

Statistical significance, statistical power, and clinical significance: How to understand your results in context

Michael B. Altman, Ph.D.
 Department of Radiation Oncology
 Physics Division

Washington University in St. Louis
 SCHOOL OF MEDICINE

Hypothesis Testing

- H_0 = Null Hypothesis
- H_1 = Alternative Hypothesis

Possible Hypothesis Test Outcomes		
Decision	Accept H_0	Reject H_0 (i.e. Accept H_1)
H_0 is true	Correct Decision (No error)	Type I Error
	Probability = $1 - \alpha$	Probability = α
H_0 is false (i.e. H_1 is true)	Type II Error	Correct Decision (No error)
	Probability = β	Probability = $1 - \beta$

<https://keydifferences.com/difference-between-type-i-and-type-ii-errors.html>

Washington University School of Medicine in St. Louis Department of Radiation Oncology Physics Division

Hypothesis Testing

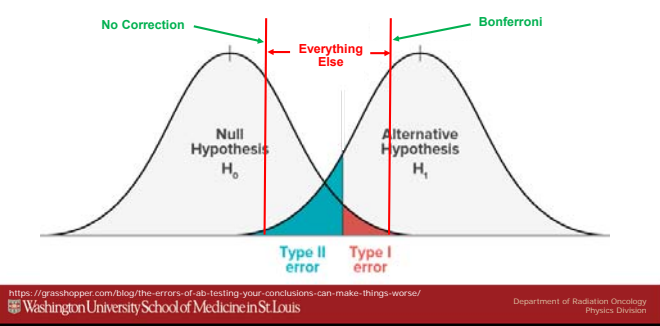
- H_0 = Null Hypothesis
- H_1 = Alternative Hypothesis

Possible Hypothesis Test Outcomes		
Decision	Accept H_0	Reject H_0 (i.e. Accept H_1)
H_0 is true	Correct Decision (No error)	Type I Error FALSE POSITIVE
	Probability = $1 - \alpha$	Probability = α
H_0 is false (i.e. H_1 is true)	Type II Error FALSE NEGATIVE	Correct Decision (No error) POWER
	Probability = β	Probability = $1 - \beta$

<https://keydifferences.com/difference-between-type-i-and-type-ii-errors.html>

Washington University School of Medicine in St. Louis Department of Radiation Oncology Physics Division

Type I / Type II Error Tradeoff



Alternative Methods

- Pick a primary comparison and make all others secondary
 - Likely impractical, inapplicable, or both for radiomics studies
- Alternative FWER corrections
 - Sidak Correction
 - Hochberg's Step-Up Procedure
 - O'Brien's Global Test
 - Many others
- False discovery rate (FDR)-based methods
 - $FDR = \frac{\# \text{ False Positives}}{\# \text{ Total Significant Features}}$
 - Ex: q - value (Storey and Tibshirani, PNAS, 2003)
 - Controlling FDR controls the percentage of features called significant that truly are not
 - 5% FDR rage means that among all features called significant, 5% of these are actually not.

Washington University School of Medicine in St. Louis Department of Radiation Oncology Physics Division

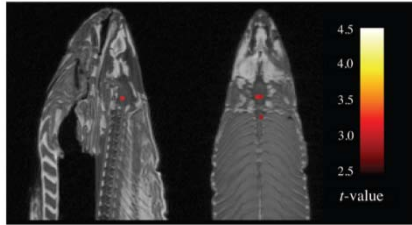
Are Multiple Comparison Corrections Necessary?

- Example study: Bennett *et al.* Neuroimage, 2009.
- EPI images of single subjects brain while being subjected to visual stimuli.
- Images analyzed for regions of brain with significant BOLD signal changes during stimuli photos relative to rest.
- Several active voxels observed in a cluster located within the brain.
 - Statistics uncorrected for multiple comparisons showed activated cluster size was significant at $p = 0.0001$ level.
 - Both Hochberg correction and setting FDR at 0.05 showed no significant voxels in dataset.

Washington University School of Medicine in St. Louis Department of Radiation Oncology Physics Division

Multiple Comparison Corrections Are Necessary

Subject Data	
Height	45.7 cm
Weight	1.73 kg
Sex	????
Species	Salmo salar
Common Name	Atlantic Salmon
Status	DEAD
Suggested Wine Pairing	Depends on the preparation, but a full bodied white like a Chardonnay or a rosé or light-tannin red

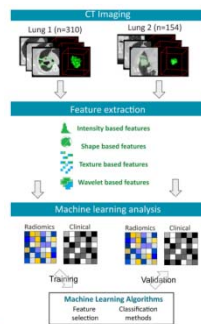


Power

- Power = $1 - \beta \rightarrow$ As power \uparrow , probability of making at Type II error (i.e. β) \downarrow
- A typical desired power is 0.8
- This means 80% chance that a feature called significant is actually significant (true positive rate)
- Increasing power typically comes at the expense of increasing Type I error rate.
- Power decreasing to 0.5 (or worse) means a test's ability to discriminate between false and true values is at or worse than random chance.

