







- If n = 1, the distribution of the average is just the distribution itself, since we have only the single data point
- If *n* is larger than one, the distribution of the mean must be narrower than the distribution of the population
 - i.e. the variance and standard deviation must be smaller
- In fact, the mean & variance of the mean of n samples is

$$\mu_{\bar{x}} = \mu \qquad \qquad \sigma_{\bar{x}} = \frac{\sigma}{\sqrt{n}}$$











- The system with the better mean can be said to be better on average with a probability better than the Confidence Level
- If the CIs do overlap
 - · Can't say that the two systems are different with this technique
 - Either:
 - 1. The two systems are equivalent
 - 2. We haven't sampled enough to discriminate between the two











So what should we do?

First test for normality

- Many such tests
- Recommended
 - Normal Probability Plot (QQ plot: sorted data vs Normal quantiles)
 - Lilliefors test (variant of the KS test)



There are 3 basic remedial measures:

- 1. Transforming data to make them normally distributed
 - also called data re-expression
 - traditional approach (required before the advent of fast computers)
- 2. Resampling techniques
- 3. Non-parametric statistics



Non-Parametric Statistics

- Basic Idea
 - Sort the data and then rank them
 - Use Ranks instead of actual values to perform statstics
- Also known as
 - order statistics,
 - ordinal statistics
- rank statistics
- Measures how interspersed the samples are from the 2 treatments
 If the result is "alternating" it is assumed that there is no difference
- Can't be affected by outliers (extrememly large or small values)
 - Just the highest or lowest rank



rank rank 0.99 0.99 1 1 0.91 2 0.91 2.5 t1 Average tied ranks 0.91 3 0.91 2.5 t1 together 0.64 0.64 4 5.5 t2 2.5 0.64 5 0.64 5.5 t2 5.5 0.64 0.64 5.5 t2 6 Sort the combined data 8.5 5.5 t2 0.64 7 0.64 12 0.27 0.27 8.5 t3 8 17 8.5 t3 0.27 9 0.27 0.16 10 0.16 12 t4 Give each data element Replace tied ranks 12 t4 0.16 11 Α 0.16 its corresponding rank with average tied ranks 12 t4 0.16 12 0.16 12 t4 0.16 13 0.16 0.16 14 0.16 12 t4 0.08 15 0.08 15 16 0.03 0.03 17 t5 0.03 17 0.03 17 t5 t5 0.03 18 0.03 17 0.02 0.02 19 19 Ranked Example Ranked Example 0.01 20 0.01 20









A Confidence Interval Around the Median: Thompson-Savur

- Find the *b* the binomial value that has a cumulative upper tail probability of $\alpha/2$
 - *b* will have a value near n/2

The lower percentile
$$l = -\frac{b}{b}$$

$$n-1$$

The upper percentile u = 1 - l

In Excel:

- To calculate *b* use CRITBINOM $(n, 1/2, \alpha/2)$
- to compute the $value_u$ use the function PERCENTILE (dataArray, u)

to compute the *value*_l use the function PERCENTILE (dataArray, *l*)

- Confidence Interval is [*value*_l,*value*_u]
 - i.e. $value_l \leq median \leq value_u$
 - With a confidence level of $1-\alpha$



ANOVA: Analysis of Variance

Part 1: Multi-Level Analysis Basic Concept



More Than 2 Treatments

- Preceding stats to be used for simple experiment designs
- More sophisticated stats needs to be done if:
 - Comparing multiple systems instead of just 2 treatments
 - E.g. comparing the effect on a Genetic Algorithm of using no mutation, low, medium and high levels of mutation
 - We say there are 4 *levels* of the mutation variable

• Need $\binom{4}{2} = 6$ possible comparisons to test all pairs of treatments

• Called a 'multi-level' analysis





Comparing Variances

- Up to now we have been comparing means
 - Student's T test
 - Difference between averages (after normalization) • see if it equals 0
- From here on we will be comparing variances
 - Won't take the difference between variances
 - Difference between variances not a nice distribution
 - Rather will take the ratio of variances
 - see if it equals 1
 - distribution known: F distribution













Polynomial Regression E.g.

