

Basic Methods

# Module 5 – Review of Basic Data Analytic Methods Using R- Part II

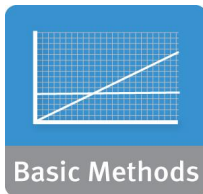
EMC<sup>2</sup> PROVEN PROFESSIONAL



Introduction



Analytics Lifecycle



Basic Methods



Adv. Methods



Tools



Lab

# Module 5: Review of Basic Data Analytic Methods Using R

## Lesson1: Statistics for Model Building and Evaluation

During this lesson the following topics are covered:

- Statistics in the Analytic Lifecycle
- Hypothesis Testing
- Difference of means
- Significance, Power, Effect Size
- ANOVA
- Confidence Intervals

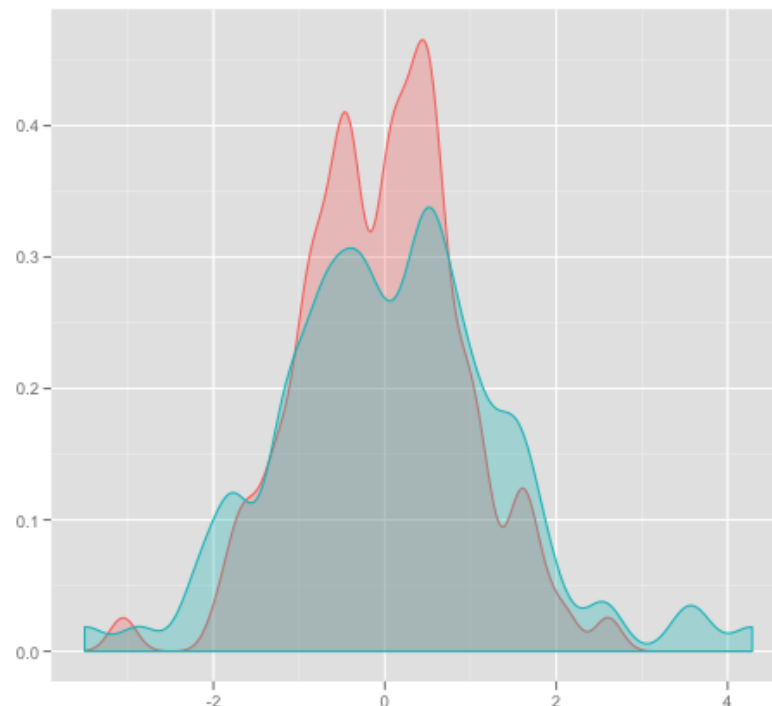


# Statistics in the Analytic Lifecycle

- Model Building and Planning
  - ▶ Can I predict the outcome with the inputs that I have?
  - ▶ Which inputs?
- Model Evaluation
  - ▶ Is the model accurate?
  - ▶ Does it perform better than "the obvious guess"?
  - ▶ Does it perform better than another candidate model?
- Model Deployment
  - ▶ Do my predictions make a difference?
    - ▶▶ Are we preventing customer churn?
    - ▶▶ Have we raised profits?