Factors Affecting Power

- In general, power impacted by the p-value used (decision threshold), sample size, magnitude of the effect
- For radiomics studies, additionally: feature selection method, classification method, # of selected features (Parmar *et al. Sci. Rep.* 2015).
- Parmar *et al.* study: compared 14 feature selection and 12 classification methods.
 - Five different feature numbers used per combination from 5-50 features (at 5 feature intervals)



Trends Toward Significance

- Verbiage in the literature: (barely) not statistically significant ($p=0.052$); a barely detectable statistically significant difference ($p=0.073$); a borderline significant trend ($p=0.09$); a certain trend toward significance ($p=0.08$); a clear tendency to significance ($p=0.052$); a distinct trend toward significance ($p=0.07$); a favorable trend ($p=0.09$); a favourable statistical trend ($p=0.09$); a little significant ($p<0.1$) ... among many others.
- There is no 'trend to significance' in any direction, and nowhere for the trend to be 'towards'.
- Think of it as pregnancy, you either are or you are not.

Clinical Significance vs. Statistical Significance

- Statistical significance is not clinical significance
- Clinical significance: Differences representing something that is relevant to a level of significantly impacting patient care.
- Example drawn from JCO 2001 (anonymous)
 - Health status questionnaire (HSQ) before / after scores on 1300 patients
 - All p-values <0.0001
 - Conclusion: all domains of QOL were statistically significantly different across treatment groups
 - Problem: 1300 patients provides 80% power to detect a change of 1 unit on 0-100 point scale

Clinical Significance vs. Statistical Significance

- EORT QLA-LC13: Clinical trial assessing change in QOL for lung cancer patients before and during TX using QOL questionnaires and targeted modules

Item	n=537	n=346	Effect Size
Coughing	46.2	44.3	small
Dyspnea	17.2	16.2	small
Pain	26.9	25.5	small

- All p-values were statistically significant

Clinical Significance Approaches

- Assessing clinical significance can be difficult and subjective
- Two Primary Method Classes:
 1. Anchor-based methods requirements
 - Independent interpretable measure (the anchor) which has appreciable correlation between anchor and target
 2. Distribution-based methods
 - Rely on expression of magnitude of effect in terms of measure of variability of results (effect size or ES)
 - Many version attempt to divide ES into different levels of magnitude: small, moderate, large

ES Level Calculations

Table 1
ES Indexes and Their Values for Small, Medium, and Large Effects

Test	ES index	Effect size		
		Small	Medium	Large
1. m_1 vs. m_2 for independent means	$d = \frac{m_1 - m_2}{s}$.20	.50	.80
2. Significance of product-moment r	r	.10	.30	.50
3. r_1 vs. r_2 for independent	$q = z_1 - z_2$ where $z = \text{Fisher's } z$.10	.30	.50
4. $P = .5$ and the sign test	$g = P - .50$.05	.15	.25
5. P_1 vs. P_2 for independent proportions	$h = P_1 - P_2$ where $\phi = \text{arcsine transformation}$.20	.50	.80
6. Chi-square for goodness of fit and contingency	$w = \sqrt{\frac{\chi^2 (P_i - P_{e_i})^2}{n}}$.10	.30	.50
7. One-way analysis of variance	$f = \frac{s_b}{s}$.10	.25	.40
8. Multiple and partial correlation	$r^2 = \frac{R^2}{1 - R^2}$.02	.15	.35

Note: ES = population effect size.

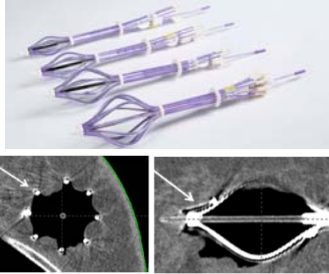
• J Cohen, 1992

Four Guidelines from QOL studies

- Sloan, Cella, Hays, JCE 2005:
 1. The method used to obtain an estimate of clinical significance should be scientifically supportable.
 2. 1/2 SD is a reasonable and scientifically supportable conservative estimate of ES that is likely to be clinically meaningful and can be used in the absence of other information.
 3. ESs < 1/2 SD may be meaningful if they are supported by data regarding the specific characteristics of a particular assessment or application, may also be meaningful.
 4. If feasible, multiple approaches to estimating a tool's clinically meaningful effect size in multiple patient groups are helpful in assessing the variability of the estimates. However, the lack of these should not preemptively restrict the application of information gained.

Radiomics-(ish) Example Study

- Altman et al. *Brachytherapy*, 2018
- Strut adjusted volume implant (SAVI) used for accelerated partial breast irradiation (APBI)
- Implant device → "Day 0" CT → Create TX plan
- CT 48-72 hours later → Replan if necessary → Treat
- Question: Can we use imaging-based features to predict which patients will need replanning from the Day 0 CT?



