September 16, 2009

Do we really need statistics in science?

"FWF Graduate Seminar"

Timothy M. Young, Ph.D. **Associate Professor Department of Forestry, Wildlife & Fisheries Forest Products Center**













"Don't worry, it will be OK....."







Overview

- Definition of Statistics
 - Variance (σ^2)
 - Random variable (sample space)
 - Probability
- Definition of Science
- 1st Law of Statistics
- Key Assumptions
- Research Program







Statistics

"The Measurement of Uncertainty"









0

Definition of Statistics?

- Many, many definitions.....most people in this room would have different definitions
- Common theme of most definitions:

"....study of variance (σ^2)....."

".....quantifying variance....."









Francis Galton 1822-1911 Variance (σ^2) Components of a System

• If X and Y are two random variables,

 $\frac{(X,Y \text{ dependent}):}{Var(X + Y) - VarX + VarY + 2Cov(X,Y)}$ or, $Var(aX + bY) = a^{2}VarX + b^{2}VarY + 2abCov(X,Y)$

(X,Y independent):

Var(X + Y) = VarX + VarY









Variance (σ^2) Components of a System

• Generalization

(dependent random variables):

$$\operatorname{var}\left[\sum_{i}^{n} X_{i}\right] = \sum_{i}^{n} \operatorname{var}\left[X_{i}\right] + 2\sum_{i < j} \operatorname{cov}\left[X_{i}, X_{j}\right]$$

(independent random variables):
$$\operatorname{var}\left[\sum_{i}^{n} X_{i}\right] = \sum_{i}^{n} \operatorname{var}\left[X_{i}\right]$$







Variance (σ^2) Components of a System

• Reduce (or increase) variance (dependent random variables): $\Psi Var(X + Y) = VarX + \Psi VarY + \Psi 2Cov(X,Y)$

(X, Y independent): $\Psi Var(X + Y) = VarX + \Psi VarY$ or $Var(X + Y) = VarX + \Lambda VarY$

• <u>What is the difficulty</u>: "quantifying σ^2 "







Random variable (*X*) is allowed to vary within a sample space for the set of real numbers

e.g., weight, height, moisture content, number of spots, distance traveled, survival rate, etc.







Two dice:



Outcomes: (2, 3, 4, 5, 6, 7, 8, 9, 10, 11, 12)

Probability of an outcome?









Two dice:













Two dice:













Sample Space

Two dice:













Sample Space

Two dice:



Lucky 7?













Probability Density Function (Discrete pdfs)

Bernoulli

$$f(x) = p^x (1-p)^{1-x}$$

Binomial

$$f(x) = \left(\frac{n}{x}\right) p^{x} (1-p)^{n-x}$$

Poisson

$$f(x) = \frac{e^{-\lambda}\lambda^x}{x!}$$

 $f(x) = p(1-p)^x$

Geometric (or Pascal)

Hypergeometric
$$f(x) = \frac{\left(\frac{K}{x}\right)\left(\frac{M-K}{n-x}\right)}{\left(\frac{M}{n}\right)}$$

etc., etc. etc.....





Probability Density Function (pdf)

Continuous









Probability Density Function (pdf)

Continuous (e.g., t distribution)



The <u>*t* statistic</u> was invented by William Sealy Gosset (1876 - 1937) for cheaply monitoring the quality of beer brews. "Student" was his pen name.











Probability Density Function (Continuous pdfs)

Normal

$$f(x) = \frac{1}{\sqrt{2\pi\sigma}} e^{-(1/2)[(x-\mu)/\sigma]^2}$$

Standard Normal

$$f(x) = \frac{1}{\sqrt{2\pi}} e^{-(1/2)[(x]^2]} (\mu = 0; \sigma = 1)$$

Student's t

$$f(x) = \frac{\Gamma(\frac{p+1}{2})}{\Gamma(\frac{p}{2})} \frac{1}{p\pi^{1/2}} \frac{1}{(1+x^2/p)^{(p+1)/2}}$$

F distribution

$$f(x) = \frac{\Gamma(\frac{p+q}{2})}{\Gamma(\frac{p}{2})\Gamma(\frac{q}{2})} \left(\frac{p}{q}\right)^{p/2} \frac{x^{(p/2)-1}}{\left[1 + (p/q)x\right]^{(p+q)/2}}$$
$$f(x) = \frac{k}{\lambda} \left(\frac{x}{\lambda}\right)^{k-1} e^{-(x/\lambda)^k}$$

Weibull

etc., etc. etc.....