# Evaluating a Model: Hypothesis Testing

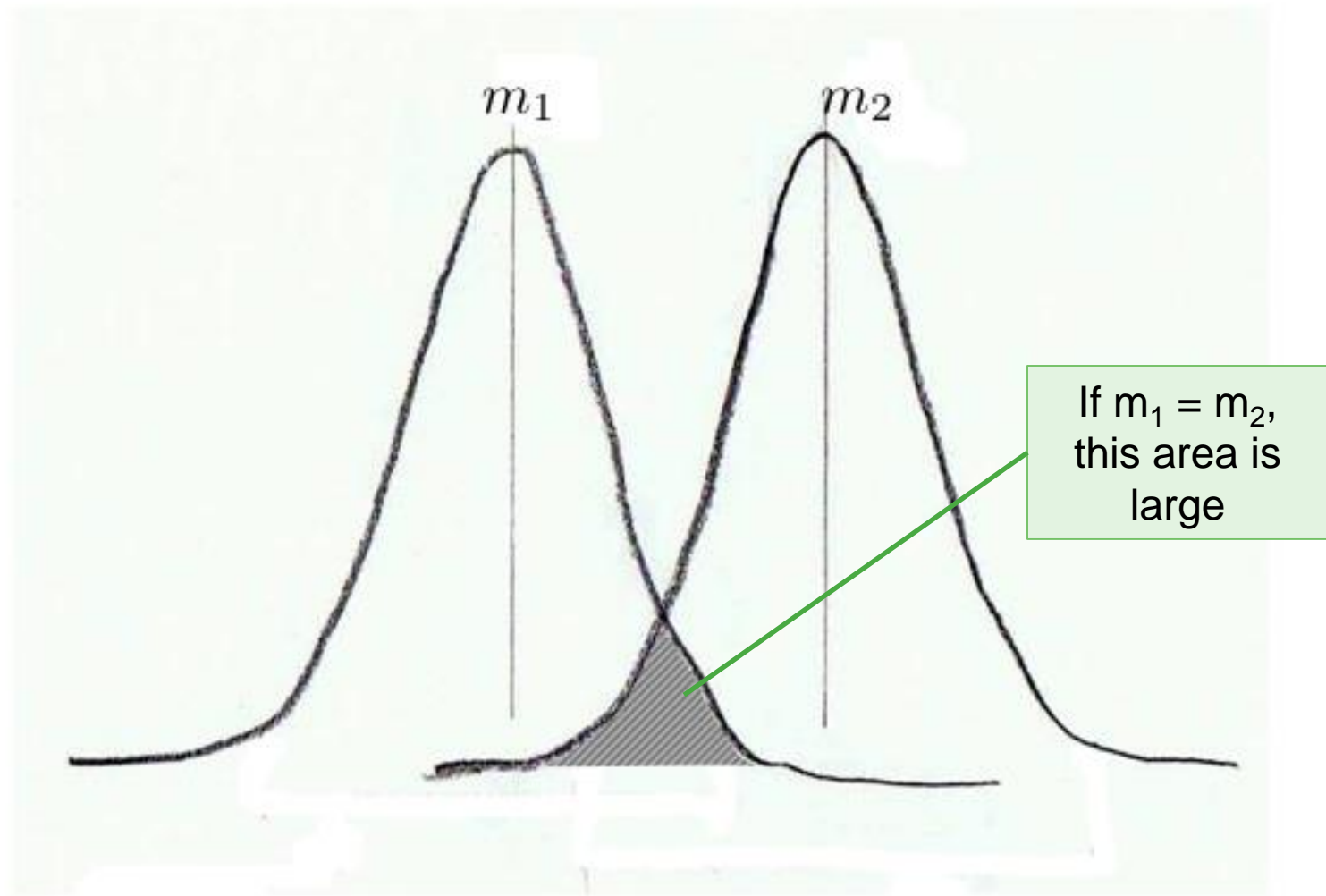
- Fundamental question: "Is there a difference?"
  - ▶ Specifically: "Would I see this value if there is no difference?"
- The baseline scenario: "There is no difference."
  - ▶ Statisticians call this the **Null Hypothesis**
  - ▶ "There is a difference." – **The Alternative Hypothesis**



# Null and Alternative Hypotheses: Examples

Null Hypothesis	Alternative Hypothesis
<p>The best estimate of the outcome is the average observed value:</p> <ul style="list-style-type: none"><li>• The mean is the "Null Model"</li></ul>	<p>The model predicts better than the null model:</p> <ul style="list-style-type: none"><li>• The average prediction error from the model is smaller than that of the null model</li></ul>
<p>This variable does not affect the outcome:</p> <ul style="list-style-type: none"><li>• The coefficient value is zero</li></ul>	<p>The variable does affect outcome:</p> <ul style="list-style-type: none"><li>• Coefficient value is non-zero</li></ul>
<p>The model predictions do not improve revenue:</p> <ul style="list-style-type: none"><li>• Revenue is the same with or without intervention</li></ul>	<p>Interventions based on model predictions improve revenue:</p> <ul style="list-style-type: none"><li>• A/B Testing, ANOVA</li></ul>

# Intuition: Difference of Means



# In Practice: t-test

t-statistic: 
$$t = \frac{\bar{X}_1 - \bar{X}_2}{\sqrt{\frac{s_1^2}{N_1} + \frac{s_2^2}{N_2}}}$$

(this is the t-statistic for the Welch t-test)

```
> x = rnorm(10) # distribution centered at 0
> y = rnorm(10,2) # distribution centered at 2
> t.test(x,y)
```