R squared = 70.2% R squared (adjusted) = 70.2% s = 0.1465 with 1000 - 4 = 996 degrees of freedom

Source Regression	Sum of Squares	df 3	Mean 16 82	Square 28	F-ratio 784	p-value < 0.0001
Residual	21.3807	996	0.021	467	704	20.0001
Variable	Coefficient	s.e. of (Coeff	t-ratio	p-value	
Constant	0.510755	0.0190	00	26.9	≤ 0.000)]
Χ	-2.17801	0.1636		-13.3	≤0.000	01
X^2	8.45358	0.3813		22.2	≤0.000	01
X^3	-6.28741	0.2515		-25.0	≤ 0.000	01

Regression model is statistically significant F-ratio = 784 >> 1

•	ANOVA: Back to Discrete Levels							
	no xover	<i>xover</i> = 1pt	<i>xover</i> = 2pt	<i>xover</i> = 3pt	xover = 4pt			
	4.3	8.8	5.0	6.3	5.4			
	3.7	7.7	5.3	6.6	5.9			
	4.7	8.3	5.1	7.2	5.4			
	3.7	8.1	5.2	7.4	5.4			
Fitness	4.2	8.1	5.5	7.4	6.2			
Values	3.6	8.0	4.9	7.3	6.7			
avg fitness	4.02	8.13	5.09	7.02	5.76			
std dev	0.451	0.313	0.424	0.478	0.471			
T test T test all pairwise T test								
Question: Do crossover settings make a difference at all?								

~~						
	٨	NOV		a tata	т	1
1	P = P		A: Dis	crete	fitness	xover
					4.3	no
\$ 200	no xover	rover = 1 nt	rover = 2nt		3.7	no
e l	1 2	00	5 0		4.7	no
	4.5		5.0		3.7	no
ť	3.7	7.7	5.3		4.2	no
ň	4.7	8.3	5.1	\rightarrow	3.6	no
е	3.7	8.1	5.2		8.8	1pt
S	4.2	8.1	5.5		7.7	1pt
5	3.6	8.0	4.9		8.3	1pt
					8.1	1pt
Most st	atistic pack		8.1	1pt		
formatt	ed as the ch	[8.0	1pt		
		5.0	2pt			
each co	olumn is a vo	5.3	2pt			
(fitness is a response variable, xover is a factor)					5.1	2pt
each row is a treatment					5.2	2pt
(i e the	settings and	5.5	2pt			
(the the	(i.e the settings and results of a single run)					2pt





ANOVA: Analysis of Variance



Part 2: Multi-Level Analysis Pairwise Comparisons Post-Hoc Analysis



- What if we want to know more detailed information?
 - Which of the means is the significantly different one?
 - Are there more than one significantly different mean?
 - If so, what are the pair-wise differences and are they statistically significant?



Pairwise Comparisons

- between Factor-Level Means
- This is determined by a series of pair-wise T tests
- However, commonly uses pooled information from the model for the variance to provide greater accuracy

Called standard error

original T test comparison

comparing level *i* with level *j* across the ANOVA model

$$t value = \frac{\overline{X}_1 - \overline{X}_2}{\sqrt{\frac{s_{X_1}^2}{n_1} + \frac{s_{X_2}^2}{n_2}}} \longrightarrow t value = \frac{\overline{X}_i - \overline{X}_j}{\sqrt{\frac{2 \cdot MSE}{n_1}}}$$
when $n_i = n_i = n$

Assumption: variances for each factor level is the same (σ^2) which is best estimated by the *MSE*

Multiple Levels: Post-hoc Analysis

- For 4 levels of mutation there are 6 comparisons possible
 - *Each one* of the comparison holds at a 95% C.L. independent of the other comparisons
 - If *all* comparisons are to hold at once the odds are 0.95 x 0.95 x 0.95 x ... x 0.95 = (0.95)⁶ = 0.735
 - So in practice we only have 73.5% C.L
 - Wrong 1/4 of the time
- For 7 levels of mutation there are 21 comparisons possible
 - C.L. = $(0.95)^{21} = 0.341$
 - Chances are better than half that at least one of the decisions may be wrong!



