



Instructor

Mark Wineberg is an Associate Professor at the University of Guelph.



He has been actively researching the field of GEC since 1993 while he was still a graduate student. Over the years he has published on various topics including: the intersection of GA and GP, enhancing the GA for improved behavior in dynamic environments through specialized multiple populations, and exploring the concept of distances and diversity in GA populations.

Prof. Wineberg also teaches an undergraduate course on computer simulation and modeling of discrete stochastic systems with an emphasis on proper statistical analysis, as well as a graduate course on experimental design and analysis for computer science, which is an outgrowth of the statistical analysis tutorial given at GECCO.





	· N	Normal Distr	ibution	S	
Env 1	-75	9	75	Variation Expected:	
	·		•	$\sigma = 5$	
Env 2	-75	· ···································	75	Variation Expected:	
				$\sigma = 10$	
Env 3	-75	an se se se se se se se	75-	Variation Expected:	
•	· · · · · · · · · · · · · · · · · · ·			$\sigma = 50$	







- For most statistical analysis for EC the question is
 - Is one way better than another way?
 - Statistically this translates into a statement about the difference between means: "Is the difference between 'my mean' and 'the other mean' greater than zero?"
- We will approach this question in 2 steps:
 - What can we say about the true mean of a *single* distribution?
 Called *point estimation*
 - 2. How can we compare the true means of *two* or more distributions?

	à	Sampling Fro Normal Distr	om Two ibution) S
Env 1	-75	Q	75	Variation Expected:
	·	-10.7 9.7		$\sigma = 5$
Env 2	-75	A CALLER STOR	75	Variation Expected:
Env 3	-75.	-9.7 10.5	75	$\sigma = 10$ Variation Expected: $\sigma = 50$
True Av	g10	+10	Number of R	tuns = 100















Confidence Intervals

- Of course, we don't know the true mean, μ , or true standard deviation, σ
- We *do* know the mean of the samples, \overline{X} , the sample size, *n*, and the sample standard deviation, s_X
- If the source distribution is *normally distributed*, the shape as well as the size of the "finger" is known exactly!
 - We can determine the odds that the true mean lies within a specified range of \overline{X}

































• Confidence Intervals can be written in 3 equivalent ways





Example:

- An experimenter runs a <u>New Evolutionary</u> <u>Algorithm on a TSP</u>
- At the end of each run, the smallest length tour that had been found during the run was recorded
- NEA is run 50 times on the same TSP problem
- On average NEA found solutions with a tour length of 272
- The standard deviation of these tours is 87
- We want to compute a Confidence Interval using a 99% Confidence level



Basic Statistical Tests



Part 2 - Comparisons: Non-Overlapping Confidence Intervals and the Student's T Test



Using Confidence Intervals to Determine Whether My Way is Better

If we have two different EC systems how can we tell if one is better than the other?

Trivial method: Find confidence intervals around both means

- If the CIs don't overlap
 - Then it is a rare occurrence when the two systems do have identical means
 - 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







- The Student *t* Test is the basic test used in statistics
 - Idea: Gain sensitivity by looking at the difference between the means of the two systems















- Everything so far has depended on the assumption of normality which in turn depends on the Central Limit Theorem holding
 - But this is not always true
 - In in many areas of CS it rarely holds
- Problems occur when
 - ...you have a non-zero probability of obtaining infinity
 Mean and standard deviation are infinite!
 - ... the sample average depends highly on a few scores
 - When the mean of your distribution is not measuring what you want, consider using the median instead (rank-based statistics)
 - ... you don't know how fast your sample series converges to normal
 - if your sample average distribution converges very slowly than the number of samples may be *insufficient to assume normality*



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



· All are effectively equivalent

00,			•		1	ranks	
00 in 1/20	А	0.03		А	0.99	1	
	А	0.91		А	0.91	2	
	А	0.64		А	0.91	3	
20	А	0.99		А	0.64	4	
	А	0.64		А	0.64	5	
	А	0.16		В	0.64	6	
Two data sata	А	0.16		В	0.64	7	
	А	0.91		А	0.27	8	Circuit 1 to the large of
combined	А	0.16	Sort	В	0.27	9	Give each data element
into a single	А	0.27	~	А	0.16	10	its corresponding rank
array	В	0.64		А	0.16	11	
	В	0.08		А	0.16	12	
	В	0.16		В	0.16	13	
	В	0.27		В	0.16	14	
	В	0.02		В	0.08	15	
	В	0.01		А	0.03	16	
	В	0.16		В	0.03	17	
	В	0.03		В	0.03	18	
	В	0.03		В	0.02	19	Ranked Example
	В	0.64		В	0.01	20	pre-





A Non-Parametric 'Mean': The Median

- Average of a data set that is not normally distributed produces a value that behaves non-intuitively
 - Especially if the probability distribution is skewed
 - Large values in 'tail' can dominate
 - Average tends to reflect the typical value of the "worst" data not the typical value of the data in general
- Instead use the Median
 - 50th percentile
 - Counting from 1, it is the value in the $\frac{n+1}{2}$ position
 - If *n* is even, (n+1)/2 will be between 2 positions,
 - average the values at that position



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 $l = \frac{b}{n-1}$
- The upper percentile u = 1 l
- Confidence Interval is [*value*_l,*value*_u]
 - i.e. $value_l \leq median \leq value_u$
 - With a confidence level of $1-\alpha$



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

₩		Box Plot: Example
	Sort Data	18 0.99 17 0.91 16 0.91 15 0.64 14 0.64 13 0.64 12 0.64 11 0.27 9 0.16 8 0.16 7 0.16 6 0.16 5 0.16 4 0.08 3 0.03 1 0.03









Does My Difference Matter?

