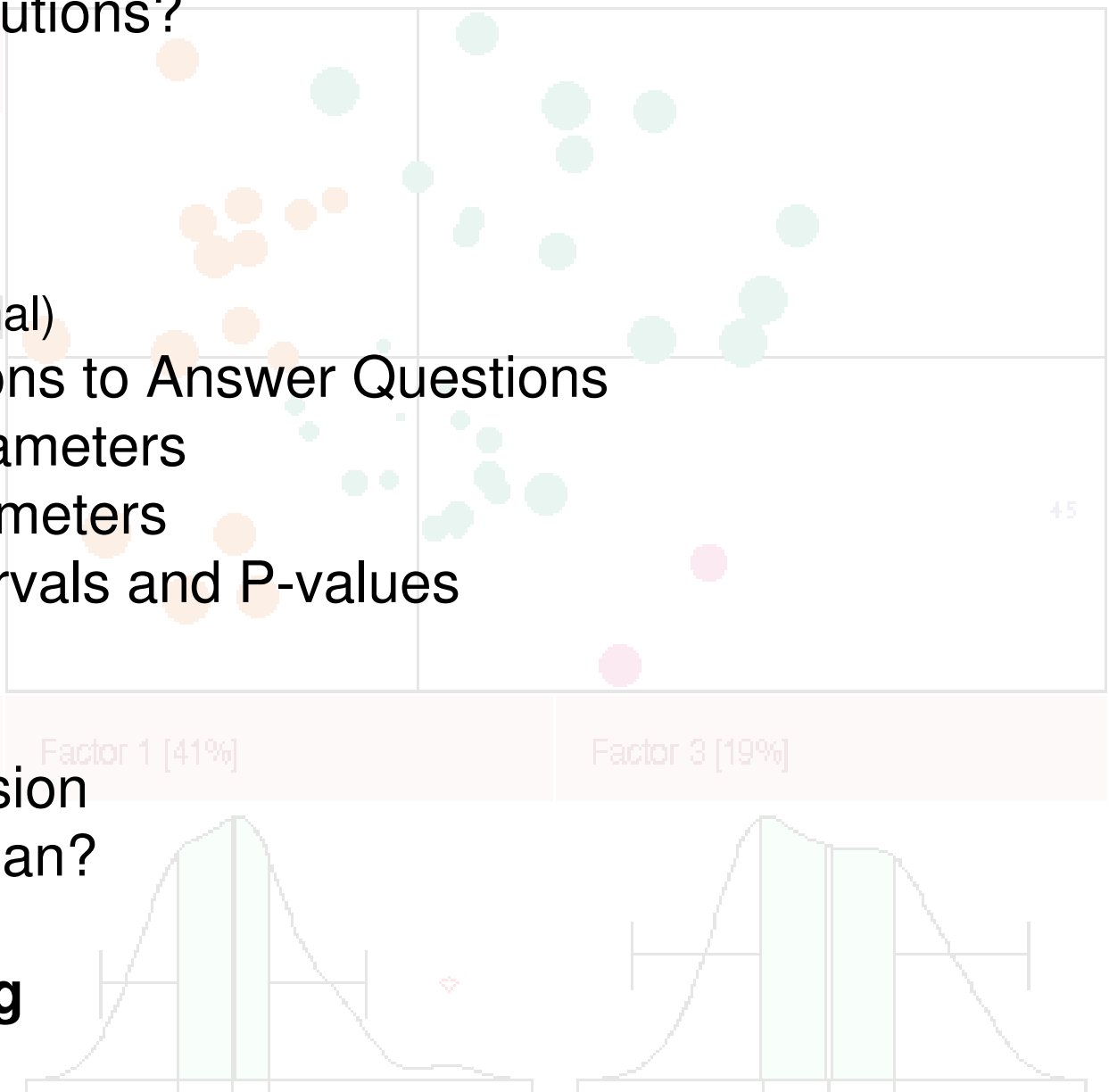
- Logistic Regression

- Why Not Gaussian?