Why are the probability density functions important?









How are your data collected from some sample space?









Statistical Methods

Parametric Methods

Analysis of variance (ANOVA) Chi-square test Correlation **Factor Analysis** Mann-Whitney U Mean Square Weighted Deviation **MSWD** Pearson product-moment correlation coefficient **Regression analysis** Logistic regression Spearman's rank correlation coefficient Student's t-test **Time Series Analysis** etc., etc., etc....

Non-Parametric Methods

Anderson-Darling test Cochran's Q Cohen's kappa Efron-Petrosian test Friedman two-way analysis Kendall's tau Kendall's W Kolmogorov-Smirnov test Kruskal-Wallis one-way analysis Kuiper's test Wilcoxon rank sum test Pitman's permutation test Rank products Siegel-Tukey test Wilcoxon signed-rank test. etc., etc., etc....









Statistical Inference

• Statistical inference (or statistical induction) is the use of statistics and random sampling to make inferences concerning some unknown aspect of the population.









Definition of Science?

























Definition of Science?

The word science comes from Latin "scientia," meaning knowledge.

"Science is an *intellectual activity carried on by humans* that is designed to discover information about the natural world in which humans live and to discover the ways in which this *information* can be *organized* into *meaningful patterns*"

"It is done through *observation of natural phenomena*, and/or through *experimentation* that tries to simulate natural processes under controlled conditions"

"A primary aim of science is to collect *facts (data)*"

"An ultimate purpose of science is to *discern the order* that exists *between and amongst* the various *facts*"

".....systematic knowledge-base or prescriptive practice that is capable of resulting in a *prediction* or predictable type of outcome....."









Scientific Reasoning or Inference











Do we really need Statistics in Science?

- Observations (data)
- Meaningful patterns
- Experimentation
- Discern order between and amongst facts
- Hypothesis testing (or generation)
- Prediction
- etc., etc., etc......







Do we really need Statistics in Science?

How will you quantify variance (σ^2) of observational data without the use of statistical methods?



"An approximate answer to the right question is worth a good deal more than the exact answer to an approximate problem"

John W. Tukey (1915-2000)









First Law of Applied Statistics

(Gleser 1996)

"....two individuals using the same statistical method on the same data should arrive at the same conclusion."









Key Assumptions

- What is the question (problem definition)?
- Sample space (bias?)
- Data quality
- Parametric (pdf) assumption?
- Most appropriate method to provide the approximate answer to a welldefined question





Research Program

(Bio-based Products Industries)











Time Ordered







Hardwood/Softwood Sawmills





Hypothesis, a priori,

"...the use of real-time statistical process control to monitor and reduce lumber thickness variation does <u>not</u> improve lumber recovery, lumber quality or financial performance."







Hardwood/Softwood Sawmills

Summary - Sawmill A Quercus rubra



±								
				Changes		Changes	Standard	Changes
			Average	in	Median	in	Deviation	in
	Month-Year	n	x		(M)	M ***	$\hat{\sigma}_{x}$	$\hat{\sigma}_{\!X}^{}$ ****
Γ	August — 1999	171	1.134"	a	1.135"	а	0.044	а
	September – 1999	230	1.124"		1.117"		0.044	ab
Γ	October – 1999	61	1.126"	a	1.123"	ab	0.027	
Γ	November - 1999	188	1.108"	ab	1.109"	c	0.040	ab d
	December - 1999	641	1.121"	а	1.121"	ab	0.045	ab d
Γ	January - 2000	*						
	February – 2000							
Γ	March - 2000							
	April – 2000	328	1.103"	ab	1.115"	abc	0.035	<u></u> cd
	May — 2000	414	1.114"	а	1.112"	abc	0.033	<u></u> cd
	June — 2000	431	1.180"	а	1.106"	c	0.035	<u>cd</u>
Γ	July — 2000	44	1.096"	Ъ	1.097"		0.024	c







Hardwood/Softwood Sawmills

Company	Investment	Return	ROI
А	\$15,000	\$180,000	12:1
В	\$27,000	\$752 , 000	28:1
C (softwood)	\$13,000	\$210,000	16:1
D	\$21,000	\$147,000	7:1
Total:	\$76,000	\$1,289,000	17:1