Welch Two Sample t-test

data: x and y

t = -7.2643, df = 15.05, **p-value = 2.713e-06**

alternative hypothesis: true difference in means is not equal to 0

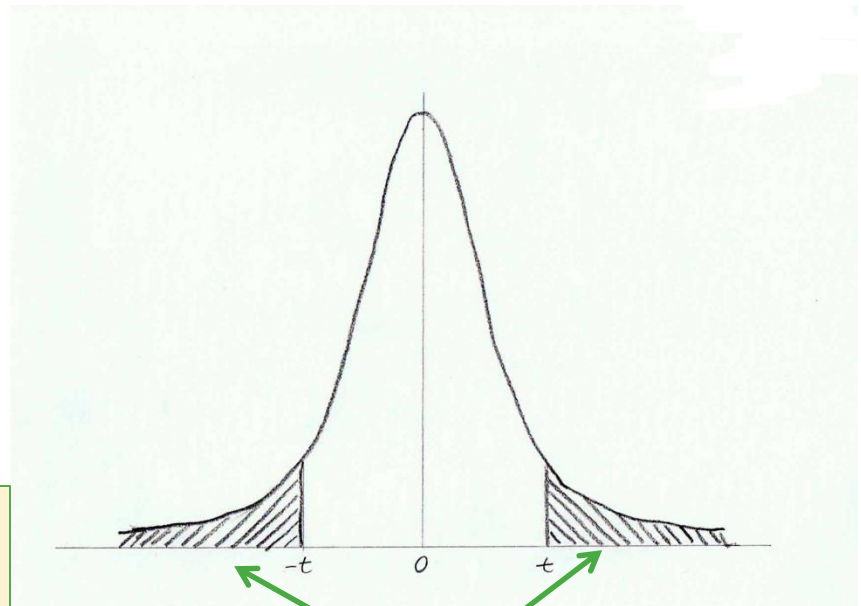
95 percent confidence interval:

-2.364243 -1.291811

sample estimates:

**mean of x mean of y**

**0.5449713 2.3729984**



**p-value:** area under the tails of the appropriate student's distribution

if p-value is small (say < 0.05), then reject the null hypothesis and assume that  $m_1 \neq m_2$

$m_1$  and  $m_2$  are "significantly different"



# In Practice: Wilcoxon Rank Sum test

- t-test assumes that the populations are normally distributed
  - ▶ Sometimes this is close to true, sometimes not
- Wilcoxon Rank Sum test
  - ▶ Makes no assumption about the distributions of the populations
  - ▶ More robust test for difference of means
  - ▶ if p-value is small: reject the null hypothesis (equal means)

```
> mean(x)
[1] 0.5449713
> mean(y)
[1] 2.372998
> wilcox.test(x, y)
```

Wilcoxon rank sum test

```
data: x and y
W = 2, p-value = 4.33e-05
alternative hypothesis: true location shift is not equal to 0
```



# Hypothesis Testing: Summary

- Calculate the **test statistic**
  - ▶ Different hypothesis tests are appropriate, in different situations
- Calculate the **p-value** on the test statistic
- If p-value is "small" then reject the null hypothesis
  - ▶ "small" is often  $p < 0.05$  by convention (95% confidence)
  - ▶ Many data scientists prefer a smaller threshold.

$$t = \frac{\bar{X}_1 - \bar{X}_2}{\sqrt{\frac{s_1^2}{n_1} + \frac{s_2^2}{n_2}}}$$

# Generating a Hypothesis: Type I and Type II Error

If $H_0$ is X, and we ...:	Null hypothesis( $H_0$ ) is <b>true</b>	Null hypothesis( $H_0$ ) is <b>false</b>
Fail to accept the Null Hypothesis → we claim something happened	<b>Type I error</b> <b>False positive</b> $\alpha$	<b>Correct Outcome</b> <b>True positive</b> <b>We reject the Null hypothesis</b>
Fail to reject the null hypothesis → we claim nothing happened.	<b>Correct outcome</b> <b>True negative</b> <b>Accept the NULL hypothesis</b>	<b>Type II error</b> <b>False negative</b> $\beta$