- Okay, so your results are significantly better than the published results. So what?
 - Statistics can answer, "is it better?", but not "does it matter?"
- You perform 100 000 runs of your classifier and 100 000 runs of the reference classifier
 - You get a *t* score of 31.6! ^(C)
 - The *p*-score is reported by Excel as 0! (Actually 2.0x10⁻²¹⁹)
 - But...your way classifies data at 91.0% accuracy, whereas the reference technique classifies at 90.8% accuracy.
 - Not much difference!
 - Especially if your technique is much slower than the reference way



Measuring Effect Size

- One statistic for effect size: Cohen's d'
 - d' is computed by $d' = \frac{t}{\sqrt{(n_1 + n_2)/2}}$
 - Measures the difference between means in terms of the pooled standard deviation
 - Cohen suggests that 0.25 is a small difference; 0.50 is a medium-sized difference; 0.75 is a large difference
 - For our example, d' is 0.10
 - Essentially an insignificant difference
- Problem: we did too many runs!



- What is the number of repetitions needed to see if there is a difference between two means or between two medians?
 - Depends on the underlying distributions
 - But underlying distributions are unknown
- Rule of thumb for t-tests...
 - Perform a minimum of 30 repetitions for each system
 - Performing 50 to 100 repetitions is usually better

ANOVA: Analysis of Variance



Part 1a: Multi-Level Analysis Basic Concept





























	Al	NOVA	table	for exa	mple om DataDesk
Source	df	SS	MS	F-ratio	Prob
const	1	3592.9	3592.9	13967	≤ 0.0001
xover	4	210.9	52.7	204.94	\leq 0.0001
Error	95	24.4	0.257		
Total	99	235.3			
$F^* = \frac{MS_{mod}}{MS_{erro}}$	$\frac{el}{r} = \frac{52}{0.2}$	$\frac{7}{57} = 204.94$	<u>F tesi</u> fdist(2	<u>t (From Exc</u> 204.94, 4, 95	e <u>el)</u> 5) = 8.19E-46





Part 1b: 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?





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

$$\alpha' =$$

- 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"

•						-
0	1/He	Pa	irwise	Comp	pariso	ons 🔘
	🕻 🖁 b	etwe	en Fac	tor-Le	evel]	Means
1	Regular Po	air-wise I	f test (with l	Bonf. Corr	rection)	
		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

	•					
@ 🐝	1	Pai	rwise	Comp	pariso	ons 🤇
	i 1	betwee	en Fac	tor-Le	evel	Means
	ANOVA I	Pair-wise T	test (with l	Bonf. Cor	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

		Pa	irwise Comp	aris	ons 🔵
		betwe	en Factor-Le	evel	Means
F	ANOVA	A Pair-wise	T test (with Bonf. Corr	ection)	
		Diff	std. err. t-value	df	p-value
Di <u>f</u> std.	$f = \overline{Y}_{i\bullet}$	$-\overline{Y}_{j\bullet}$ $= \sqrt{\frac{MS_{error}}{n_i}} +$ $= 0.1604$	$\frac{df = n_T - r = rn - r}{= 95}$ $\frac{\overline{MS_{error}}}{n_j} = \sqrt{\frac{2 \cdot MS_{error}}{n}} = $	$= 5 * 20 - \frac{2 * 0}{20}$	- 5 <u>257</u>
t - 1	value =	= Diff stdError	Student-T with Bor p-value = m * tdist(= 10 * tdist(u<u>f. Corr</u> t-value, (t-value,	<u>ection</u> df, two-sided) .95, 2)











Other Post-Hoc Corrections

- Tukey
 - Same as T test except uses the q distribution instead of the t distribution
 - q(1 α, r, n_T r) value is the cut off value where the difference observed would be less than this value with a probability of 1 - α
 if r values are sampled from a normal distribution N(0,1)
 - $DofF = n_T r$
 - q distribution is called the studentized range distribution
 - q "broader" than t,
 - *q* is not as "broad" as *t* after Bonferroni correction
 - *q* distribution is not in Excel, but it is in most other stats packages including R









Multiple Factors: Factorial Design

E.g. if we have 2 EC systems, new and standard (New and Std) and we want to see their behavior under

- crossover and no crossover (x and x)
- 3 different selection pressures (p1, p2 and p3)