- Bootstrapping

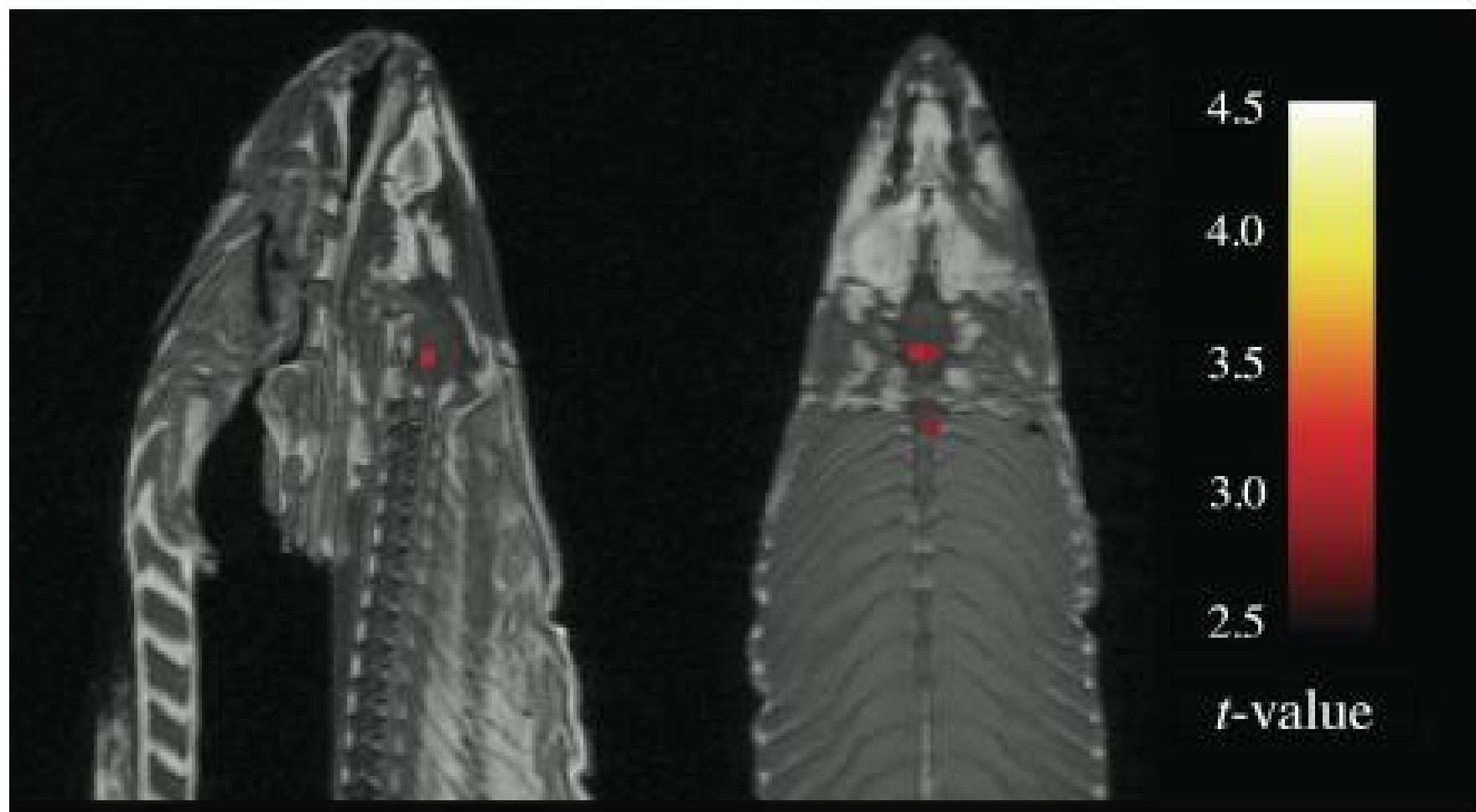
- **Multiple Testing**

- Useful Tools





Creative Commons licensed, from CharlotteKinzie on flickr 55/61



| Probability of False Detection (per test) | Number of Tests | Probability of False Detection (in at least 1 test) |
|---|-----------------|---|
| 0.01 | 1 | 0.01 |
| 0.01 | 10 | 0.10 |
| 0.01 | 100 | 0.63 |
| 0.05 | 1 | 0.05 |
| 0.05 | 10 | 0.40 |
| 0.05 | 100 | 0.99 |

- What are Distributions?