The Bonferroni Correction

• To correct, choose a smaller α

- Where *m* is the number of comparisons
- So for 95% CL use $\alpha = 0.025/6 = 0.004167$
- For a Z test the critical value changes from 1.96 to 2.64
- You should apply the Bonferroni (etc.) correction:
 - To *t* tests (*t* tests and ranked *t* tests)
 - To Confidence Intervals and Error Bounds
 - Whenever you mean "all the significant results we found hold at once"

@	ite	Pairwise Comparisons							
0	🍹 b	etwe	en Fac	tor-Le	evel	Means			
Regular Pair-wise T test (with Bonf, Correction)									
1.7	0	Diff	std. err.	t-value	df	p-value			
	n - 1	-4.04	0.15	-27.5	18	3.6E-15			
	n - 3	-3.18	0.16	-20.5	18	6.3E-13			
	2 - 1	-3.04	0.16	-20.2	18	8.4E-13			
	3 - 2	2.16	0.17	13.7	18	5.5E-10			
	4 - 1	-2.09	0.17	-12.7	18	2.0E-09			
	n - 4	-1.95	0.17	-11.4	18	1.1E-08			
	4 - 3	-1.22	0.18	-7.1	18	1.3E-05			
	n - 2	-1.00	0.16	-6.3	18	5.8E-05			
	4 - 2	0.95	0.16	5.6	18	2.6E-04			
	3 - 1	-0.86	0.15	-5.6	18	2.6E-04			

Pairwise Comparisons								
	b	etwe	en Fac	tor-Le	evel	Means		
	rection)							
		Diff	std. err.	t-value	df	p-value		
	n - 1	-4.04	0.16	-25.2	95	7.7E-43		
	n - 3	-3.18	0.16	-19.8	95	1.7E-34		
	2 - 1	-3.04	0.16	-19.0	95	4.8E-33		
	3 - 2	2.16	0.16	13.6	95	6.0E-23		
	4 - 1	-2.09	0.16	-13.0	95	7.5E-22		
	n - 4	-1.95	0.16	-12.2	95	4.4E-20		
	4 - 3	-1.22	0.16	-7.6	95	1.8E-10		
	n - 2	-1.00	0.16	-6.2	95	1.2E-07		
	4 - 2	0.95	0.16	5.9	95	4.8E-07		
	3 - 1	-0.86	0.16	-5.4	95	5.1E-06		















- Many others
 - Scheffé
 - used when comparing pairs, and triples and quadruples etc., not just pairs
 - many many others
 - Duncan's multiple range test
 - The Nemenyi test
 - The Bonferroni–Dunn test
 - Newman-Keuls post-hoc analysis

Important Topics Not Covered



No time

Important Topics Not Covered

- Multifactor ANOVA (MANOVA)
 - Main Effects vs Interaction Terms
 - F tests to determine which factors are statistically significant (validating the model
 - T tests to compare treatments
- Regression
 - Multivariate regression, Polynomial Regression
 - Confidence Intervals around model parameters
 - Statistical Testing for factor relevance
 - Correlation Coefficients: *r*, *r*², adjusted *r*²
- How to perform ANOVA as a multivariate regression
 - Indicator Variables



Important Topics Not Covered

- Testing for equality (homogeneity) of variance across different factor-levels / treatments
 - Levene's Test
- Correcting for inequality of variance
 - · Convert to multivariate regression using indicator variables
 - Perform Weighted Least Squares
- · How to perform ANOVA when using different test functions
 - Test functions as *blocking variables*
 - Non-parametric blocking
- What if one EC system has parameters the other EC system doesn't?
 - Nesting factor analysis



- Mathematical statistics with applications
 - Dennis D. Wackerly, William Mendenhall, Richard L. Scheaffer.
 - Boston : Duxbury Press, (6th Ed.)
 - · Introductory material probability distributions, simple sample statistics
 - · Easy to understand concrete proofs and examples good exercises
- Applied linear statistical models
 - Michael H. Kutner, Christopher J. Nachtsheim, John Neter, William Li
 - Boston: McGraw-Hill Irwin, 2005. (5th Ed.)
 - Advanced Regression techniques, ANOVA, and GLM
- Nonparametric statistical methods
 - Myles Hollander and Douglas A. Wolfe.
 - New York: Wiley, 1973
 - Classic nonparametric statistics textbook (very practical)



Online Resources

Websites

- Wikipedia (various pages)
 - http://en.wikipedia.com
- HyperStat Online
 - http://davidmlane.com/hyperstat
- Mathworld
 - <u>http://mathworld.wolfram.com/</u>