Staves for Bourbon Barrels







<u>Hypothesis</u>: A reduction in stave width variability and bilge variability in any of the 15 jointer wheels will reduce both within and between barrel circumference variability













Staves for Bourbon Barrels











Staves for Bourbon Barrels











Statistical Process Control Staves for Bourbon Barrels

- Barrel circumference variation was reduced
- Allowed for increase in barrel target circumference size
- Yield per barrell (<u>Official Proof Gallons</u>) at Jack Daniels improved by 0.3 OPG after SPC
- Additional 938 barrels of Jack Daniels and approximately \$300,000 of cost savings over the sixmonth study period (estimated by Brown Forman)











<u>Question</u>: Can improved methodologies for real-time process modeling improve the scientific understanding of undiscovered correlations in bio-based products manufacturing (facilitate improved causation investigation) ?















<u>Idea</u>: Reduce generalized error of prediction by combining predictions from several models and various types of

algorithms into an "ensemble"

- MLR
- Regression Trees
- Partial Least Squares
- Ridge Regression
- Neural Networks











Several large projects with USDA SBIR competitive grants and private industry (T: \$1.6M)





als IB Graph Rese	duals Graph Data Sets	Server Settings and	Client Control Common	n Parameter History				
			Dece	sace in Parameter sces IB	Increace in Paramet Increaces I®	e		
	Hin.	Value.	Max.	Scaled	Estimate.	Coefficient	p-value.	VIE.
						1175	0.0009	
	145.0	151.6	157.9		-10.20	-3.021	<0.0001	2.090
	65.70	29.96	90.98	18.25		1.444	<0.0001	2.020
	17.38	25.41	84.58		-17.58	-0.523	<0.0001	1.222
	211.8	266.3	268.1	15.52		0.550	<0.0001	2.438
	24.68	64.62	95.35		-11.08	0.314	<0.0001	2.300
Distance	0.097	0.090	0.140		-0.207	-204.2	0.0002	2.471
n Amps	22.45	31.15	34.20		-7.045	-1.335	0.0004	2.391
	23.76	80.03	97.38	7.837		0.213	0.0019	4.360
	440.0	490.0	540.0		6.957	-0.139	0.0002	1.526
stance	0.056	0.095	0.138		6.663	-163.3	0.0285	1.773
circl.	28.00	30.00	50.00	5 168		0.470	0.0031	4.511
	2.409	2.499	2.514		-3.501	-200.2	0.0479	1.345

TOP SOURCES OF VARIATION FOR CU	IBRENT PRODUCT FOR LAS	T 30 DAVS			ck		
Top Five Sources of Variation	Number of Occurences	Change in Parameter to Increase IB by 1	Scaled Estimate.	Proportion Positive	Proportion Negative		
1. Flap. #2 Feedback	12	-2.940	-12.037711	0	100		
2. Press Roll #4 Operator Distance	10	-0.00248	-10.12573455	0	100		
3 Board Thickness Calculated	9	0.614	9.3848695	0	100		
4. Dryer Dutlet Temp	9	-0.790	-13.787005	0	100		
5. Dam Drive Pressure Roll #2	8	15.93	13.23122795	100	0		
Model Statur: No Enors				Model Source Files: 450	model		
				C	urrent Product		
	ş	U Number = 450					
EWMA OF FIEdicion = 50.0							
Enable IB Prediction 2 minutes	Prediction =	97.9			See		







<u>"BioSAT" Model</u>: Modeling system for determining optimal locations for biomass using facilities in the eastern U.S.









Regression Trees



- Regression Tree:
 - Piecewise estimate
 of a regression function
 - Constructed by recursively partitioning the data and sample space















Training "Statistical Thinking"







- Advanced statistical seminars for the bio-based products industries
- Applied design of experiments for the bio-based products industries









Conclusion

- Statistics is a key foundation of science
- <u>Bottom-line</u>: Statistics helps minimize the risk of being wrong

"What makes a scientist great is the care that he/she takes in telling



you what is wrong with his/her results, so that you will not misuse them"

W. Edwards Deming (1900-1993)





Questions & Discussion



"I think you should be more explicit here in step two."

http://www.nytimes.com/2009/08/06/technology/06stats.html "For Today's Graduate, Just One Word: Statistics"