Γ		t1	t2	t3	t4	t5	t6	t7	t8	t9	t10	t11	t12
	S	new	new	new	new	new	new	std	std	std	std	std	std
Γ	Χ	х	х	х	x	x	x	х	х	Х	x	×	×
	Р	p1	p2	p3	p1	p2	p3	p1	p2	p3	p1	p2	p3







1	$a_T = 180$ a = 2	
	$b = \frac{1}{2}$	8
	c = 3 $n = 15$	្ន

Multi-Factor ANOVA: Results Report

Source	df	SS	MS	F-ratio	p-value
Const	1	16970	16970	12930	≤ 0.0001
S	1	113	113	86.5	≤ 0.0001
Х	1	775	775	591.0	≤ 0.0001
Р	2	939	469.5	357.7	\leq 0.0001
S*X	1	4.05	4.05	3.1	0.0809
S*P	2	307	153.5	116.8	≤ 0.0001
X*P	2	0.570	0.285	0.217	0.8049
S*X*P	2	0.308	0.154	0.117	0.8892
Error	168	220.5	1.312		
Total	179	2360.12			

































What are the distributions
of
$$b_1$$
 and b_0 ?
 b_1 can be rewritten as
 $b_1 = \sum_{i=1}^n k_i Y_i$ where $k_i = \frac{(x_i - \overline{x})}{\sum (x_i - \overline{x})^2}$
and $b_0 = \overline{Y} - b_1 \overline{x}$
• since the x_i are constant
 b_1 is a linear combination of Y_i 's
• linear combinations of normally distributed
random variables are normally distributed
• S0 ...

What are the distributions
of b_1 and b_0 ? b_1 can be rewritten as
 $b_1 = \sum_{i=1}^n k_i Y_i$ If Y is normally distributed,
 b_1 is normally distributed,
 b_1 is normally distributedand $b_0 = \overline{Y} - b_1 \overline{x}$ Same for b_0 • since the x_i are constant
 b_1 is a linear combination of Y_i 's• linear combination of Y_i 's• linear combinations of normally distributed
random variables are normally distributed• S0 ...35











T test to see if a the slope is statistically significant

- To see if the slope b_1 is statistically different from 0
 - use the T test

$$T = \frac{(b_1 - 0)}{S_{b_1}} = \frac{b_1}{S_{b_1}}$$

- and find the corresponding p-value
- because we we originally estimated 2 parameters use

$$df = n - 2 - 1 = n - 3$$

43



T test to see if a y intercept is statistically significant

43

• To see if the regression line goes through the origin check if b_0 is statistically different from 0

• use the T test

$$T = \frac{(b_0 - 0)}{S_{b_0}} = \frac{b_0}{S_{b_0}}$$

- and find the corresponding p-value
- again because we originally estimated 2 parameters use

df = n - 2 - 1 = n - 3



T test to see if a y intercept is statistically significant

• To see if the regression line goes through the origin check if b_0 is statistically different from 0

These confidence intervals and tests are very important to perform.

Yet they are not commonly done!

• again because we originally estimated 2 parameters use

$$df = n - 2 - 1 = n - 3$$

Part 4 Multi-factor and Polynomial Regression

43













Polynomial Regression E.g.

R squared = 70.2% R squared (adjusted) = 70.1% s = 0.1466 with 1000 - 5 = 995 degrees of freedom

Source Regression	Sum of Squares 50 4708	df 4	Mean S 12 6177	quare 7	F-ratio 587
Residual	21.3783	995	0.0215		507
Variable	Coefficient	s.e. of C	Coeff	t-ratio	p-value
Constant	0.515460	0.0236		21.9	≤ 0.000
Χ	-2.27114	0.3210		-7.07	≤0.000
X^2	8.87396	1.303		6.81	≤0.000
X^3	-6.94563	1.968		-3.53	0.0004
X^4	0.331472	0.9828		0.337	0.7360



R squared = 70.2% R squared (adjusted) = 70.1% s = 0.1466 with 1000 - 5 = 995 degrees of freedom

Source Regression	X ^A 4 is not statistically significant reduce the number of terms by one			F-ratio 587
Residual	21.3783	995 0.021	5	
Variable	Coefficient	s.e. of Coeff	t-ratio	p-value
Constant	0.515460	0.0236	21.9	
X	-2.27114	0.3210	-7.07	≤ 0.0001
X^2	8.87396	1.303	6.81	≤ 0.0001
X^3	-6.94563	1.968	-3.53	0.0004
X^4	0.331472	0.9828	0.337	0.7360





 $\begin{array}{ll} R \ squared = 70.2\% & R \ squared \ (adjusted) = 70.2\% \\ s = \ 0.1465 \ with \ 1000 - 4 = 996 \ degrees \ of freedom \\ \end{array}$

Regress Residua.	regression fur	nction is a cubic	e polynon	nial
Variable	Coefficient	s.e. of Coeff	t-ratio	p-value
Constant	0.510755	0.0190	26.9	≤ 0.000
Χ	-2.17801	0.1636	-13.3	≤ 0.000
X^2	8.45358	0.3813	22.2	≤ 0.000
X^3	-6.28741	0.2515	-25.0	≤ 0.000