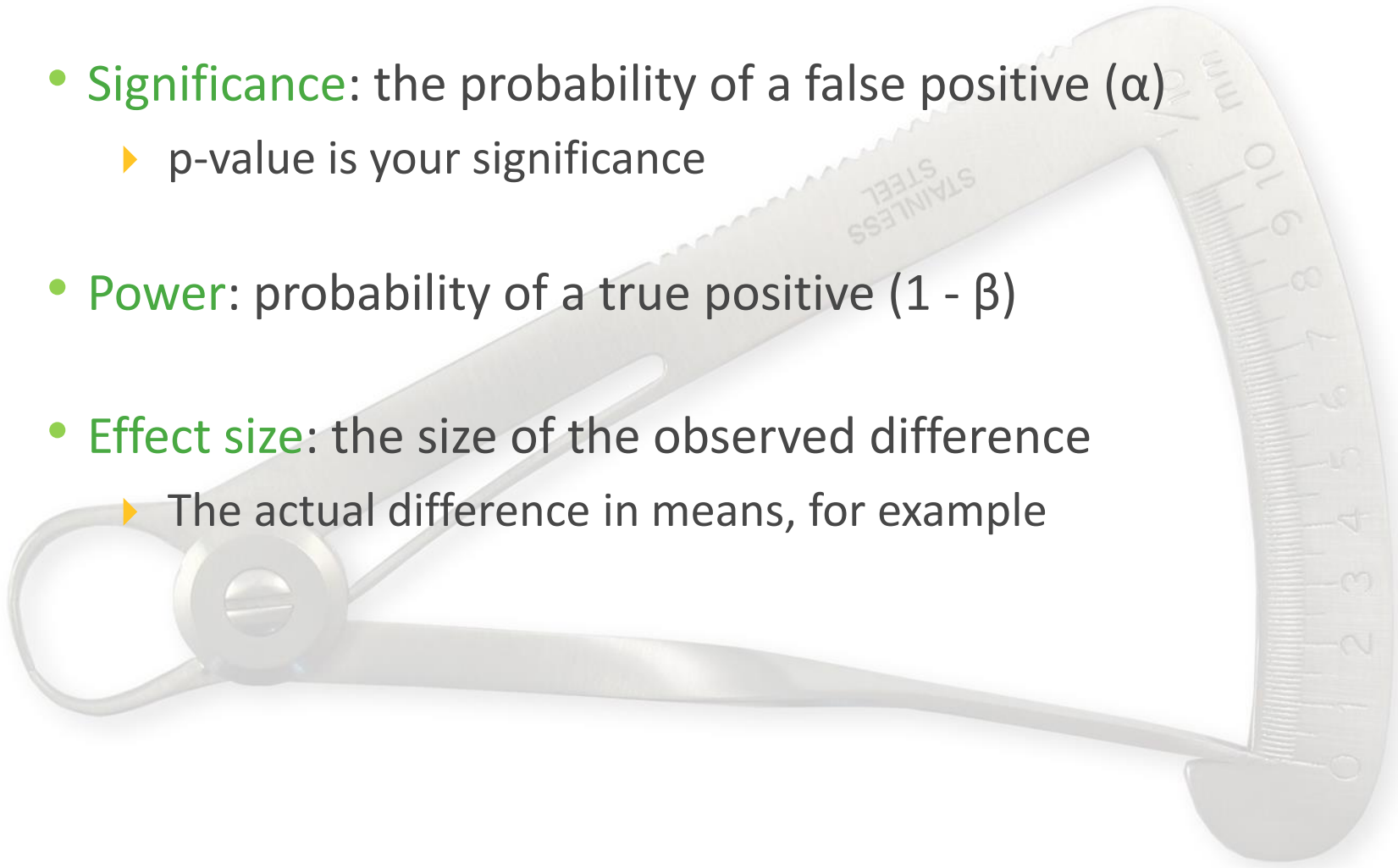
**Example: Ham or Spam?  $H_0$ : it's ham  $H_A$ : it's spam**

If it's ↓, and we say it's →	SPAM	HAM
HAM	<b>Type I – false positive</b>	OK – true positive
SPAM	OK – true negative	<b>Type II – false negative</b>

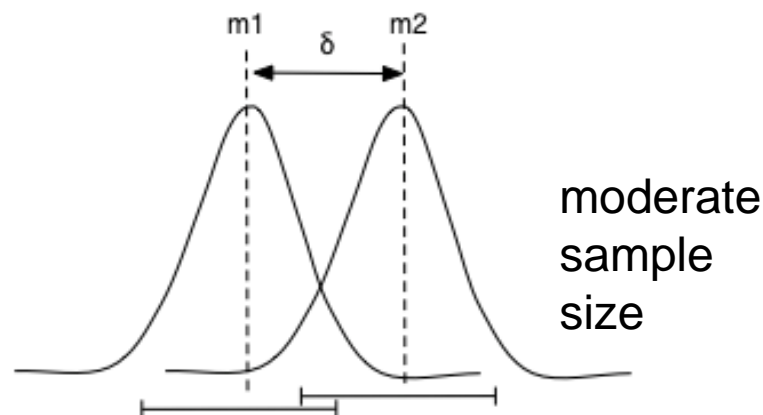
- **Goal: Identify spam**
- **Which error is worse?**