- Models

- Binomial
- Poisson
- Uniform
- Gaussian (normal)

- Using Distributions to Answer Questions

- Distribution Parameters

- Estimating Parameters

- Confidence Intervals and P-values

- Why Gaussian?

- Regression

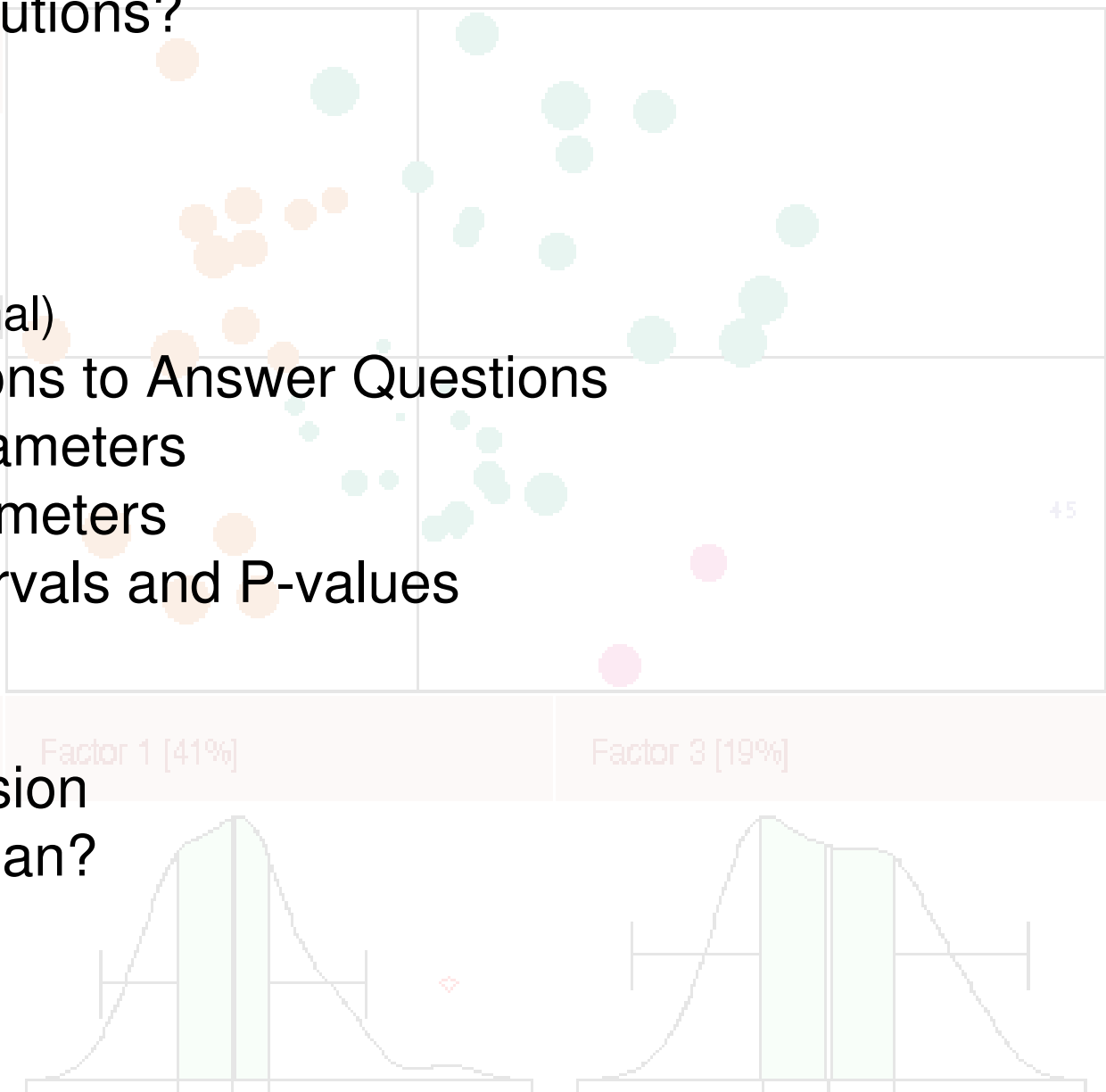
- Logistic Regression

- Why Not Gaussian?