Washington University School of Medicine in St. Louis

Department of Radiation Oncology
Physics Division

Methods and Metrics

- 62 consecutive patients, 36% requiring replans

Table 1
Summary of metrics used in the retrospective analysis

Metric	Definition	Variable type
Plan adapted	Yes or no: was treatment plan replanned using Verity CT?	Binary
Applicator size	Use of which size SAVI applicator: 6-1 mini, 6-1, 8-1, or 10-1	Categorical
Attending physician	Which one of two attending physicians implanted the device and evaluated the need to replan the DOI CT plan?	Binary
Splayed struts	Yes or no: were any of the struts in the applicator splayed?	Binary
Splay toward skin	Yes or no: was the direction of the splay in the applicator directed toward the skin?	Binary
Absolute volume of air	Absolute volume of air in the PTV_EVAL	Continuous
Percent volume of air	Percent volume of air in the PTV_EVAL	Continuous
Minimum distance to skin or ribs (SAVI + 1)/(SAVI + 1 - SKIN)	Minimum distance from any part of the applicator volume to the skin or ribs	Continuous
	Ratio of volume of 1-cm thick shell around the applicator volume to that same shell cropped from a margin of 0.5 cm from the outer surface of the skin	Continuous

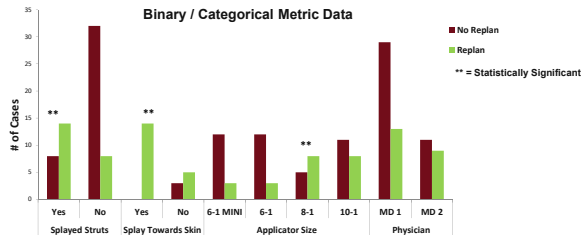
DOI = day-of-implant; SAVI = strut-adjusted volume implant.
Definitions of Verity CT, PTV_EVAL, and splay discussed in the text.

- Evaluation: Fisher's exact test, Wilcoxon rank-sum, 2 tailed p-test for categorical, binary, and continuous variables, respectively.

Washington University School of Medicine in St. Louis

Department of Radiation Oncology
Physics Division

Results



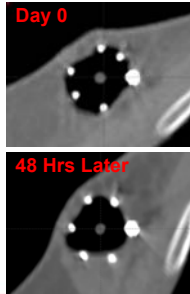
- No continuous variables showed statistical significance

Washington University School of Medicine in St. Louis

Department of Radiation Oncology
Physics Division

Statistical vs. Clinical Significance

- Splay (Towards Skin)
 - Flexible struts can be pushed around by patient's anatomy
 - Anatomy (and device) can shift/change over the 24 hours after implant → altered dose distribution vs. Day 0 plan
 - Statistical significance seems to correlate to a clinical result → clinically significant result → change in clinical workflow
- Use of 8-1 SAVI
 - Use of only 1 (out of 4) middle-sized device was statistically significant
 - Clinically intuitive reason why: ??? → Not deemed clinically significant



Washington University School of Medicine in St. Louis

Department of Radiation Oncology
Physics Division

Summary

- For “big data” type studies as in radiomics with multiple hypotheses tested simultaneously, corrections must be made to statistical significance thresholds to limit Type I errors.
- Too conservative corrections can result in a rise in Type II errors (i.e. limited power).
- Other factors such as feature selection methods, classification methods, and number of features used can also impact power.
- Clinical significance statistical significance.
- Clinical significance should be assessed case-wise using accepted methods and/or guidelines in concert with input from clinicians
- Don't go it alone – statisticians are your friends!

Washington University School of Medicine in St. Louis

Department of Radiation Oncology
Physics Division

Bibliography



1. Altman *et al.* *Brachytherapy*. 17, 2008; 40-49
2. Bennett CM *et al.* *Neuroimage*, 47, 2009; S125
3. Cohen J. *Psycholog Bulletin*, 112(1), 1992; 155-159
4. Kumar *et al.* *MRI* 30, 2012; 1234-1248
5. Parmar *et al.* *Sci Rep*, 5(13087), 2015
6. Sloan JA, Cella D, Hays RD. *J Clin Epid* 58(12), 2005; 1217-1219.
7. Storey JD and Tibshirani R., *Proc Natl Acad Sci USA*, 100 (16), 2003; 9440-9445
8. Yip SSF and Aerts HJWL. *Phys Med Bio.* 61, 2016. R150-166

Washington University School of Medicine in St. Louis

Department of Radiation Oncology
Physics Division

Michael B. Altman
Assistant Professor of Radiation Oncology
Campus Box XXXX
Street address
St. Louis, MO 63110
(314) 123-4567
m.altman@wustl.edu

©2018

Washington University School of Medicine in St. Louis
Department of Radiation Oncology
Physics Division