# Significance, Power and Effect Size

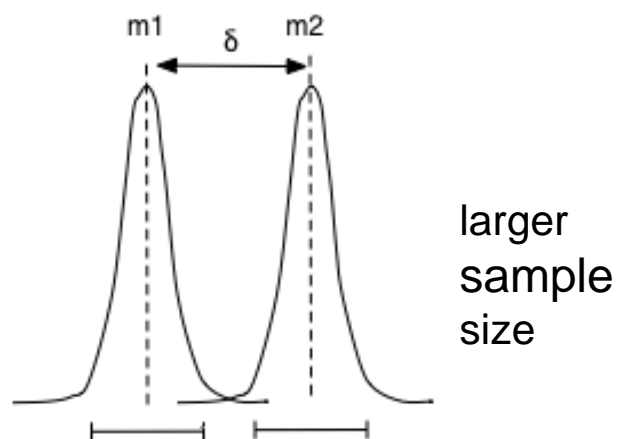
- **Significance**: the probability of a false positive ( $\alpha$ )
  - ▶ p-value is your significance
- **Power**: probability of a true positive ( $1 - \beta$ )
- **Effect size**: the size of the observed difference
  - ▶ The actual difference in means, for example



# Always Keep Effect Size in Mind!



Both power and significance increase with larger sample sizes.



So you can observe an effect size that is *statistically* significant, but *practically* insignificant!

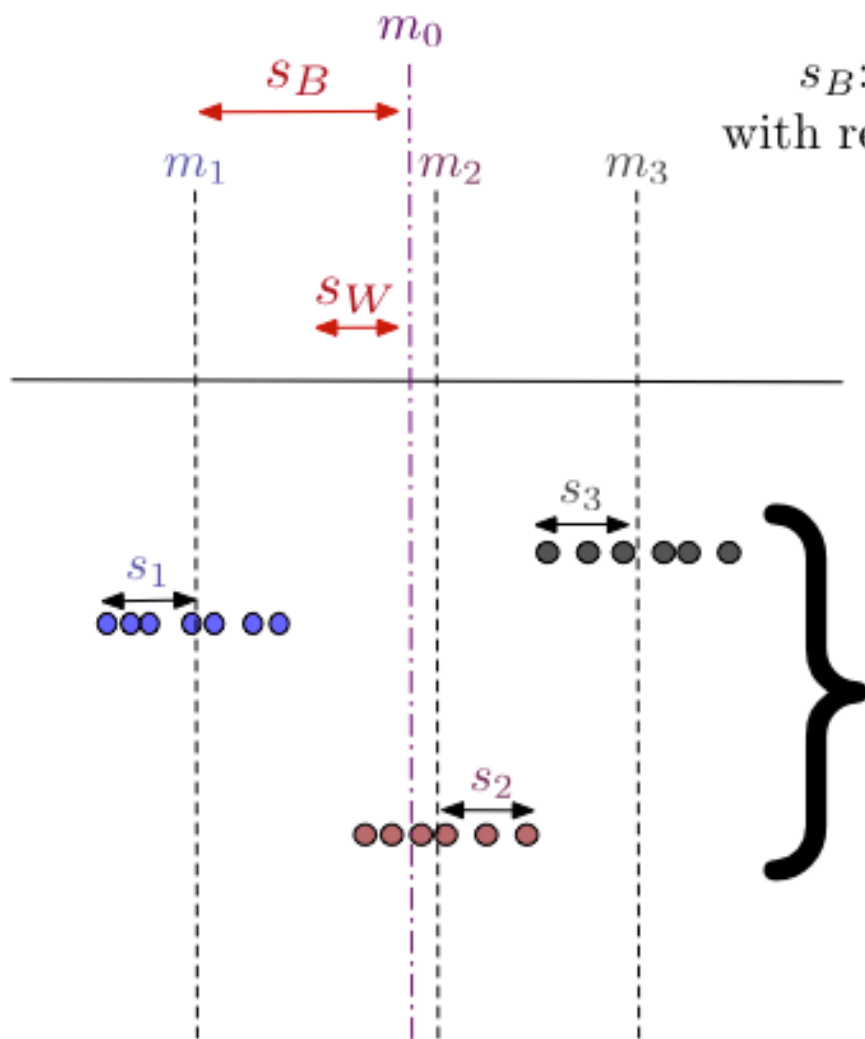
# Hypothesis Testing: ANOVA

ANOVA is a generalization of the difference of means

- One-way ANOVA
  - ▶ k populations ("treatment groups")
  - ▶  $n_i$  samples each – total N subjects
  - ▶ Null hypothesis: ALL the population means are equal

Population	$n_i$ : # offers made	$m_i$ : avg purchase size
Offer 1	100	\$55
Offer 2	102	\$40
No intervention	99	\$25