- Bootstrapping

- Multiple Testing

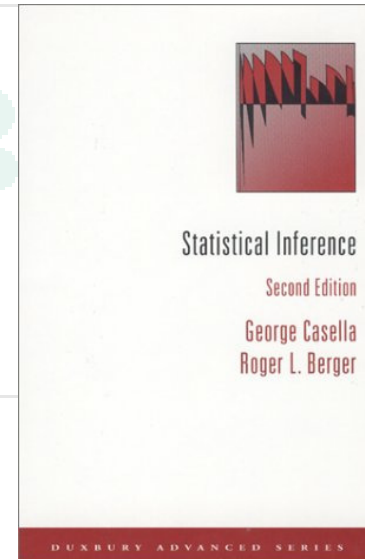
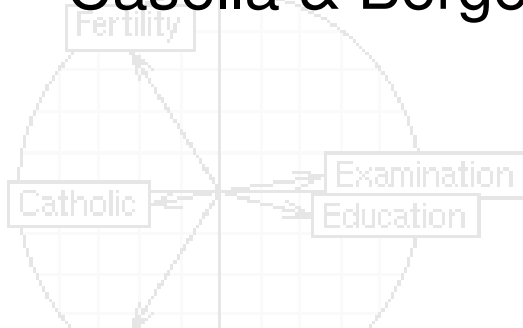
- **Useful Tools**



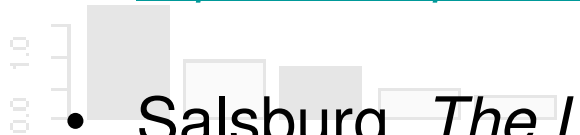
- R: <http://cran.r-project.org/>

PCA 5 vars
princomp(x = data, cor = cor)

- Casella & Berger, *Statistical Inference*



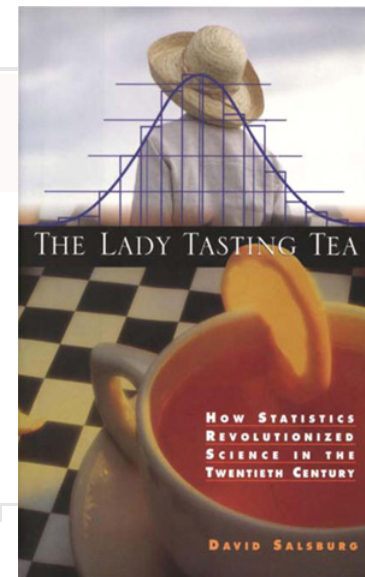
- Wikipedia:
http://en.wikipedia.org/wiki/List_of_probability_distributions



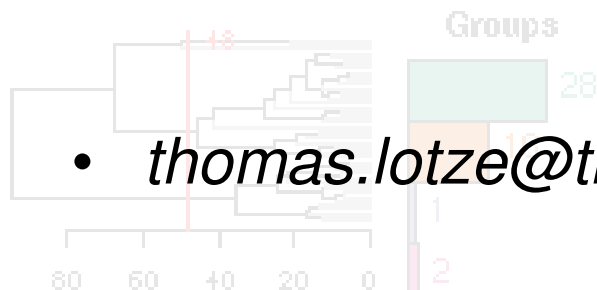
- Salsburg, *The Lady Tasting Tea*

Clustering 4 groups

Factor 1 [41%]



- thomas.lotze@thomaslotze.com



1. Everything is a Distribution
2. Many Kinds of “Random” (many Distributions)
3. Estimated Parameters are Random
They have Distributions!
4. Statistical Decisions come from Distribution Estimates
5. Be Skeptical of Normality
Mean and Variance are not sufficient!
6. Be Skeptical of Multiple Testing

- Practical Take-home

- Normality test
- T-test
- Wilcoxon rank-sum test

- Other Distributions

- Student's T
- F
- Lognormal
- Geometric
- Levy
- Weibull
- Benford

- Goodness of fit

- Chi-squared test
- Q-Q plots

- Distribution Connections

- Multivariate Distributions

- Bias/Variance Tradeoff

- Nonparametric Distributions

- Model Comparison:
Parameters and Fit (AIC, BIC)

- Bayesian Statistics

- Bayes' Law

- Cognitive Biases

- Time series

- Counterintuitive Probability

- Monty Haul

- Two Aces

- Poisson Waiting Times

- Markov Chain Monte Carlo

- Extreme Value Statistics