# ANOVA: Understanding the F statistic



$s_B$ : how the population means vary with respect to the total mean  $m_0$

$$s_B^2 = \frac{1}{k-1} \sum_i n_i \cdot (m_i - m_0)^2$$

$$s_W^2 = \frac{1}{N-k} \sum_i^k \sum_j^{n_i} x_{ij} - m_i^2$$

$s_W$ : the "average" of the  $s_i$

$$\text{Test statistic: } F = s_B^2 / s_W^2$$

# R Example: ANOVA

3 different offers, and their  
outcomes

Use `lm()` to do the ANOVA

offer1-nooffer  
offer2-nooffer

**F-statistic: reject the null hypothesis**

Tukey's test: all pair-wise tests for difference of means

**95% confidence intervals** for difference between means

**.No appreciable difference between offer1 and offer2**

```
>offers = sample(c("nooffer", "offer1", "offer2"),
  size=500, replace=T)
>purchasesize = ifelse(offers=="nooffer", rlnorm(500,
meanlog=log(25)), ifelse(offers=="offer1", rlnorm(500,
meanlog=log(50)), rlnorm(500, meanlog=log(55))))
>offertest = data.frame(offer=as.factor(offers),
  purchase_amt=purchasesize)
> model = lm(log10(purchase_amt) ~ as.factor(offers),
  data=offertest)
>summary(model)
```

Residuals:

Min	1Q	Median	3Q	Max
-1.1940	-0.2837	0.0135	0.2863	1.3374

Coefficients:

	Estimate	Std. Error	t value	Pr(> t )
(Intercept)	1.49092	0.03240	46.011	<b>&lt; 2e-16 ***</b>
as.factor(offers)offer1	0.20424	0.04706	4.340	<b>1.73e-05 ***</b>
as.factor(offers)offer2	0.22371	0.04596	4.867	<b>1.52e-06 ***</b>

---

Signif. codes: 0 '\*\*\*' 0.001 '\*\*' 0.01 '\*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 0.4262 on 497 degrees of freedom

Multiple R-squared: 0.05479, Adjusted R-squared: 0.05098

**F-statistic: 14.4 on 2 and 497 DF, p-value: 8.304e-07**

```
> TukeyHSD(aov(model))
```

Tukey multiple comparisons of means

95% family-wise confidence level

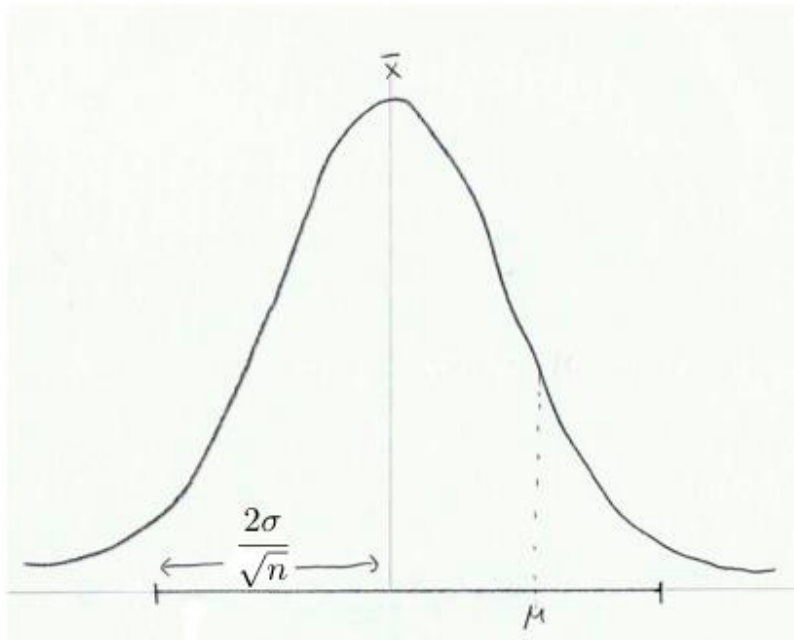
Fit: aov(formula = model)

\$offers

	diff	lwr	upr	p adj
offer1-nooffer	0.20424099	0.09361976	0.3148622	0.0000512
offer2-nooffer	0.22370761	0.11566775	0.3317475	0.0000045
offer2-offer1	0.01946663	-0.09146092	0.1303942	0.9104871



# Confidence Intervals



## Example:

- Gaussian data  $N(\mu, \sigma)$
- $\bar{x}$  is the estimate of  $\mu$ 
  - based on  $n$  samples

$\mu$  falls in the interval

$$\bar{x} \pm 2\sigma/\sqrt{n}$$

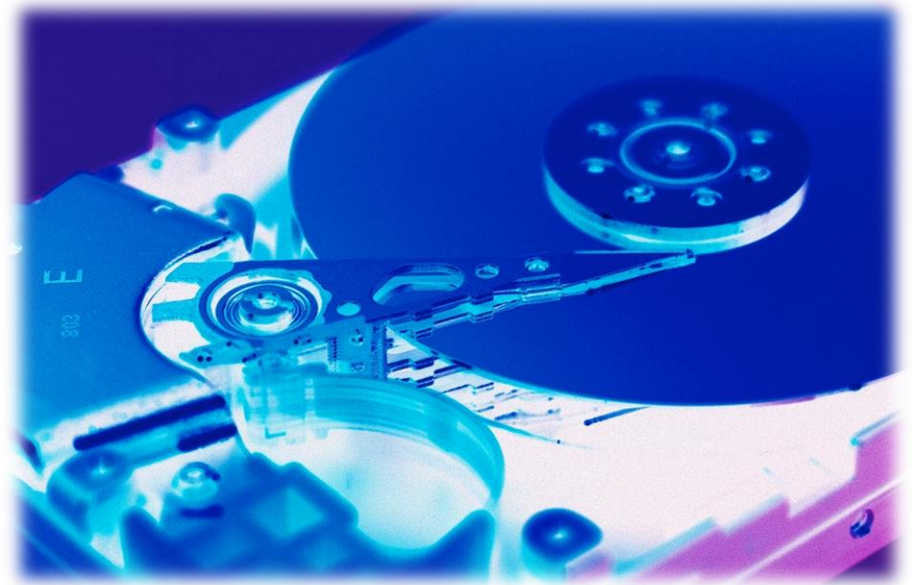
with approx. 95% probability  
("95% confidence")

If  $\bar{x}$  is your estimate of some unknown value  $\mu$ ,  
the  $P\%$  confidence interval  
is the interval around  $\bar{x}$  that  $\mu$  will fall in, with  
probability  $P$ .

# Example

The defect rate of a disk drive manufacturing process is within 0.9% - 1.7%, with 98% confidence. We inspect a sample of 1000 drives from one of our plants.

- We observe 13 defects in our sample.
  - Should we inspect the plant for problems?
- What if we observe 25 defects in the sample?



# Check Your Knowledge

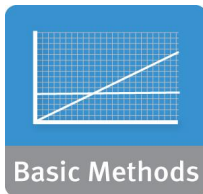
- Refer back to the Anova example on an earlier slide. What do you think? Does the difference between *offer1* and *offer2* make a practical difference? Should we go ahead and implement one of them?
- If yes, and the costs were US \$25 for each *offer1* and US \$10 for *offer2*, would you still make the same decision?
- In our manufacturing plant example, assuming you would check the plant for problems in the manufacturing process, how might you justify this decision financially?



Introduction



Analytics Lifecycle



Basic Methods



Adv. Methods



Tools



Lab

# Module 5: Review of Basic Data Analytic Methods Using R

## Lesson 1: Summary

During this lesson the following topics were covered:

- The role of Statistics in the Analytic Lifecycle
- Developing a model and generating the null and the alternative hypothesis
- Difference between means
- Difference between significance, power and effect size, and how they relate to Type I and Type II errors
- Applying ANOVA and determining whether the results are significant
- Defining confidence intervals and applying them

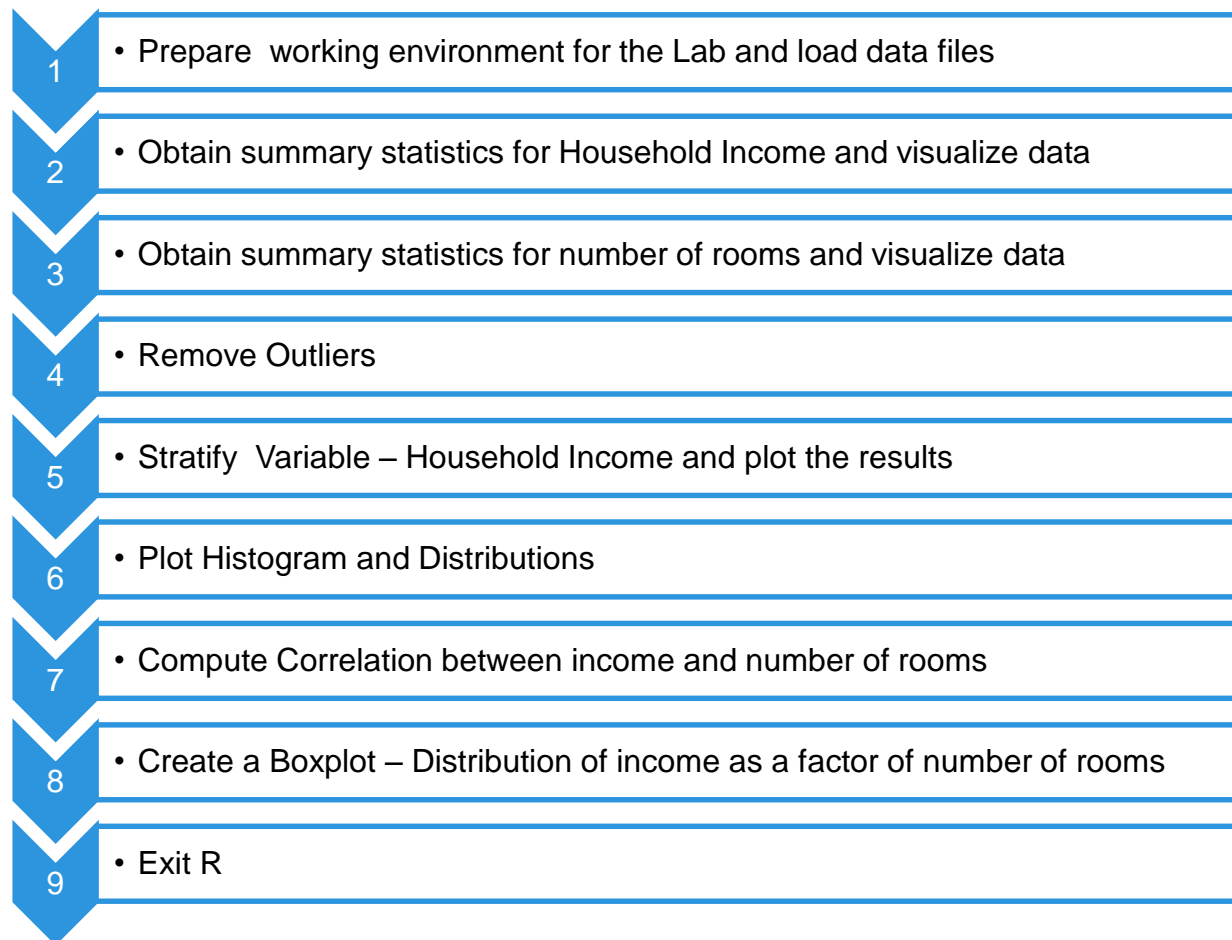
# Lab Exercise 3: Basic Statistics, Visualization and Hypothesis Tests



This lab is designed to investigate and practice using R to perform basic statistics and visualization on data and to perform hypothesis testing.

- After completing the tasks in this lab you should be able to:
  - Perform basic data analysis
  - Visualize data with R
  - Create and test a hypothesis

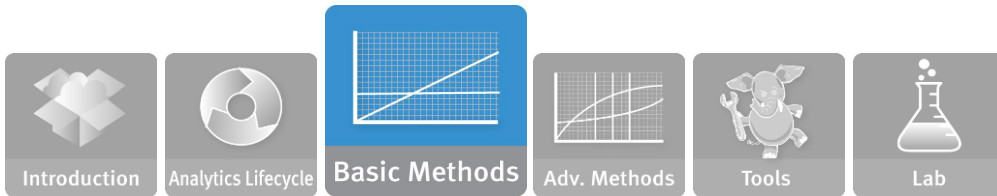
# Lab Exercise 3: Basic Statistics, Visualization and Hypothesis Tests– Part1 - Workflow



# Lab Exercise 3: Basic Statistics, Visualization and Hypothesis Tests - Part 2 - Workflow

- 1 • Define problem – Analysis of Variance (ANOVA)
- 2 • Generate the Data
- 3 • Examine the Data
- 4 • Plot and determine how purchase size varies within the three groups
- 5 • Use `lm()` to do the ANOVA
- 6 • Use Tukey's test to check all the differences of means
- 7 • Use the lattice package for density plot
- 8 • Plot the Logarithms of the Data
- 9 • Use `ggplot()` package
- 10 • Generate the example data to perform a Hypothesis Test with manual calculations
- 11 • Create a function to calculate the pooled variance, which is used in the Student's t statistic
- 12 • Examine the Data
- 13 • Calculate the t statistic for Student's t-test
- 14 • Calculate the degrees of freedom
- 15 • Compute the area under the curve
- 16 • Perform Student's t-test directly and compare the results





## Module 5: Summary

Key points covered in this module:

- How to use basic analytics methods such as distributions, statistical tests and summary operations to investigate a data set
- How to use R to apply visualization patterns to better understand the data, help develop a model and derive hypotheses, and determine if our actions had a practical affect.